

SCIENTIFIC AMERICAN

JUNE 2007

WWW.SCIAM.COM

Bring
Back
America's
**Prehistoric
Beasts**



Did this molecule start life?



FORGET DNA AND RNA. MAYBE IT
ALL BEGAN WITH SOMETHING
MUCH SIMPLER

The Mysteries of **Anesthesia**

When It Pays to **Be Irrational**

How Particles **Shape the Cosmos**

june 2007
contents

SCIENTIFIC AMERICAN Volume 296 Number 6

features

BIOLOGY

46 A Simpler Origin for Life

BY ROBERT SHAPIRO

Energy-driven networks of small molecules may be more likely first steps for life than the commonly held idea of the sudden emergence of large self-replicating molecules such as RNA.

MEDICINE

54 Lifting the Fog around Anesthesia

BY BEVERLEY A. ORSER

Learning why current anesthetics are so potent and sometimes dangerous will lead to a new generation of safer targeted drugs.

PARTICLE COSMOLOGY

62 When Fields Collide

BY DAVID KAISER

The history of particle cosmology shows that science can benefit from wrenching changes.

ECOSYSTEMS

70 Restoring America's Big, Wild Animals

BY C. JOSH DONLAN

Pleistocene rewilding—a proposal to bring back animals that disappeared from North America 13,000 years ago—offers an optimistic agenda for 21st-century conservation.

INFORMATION TECHNOLOGY

78 Breaking Network Logjams

BY MICHELLE EFFROS, RALF KOETTER AND MURIEL MÉDARD

Network coding could dramatically enhance the efficiency of communications networks.

INNOVATIONS

86 Seeing Triple

BY STUART F. BROWN

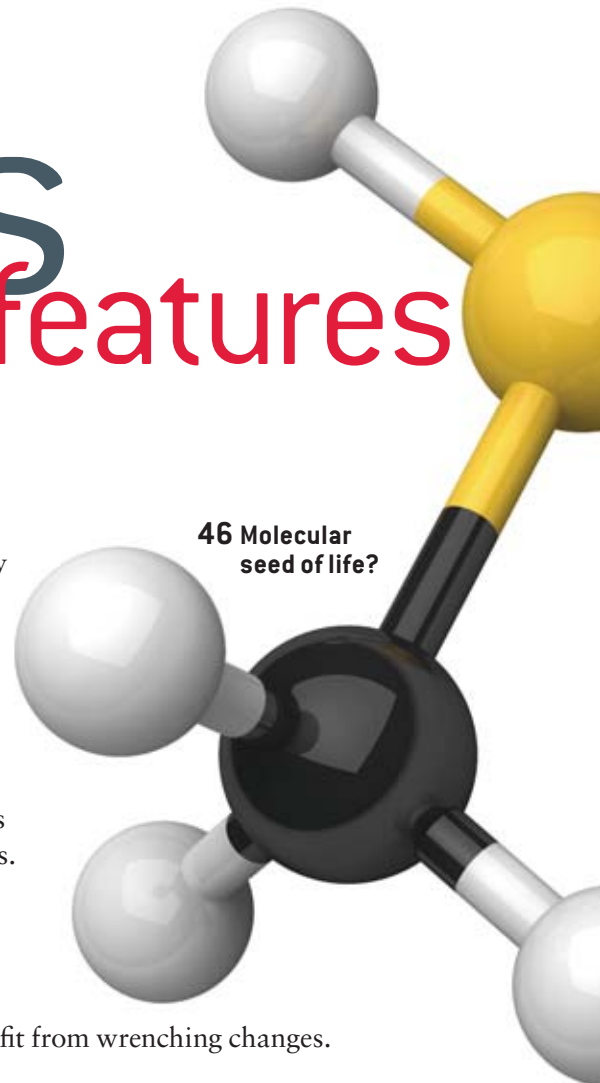
Anticipated for decades, machines are finally displaying objects in three true dimensions.

GAME THEORY

90 The Traveler's Dilemma

BY KAUSHIK BASU

People playing this simple game consistently reject the rational choice. In fact, by acting illogically, they end up reaping a larger reward—an outcome that demands a new kind of formal reasoning.



46 Molecular seed of life?

departments



- 8 SA Perspectives**
Serengeti in the Dakotas.
- 10 How to Contact Us**
- 10 On the Web**
- 12 Letters**
- 18 50, 100 & 150 Years Ago**
- 20 News Scan**
 - Planting a better biodiesel.
 - Fingerprinting individual atoms.
 - Meet the Memjet.
 - Less lethal turbines for dams.
 - Solid-state laser blaster?
 - New fish food: psoriasis patients.
 - Computers can finally pass Go.
 - Data Points: Large Hadron Collider.

- 40 Insights**
New genetic studies, Eric Vilain says, should force a rethinking about mixed-sex babies and gender identity.
- 96 Working Knowledge**
Optical character recognition finds the write type.
- 98 Reviews**
Books on dirt, flowers, uncertainty, perishable artifacts.



Eric Vilain,
University of California, Los Angeles

40

columns

- 39 Skeptic** BY MICHAEL SHERMER
The inverse square law trumps the law of attraction.
- 43 Sustainable Developments**
BY JEFFREY D. SACHS
As global warming tightens the availability of water, prepare for forced mass migrations.
- 102 Anti Gravity** BY STEVE MIRSKY
Coffee cup full of beans.
- 104 Ask the Experts**
How do itches come about?
Why is the sun in the middle of the solar system?

Cover image and illustration of molecule on preceding page by Ken Eward, BioGrafx; cover photograph of cheetah by Shani; photograph at left by Brad Swonetz, Redux.

Scientific American (ISSN 0036-8733), published monthly by Scientific American, Inc., 415 Madison Avenue, New York, N.Y. 10017-1111. Copyright © 2007 by Scientific American, Inc. All rights reserved. No part of this issue may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying and recording for public or private use, or by any information storage or retrieval system, without the prior written permission of the publisher. Periodicals postage paid at New York, N.Y., and at additional mailing offices. Canada Post International Publications Mail (Canadian Distribution) Sales Agreement No. 40012504. Canadian BN No. 127387652RT; QST No. Q1015332537. Publication Mail Agreement #40012504. Return undeliverable mail to Scientific American, P.O. Box 819, Stn Main, Markham, ON L3P 8A2. Subscription rates: one year \$34.97, Canada \$49 USD, International \$55 USD. Postmaster: Send address changes to Scientific American, Box 3187, Harlan, Iowa 51537. Reprints available: write Reprint Department, Scientific American, Inc., 415 Madison Avenue, New York, N.Y. 10017-1111; (212) 451-8877; fax: (212) 355-0408. Subscription inquiries: U.S. and Canada (800) 333-1199; other (515) 248-7684. Send e-mail to sacust@sciam.com Printed in U.S.A.



SA Perspectives

Serengeti in the Dakotas

Some people love science for the crazy ideas, the ones that transport you beyond the everyday grind: black holes, alien life, anything with the word “quantum.” Others prefer the not-crazy ideas, the practical solutions: zippier computers, 100-mpg cars, cures for cureless diseases.

So what do you make of an idea like Pleistocene rewilding? It manages to be both crazy and not crazy at the same time. As the article by C. Josh Donlan

beginning on page 70 describes, a team of biologists has proposed a decades-long project to restock North America with large mammal species like those that roamed the continent before humans crossed the Bering Strait—species such as camels, lions and elephants (the nearest thing to mammoths). The undertaking would culminate in a vast national park—1,000 square miles or more—stretching across the Great Plains. The plains states are depopulating anyway, whereas Africa and Asia are filling up. So the project

would transplant wildlife from where it gets in the way to where it would have plenty of room.

To be sure, Midwesterners might not see it that way. Elephant families running free under big skies sounds romantic—unless you have to dodge them on your morning commute. Lion cubs are so very cute—except when they wander into your backyard. Farmers worry about rampaging rogues, cattle ranchers about novel diseases. Proponents have addressed some of the concerns but clearly have a lot more work to do.

Whether or not cheetahs ever chase pronghorn across the continent again, the rewilding concept has drawn attention to the fact that the loss of biodiversity is not just a problem for the rain forest; it affects less exotic locations, too. The demise of large animals has thrown entire ecosystems out of balance. Even if humans decided now to leave these ecosystems alone, they are too far gone to recover on their own. The prairie would revert not to its Pleistocene glory but to a scraggly weedland.

Instead of merely bemoaning nature’s plight, the proponents of rewilding are doing something about it. The reintroduction of wolves in Yellowstone a decade ago brought the population of moose and elk under control, with a cascade of benefits for vegetation, birds and beavers. Other projects have reintroduced tortoises, bison and falcons into their old haunts. Wildlife-stocked private ranches let you go on safari within half an hour’s drive of the Alamo. More broadly, biologists are also working to restore fisheries, forests and wetlands.

Most of these efforts are too piecemeal to nurse whole ecosystems back to health. They need a broader framework, and a grand rewilding project can provide it. Visionary schemes have a checkered history. The ones that work combine big ideas with baby steps. They need both: for lack of attention to detail, a grand project can easily fail; equally so, for lack of an overarching plan, incremental steps can seem like spitting in the sea. Recent history offers many examples of enterprises that, deprived of a real goal, lost their way. One that comes to mind is the human space program. After the moon landings, NASA became little more than a delivery van service.

A thriving ecosystem can underpin economic prosperity and enhance our quality of life, but it won’t happen on its own. We have to make it happen.



NOT QUITE a woolly mammoth, but close.

THE EDITORS editors@sciam.com

I How to Contact Us

EDITORIAL

For Letters to the Editors:

Letters to the Editors
Scientific American
415 Madison Ave.
New York, NY 10017-1111

or

editors@sciam.com

Please include your name
and mailing address,
and cite the article
and the issue in
which it appeared.

Letters may be edited
for length and clarity.

We regret that we cannot
answer all correspondence.

For general inquiries:

Scientific American
415 Madison Ave.
New York, NY 10017-1111

212-451-8200

fax: 212-755-1976

or

editors@sciam.com

SUBSCRIPTIONS

For new subscriptions,
renewals, gifts, payments,
and changes of address:

U.S. and Canada

800-333-1199

Outside North America

515-248-7684

or

www.sciam.com

or

Scientific American

Box 3187

Harlan, IA 51537

REPRINTS

To order reprints of articles:

Reprint Department

Scientific American

415 Madison Ave.

New York, NY 10017-1111

212-451-8877

fax: 212-355-0408

reprints@sciam.com

PERMISSIONS

For permission to copy or reuse
material from SA:

www.sciam.com/permissions

or

212-451-8546 for procedures

or

Permissions Department

Scientific American

415 Madison Ave.

New York, NY 10017-1111

Please allow three to six weeks

for processing.

ADVERTISING

www.sciam.com has electronic contact
information for sales representatives
of Scientific American in all regions of
the U.S. and in other countries.

New York

Scientific American

415 Madison Ave.

New York, NY 10017-1111

212-451-8893

fax: 212-754-1138

West/Northwest

818-827-7138

fax: 818-827-7139

Detroit

Karen Teegarden & Associates

248-642-1773

fax: 248-642-6138

Midwest

Derr Media Group

847-615-1921

fax: 847-735-1457

Southeast and Southwest

Publicitas North America, Inc.

972-386-6186

fax: 972-233-9819

Direct Response

Special Aditions Advertising, LLC

914-461-3269

fax: 914-461-3433

Australia

IMR Pty Ltd.

+612-8850-2220

fax: +612-8850-0454

Belgium

Publicitas Media S.A.

+32-(0)2-639-8420

fax: +32-(0)2-639-8430

Canada

Derr Media Group

847-615-1921

fax: 847-735-1457

France and Switzerland

PEM-PEMA

+33-1-46-37-2117

fax: +33-1-47-38-6329

Germany

Publicitas Germany GmbH

+49-211-862-092-0

fax: +49-211-862-092-21

Hong Kong

Hutton Media Limited

+852-2528-9135

fax: +852-2528-9281

India

Convergence Media

+91-22-2414-4808

fax: +91-22-2414-5594

Japan

Pacific Business, Inc.

+813-3661-6138

fax: +813-3661-6139

Korea

Biscom, Inc.

+822-739-7840

fax: +822-732-3662

Middle East

Peter Smith Media & Marketing

+44-140-484-1321

fax: +44-140-484-1320

The Netherlands

Insight Publicitas BV

+31-35-539-5111

fax: +31-35-531-0572

Scandinavia and Finland

M&M International Media AB

+46-8-24-5401

fax: +46-8-24-5402

U.K.

The Powers Turner Group

+44-207-592-8331

fax: +44-207-630-9922

I On the Web

WWW.SCIAM.COM

UPDATED EVERY WEEKDAY

Visit www.sciam.com/ontheweb

to find these recent additions to the site:

NEWS

Gene Makes Small Dogs Small



In big news for small dogs everywhere, researchers have found a tie that binds the small breeds, from Chihuahua to Pomeranian to Pekingese: they all share the same version of a gene for a growth hormone called insulinlike growth factor 1 (IGF1).

SPECIAL REPORT

The Poisoning of Our Pets

Scientists and government agencies home in on the cause of more than 100 pet deaths from tainted food.

FACT OR FICTION?

Waking Sleepwalkers May Kill Them

On the contrary, rousing a sleepwalker could save the person's life.

BLOG

When the Sky Falls, Where Will NASA Be?

Everybody duck and cover. In March, NASA shrugged and told Congress that it neither has the funding nor the resources to meet its goal of identifying 90 percent of near-Earth asteroids 150 yards or more in diameter by the year 2020.

ASK THE EXPERTS

Why is there an ozone hole in the atmosphere, whereas there is too much ozone at ground level?

Why doesn't ground level ozone rise to fill the hole?

Ross J. Salawitch, a senior research scientist at the NASA Jet Propulsion Laboratory in Pasadena, Calif., explains.

NEW: VIDEO

SciAm.com now features the latest video science news every weekday.

Subscribe to Scientific American Digital
Visit www.sciamdigital.com

EDITOR IN CHIEF: John Rennie
 EXECUTIVE EDITOR: Mariette DiChristina
 MANAGING EDITOR: Ricki L. Rusting
 CHIEF NEWS EDITOR: Philip M. Yam
 SPECIAL PROJECTS EDITOR: Gary Stix
 SENIOR EDITOR: Michelle Press
 EDITORS: Mark Alpert, Steven Ashley, Graham P. Collins, Mark Fischetti, Steve Mirsky, George Musser, Christine Soares
 CONTRIBUTING EDITORS: W. Wayt Gibbs, Marguerite Holloway, Michael Shermer, Sarah Simpson

EDITORIAL DIRECTOR, ONLINE: Kate Wong
 NEWS EDITOR, ONLINE: Lisa Stein
 ASSOCIATE EDITORS, ONLINE: David Biello, Christopher Mims
 NEWS REPORTER, ONLINE: JR Minkel

ART DIRECTOR: Edward Bell
 SENIOR ASSOCIATE ART DIRECTOR: Jana Brenning
 ASSOCIATE ART DIRECTOR: Mark Clemens
 ASSISTANT ART DIRECTOR: Johnny Johnson
 PHOTOGRAPHY EDITOR: Emily Harrison
 PRODUCTION EDITOR: Richard Hunt

COPY DIRECTOR: Maria-Christina Keller
 COPY CHIEF: Molly K. Frances
 COPY AND RESEARCH: Daniel C. Schlenoff, Michael Battaglia, Smitha Alampur, Michelle Wright, John Matson, Aaron Shattuck

EDITORIAL ADMINISTRATOR: Jacob Lasky
 SENIOR SECRETARY: Maya Hartly

ASSOCIATE PUBLISHER, PRODUCTION: William Sherman
 MANUFACTURING MANAGER: Janet Cermak
 ADVERTISING PRODUCTION MANAGER: Carl Cherebin
 PREPRESS AND QUALITY MANAGER: Silvia De Santis
 PRODUCTION MANAGER: Christina Hippeli
 CUSTOM PUBLISHING MANAGER: Madelyn Keyes-Milch

ASSOCIATE PUBLISHER, CIRCULATION: Simon Aronin
 CIRCULATION DIRECTOR: Christian Dorbandt
 RENEWALS MANAGER: Karen Singer
 FULFILLMENT AND DISTRIBUTION MANAGER: Rosa Davis

VICE PRESIDENT AND PUBLISHER: Bruce Brandfon
 SALES DEVELOPMENT MANAGER: David Tirpack
 SALES REPRESENTATIVES: Jeffrey Crennan, Stephen Dudley, Stan Schmidt

ASSOCIATE PUBLISHER, STRATEGIC PLANNING: Laura Salant

PROMOTION MANAGER: Diane Schube
 RESEARCH MANAGER: Aida Dadurian
 PROMOTION DESIGN MANAGER: Nancy Mongelli
 GENERAL MANAGER: Michael Florek
 BUSINESS MANAGER: Marie Maher
 MANAGER, ADVERTISING ACCOUNTING AND COORDINATION: Constance Holmes

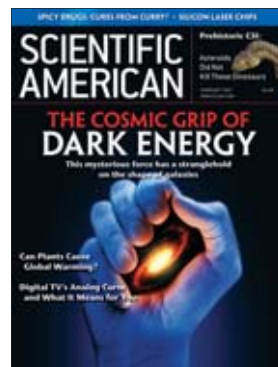
DIRECTOR, SPECIAL PROJECTS: Barth David Schwartz

MANAGING DIRECTOR AND VICE PRESIDENT, ONLINE: Mina C. Lux
 DIRECTOR, WEB TECHNOLOGIES, ONLINE: Vincent Ma
 SALES REPRESENTATIVE, ONLINE: Gary Bronson

DIRECTOR, ANCILLARY PRODUCTS: Diane McGarvey
 PERMISSIONS MANAGER: Linda Hertz

CHAIRMAN EMERITUS: John J. Hanley
 CHAIRMAN: Brian Napack
 PRESIDENT AND CHIEF EXECUTIVE OFFICER: Gretchen G. Teichgraber
 VICE PRESIDENT AND MANAGING DIRECTOR, INTERNATIONAL: Dean Sanderson
 VICE PRESIDENT: Frances Newburg

UNSURPRISINGLY, Frank Keppler and Thomas Röckmann's February article on the discovery of methane emissions in plants elicited strong reactions. Although some correspondents seemed eager to jump to the exact global warming-denying conclusions Keppler and Röckmann had cautioned against, others expressed a more thoughtful attitude as to how this discovery might affect our understanding of global climate change. Readers were also drawn to the more cosmic implications of Christopher J. Conselice's overview of how scientists believe dark energy shapes the universe. Many were unsatisfied with our inability to detect such a force from anything but its effects and offered their own theories as to what dark energy might tangibly represent or how the universe could function without it.



DENTAL DELIMITATION

"License to Work" [News Scan], by Roger Doyle, suggests that the reason the number of dentists in the U.S. has not grown substantially compared with other professions is restrictive licensing practices. Doyle has the cart before the horse. The license to practice dentistry is obtained after the completion of educational requirements and is typically passed by most dentists, although it sometimes calls for more than one attempt. The restriction on numbers is at the beginning of the road, where the educational system has not changed the number of dentists: it is capable of training on average since the 1970s. This restriction is not caused by the licensing board but by the cutting of direct and indirect federal and state support for dental education (number of schools, class size, faculty numbers, student loans, and so on).

Keith J. Lemmerman
 Graduate Periodontics
 University of Kentucky
 College of Dentistry

DIGITAL DOMINION

In "Digital TV at Last?" Michael Antonoff writes about consumer issues related to switching from analog broadcasts to digital with no mention of digital-rights management (DRM). How many know that there is a serious move afoot to allow content owners to limit time shifting of their programs to 90

minutes after we make this switch (effectively reversing the principle of fair use forged in the early days of videotape)? Or that we could be charged every time we view their programming with a premium for every different device we view it on?

Thomas Phelan
 Hightstown, N.J.

ANTONOFF REPLIES: Digital TV potentially affords more control by broadcasters and studios over their products and can restrict what consumers are allowed to do with programs once they enter their homes. But to do this subject justice would have required an in-depth discussion of the various aspects of DRM (many of which may or may not be implemented). Recording rights, home networking, down conversion of digital signals to analog devices and legal fair-use issues would all need to be examined in detail, requiring a full article unto itself.

ILLUSORY IQ

In "Unsettled Scores: Has the Black-White IQ Gap Narrowed?" [News Scan], Marina Krakovsky cites so-called hereditarians who believe that lower IQ test scores among blacks are the result of inherent genetic factors. The presupposition that what is being measured by such tests is an innate intelligence that can be described as biological and unified is one that has never been proved. Psychologist Alfred Binet, who created the first standard for testing, insisted

himself that his tests merely identified where below-average performers lacked skills in the classroom and could best be tutored to bring them up to normal performance standards. The American twist on IQ testing, which included the conception of a biological and unitary intelligence, was based on assumptions that all character traits are Mendelian in nature.

Aaron Method
Indianapolis

RUMINANT RUMINATION

In “Methane, Plants and Climate Change,” by Frank Keppler and Thomas Röckmann, two graphs compare sources of methane in the atmosphere during

preindustrial times with those of today. Ruminants are listed as a major source of current emissions but are not included in the preindustrial chart. Did as many as 70 million bison really produce that much less methane than today’s cattle?

Ed Miller
Falls City, Ore.

KEPPLER REPLIES: *Although wildlife certainly produced methane in preindustrial times, this output was just a minor fraction of the 233 million metric tons of yearly global methane emissions. According to estimates made by environmental scientist Susan Subak and her colleagues in a 1994 article for Chemosphere, the total production of methane by wild ruminants in that period was no more than 10 million metric tons a year—a figure that takes into account the North American bison population (which Subak estimates to have comprised 60 million animals) and the natural ruminants of Africa and other continents. An estimated 1.4 billion head of cattle populate the world today—far more ruminants than existed in preindustrial times.*

Furthermore, modern cattle are bred for productivity, which probably leads them to emit more methane than their wild relatives did. Estimates put their methane production at 115 million metric tons a year.

PHYS ED PHYSICS

In “Eat, Drink and Be Merry” [Skeptic], Michael Shermer suggests that the mechanics of a cyclist taking longer to climb a 5 percent grade with 10 pounds of extra weight can be described in terms of Newton’s law of $F = MA$ (Force equals Mass times Acceleration): “The Force needed to turn the pedals equals Acceleration times that Mass on the saddle.”

Shermer cites the wrong Newtonian concept. The correct equation would be $P = FV$ (Power equals Force times Velocity). Assuming the Power (the rate at which the rider converts internal energy into mechanical energy to propel the bicycle) remains the same with or without the extra 10 pounds, the increase in Force required to “lift” that weight must result in a decrease in the Velocity of the bicycle and rider. The increase in Force is not equal to the full 10 pounds, because the rider is on an inclined plane, but this is a separate issue.

Allen Zimmerman
Ohio State University

ERRATA In “The Universe’s Invisible Hand,” by Christopher J. Conselice, a graph in the box “Dark Energy Takes Charge” represents the fraction of galaxies taking a spiral shape with a pink line and those settling into an elliptical shape with a yellow line. The labels on both of these should be swapped.

“Making Silicon Lase,” by Bahram Jalali, incorrectly states that a data-transfer rate of 10 gigabits a second would constitute a 10,000-fold improvement over a rate of one megabyte a second. Because one byte is equal to eight bits, the improvement would be 1,250-fold.

CLARIFICATION “Methane Flatline,” by David Biello [News Scan], reports that methane traps 23 times more heat per molecule than carbon dioxide. This figure refers to the average amount of infrared radiation absorbed by a methane molecule during its 10-year life span in the atmosphere as compared with the amount absorbed by a CO_2 molecule over a century. It is a separate figure from the Intergovernmental Panel on Climate Change’s estimates of the global-warming potential [GWP] of both substances [methane’s GWP is 23 times that of CO_2 per kilogram].

Send letters to editors@sciam.com or to *Scientific American*, 415 Madison Ave., New York, NY 10017. Letters become the property of *SCIENTIFIC AMERICAN* and may be edited for space and clarity.



GROWING CONCERN: Estimates of the relative contributions to global methane levels from natural sources now need to take emissions from living plants into consideration.

Lunar Idea ■ Milk Agenda ■ Balloon Dismay

JUNE 1957

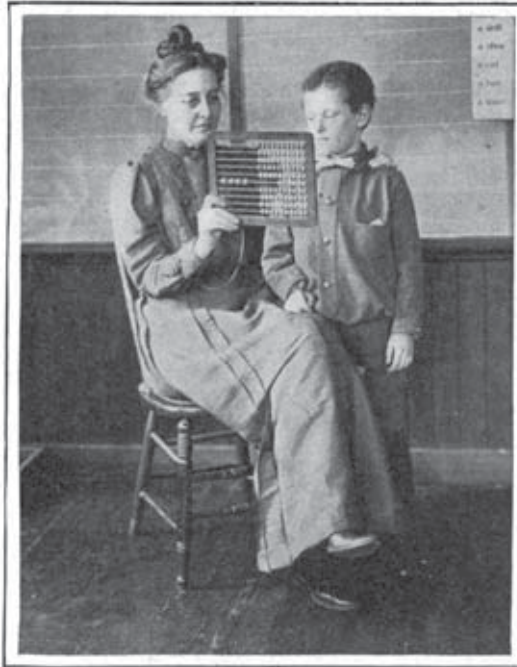
MOON DUST—“The possibility of actually bringing back some of the moon’s material is a scientific bonanza so alluring that ingenious schemes have been proposed to accomplish it, even without landing on the moon. We might, for example, send a pair of rockets, one trailing the other closely by means of a homing device. The first rocket would drop a small atomic bomb on the moon. Since the moon has no atmosphere and comparatively little gravity, the bomb cloud would rise very high. The second rocket could dive into the cloud, collect some of the spray and emerge from its dive by means of an auxiliary jet. Of course, such a maneuver would require a miracle of electronic guidance.—Krafft A. Ehrlicke and George Gamow”

JAWS—“The plain fact is that the ferociousness of the shark is too well documented by long experience to be dismissed as a legend. Jacques-Yves Cousteau, the greatest living expert on modern deep-sea diving, has reported harrowing encounters with sharks. Off the Cape Verde Islands of Portugal, Cousteau and his fellow diver Frédéric Dumas found themselves closely pressed by an eight-foot gray shark, later joined by two five-footers and several blue sharks. The pack resisted all attempts to drive them off: they would shy away only to come back almost at once. The two men used every method in the diver’s book. The sharks were not scared off by any of these tricks, or by the shark ‘repellent’ (copper acetate) strapped to the divers’ legs or by a blow Cousteau gave one approaching shark on the nose with his heavy undersea camera. The two backtracking divers got out just in time.”

JUNE 1907

GOT BOTTLED MILK?—“Milk, the article of food most susceptible to contamination, is served in bottles which are used again and again. This is vitally wrong. The solution is the abandonment of it, and the substitution of the single-service paper package of milk. These are now just beginning to be available. One of the first paper bottles to be placed on the market is a plain paper cylinder, made of new spruce-wood paper in clean, sanitary surroundings. After the bottom is

SCIENTIFIC AMERICAN



TEACHING a deaf pupil how to count, 1907

put in, the bottle is dipped in hot paraffin, the same paraffin the housewife pours over her jellies to keep out the air, moisture, and dust.”

FORCING SPEECH—“The loss of the sense of hearing should not necessarily mean deprivation of the power of speech also. It is only within recent years that we

have come to realize this fact, and in up-to-date institutions the old-fashioned finger alphabet is now unknown. Every child is taught to speak in the natural way by means of the vocal organs. The four or five years of the primary course are devoted almost exclusively to the acquirement of language and numbers [see illustration].”

JUNE 1857

ELECTRIC MOTOR—“We admit that no explosion can take place in the batteries of Prof. Vergnes’ electro-magnetic engine, as with a steam boiler; but his engine is neither as simple nor compact as a steam engine, taking the latter with all its appurtenances. His large electro-magnetic engine in the Crystal Palace, which is not claimed to exert more than 10 horse power (and which, we believe, from mere inspection of its operation, is less than five), with 128 cups of electric battery, occupies more space—engine and battery—than many steam engines working up to 20 horse power.”

LOOSE GAS—“The miniature toy balloons introduced with such success into Paris during the holidays of last winter [made from ox intestine membranes], we suggested as pleasing and beautiful toys to delight Young America. Quite recently they have come into pretty general circulation, and may be seen in various store windows in our city. A number of our boys who have purchased such balloons have been rather astonished to find them daily growing beautifully less in size, and prematurely old in wrinkles, and at last ceasing to be balloons at all. The cause of this is the percolation of the inflating gas (hydrogen) through the pores of the balloon.”

Green Gold in a Shrub

ENTREPRENEURS TARGET THE JATROPHA PLANT AS THE NEXT BIG BIOFUEL **BY REBECCA RENNER**

A woody shrub with big oily seeds could be the ideal source for biofuel. For hundreds of years, Africans in places such as Tanzania and Mali have used *Jatropha curcas* (jatropha) as a living fence. Now biodiesel entrepreneurs in tropical zones in Africa and India are buying up land, starting plantations and looking forward to making fuel from the seeds, which, they ar-

gue, will be better for the global environment and economy than conventional biofuel crops grown in temperate climates.

Ethanol from corn or sugarcane and biodiesel from canola, soy or palm oil have become major players in renewable energy. In principle, biofuels do not increase the amount of carbon dioxide in the air, because as the plants grow they trap the CO₂ that is released when the biofuels are burned.

Still, biofuels face a great deal of criticism. Food commodities such as corn, canola and soy all yield oil, but they are expensive, require intensive agricultural practices and threaten food supplies. Jatropha seems to offer the benefits of biofuels without the pitfalls. The plants favor hot, dry conditions and hence are unlikely to threaten rain forests. There is no trade-off between food and fuel either, because the oil is poisonous. John Mathews, a professor of strategic management at Macquarie University in Australia, notes that many tropical developing countries have huge swaths of degraded and semiarid land that can be utilized for fuel crops. The cost of labor there is cheap, too. Biofuels made from plants such as jatropha, he argues, “represent the best bet for a last-ditch effort to industrialize the poor south and end poverty.” He advocates large-scale plantings to aid energy independence in expanding economies such as China and India



JATROPHA SEEDLINGS are planted in Zambia for U.K. biodiesel firm D1 Oils—part of an increasing effort to harvest the shrub, which favors hot, dry climates, as a source of biofuel.

and to boost exports in the less developed countries of Africa.

Mathews's vision may be coming true. U.K. biodiesel company D1 Oils has planted 150,000 hectares of jatropha in Swaziland, Zambia and South Africa, as well as in India, where it is part of a joint venture. The firm plans to double its crop sizes this year. Dutch biodiesel equipment manufacturer BioKing is developing plantings in Senegal, and the government of China has embarked on a massive project. "People aren't making much jatropha oil right now, because everyone wants seeds for planting," says Reinhard Henning, a German technology transfer consultant and expert in jatropha.

In addition to establishing plantations, jatropha boosters are starting to identify, select and propagate the best varieties for biodiesel production. Henning has found Brazilian jatropha seeds that contain 40 percent oil—about the same as canola and more than twice the 18 percent contained in soybeans. Indonesia has a dwarf variety that is especially easy to harvest.

Finding the variety best suited to particular growing conditions is crucial, explains D1 Oils agronomy director Henk Joos, because right now not much hard scientific information exists about jatropha—just lots of stories. "We know that this plant is environmentally elastic and drought-tolerant. But the aura that this is a wonder crop that you can plant in the desert and harvest gold" is a dangerous notion that threatens social and economic sustainability, Joos says, adding that jatropha needs to be managed like any other crop. He notes that at D1 Oils plantations, farmers plant in land that is as good as possible without replacing food

crops, then apply first-rate farming practices: prune branches, apply manure and provide water.

But the realization that successful large-scale operations have to function like well-run farms raises the issue of competition with food crops for water and land, says agronomist Raymond Jongschaap of Wageningen University in the Netherlands. Jongschaap is spearheading one of the research projects looking for different types of jatropha with the goal of matching plants to growing conditions and maximizing oil yields. He has the most faith in small-scale efforts based on hedges or intercropping jatropha with other plants—a method used in projects in Kenya and Madagascar, where jatropha is planted alongside vanilla.

Henning agrees that it is smart for jatropha growers to start small. Biodiesel cannot compete with current petroleum prices, which are relatively low, so jatropha would be better suited for local projects that improve rural livelihoods and basic energy services. These small projects have already started to build a framework of familiarity and expertise—in parts of Tanzania, kids learn about jatropha in school. Then, as fuel prices increase, jatropha cultivation can go to a larger scale. The wild shrub could then become a "sustainable cash crop," Joos believes, and a fuel for the future.

Rebecca Renner is an environmental writer based in Williamsport, Pa.



NUTTY GOODNESS: Jatropha fruit may contain up to 40 percent oil, depending on the variety.

NEED TO KNOW: BIOFUEL BLUES

The jatropha plant has excited some energy experts because it does not have the same negative impacts as conventional biofuel crops. Touted as a renewable resource and a means to promote energy independence, today's biofuel plants are expensive, require intensive farming and threaten food supplies. The recent U.S. push for ethanol has already contributed to rising corn prices.

Biofuel plantations can also harm the environment. In Europe a biofuel quota backfired when it increased demand for palm oil from Southeast Asia; as a result, farmers carved new plantations out of dwindling rain forest and released millions of tons of carbon dioxide when they cut into carbon-rich peat soils.

D1 OILS PLC

PHYSICS

Atomic Fingerprinting

MICROSCOPE DISCERNs AN ATOM'S CHEMICAL IDENTITY BY LUIS MIGUEL ARIZA

Deciding whether a substance is, say, steel, brick, wood or plastic is easy—but not on the atomic scale, which lacks information about such everyday characteristics. Using an atomic-force microscope (AFM), however, an international

team of physicists has developed a method of atomic "fingerprinting" that can determine the chemical identity of individual atoms on a surface mixed with all of them.

"Until now, there was not any technique that would allow us to identify atom by atom

Get
MORE
Science.

▶ www.sciam.com

▶ www.sciamdigital.com

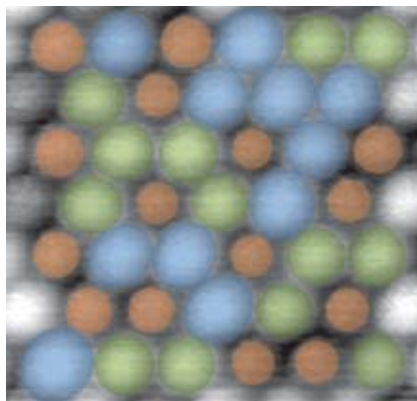
▶ www.sciammind.com

SCIENTIFIC
AMERICAN.COM

SCIENTIFIC
AMERICAN | DIGITAL

SCIENTIFIC AMERICAN
MIND

©1996-2007 Scientific American, Inc. All Rights Reserved



SEEING CHEMICALS: Physicists used an atomic-force microscope to identify the surface composition of an alloy made of silicon (red), tin (blue) and lead (green) atoms.

and see them at the same time,” says Rubén Pérez of the Autonomous University in Madrid. Using their AFM approach, Óscar Custance and his collaborators at Osaka University, along with Pérez and his colleagues and others, could discern tin, silicon and lead, which are all chemically similar. The resulting image of the atoms resembles a granulated painting, where the “grains”—the individual atoms—are distinguishable in false color.

The ability to identify and manipulate atoms first gained prominence in 1989, when IBM scientists spelled out their company logo with xenon atoms. Back then, the physicists relied on a scanning tunneling microscope (STM), which detects atoms by a slight flow of electrons between the microscope tip and an atom. The STM, however, can identify only atoms of materials that conduct electricity.

In contrast, the fingerprinting AFM technique works for conductor and insulator alike. Like a phonograph, the AFM employs an ultrafine needle mounted on a flexible cantilever. As the needle gets dragged across a surface, it jogs up and down as it encounters atoms on that surface. This oscillation actually occurs because of the attractive forces associated with the onset of chemical bonding between the silicon in the tip and the atoms on the surface.

The Japanese-Spanish team showed that the oscillation frequency depends on the atom’s chemical nature. It thus enabled the researchers to identify different atomic species even if they exist in equivalent abundance, Pérez says, “like distinguishing a tree in a noisy, fuzzy forest.”

In previous work, Custance and his team had demonstrated that they could use the AFM to move tin atoms strongly attached to a germanium surface, writing the letters “Sn” (tin’s chemical symbol). Combining the method with atomic fingerprinting opens exciting possibilities for the AFM—researchers might be able to “visualize reactions with atomic resolution,” Custance remarks. And, he adds, as microelectronics shrink into the nanoscale realm—2,000 of today’s transistors can fit across the width of a human hair—then “just by arranging a few atoms in predefined patterns, it could be possible to enhance the performance of the devices.”

Luis Miguel Ariza is based in Madrid.

PRINTERS

“Memjet” Momentum

INK-JET PRINTING AT RAMJET SPEED

BY CHARLES Q. CHOI

Stealthily, over more than a decade, a new kind of printer has been under development in Australia. The original vision was to create a printer small enough to fit inside a digital camera. Instead the research has

yielded an ink-jet printer, dubbed the Memjet, that can print color photographs up to 30 times faster than any other printer.

Ink-jet printers, which dominate the desktop market, work by spraying

ÓSCAR CUSTANCE Osaka University

NEW PODCAST

60-Second Science

Quick science that informs and entertains

SCIENTIFIC AMERICAN

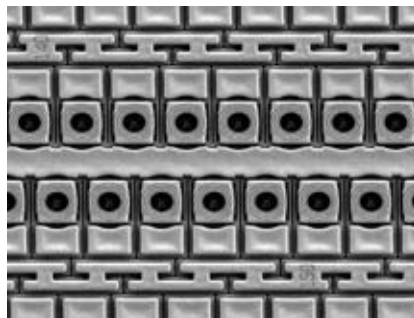
www.SciAm.com/Podcast

drops of ink from nozzles onto a sheet. Conventional ink-jets shuffle their nozzles back and forth across paper like a shuttle would between threads in a loom.

The Memjet, from Silverbrook Research in Sydney, has fixed rows of nozzles stretching from one edge of the page to the other that all fire simultaneously. Besides cutting down on moving parts, noise and vibration, the approach enables the device to finish one high-quality color page a second, roughly five times faster than consumer ink-jets, and to generate 30 four-by-six-inch color photographs a minute. "Every other printer seems glacial by comparison," says Steve Hoffenberg, an analyst for Lyra Research in Newtonville, Mass., which tracks the imaging industry.

Like many ink-jet devices, the Memjet (a name derived from microelectromechanical systems, or MEMS) relies on heating ink to a boil, which generates steam explosions that fire ink droplets out through nozzles. But whereas conventional ink-jets heat ink by using electronics located in the nozzle walls, the Memjet suspends its heaters in the fluid.

This setup cools the heaters off, which in turn allows the Memjet to pack nozzles closer together, amounting to 320 nozzles per square millimeter. That density is 17 times more than what leading printer makers can do now. Closer nozzle packing means greater efficiency in the amount of material used and lower costs. The aim is for Memjet printers to cost



NOVEL NOZZLES: The Memjet uses fixed rows of nozzles that span the entire width of a page, enabling rapid printing.



roughly as much as other consumer-grade ink-jets.

The Memjet fires droplets one-fifth the size of those sprayed in most ink-jet printers. Such picoliter volumes enable it to print in sharp detail at 1,600 dots per inch, better than the eye can see and comparable with leading printers. It would not be hard to design conventional ink-jets to print with drops as small as the Memjet's, but they would then have five times as many droplets to fire and thus "would be five times slower," explains lead researcher Kia Silverbrook. "We

have plenty of speed," he adds. The company also asserts that its device should be no more prone to clogging than existing printers.

Silverbrook unveiled its new technology March 21 at the Global Ink Jet Printing Conference in Prague. It should be available in photograph and label printers later this year and in home and office printers in 2008. And as for a printer that can be embedded inside a camera, the company says that, too, still remains a possibility.

Charles Q. Choi is a frequent contributor.

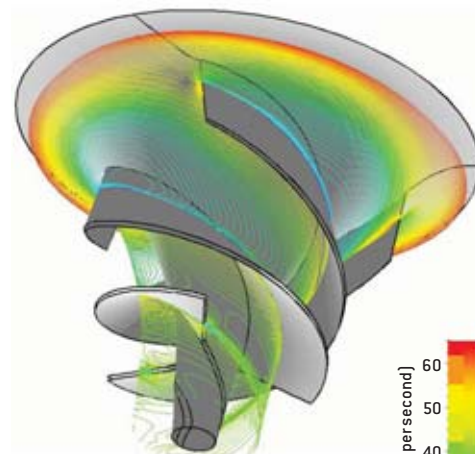
TECNOLOGY

A Good Turn

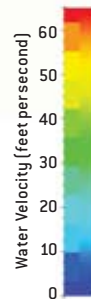
A FISH-FRIENDLY HYDROELECTRIC TURBINE GETS NEW LIFE BY MADELINE BODIN

Nearly two years ago Alden Research Laboratory in Holden, Mass., hauled the scale model of a promising hydro-power turbine out of its massive test flume and set it in a dim corner of the company's hydraulics laboratory building. As an innovation developed in the 1990s, the device proved quite promising in reducing one of hydropower's drawbacks: the turbines kill creatures that pass through them. The novel design enabled at least 98 percent of fish to survive. But orphaned by federal budget cuts, it has sat gathering dust. Now a new push has begun to retool the turbine for potential commercial use.

Conventional turbines, which resemble the blades of an electric fan, kill as many as 40 percent of the fish that are swept through them. Working from U.S. Department of Energy funds first granted in 1994, Alden Lab teamed up with Concepts NREC in White River Junction, Vt., to develop a fish-friendly turbine. The design features three rotor blades wrapped around a conical hub to create a kind of helix. A rotating case covers the rotor blades, so that only a fraction of their edges are exposed. The turbine has no gaps between the blades, fewer blades and a slower spin rate. All these features lower the chance of a fish being injured by moving parts. Moreover, the flow of water through the turbine is smooth, creating less potentially harmful shearing force.



FISH-FRIENDLY TURBINE relies on sheathed rotor blades that wrap around the hub to form a kind of helix [top]. The design keeps water flowing smoothly and at a constant rate through most of the turbine [bottom].



ALDEN RESEARCH LABORATORY, INC.

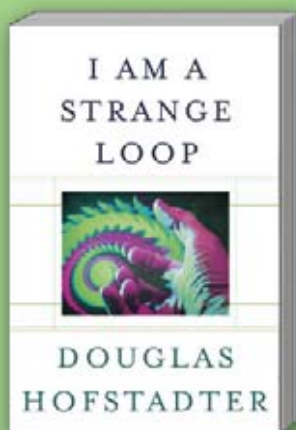
A TURBINE DUSTED OFF

The scale model of the fish-friendly Alden/Concepts NREC turbine had been mothballed until this past February, when Canadian energy firm Brookfield Power began installing the device at its School Street hydropower facility in Cohoes, N.Y. There the turbine will be located in what was once a spillway that allowed fish safe passage around the dam. It probably will not generate a lot of power, but nonetheless Brookfield considers the installation a win-win situation for the firm and the fish.

The hardest part, says Dave Culligan, manager of licensing for Brookfield's U.S. Development Group, was getting state and federal wildlife agencies "over the philosophical hurdle of letting fish go through a turbine."

**“VINTAGE
HOFSTADTER:
earnest, deep, overflowing
with ideas, building its
argument into the
experience of reading it”**

— *Los Angeles Times Book Review*



**“Fascinating...original and
thought-provoking.”**

— *Wall Street Journal*

“Nearly thirty years after his best-selling book *Gödel, Escher, Bach*, cognitive scientist and polymath Douglas Hofstadter has returned to his extraordinary theory of self in his latest book.”

— *New Scientist*



Also by Douglas Hofstadter

Gödel, Escher, Bach is an entire humanistic education between the covers of a single book.”

— John L. Casti, *Nature*

BASIC BOOKS
www.basicbooks.com

The design is particularly friendly to eel and sturgeon—long, thin swimmers whose populations are declining and for whom traditional hydropower turbines are decidedly deadly. According to initial tests, all the eel and sturgeon that pass through the new turbine can be expected to survive.

But about three years ago federal funding for the turbine was cut completely. Research stalled, and eventually the scale model was put in storage.

The Electric Power Research Institute (EPRI), based in Palo Alto, Calif., has taken up the fish-friendly turbine's cause. It has raised about \$300,000 of the \$500,000 needed to make the turbine commercially feasible and to replace conventional models as hydroelectric plants are renovated. The first step will be to increase the new turbine's power output so that it can compete with existing turbines, comments Doug Dixon, senior project manager at EPRI. The original design generated half the power of a commercial turbine of the same size. “The more efficient it is, the more attractive it will be to industry,” Dixon says.

Boosting the power output, though,

will be tough. “As a rule, what's good for engineering is not good for fish, and vice versa,” says Ned Taft, president of Alden Lab. The solution, Alden engineers believe, is to increase the amount of water flowing through the turbine. They recently figured out a way to double the volume of the spiral pipe that feeds water into the blades while increasing the turbine's diameter by only 2 percent. Besides reshaping and reangling the blades to handle the increased flow, Concepts NREC plans to broaden the leading edges of the blade to a width that is closer to the length of the typical fish moving through the turbine—a size correlation that seems to boost survival.

The design is on the brink of commercial viability, Dixon believes. “But we're struggling to find the resources to make it happen.” Supporters say that more is at stake than the success of this particular turbine design. “A lot of people don't see hydro as green power,” Taft states. “This turbine could change that.”

— *Madeline Bodin is based in Andover, Vt.*

DEFENSE SYSTEMS

Beam Weapons Get Real

SOLID-STATE LASERS NEAR BATTLEFIELD DEPLOYMENT **BY STEVEN ASHLEY**

The ray gun, long a mainstay of science-fiction tales, may actually enter the American war-fighting arsenal in a few years. Engineers at several defense firms have conducted successful tests of key prototype components of truck-size “laser cannon” systems capable of firing a beam from aircraft, naval ships or armored vehicles to zap targets many kilometers away—even through intervening dust or fog.

High-power lasers—measured in

hundreds to thousands of kilowatts—offer several advantages over conventional projectile weapons, according to Mark Neice, director of the Department of Defense's High-Energy Laser Joint Technology Office in Albuquerque, N.M. “They could provide ultra-precise, speed-of-light strike capabilities that could leave little or no collateral damage,” he says.

Although overly optimistic predictions in the past have led skeptics to quip that “lasers are the weapons of the

future, and they always will be,” this time beam weapons finally look real. “Industry is just on the cusp of delivering practical directed-energy weapons for offensive and defensive military purposes,” Neice states.

In the past year or so, with funding from the U.S. Air Force, Army and Navy, researchers at Northrop Grumman, Textron, Raytheon and Lawrence Livermore National Laboratory have achieved marked progress on the solid-state bulk laser, which runs directly on electricity. Coupled to a ground vehicle’s electric generator, fuel cell or battery bank, a solid-state laser with an average power capacity greater than 100 kilowatts would feature a nearly “infinite magazine” of low-cost shots with which to blast incoming mortars, artillery shells, rockets and missiles in flight from five to eight kilometers away. Such a system could also blind electro-optical and infrared battlefield sensors and enable troops to neutralize mines and improvised explosive devices from a safe distance.

The key to this high-energy device is the gain medium, the material that amplifies the laser photons. In laser diodes used in DVD players and other consumer electronics, semiconductor layers amplify light after a jump start from an electrical charge. In the solid-state bulk laser, the gain medium takes the form of a few-centimeter square (or rectangular) slab, explains Jackie Gish, Northrop Grumman Space Technology’s director of directed-energy technology and products. The slabs consist of a tough ceramic material, such as yttrium-aluminum-garnet doped with the rare-earth element neodymium. But instead of an electrical priming, stacks of laser diodes optically pump, or excite, the gain medium. In general, the larger the slab, the greater the power output.

Each research team has different ways to link several slabs together to create laser “chains” that produce high power levels in the tens of kilowatts, says John Boness, vice president of applied technology at Textron Systems. In the coming year, engineers expect to gang these chains coherently in series or in parallel to achieve 100 kilo-

watts, the so-called entry-level average power benchmark for military laser applications. Other key performance targets, Neice says, are an operational run time of 300 seconds (sufficient for multiple laser shots), an electrical-to-optical system energy conversion efficiency of 17 percent or more, and, crucially, adequate “beam quality” (essentially, focusing) to ensure enough photons reach the target to heat its exterior sufficiently to



Kuwait University

Faculty of Science
Kuwait

The Department of **Biological Sciences** at the Faculty of Science of Kuwait University invites applications for appointment at the rank of Full Professor, starting September 2007 in the following area:

Molecular Biology
(Human genetics with special interest in cytogenetics and/or molecular genetics)

Required Qualifications:

- Ph.D. degree in the area of specialization from a reputable university
- The applicant’s bachelor’s degree GPA should be 3 out of 4 or equivalent
- Research experience and publications in refereed international journals
- Full command of teaching in English
- University teaching experience in the specified field

Benefits include attractive tax-free salary according to rank and teaching experience (Professor’s monthly salary varies from 2950 to 3192 KD., 1 KD. = \$3.40), annual air tickets for the faculty member and his/her family (spouse and up to three children under the age of 20), a one time settling-in allowance, housing allowance, free national health medical care, paid mid-term holidays and summer vacations, and end-of-contract gratuity. The University also offers an excellent academic environment and financial support for research projects.

To apply send by express mail/courier service or e-mail, within six weeks of the date of announcement, a completed application form, updated curriculum vitae (including mailing address, phone and fax numbers, e-mail address, academic qualifications, teaching and research experience, and a list of publications in professional journals), three copies of Ph.D., Master and Bachelor certificates and transcripts (an English translation of all documents in other languages should be enclosed), a copy of your passport, three letters of recommendation, and names and addresses of three persons well-acquainted with your academic and professional work to the following address:

The Dean,
Faculty of Science
Kuwait University, P.O. Box 5969,
Safat, 13060, Kuwait

For inquiries:
Fax: +965 4836127
E-mail: biosc@kuc01.kuniv.edu.kw

LETHAL REACTIONS

The U.S. military has already developed one kind of mighty laser—a directed-energy device powered by chemical reactions. Much more potent than their solid-state cousins, these megawatt-class devices are called chemical oxygen iodine lasers (COILs). But they are large and can run only as long as their stores of reactants hold out. Nevertheless, defense contractors are now preparing to install COILs on aircraft. A Boeing 747 airliner will host a COIL system—called the YAL-1A Airborne Laser—for standoff strikes against air vehicles, especially ballistic missiles, and an AC-130 cargo plane-cum-gunship will wield the Advanced Tactical Laser, which is intended for precision air-to-ground attacks.



SOLID-STATE LASER is tested at Northrop Grumman Space Technology. A high-power version of such a device—capable of generating 100 kilowatts or more—could enable ground troops to destroy incoming mortars or rockets.

destroy it, detonate it or send it off course.

Assuming that slab lasers meet these goals, their making it into actual combat depends on their successful integration into a functioning weapons system that is compact enough to fit in a vehicle, Boness explains. Besides a renewable electrical power source of 1,000 kilowatts or more, a solid-state bulk laser weapon needs a thermal chiller to keep the slabs from heating up enough to distort the beam.

It would also require a beam director to place the photons on target—probably a large, mobile mirror equipped with adaptive or deformable optics to compensate for atmospheric distortion, which would be detected by a low-power “sensing” laser beam. Finally, aiming such a system would rely on a radar-based or optical cueing system to find and track the intended target.

Effective beam weapons would initiate nothing short of a revolution in warfare. But packing all that technology into something the size of Captain Kirk’s handheld phaser still clearly lies in the realm of science fiction.

NORTHROP GRUMMAN PHOTO

We don't know when childhood ends,



Retirement. Time once again to play. With a full family of investment options, The Hartford could help make it all you've hoped for. Call your broker, or visit hartfordinvestor.com. Prepare to Live™

•Mutual Funds •401(k)
•Annuities •Life Insurance

Fish That Go Skin-Deep

news

SCAN

TRAPPED FISH ADAPT TO A LIFE OF NIBBLING ON HUMANS BY MATT MOSSMAN

Kangal, Turkey—Tucked between brown hills in central Turkey is a natural hot spring where, for a fee, you can become fish food. Dip in a hand or foot, and within seconds small fish will swarm, bump and nibble it. Stand above the pools, and the fish will gather below, waiting. The scaly swimmers—the “Doctor Fish of Kangal”—supposedly have curative powers. But in this unusual case of adaptive ecology, the human visitors may be helping the fish more than themselves.

These fish have acquired a taste for humans largely because they have little choice. The spring is too hot to sustain enough algae and plankton to feed them all. In the past, the fish were able to move between the spring and a creek that runs nearby. But after learning of a story about a local shepherd whose wounded leg healed after being dipped into the spring in 1917, builders

walled off the spring from the creek in the 1950s to preserve a captive school. A Turkish family has now constructed a hotel, villas and a playground and markets the resort to psoriasis patients. Some 3,000 people every year pay for the privilege of sitting in the spring and allowing these omnivores to eat their dead skin, a process that may stimulate new skin growth or relax patients and thereby ease stress-triggered psoriasis.

Unquestionably for the fish, “the human skin is a big help,” remarks Fevzi Bardakci, a biologist at Turkey’s Adnan Menderes University. “It’s like meat for them.” In 2000 Bardakci published a paper in the *World Wide Web Journal of Biology* on *Garra rufa*, one of the two species found in the hot spring. He discovered that members of the same species that swim in a nearby creek grow to an average 97 millimeters and about 11 grams. In the hot spring, the fish are three-quarters



but it starts again at retirement.

Official Corporate Provider of the NCAA® NCAA® is an equal opportunity organization. "You are here" is the National Collegiate Athletic Association's slogan. © 2007 The Associated Press. All rights reserved. NCAA® is a registered trademark of the National Collegiate Athletic Association.





BITE ME: Fish in hot-spring pools at a Turkish resort survive by feasting on human skin. The water does not sustain enough plankton and algae that would otherwise serve as food.

the length and weigh one quarter as much. Moreover, during the summer spawning season, the trapped females grow fewer and smaller oocytes, the cells that develop into eggs. In the creek, the gonads balloon from 3 percent of body weight to almost 8 percent. In the hot spring, organs increase from 1 percent of body weight to 2 percent. They would grow even less without submerged skin to nibble, Bardakci concluded. Some 90 percent of visitors arrive in the summer, providing a nutritional supplement at the perfect time.

The fish come from the carp and minnow family, which is known for adaptabil-

ity, says Richard Londraville, a biologist at the University of Akron. He adds that those in the hot spring may eventually evolve into a separate species in a few thousand years.

Other fish survive in waters as hot as or hotter than this spring, which hovers near 34 degrees Celsius. None are widely known to feed on skin, which may explain why *G. rufa* are catching on elsewhere. A Chinese spa-building company, which claims to have invented the concept, says on its Web site that it trained its own doctor fish and built 10 spas in China, including in Beijing. Some of the Turkish fish were scooped up and enlisted in springs in Japan, where at several spas they are now also performing fish pedicures.

Matt Mossman is based in Istanbul.

MATT MOSSMAN

COMPUTING

Silicon Smackdown

NEW GO ALGORITHM AIMS TO DEPOSE HUMANS BY KAREN A. FRENKEL

A decade ago IBM's chess program, Deep Blue, beat world champion Garry Kasparov in a six-game match. The event marked a milestone, forcing humans to yield dominance of yet another strategic diversion. Only the Asian board game Go seemed to be computer science's Achilles' heel: humans could soundly beat the machines. A new algorithm can now take on strong human players—and win.

Go has proved enormously difficult for computer programmers because of the game's deceptive complexity. The objective of Go is to stake out territory and surround

an opponent by placing black or white stones on the intersections of a nine-by-nine or 19-by-19 line grid. Especially on the large board, the number of possible moves per turn is huge—200 on average for each midgame position compared with the several dozen possible in chess. There are also enormous branching factors. Given N positions on the board, the total number of possible game positions is 3^N , because every position can be occupied by a black or white piece, or it can be empty. The total number of legal positions on the small board is about 10^{38} ; on the large board, about 10^{170} . Additionally,

more stones do not ensure victory, and players must be able to consider local positions and the board as a whole.

To cope with such an enormous number of options, artificial-intelligence experts have designed algorithms to limit searches, but the programs have not been able to beat the better human players on large boards. Last fall two Hungarian researchers reported that their algorithm outdid the win rates of the best Go programs by 5 percent and could compete with professional Go players on small boards. Levente Kocsis of the Computer and Automation Research Institute at the Hungarian Academy of Sciences in Budapest and Csaba Szepesvári, now at the University of Alberta in Edmonton, developed the algorithm, called UCT (for *u*pper *c*onfidence bounds applied to *t*rees). It extends the well-known Monte Carlo method.

First incorporated into Go programs in the 1970s, Monte Carlo works like a political poll: it performs statistical sampling to predict the behavior or characteristics of a large group. When applied to Go, the algorithm evaluates and ranks candidate moves by playing a large number of random games. But playing the move with the highest score in each position does not guarantee that the player will win the game. Instead this type of search merely restricts

the number of relevant potential moves.

UCT takes Monte Carlo further by focusing the search on the most promising moves. “The main idea is to sample actions selectively,” Kocsis says. The algorithm must strike a balance, testing alternatives that look the best at the moment to find possible weaknesses and exploring “less optimal-looking alternatives, to ensure that no good alternatives are missed because of early estimation errors,” he explains.

UCT calculates indices for moves and selects the move that has the highest index. The algorithm computes the index from the win rate, which describes how often that position leads to a win, as well as the number of times the position has been visited but not played. UCT grows a decision tree in memory and uses it to track these statistics. When UCT encounters a move that has not been visited previously, it adds the move to the tree and plays the rest of the game randomly.

UCT decides if the finished random game is a win or loss, then updates the statistics of all moves made during the game. If the index equals the win rate of the move, the algorithm quickly focuses on the most promising path. But nothing guarantees that an initially successful path will eventually yield a winning move. So when selecting moves, UCT inflates win rates by weighting less visited candidate moves more heavily. The researchers borrowed this idea from bandit problems—selective weighting yields the maximum gain for a gambler playing several slot machines with unknown average payoffs.

Mathematicians Sylvain Gelly of the University of Paris-South and Yizao Wang of the Polytechnic School outside Paris have incorporated UCT into a program they call MoGo. It has a 95 percent better win rate as compared with a previous state-of-the-art Monte-Carlo extension algorithm. Now the top-ranked Go algorithm, MoGo demonstrated its abilities this past spring, vanquishing strong amateur players on nine-by-nine boards and beating weaker ones on large boards. Gelly says that UCT is simple to implement and can be improved. So turning the ultimate corner—ending the reign of professional human Go players—could occur in 10 years, Kocsis states.

Karen A. Frenkel is based in New York City.

GO BEYOND GAMES

The new Go-playing algorithm, called *upper confidence bounds applied to trees* (UCT), is not limited to games. It applies to any problem that involves choosing the best option, as long as alternatives have an internal treelike structure (that is, a cascading set of choices) and their values can be recursively computed. UCT may prove useful for targeting advertisements on the Web, finding the best settings for an industrial plant or optimizing channel allocation in cellular systems.



GO ATTACK: Humans still rule when it comes to Go, but a new algorithm can topple strong players.



DATA POINTS: COLLISION DECISION

Design flaws caused a support structure for a magnet to fail during a stress test of the Large Hadron Collider (LHC) on March 27.

Built to be the world's most powerful particle accelerator, the LHC will smash lead ions into one another with energies of trillions of electron volts (TeV). Repairs may delay the LHC's target start date of November.

Circumference of LHC
in meters: **26,659**

Energy imparted to each
proton in TeV: **7**

Total collision energy of lead-ion
beams (consisting of many
protons) in TeV: **1,150**

Number of particles in each beam:
300 trillion

Odds that two particles will
actually collide: **1 in 10 billion**

Approximate power consumption
in watts of:

A household microwave oven:
1,100

The LHC: **120 million**

Time it takes to defrost a pizza in:

A household microwave oven:
6 minutes

The LHC: **7.6 femtoseconds**
(7.6×10^{-15} second)

SOURCES: CERN; DiGiorno's Microwave Rising Crust Four-Cheese Pizza. Pizza calculation is theoretical and assumes that the collision energy can be evenly distributed over the pizza.

NEUROSCIENCE

Brain Damage for Easier Moral Choices

How long would you hesitate before pushing someone in front of a runaway train to keep it from killing five other people? The answer may be no time at all, if you have damage to the ventromedial prefrontal cortex (VMPC)—a region in the forebrain associated with emotional response. Researchers confronted volunteers with such scenarios and found that those with VMPC injury were three times more likely than healthy people to advocate throwing the person to certain death for the good of the

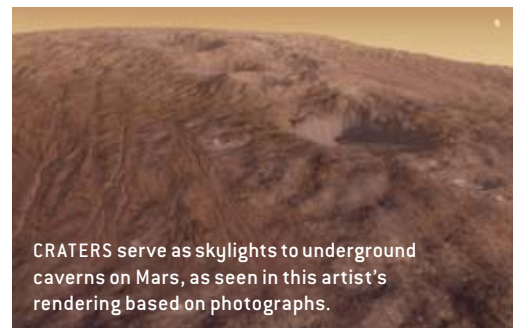
many. In a similar scenario, VMPC patients were five times more likely to advocate smothering one's baby to save others. Senior author Antonio Damasio of the University of Southern California says that the patients are not amoral but seem to lack the natural conflict between emotion and reason. The study, in the March 22 *Nature*, also shows that such decisions result not from a single moral faculty but from two different processes that can compete with each other.

—Nikhil Swaminathan

PLANETARY SCIENCE

Martian Cave Dwellings

Seven football field–size caves may have been discovered on Mars. Analysis of photographs from NASA's Mars Odyssey orbiter revealed black spots near the massive Martian volcano Arsia Mons that do not look like impact craters because they lack blast patterns and raised rims. Scientists at Northern Arizona University and their colleagues say the possible caverns range from 330 to 825 feet wide and are 425 feet deep and have named them after their loved ones: Dena, Chloe, Wendy, Annie, Abbey, Nikki and Jeanne. Caves would serve as havens from radiation on the surface and so would be the most likely areas to harbor life. They could also accumulate ice, which could help to support future human exploration. NASA's Mars Reconnaissance Orbiter could take sidelong glances at the putative caves, a view that might show whether wider chambers exist underneath. The findings were unveiled during a March meeting of the Lunar and Planetary Science Conference in League City, Tex. —Charles Q. Choi



CRATERS serve as skylights to underground caverns on Mars, as seen in this artist's rendering based on photographs.

EPIDEMICS

Stick It to the Kids

When flu epidemics loom, the long-standing recommendation of the U.S. Centers for Disease Control and Prevention is to vaccinate the elderly first, because they are at greater risk of dying if they contract the virus. New evidence suggests that youngsters should take priority. Work at Yale University and Rutgers University underlines that children are the group most responsible for spreading the flu: they carry the virus into the home and infect adults, who then bring the flu into the workplace. Vaccinating most young people would virtually eliminate the flu, the researchers calculate, thereby cutting down on the mortality of the elderly, the young and people overall. The details are in the March 27 *Proceedings of the National Academy of Sciences USA*. —Charles Q. Choi



TARGETED for flu shots: children are the main vectors of the illness.

BRIEF
POINTS

■ **Blood for all:** researchers have discovered two bacterial enzymes that can efficiently remove immune-triggering sugar molecules from red blood cells, thereby turning A, B and AB blood types into the universal type O.

Nature Biotechnology, April

■ **New car smell**—a stew of volatile organic compounds—is not toxic, at least not to human cells in culture. The chemicals did aggravate the cells' immune response, suggesting that people with allergies should beware.

Environmental Science & Technology, April 1

■ **Lower-dose chemo:** a fast-screening process using small interfering RNA molecules has uncovered 87 genes that affect chemotherapy sensitivity; silencing some of those genes, for instance, made lung cancer cells 10,000 times more sensitive to the drug Taxol.

Nature, April 12

■ **Planets of other solar systems** might have plants whose dominant colors are yellow, red or even "infrared." NASA scientists say they can predict the foliage color by the type of light emitted by the parent star.

Astrobiology, March

PALEONTOLOGY

Collagen from *T. Rex* 

Researchers have extracted collagen protein from a 68-million-year-old *Tyrannosaurus rex* femur, which two years ago was revealed to have soft tissue. Chemical analysis of the protein yielded seven sequences of about 10 to 20 amino acids in length. Three sequences matched collagen peptide scripts from chickens, one matched a frog and another a salamander; the other two matched multiple organisms, including chickens and salamanders. The results strengthen the bird-dinosaur connection and jettison the belief that ancient fossils could not provide protein samples for study (sorry, *Jurassic Park* fans—any genetic material degraded long ago). The findings and the techniques used to uncover them should clarify the relations between extinct species and modern-day animals and reveal more about “patterns of molecular change and the rates and directions of molecular evolution,” says Mary Schweitzer of North Carolina State University, who reports the work with her colleagues in the April 13 *Science*.

—Nikhil Swaminathan



BEHAVIOR

Brain Brakes

So what stops you from pressing that send button for an e-mail that tells off your boss? Three distant brain regions connected by “hyperdirect cables,” believe scientists at the University of California, San Diego. They asked participants to plan an action, listen for a stop signal and decide whether to obey or continue as planned. Brain scans revealed that a neural braking network kicks in for a few milliseconds—just long enough for participants to make a decision. The inferior frontal cortex sends the braking signal to the midbrain’s subthalamic nucleus, which stops motor movement; a third region, the presupplementary motor area, initiates the plan to halt or continue the action. No synapses lie between the areas, enabling direct and fast communication. Understanding this network could explain neurological disorders such as stuttering, which may arise from the brain’s inability to coordinate stop signals. The work appears in the April 4 *Journal of Neuroscience*. —Thania Benios

VISION

Rodent Roy G. Biv 

Mice, like most mammals, normally view the world in colors limited to yellows, blues and grays, similar to what people with red-green color blindness see. By introducing a single human gene into mice, scientists endowed the animals with full-color vision. Humans and closely related primates possess an extra light-sensitive pigment permitting them to see red. (Color-seeing mammals have at least two pigments, for blue and green.) Researchers at the University of California, Santa Barbara, and their colleagues inserted the gene for this extra pigment into the mouse X chromosome. Even though the rodent brains had not evolved to use these signals, they were able to rewire themselves to handle the upgrade, correctly discriminating between colored lights to win soy-milk rewards. The investigators, who detail the work in the March 23 *Science*, say this finding could help explain how color vision evolved in humans.

—Charles Q. Choi



News Scan briefs with this icon have extended coverage at www.sciam.com/ontheweb



The (Other) Secret

The inverse square law trumps the law of attraction By MICHAEL SHERMER

An old yarn about a classic marketing con game on the secret of wealth instructs you to write a book about how to make a lot of money and sell it through the mail. When your marks receive the book, they discover the secret—write a book about how to make a lot of money and sell it through the mail.

A confidence scheme similar to this can be found in *The Secret* (Simon & Schuster, 2006), a book and DVD by Rhonda Byrne and a cadre of self-help gurus that, thanks to Oprah Winfrey's endorsement, have now sold more than three million copies combined. The secret is the so-called law of attraction. Like attracts like. Positive thoughts sally forth from your body as magnetic energy, then return in the form of whatever it was you were thinking about. Such as money. "The only reason any person does not have enough money is because they are blocking money from coming to them with their thoughts," we are told. Damn those poor Kenyans. If only they weren't such pessimistic sourpusses. The film's promotional trailer is filled with such vainglorious money mantras as "Everything I touch turns to gold," "I am a money magnet," and, my favorite, "There is more money being printed for me right now." Where? Kinko's?

A pantheon of shiny, happy people assures viewers that *The Secret* is grounded in science: "It has been proven scientifically that a positive thought is hundreds of times more powerful than a negative thought." No, it hasn't. "Our physiology creates disease to give us feedback, to let us know we have an imbalanced perspective, and we're not loving and we're not grateful." Those ungrateful cancer patients. "You've got enough power in your body to illuminate a whole city for nearly a week." Sure, if you convert your body's hydrogen into energy through nuclear fission. "Thoughts are sending out that magnetic signal that is drawing the parallel back to you." But in magnets, opposites attract—positive is attracted to negative. "Every thought has a frequency.... If you are thinking that thought over and over again you are emitting that frequency."

The brain does produce electrical activity from the ion currents flowing among neurons during synaptic transmission,

and in accordance with Maxwell's equations any electric current produces a magnetic field. But as neuroscientist Russell A. Poldrack of the University of California, Los Angeles, explained to me, these fields are minuscule and can be measured only by using an extremely sensitive superconducting quantum interference device (SQUID) in a room heavily shielded against outside magnetic sources. Plus, remember the inverse square law: the intensity of an energy wave radiating from a source is inversely proportional to the square of the distance from that source. An object twice as far away from the source of energy as another object of the same size receives only one-fourth the energy that the closer object receives. The brain's magnetic field of 10^{-15} tesla quickly dissipates from the skull and is promptly swamped by other magnetic sources, not to mention the earth's magnetic field of 10^{-5} tesla, which overpowers it by 10 orders of magnitude!

Ceteris paribus, it is undoubtedly better to think positive thoughts than negative ones. But in the real world, all other things are never equal, no matter how sanguine your outlook. Just ask the survivors of Auschwitz. If the law of attraction is true, then the Jews—along with the butchered Turkish-Armenians, the raped Nanking Chinese, the massacred Native Americans and the enslaved African-Americans—had it coming. The latter exemplar is especially poignant given Oprah's backing of *The Secret* on her Web site: "The energy you put into the world—both good and bad—is exactly what comes back to you. This means *you* create the circumstances of your life with the choices you make every day." Africans created the circumstances for Europeans to enslave them?

Oprah, please, withdraw your support of this risible twaddle—as you did when you discovered that James Frey's memoir was a million little lies—and tell your vast following that prosperity comes from a good dollop of hard work and creative thinking, the way you did it. ■

Michael Shermer is publisher of Skeptic (www.skeptic.com). His new book is Why Darwin Matters.

A pantheon of shiny, happy people assures viewers that *The Secret* is grounded in science.

Going beyond X and Y

Babies born with mixed sex organs often get immediate surgery. New genetic studies, Eric Vilain says, should force a rethinking about sex assignment and gender identity **By SALLY LEHRMAN**

When Eric Vilain began his medical school rotation two decades ago, he was assigned to France's reference center for babies with ambiguous genitalia. He watched as doctors at the Paris hospital would check an infant's endowment and quickly decide: boy or girl. Their own discomfort and social beliefs seemed to drive the choice, the young Vilain observed with shock. "I kept asking, 'How do you know?'" he recalls. After all, a baby's genitals might not match the reproductive organs inside.

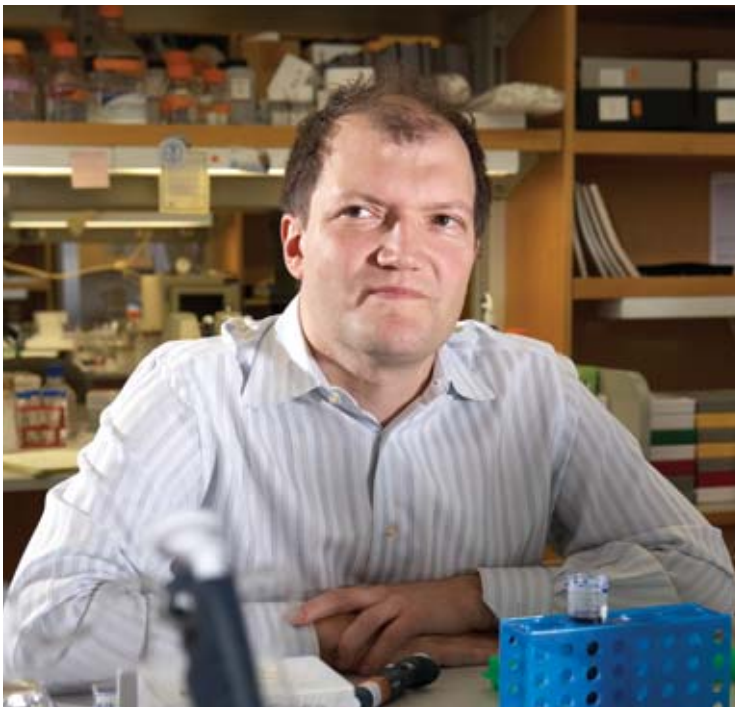
By coincidence, Vilain was also reading the journals of Herculine Barbin, a 19th-century hermaphrodite. Her story of love and woe, edited by famed social constructionist Michel Foucault, sharpened his questions. He set on a path to find out what sexual "normality" really meant—and to find answers to the basic biology of sex differences.

Today the 40-year-old French native is one of a handful of geneticists on whom parents and doctors rely to explain how and why sex determination in an infant may have taken an unusual route. In his genetics laboratory at the University of California, Los Angeles, Vilain's findings have pushed the field toward not only improved technical understanding but more thoughtful treatment as well. "What really matters is what people feel they are in terms of gender, not what their family or doctors think they should be," Vilain says. Genital ambiguity occurs in an estimated one in 4,500 births, and problems such as undescended testes happen in one in 100. Altogether, hospitals across the U.S. perform about five sex-assignment surgeries every day.

Some of Vilain's work has helped topple ancient ideas about sex determination that lingered until very recently. Students have long learned in developmental biology that the male path of sex development is "active," driven by the presence of a Y chromosome. In contrast, the female pathway is passive, a default route. French physiologist Alfred Jost seemed to prove this idea in experiments done in the 1940s, in which castrated rabbit embryos developed into females.

In 1990, while at the University of Cambridge, Peter Goodfellow discovered *SRY*, a gene on the Y chromosome hailed as the "master switch." Just one base pair change in this sequence would produce a female instead of a male. And when researchers integrated *SRY* into a mouse that was otherwise chromosomally female, an XX fetus developed as a male.

But studies by Vilain and others have shaped a more complex picture. Instead of turning on male de-



ERIC VILAIN: SEX BIOLOGY IN THE MIX

- Studies genetics of sex determination; advises on intersex diagnoses.
- Birth rate in which genitals are ambiguous: 1 in 4,500.
- Discourages hasty sex-assignment surgery based on distress of parents or physicians: "That should appropriately be treated by a psychologist."

velopment directly, *SRY* works by blocking an “antitestis” gene, he proposes. For one, males who have *SRY* but two female chromosomes range in characteristics from normal male to an ambiguous mix. In addition, test-tube studies have found that *SRY* can repress gene transcription, indicating that it operates through interference. Finally, in 1994, Vilain’s group showed that a male could develop without the gene. Vilain offers a model in which sex emerges out of a delicate dance between a variety of promale, antimale, and possibly profemale genes.

Because researchers have long viewed the development of females as a default pathway, the study of profemale genes has taken a backseat. Over the past few years, though, geneticists have uncovered evidence for active female determination. *DAX1*, on the X chromosome, seems to start up the female pathway while inhibiting testis formation—unless the gene has already been blocked by *SRY*. With too much *DAX1*, a person with the XY complement is born a female. Vilain’s group found that another gene, *WNT4*, operates in a similar way to promote the formation of a female. The researchers discovered that these two work together against *SRY* and other promale factors. “Ovary formation may be just as coordinated as testis determination, consistent with the existence of an ‘ovarian switch,’” report geneticist David Schlessinger and his collaborators in a 2006 review in the journal *Bioessays*.

Lately Vilain has been exploring molecular determinants of sex within the brain and whether they may be linked to gender identity. Despite classic dogma, he is certain that sex hormones do not drive neural development and behavioral differences on their own. *SRY* is expressed in the brain, he points out, suggesting that genes influence brain sexual differentiation directly. His lab has identified in mice 50 new gene candidates on multiple chromosomes for differential sex expression. Seven of them begin operating differently in the brain before gonads form. Vilain’s group is testing these findings using mice and is collaborating with a clinic in Australia to study expression patterns of the sex-specific genes in transsexual people.

This work, like much of Vilain’s efforts, treads on fairly touchy ground. He copes by sticking to his findings conservatively. “You also have to be aware of the social sensibilities,” he explains. Accordingly, he has come to agree with some gender activists that it is time to revamp the vocabulary used to describe ambiguously sexed babies.

At the 2005 Intersex Consensus Meeting in Chicago, he stood before a group of 50 geneticists, surgeons, psychologists and other specialists and argued that terms such as “hermaphrodite,” male or female “pseudohermaphrodite” and “intersex” were vague and hurtful. Instead of focusing on a newborn’s confusing mix of genitals and gonads, he urged his

colleagues to let the explosion of new genetic findings point toward a more scientific approach. Rather than using “hermaphrodite,” for instance, he recommended referring to a “disorder of sexual development” (DSD) and applying the more precise term of “ovotesticular DSD.”

Although the attendees eventually concurred, not everyone likes the new terminology. Some who prefer “intersex” feel that a “disorder” is demeaning. Milton Diamond, who studies sex identity at the University of Hawaii, complains that it stigmatizes people who have nothing wrong with their bodies.

But the decision to change nomenclature realizes a 15-year dream for Cheryl Chase, executive director of the Intersex Society of North America (ISNA). Chase has fought for years against secret, rushed surgeries intended to comfort parents and adjust anatomy to match an assigned social gender. Recalling how a doctor once called her “formerly intersex,” she hopes physicians will begin to see mixed sex characteristics as a lifelong medical condition instead of a problem to be quickly fixed. “Now that we’ve accomplished the name change, culture can accomplish a little magic for us,” she predicts.

For her, Vilain has been a valued ally in the process as a member of ISNA’s medical advisory board. The job, he admits, forces him to listen to patients, a practice he considers unusual for the field. He expects the new, medicalized terminology for DSD to have what he describes wryly as “an interesting side effect,” in which “medical science should apply” to clinical decisions about ambiguous sex.

Indeed, the new consensus statement on managing intersex disorders encourages physicians to see beyond a patient’s sex organs, agrees conference co-organizer Peter A. Lee. The statement, released last fall, recommends speedy gender assignment but a more cautious approach to surgery. The family should participate in decision making, along with a multidisciplinary team of caregivers in specialties that include psychology and ethics. But Lee, a pediatric endocrinologist at the Penn State College of Medicine, cautions that much more work lies ahead to fill in data gaps. For instance, physicians have not measured how their choices affect patients over a lifetime.

On one Friday afternoon Vilain’s white coat and stethoscope lay tossed amid the papers on his desk, a reminder that his discoveries have more than philosophical meaning. He sees six to eight patients in the U.C.L.A. intersex clinic every month, and in his on-call capacity, he receives two calls about babies in the hospital within the space of a couple of hours. Even while immersed in the workings of DNA transcription, Vilain stays grounded in what his findings mean for people’s lives. ■

Sally Lehrman is based in the San Francisco Bay Area.



Climate Change Refugees

As global warming tightens the availability of water, prepare for a torrent of forced migrations By JEFFREY D. SACHS

Human-induced climate and hydrological change is likely to make many parts of the world uninhabitable, or at least uneconomic. Over the course of a few decades, if not sooner, hundreds of millions of people may be compelled to relocate because of environmental pressures.

To a significant extent, water will be the most important determinant of these population movements. Dramatic alterations in the relation between water and society will be widespread, as emphasized in the new report from Working Group II of the Intergovernmental Panel on Climate Change. These shifts may include rising sea levels, stronger tropical cyclones, the loss of soil moisture under higher temperatures, more intense precipitation and flooding, more frequent droughts, the melting of glaciers and the changing seasonality of snowmelt.

Impacts will vary widely across the world. It will be important to keep our eye on at least four zones: low-lying coastal settlements, farm regions dependent on rivers fed by snowmelt and glacier melt, subhumid and arid regions, and humid areas in Southeast Asia vulnerable to changes in monsoon patterns.

A significant rise in sea levels, even by a fraction of a meter, could wreak havoc on tens or even hundreds of millions of people. One study found that although coastal areas less than 10 meters above sea level constitute only 2 percent of the world's land, they contain 10 percent of its population. These coastal zones are vulnerable to storm surges and increased intensity of tropical cyclones—call it the New Orleans Effect.

Regions much farther inland will wither. Hundreds of millions of people, including many of the poorest farm households, live in river valleys where irrigation is fed by melting glaciers and snow. The annual snowmelt is coming earlier every year, synchronizing it less and less well with the summer growing season, and the glaciers are disappearing altogether.


Thus, the vast numbers of farmers in the Indo-Gangetic

Plain and in China's Yellow River Basin will most likely face severe disruptions in water availability. Yet those regions are already experiencing profound water stress because of unsustainable rates of groundwater pumping performed to irrigate large expanses of northern China and northern India.

In Africa, all signs suggest currently subhumid and arid areas will dry further, deepening the food crisis for many of the world's poorest and most vulnerable people. The severe decline in precipitation in the African Sahel during the past 30 years seems to be related to both anthropogenic warming and aerosol pollutants. The violence in Darfur and Somalia is fundamentally related to food and water insecurity. Ivory Coast's civil war stems, at least in part, from ethnic clashes after people fled the northern drylands of Burkina Faso for the coast. Worse chaos could easily arise.

Each El Niño cycle brings drying to thousands of islands in the Indonesian archipelago, with attendant crop failures, famine and peat fires. Some climatologists hypothesize that global warming could induce a more persistent El Niño state; if so, the 200 million people in Indonesia and neighboring areas could experience lasting drought conditions.

Until now, the climate debate has focused on the basic science and the costs and benefits of reducing greenhouse gas emissions. Attention will now increasingly turn to the urgent challenge of adapting to the changes and helping those who are most affected.

Some hard-hit places will be salvaged by better infrastructure that protects against storm surges or economizes on water for agriculture. Others will shift successfully from agriculture to industry and services. Yet some places will be unable to adjust altogether, and suffering populations will most likely move. We are just beginning to understand these phenomena in quantitative terms. Economists, hydrologists, agronomists and climatologists will have to join forces to take the next steps in scientific understanding of this human crisis. 

Hundreds of millions of people may be compelled to relocate.

Jeffrey D. Sachs is director of the Earth Institute at Columbia University (www.earth.columbia.edu).



A Simpler
ORIGIN
for **LIFE**



The sudden appearance of a large self-copying molecule such as RNA was exceedingly improbable. Energy-driven networks of small molecules afford better odds as the initiators of life

Extraordinary discoveries inspire extraordinary claims. Thus, James Watson reported that immediately after he and Francis Crick uncovered the structure of DNA, Crick “winged into the Eagle (pub) to tell everyone within hearing that we had discovered the secret of life.” Their structure—an elegant double helix—almost merited such enthusiasm. Its proportions permitted information storage in a language in which four chemicals, called bases, played the same role as 26 letters do in the English language.

Further, the information was stored in two long chains, each of which specified the contents of its partner. This arrangement suggested a mechanism for reproduction: The two strands of the DNA double helix parted company, and new DNA building blocks that carry the bases, called nucleotides, lined up along the separated strands and linked up. Two double helices now existed in place of one, each a replica of the original.

The Watson-Crick structure triggered an avalanche of discoveries about the way living cells function today. These insights also stimulated speculations about life’s origins. Nobel laureate H. J. Muller wrote that the gene material was “living material, the present-day representative of the first life,” which Carl Sagan visualized as “a primitive free-living naked gene situated in a dilute solution of organic matter.” (In this context, “organic” specifies compounds containing bound carbon atoms, both those present in life and those playing no

part in life.) Many different definitions of life have been proposed. Muller’s remark would be in accord with what has been called the NASA definition: life is a self-sustained chemical system capable of undergoing Darwinian evolution.

Richard Dawkins elaborated on this image of the earliest living entity in his book *The Selfish Gene*: “At some point a particularly remarkable molecule was formed by accident. We will call it the *Replicator*. It may not have been the biggest or the most complex molecule around, but it had the extraordinary property of being able to create copies of itself.” When Dawkins wrote these words 30 years ago, DNA was the most likely candidate for this role. Later, researchers turned to other possible molecules as the earliest replicator, but I and others think that this replicator-first model of the origin of life is fundamentally flawed. We prefer an alternative idea that seems much more plausible.

When RNA Ruled the World

COMPLICATIONS to the DNA-first theory soon set in. DNA replication cannot proceed without the assistance of a number of proteins—members of a family of large molecules that are chemically very different from DNA. Both are constructed by linking subunits together to form a long chain, but whereas DNA is made of nucleotides, proteins are made of amino acids. Proteins are the handymen of the living cell. Enzymes, proteins’ most famous subclass, act as expeditors, speeding up chemical pro-

An earlier, longer version of this story was posted on www.sciam.com. Feedback about that version helped to shape the article that appears here.

cesses that would otherwise take place too slowly to be of use to life. Proteins used by cells today are built following instructions encoded in DNA.

The above account brings to mind the old riddle: Which came first, the chicken or the egg? DNA holds the recipe for protein construction. Yet that information cannot be retrieved or copied without the assistance of proteins. Which large molecule, then, appeared first—proteins (the chicken) or DNA (the egg)?

plate an RNA world, containing only RNA molecules that serve to catalyze the synthesis of themselves.... The first step of evolution proceeds then by RNA molecules performing the catalytic activities necessary to assemble themselves from a nucleotide soup.” In this vision, the first self-replicating RNA that emerged from nonliving matter carried out the various functions now executed by RNA, DNA and proteins.

A number of additional clues support

tinguish it from the assertion that RNA merely arose before DNA and proteins.

The Soup Kettle Is Empty

THE RNA-FIRST HYPOTHESIS faces a tremendously challenging question: How did that first self-replicating RNA arise? Enormous obstacles block Gilbert’s picture of RNA forming in a non-living nucleotide soup.

RNA’s building blocks, nucleotides, are complex substances as organic molecules go. Each contains a sugar, a phosphate and one of four nitrogen-containing bases as sub-subunits. Thus, each RNA nucleotide contains nine or 10 carbon atoms, numerous nitrogen and oxygen atoms and the phosphate group, all connected in a precise three-dimensional pattern. Many alternative ways exist for making those connections, yielding thousands of plausible nucleotides that could readily join in place of the standard ones but that are not represented in RNA. That number is itself dwarfed by the hundreds of thousands to millions of stable organic molecules of similar size that are not nucleotides.

The idea that suitable nucleotides might nonetheless form draws inspiration from a well-known experiment published in 1953 by Stanley L. Miller. He applied a spark discharge to a mixture of simple gases that were then thought to represent the atmosphere of the early earth and saw that amino acids formed. Amino acids have also been identified in the Murchison meteorite, which fell in Australia in 1969. Nature has apparently been generous in providing a supply of these particular building blocks. By extrapolation of these results, some writers have presumed that *all* life’s building blocks could be formed with ease in Miller-type experiments and were present in meteorites. This is not the case.

Amino acids, such as those produced in experiments like Miller’s, are far less complex than nucleotides. Their defining features are an amino group (a nitrogen and two hydrogens) and a carboxylic acid group (a carbon, two oxygens and a hydrogen), both attached to the same carbon. The simplest of the 20 amino acids used to build natural proteins contains

Inanimate nature provides us with a variety of mixtures of small molecules as potential incubators for life.

A possible solution appeared when attention shifted to a new champion—RNA. This versatile class of molecule is, like DNA, assembled of nucleotide building blocks but plays many roles in our cells. Certain RNAs ferry information from DNA to ribosomes, structures (which themselves are largely built of other kinds of RNA) that construct proteins. In carrying out its various duties, RNA can take the form of a double helix that resembles DNA or of a folded single strand, much like a protein.

In the early 1980s scientists discovered ribozymes, enzymelike substances made of RNA. A simple solution to the chicken-and-egg riddle now appeared to fall into place: life began with the appearance of the first self-copying RNA molecule. In a germinal 1986 article, Nobel laureate Walter Gilbert wrote in the journal *Nature*: “One can contem-

the idea that RNA appeared before proteins and DNA in the evolution of life. For example, many small molecules, called co-factors, play a role in enzyme-catalyzed reactions. These co-factors often carry an attached RNA nucleotide with no obvious function. Such structures have been considered “molecular fossils,” relics descended from the time when RNA alone, without DNA or proteins, ruled the biochemical world.

This clue and others, however, support only the conclusion that RNA preceded DNA and proteins; they provide no information about the origin of life, which may have involved stages prior to the RNA world in which other living entities ruled supreme. Confusingly, researchers use the term “RNA world” to refer to both notions. Here I will use the term “RNA first” for the claim that RNA was involved in the origin of life, to dis-

Overview/Origin of Life

- Theories of how life first originated from nonliving matter fall into two broad classes—replicator first, in which a large molecule capable of replicating (such as RNA) formed by chance, and metabolism first, in which small molecules formed an evolving network of reactions driven by an energy source.
- Replicator-first theorists must explain how such a complicated molecule could have formed before the process of evolution was under way.
- Metabolism-first proponents must show that reaction networks capable of growing and evolving could have formed when the earth was young.

only two carbon atoms. Seventeen of the set contain six or fewer carbons. The amino acids and other substances that were prominent in the Miller experiment contained two and three carbon atoms.

In contrast, no nucleotides of any kind have been reported as products of spark-discharge experiments or in studies of meteorites. Apparently inanimate nature has a bias toward the formation of mol-

ecules made of fewer rather than greater numbers of carbon atoms and thus shows no partiality in favor of creating the nucleotides required by our kind of life.

To rescue the RNA-first concept from this otherwise lethal defect, its advocates have created a discipline called prebiotic synthesis. They have attempted to show that RNA and its components can be prepared in their laboratories in a sequence of carefully controlled reactions, using what they consider to be relevant conditions and starting materials.

The Web version of this article, available at www.sciam.com/ontheweb, goes into more detail about the shortcomings of prebiotic synthesis research. The problems bring the following analogy to mind: Consider a golfer who, having played a ball through an 18-hole course, then assumes that the ball could also play itself around the course in his absence. He had demonstrated the possibility of the event; it was only necessary to presume that some combination of natural forces (earthquakes, winds, tornadoes and floods, for example) could produce the same result, given enough time. No physical law need be broken for spontaneous RNA formation to happen, but the chances against it are immense.

Some chemists have suggested that a simpler replicator molecule similar to RNA arose first and governed life in a "pre-RNA world." Presumably this first replicator would also have the catalytic capabilities of RNA. Because no trace of this hypothetical primal replicator and catalyst has been recognized so far in modern biology, RNA must have completely taken over all its functions at some point after its emergence.

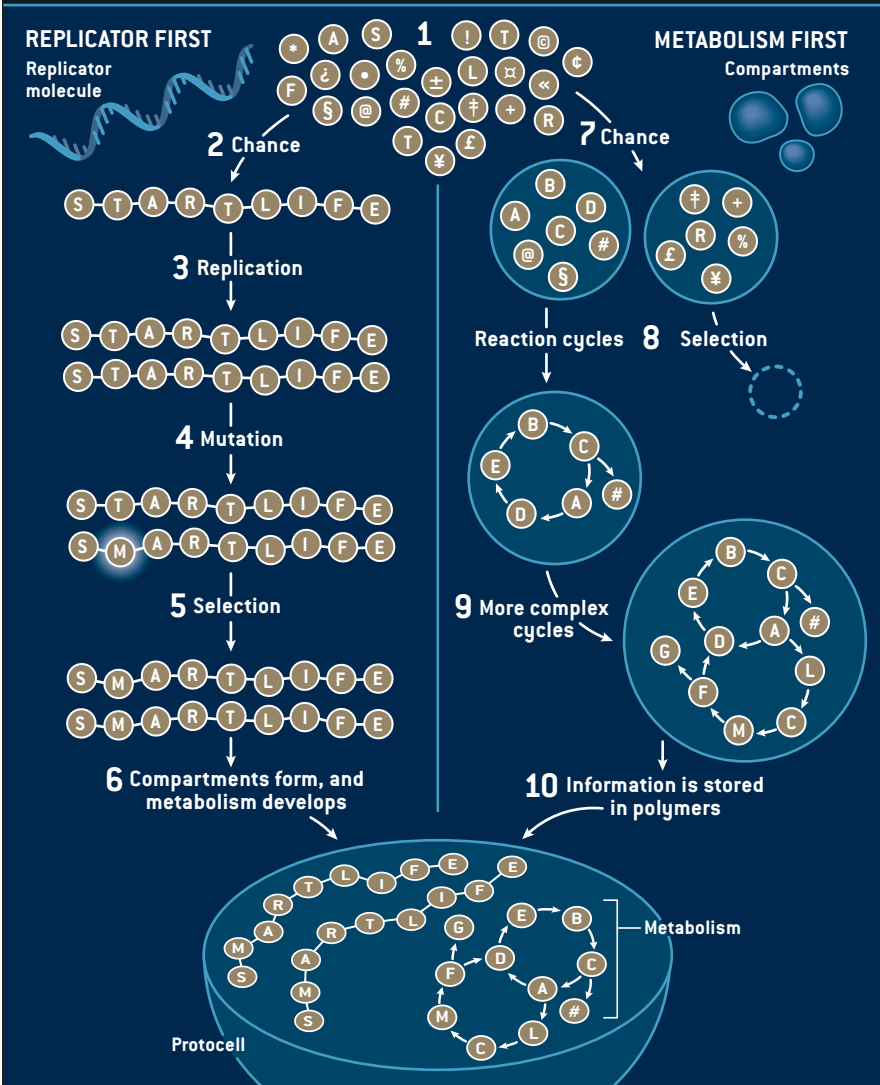
Yet even if nature could have provided a primordial soup of suitable building

REPLICATOR VS. METABOLISM

Scientific theories of the origin of life largely fall into two rival camps: replicator first and metabolism first. Both models must start from molecules formed by nonbiological chemical processes, represented here by balls labeled with symbols (1).

In the replicator-first model, some of these compounds join together in a chain, by chance forming a molecule—perhaps some kind of RNA—capable of reproducing itself (2). The molecule makes many copies of itself (3), sometimes forming mutant versions that are also capable of replicating (4). Mutant replicators that are better adapted to the conditions supplant earlier versions (5). Eventually this evolutionary process must lead to the development of compartments (like cells) and metabolism, in which smaller molecules use energy to perform useful processes (6).

Metabolism first starts off with the spontaneous formation of compartments (7). Some compartments contain mixtures of the starting compounds that undergo cycles of reactions (8), which over time become more complicated (9). Finally, the system must make the leap to storing information in polymers (10).



ROBERT SHAPIRO is professor emeritus of chemistry and senior research scientist at New York University. He is author or co-author of more than 125 publications, primarily in the area of DNA chemistry. In particular, he has studied the ways in which environmental chemicals can damage our hereditary material, causing changes that can lead to mutations and cancer.

blocks, whether nucleotides or a simpler substitute, their spontaneous assembly into a replicator involves implausibilities that dwarf those required for the preparation of the soup. Let us presume that the soup of building blocks has somehow been assembled, under conditions that favor their connection into chains. They would be accompanied by hordes of defective units, the inclusion of which in a nascent chain would ruin its ability to act as a replicator. The simplest kind of flawed unit would have only one “arm” available for connection to a building block, rather than the two needed to support further growth of the chain.

An indifferent nature would theoretically combine units at random, producing an immense variety of short, terminated chains, rather than the much longer one of uniform backbone geometry needed to support replicator and catalytic functions. The probability of this latter process succeeding is so vanishingly small that its happening even once anywhere in the visible universe would count as a piece of exceptional good luck.

Life with Small Molecules

NOBEL LAUREATE Christian de Duve has called for “a rejection of improbabilities so incommensurably high that they can only be called miracles, phenomena

that fall outside the scope of scientific inquiry.” DNA, RNA, proteins and other elaborate large molecules must then be set aside as participants in the origin of life. Inanimate nature instead provides us with a variety of mixtures of small molecules with which to work.

Fortunately, an alternative group of theories that can employ these materials has existed for decades. The theories use a thermodynamic, rather than a genetic, definition of life, under a scheme put forth by Sagan in the *Encyclopedia Britannica*: a localized region that increases in order (decreases in entropy) through cycles driven by an energy flow would be considered alive. This small-molecule approach is rooted in the ideas of Soviet biochemist Alexander Oparin. Origin-of-life proposals of this type differ in specific details; here I will list five common requirements (and add some ideas of my own).

1. A boundary is needed to separate life from nonlife. Life is distinguished by its great degree of organization, yet the second law of thermodynamics requires that the universe move in a direction in which disorder, or entropy, increases. A loophole, however, allows entropy to decrease in a limited area, provided that a greater increase occurs outside the area. When living cells grow and multiply, they convert chemical energy or radia-

tion to heat. The released heat increases the entropy of the environment, compensating for the decrease in living systems. The boundary maintains this division of the world into pockets of life and the nonliving environment in which they must sustain themselves.

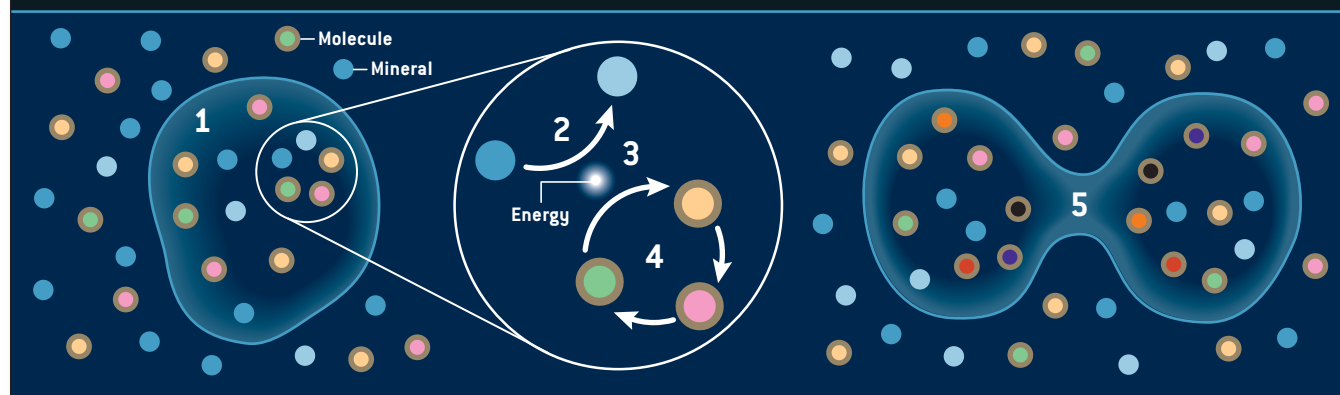
Today sophisticated double-layered cell membranes, made of chemicals classified as lipids, separate living cells from their environment. When life began, some natural feature probably served the same purpose. In support of this idea, David W. Deamer of the University of California, Santa Cruz, has observed membranelike structures in meteorites. Other proposals have suggested natural boundaries not used by life today, such as iron sulfide membranes, rock surfaces (in which electrostatic interactions segregate selected molecules from their environment), small ponds and aerosols.

2. An energy source is needed to drive the organization process. We consume carbohydrates and fats, combining them with oxygen that we inhale, to keep ourselves alive. Microorganisms are more versatile and can use minerals in place of the food or the oxygen. In either case, the transformations that are involved are called redox reactions. They entail the transfer of electrons from an electron-rich (or reduced) substance to

FIVE REQUIREMENTS FOR METABOLISM FIRST

At least five processes must occur for small molecules to achieve a kind of life—here defined as the creation of greater order in localized regions by chemical cycles driven by an energy flow. First, something must create a boundary to separate the living region from the nonliving environment (1). A source of energy must be available, here depicted as a mineral (blue) undergoing a heat-producing reaction (2). The released energy

must drive a chemical reaction (3). A network of chemical reactions must form and increase in complexity to permit adaptation and evolution (4). Finally, the network of reactions must draw material into itself faster than it loses material, and the compartments must reproduce (5). No information-storing molecule (such as RNA or DNA) is required; heredity is stored in the identity and concentration of the compounds in the network.



an electron-poor (or oxidized) one. Plants can capture solar energy directly and adapt it for the functions of life. Other forms of energy are used by cells in specialized circumstances—for example, differences in acidity on opposite sides of a membrane. Yet others, such as radioactivity and abrupt temperature differences, might be used by life elsewhere in the universe.

3. A coupling mechanism must link the release of energy to the organization process that produces and sustains life. The release of energy does not necessarily produce a useful result. Chemical energy is released when gasoline is burned within the cylinders of an automobile, but the vehicle will not move unless that energy is used to turn the wheels. A mechanical connection, or coupling, is required. Every day, in our own cells, each of us degrades pounds of a nucleotide called ATP. The energy released by this reaction serves to drive processes necessary for our biochemistry that would otherwise proceed too slowly or not at all. Linkage is achieved when the reactions share a common intermediate, and the process is sped up by the intervention of an enzyme. One assumption of the small-molecule approach is that coupled reactions and primitive catalysts sufficient to get life started exist in nature.

4. A chemical network must be formed to permit adaptation and evolution. We come now to the heart of the matter. Imagine, for example, that an energetically favorable redox reaction of a mineral drives the conversion of an organic chemical, A, to another one, B, within a compartment. I call this key transformation a driver reaction, because it serves as the engine that mobilizes the organization process. If B simply reconverts back to A or escapes from the compartment, we would not be on a path that leads to increased organization. In contrast, if a multistep chemical pathway—say, B to C to D to A—reconverts B to A, then the steps in that circular process (or cycle) would be favored to continue operating because they replenish the supply of A, allowing the continuing useful discharge of energy by the mineral reaction [see box on page 53].

What Readers Want to Know

In Scientific American's blog, Robert Shapiro answered questions raised by readers of the Web version of this article. An edited selection follows.

Does the metabolism-first hypothesis point to a single origin or multiple independent origins of life? —JR

A: Multiple origins seem more viable with the metabolism-first scenario. Gerald Feinberg and I discussed the possibility of alien life (life not based on RNA, DNA and other biochemistry familiar to us) in our 1980 book, *Life beyond Earth*. Researchers at a conference hosted by Paul Davies at Arizona State University in December 2006 concluded that alien life may even exist, undetected, on this planet. The great majority of microorganisms that can be observed under a microscope cannot be grown in conventional culture media and remain uncharacterized. Alien microbes may also exist in habitats on the earth that are too extreme for even the hardest forms of our familiar life.

Why do we have to demonstrate metabolism first in a reaction vessel? Can't we simulate it in software? —Dave Evanoff

A: Stuart Kauffman, Doron Lancet and others have used computer simulations to illustrate the feasibility of self-sustaining reaction cycles. Such simulations have not specified the exact chemical mixtures and reaction conditions needed to establish self-sustaining chemical networks. We do not yet know all the reaction pathways open to mixtures of simple organic compounds, let alone their thermodynamic constants. Even if such data were available, most chemists would not be convinced by a computer simulation but would demand an experimental demonstration.

The fact that all biological molecules are of one handedness needs some explanation. —John Holt

A: If the mineral transformation that powered the reaction cycle I discuss in my article were selective for only one mirror-image form of chemical A, then the product B and other members of the cycle might also occur in only one mirror-image form. Control of handedness, or chirality, becomes crucial when small chiral molecules are linked together to form larger ones. A modern enzyme may contain 100 linked amino acids, all of the same handedness (so-called L-amino acids). If a D-amino acid were substituted for its mirror-image L-form at a sensitive site within the enzyme, then the enzyme's shape would change and its function might be lost.

Branch reactions will occur as well, such as molecules converting back and forth between D and another chemical, E, that lies outside the ABCD cycle. Because the cycle is driven, the E-to-D reaction is favored, moving material into the cycle and maximizing the energy release that accompanies the driver reaction.

The cycle could also adapt to changing circumstances. As a child, I was fascinated by the way in which water, released from a leaky hydrant, would find a path downhill to the nearest sewer. If falling leaves or dropped refuse blocked that path, the water would back up until another route was found around the obstacle. In the same way, if a change in the acidity or in some other environmental circumstance should hinder a step in the

pathway from B to A, material would back up until another route was found. Additional changes of this type would convert the original cycle into a network. This trial-and-error exploration of the chemical "landscape" might also turn up compounds that could catalyze important steps in the cycle, increasing the efficiency with which the network used the energy source.

5. The network must grow and reproduce. To survive and grow, the network must gain material faster than it loses it. Diffusion of network materials out of the compartment into the external world is favored by entropy and will occur to some extent. Some side reactions may produce gases, which escape, or form tars, which will drop out of so-

lution. If these processes together should exceed the rate at which the network gains material, then it would be extinguished. Exhaustion of the external fuel would have the same effect. We can imagine, on the early earth, a situation where many start-ups of this type occur, involving many alternative driver reactions and external energy sources. Finally, a particularly hardy one would take root and sustain itself.

A system of reproduction must eventually develop. If our network is housed in a lipid membrane, physical forces may split it after it has grown enough. (Freeman Dyson of the Institute for Advanced Study in Princeton, N.J., has described such a system as a “garbage bag world” in contrast to the “neat and beautiful scene” of the RNA world.) A system that

functions in a compartment within a rock may overflow into adjacent compartments. Whatever the mechanism may be, this dispersal into separated units protects the system from total extinction by a local destructive event. Once independent units were established, they could evolve in different ways and compete with one another for raw materials; we would have made the transition from life that emerges from nonliving matter through the action of an available energy source to life that adapts to its environment by Darwinian evolution.

Changing the Paradigm

SYSTEMS OF THE TYPE I have described usually have been classified under the heading “metabolism first,” which implies that they do not contain a mecha-

nism for heredity. In other words, they contain no obvious molecule or structure that allows the information stored in them (their heredity) to be duplicated and passed on to their descendants. Yet a collection of small items holds the same information as a list that describes the items. For example, my wife gives me a shopping list for the supermarket; the collection of grocery items that I return with contains the same information as the list. Doron Lancet of the Weizmann Institute of Science in Rehovot, Israel, has given the name “compositional genome” to heredity stored in small molecules, rather than a list such as DNA or RNA.

The small-molecule approach to the origin of life makes several demands on nature (a compartment, an external energy supply, a driver reaction coupled to that supply, a chemical network that includes that reaction, and a simple mechanism of reproduction). These requirements are general in nature, however, and are immensely more probable than the elaborate multistep pathways needed to form a molecule that is a replicator.

Over the years, many theoretical papers have advanced particular metabolism-first schemes, but relatively little experimental work has been presented in support of them. In those cases where experiments have been published, they have usually served to demonstrate the plausibility of individual steps in a proposed cycle. The greatest amount of new data has perhaps come from Günter Wächtershäuser and his colleagues at Munich Technical University. They have demonstrated parts of a cycle involving the combination and separation of amino acids in the presence of metal sulfide catalysts. The energetic driving force for the transformations is supplied by the oxidation of carbon monoxide to carbon dioxide. The researchers have not yet demonstrated the operation of a complete cycle or its ability to sustain itself and undergo further evolution. A “smoking gun” experiment displaying those three features is needed to establish the validity of the small-molecule approach.

The principal initial task is the identification of candidate driver reactions—small-molecule transformations (A to B

An RNA-First Researcher Replies

Steven A. Benner of the Westheimer Institute for Science and Technology in Gainesville, Fla., argues that RNA-first models are alive and well.

Even as some declare that the RNA-first model of life's origin is dead because RNA arising spontaneously is fantastically improbable, research is lending support to the model.

Let me first acknowledge that most organic molecules when hit with energy (such as lightning or heat from volcanoes) become something resembling asphalt, more suitable for paving roads than sparking life. But metabolism-first models, to the extent that they have been supported with *any* real chemicals, must also deal with this paradox: molecules reactive enough to participate in metabolism are also reactive enough to decompose. There are no easy solutions.

Like many others, my research group has returned to the scientific imperative: actually do laboratory research to learn about how RNA might have arisen on the earth.

The sugar ribose, the “R” in RNA, provides an object lesson in how a problem declared “unsolvable” may instead merely be “not yet solved.” Ribose long remained “impossible” to make by prebiotic synthesis (reactions among mixtures of molecules that could plausibly have existed on a prebiotic earth) because it contains a carbonyl group—a carbon atom twice bonded to an oxygen atom. The carbonyl group confers both good reactivity (the ability to participate in metabolism) and bad reactivity (the ability to form asphalt). A decade ago Stanley L. Miller concluded that the instability of ribose stemming from its carbonyl group “preclude[s] the use of ribose and other sugars as prebiotic reagents.... It follows that ribose and other sugars were not components of the first genetic material.”

But prebiotic soups need soup bowls made of appropriate minerals, not Pyrex beakers. One attractive “bowl” is found today in Death Valley. In a primordial Death Valley, the environment was alternately wet and dry, rich in organic molecules from planetary accretion and (most important) full of minerals containing boron. Why care about boron? Because boron stabilizes carbohydrates such as ribose. Further, if borate (an oxide of boron) and organic compounds abundant in meteorites are mixed and hit with lightning, good quantities of ribose are formed from formaldehyde and the ribose does not decompose.

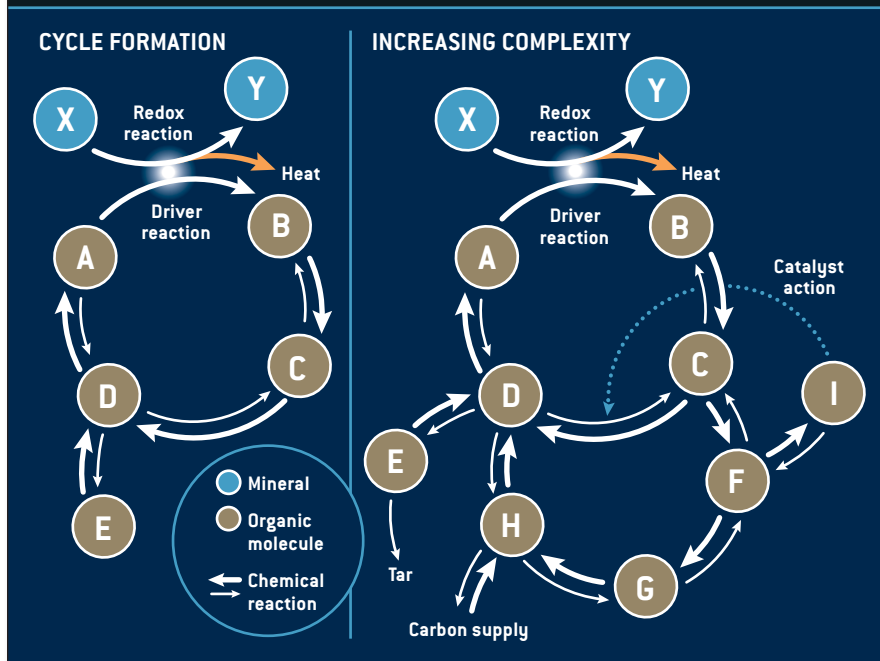
The fact that such a simple solution can be found for a problem declared “unsolvable” does not mean that the first form of life definitely used RNA to do genetics. But it should give us pause when advised to discard avenues of research simply because some of their problematic pieces have not yet been solved.

EVOLUTION OF CHEMICAL NETWORKS

The metabolism-first hypothesis requires the formation of a network of chemical reactions that increases in complexity and adapts to changes in the environment.

CYCLE FORMATION: An energy source (here the so-called redox reaction converting mineral X to mineral Y) couples to a reaction that converts the organic molecule A to molecule B. Further reactions (B to C, C to D...) form a cycle back to A. Reactions involving molecular species outside the cycle (E) will tend to draw more material into the cycle.

INCREASING COMPLEXITY: If a change in conditions inhibits a reaction in the cycle (for example, C to D), then other paths can be explored. Here a bypass has been found by which C is converted to D through intermediates F, G and H. Another solution would be the incorporation into the reaction network of a catalyst (I) whose action (dotted line) unblocks the C to D transformation. To survive, the evolving network must draw in carbon-containing materials from the environment more rapidly than it loses them by diffusion and side reactions, such as the formation of tars that settle out of the solution.



in the preceding example) that are coupled to an abundant external energy source (such as the oxidation of carbon monoxide or a mineral). Once a plausible driver reaction has been identified, there should be no need to specify the rest of the system in advance. The selected components (including the energy source), plus a mixture of other small molecules normally produced by natural processes (and likely to have been abundant on the early earth), could be combined in a suitable reaction vessel. If an evolving network were established, we would expect the concentration of the participants in the network to increase and alter with time. New catalysts that increased the rate of key reactions might appear, whereas irrelevant materials would de-

crease in quantity. The reactor would need an input device (to allow replenishment of the energy supply and raw materials) and an outlet (to permit removal of waste products and chemicals that were not part of the network).

In such experiments, failures would be easily identified. The energy might be dissipated without producing any significant changes in the concentrations of the other chemicals, or the chemicals might be converted to a tar, which would clog the apparatus. A success might demonstrate the initial steps on the road to life. These

steps need not duplicate those that took place on the early earth. It is more important that the general principle be demonstrated and made available for further investigation. Many potential paths to life may exist, with the choice dictated by the local environment.

An understanding of the initial steps leading to life would not reveal the specific events that led to the familiar DNA-RNA-protein-based organisms of today. Still, because we know that evolution does not anticipate future events, we can presume that nucleotides first appeared in metabolism to serve some other purpose, perhaps as catalysts or as containers for the storage of chemical energy (the nucleotide ATP continues to serve this function today). Some chance event or circumstance may have led to the connection of nucleotides to form RNA. The most obvious function of modern RNA is to serve as a structural element that assists in the formation of bonds between amino acids in the synthesis of proteins. The first RNAs may have served the same purpose, but without any preference for specific amino acids. Many further steps in evolution would be needed to “invent” the elaborate mechanisms for replication and specific protein synthesis that we observe in life today.

If the general small-molecule paradigm were confirmed, then our expectations of the place of life in the universe would change. A highly improbable start for life, as in the RNA-first scenario, implies a universe in which we are alone. In the words of biochemist Jacques Monod, “the universe was not pregnant with life nor the biosphere with man. Our number came up in the Monte Carlo game.”

The small-molecule alternative, however, is in harmony with the views of biologist Stuart Kauffman: “If this is all true, life is vastly more probable than we have supposed. Not only are we at home in the universe, but we are far more likely to share it with as yet unknown companions.”



Additional coverage—including commentaries, answers to questions, links to further reading and the opportunity to post your own comments—can be found at www.sciam.com/ontheweb

Lifting THE Fog AROUND *Anesthesia*

Learning why current anesthetics are so potent and sometimes dangerous will lead to a new generation of safer targeted drugs without unwanted side effects

By Beverley A. Orser

A Hollywood thriller due out this year

centers on a young man who awakens while undergoing open-heart surgery but is unable to move or cry out. The film's plot will undoubtedly take many more dramatic turns from there, but its early premise is, sadly, not entirely far-fetched. Episodes of intraoperative awareness while under general anesthesia are reported by one or two of every 1,000 patients. In reality, such incidents are usually brief and generally do not involve pain or distress, but they do highlight one of several ways that even the newest generation of anesthetic drugs can sometimes leave much to be desired. Indeed, the medical specialty of anesthesiology has evolved into a sophisticated art form because scientific understanding of how anesthetic drugs actually work, and how to make them better, has lagged behind most other areas of drug development.

Many of the modern anesthetics, in fact, share structural properties and clinical effects with ether, whose application as an anesthetic was first successfully demonstrated in public by Boston dentist William Morton in 1846. Since then, the use of general anesthesia has expanded to 40 million patients each year in North America alone. Yet advances in anesthetic care since Morton's day have come largely from the development of complex drug delivery systems and strategies for managing anesthesia's dangers and side effects.

JACK HOLLINGSWORTH Corbis



Today's general anesthetics are the most potent depressors of nervous system activity used in medicine. They even affect regulation of breathing and heart function. As a result, the drugs have a fairly narrow margin of safety, which is the difference between the therapeutic dose and a dose that is toxic, even lethal. That is one reason why individuals whose lung or cardiovascular function is already unstable—such as trauma victims undergoing emergency operations or patients in the midst of heart surgery—must receive a lighter than normal dose of anesthesia, which could make them susceptible to brief awareness incidents, as in the movie.

Although radical improvements in the care of people under general anesthesia have laid the foundation for complicated procedures such as organ transplants and open-heart surgery, the powerful neurodepressive effects of these drugs make them more likely to cause

death during an operation than the surgical procedure itself. And because anesthesia-related mortality has plateaued at a rate of approximately one patient in 13,000 for the past 15 years, it appears that anesthesiologists may have reached the limits of our ability to deliver these toxins safely. Moreover, severe side effects—ranging from loss of airway control to memory and cognitive problems after general anesthesia—may also stem from the broad yet poorly understood influence that current anesthetics exert on the central nervous system.

Science should be able to do much better, and very recent research is beginning to reveal how it can.

Pulling Out the Plugs

ALL OF TODAY'S general anesthetic drugs were developed empirically, which is to say, they were tested for their ability to produce the desirable effects that define the anesthetized state. Anes-

thesia's main components are sedation, unconsciousness (also sometimes called hypnosis), immobility, absence of pain (analgesia) and absence of memory for the anesthetized period (amnesia). By studying the mechanisms through which anesthetics achieve these end points, many groups, including my own at the University of Toronto, are beginning to tease those effects apart. Such studies are revealing that the activity of these potent drugs involves highly specific interactions with subpopulations of nervous system cells to create each of the separate properties of anesthesia.

Armed with this knowledge, we will be poised to finally move beyond the ether era and develop a new generation of highly specific drugs that can be used in combinations to deliver only the desired results without the dangers. As a bonus, this research is also yielding insights that can improve related therapies, such as sedatives and sleep aids, that

share some of anesthesia's mechanisms.

Anesthetics fall into two main categories based on whether they are delivered by inhalation, such as isoflurane, or intravenously, such as propofol. These drugs may appear to induce a deep sleep, but the state produced by most modern general anesthetics is more of a pharmacological coma. In a step toward clarifying the mechanisms underlying their effects, technologies such as magnetic

imaging anesthesia. One leading theory holds that it is simply the result of "cognitive unbinding"—a severing of communication between the many brain regions that usually cooperate in higher cognitive processing. Even at the local level, if one imagines groups of neurons as forming lines in a vast telephone network, the effect of general anesthesia is analogous to pulling out plugs at the switchboard. Researchers are, however, making en-

contain subtly different versions, however, which tend to predominate in different areas of the central nervous system. The presence of particular receptor subtypes on only certain subpopulations of cells will thus determine which cells are influenced by an anesthetic.

Contemporary studies are therefore focusing on identifying which receptor variants are the targets of current anesthetic drugs, understanding how the

The effect of anesthesia is analogous to pulling out plugs at the switchboard.

resonance imaging (MRI) and positron-emission tomography (PET) have helped identify some of the discrete brain regions and neural circuits involved in specific components of the anesthetized state. For instance, anesthetic action on the spinal cord accounts for the immobility produced by the drugs, whereas drug-induced changes to the hippocampus, a brain structure involved in memory formation, have been linked to amnesia. Chronic memory impairment following surgery, one of the undesirable side effects suffered by some patients, may also represent a leftover influence of the drugs on the hippocampus.

Because consciousness is a complex experience whose defining properties are still hotly debated by neuroscientists, it is not as easy to pinpoint a single anatomical source of unconsciousness dur-

ing progress in discovering details about the ways that anesthetic drugs physically act on the individual cells of the nervous system to block their transmissions.

During most of the 20th century, anesthetics were widely thought to work by disrupting the lipid components of cell membranes. Most anesthetics are highly fat-soluble compounds with widely differing chemical structures that range from simple inert gases to complex steroids. Their great physical and chemical diversity supported the idea that anesthetics must work in some nonspecific way to depress neuronal functioning. Modern research has shown, however, that anesthetics actually interact with multiple varieties of specific proteins, known as receptors, found on the surface of nerve cells. Families of receptors

drugs interact with the receptors to change the cell's function and how those cellular changes produce both the desired and unwanted "symptoms" of anesthesia.

Signaling Silence

MANY CATEGORIES of receptor proteins are found on the surface of neurons, but those activated by natural neurotransmitter chemicals have garnered the most interest in anesthesia research because they critically regulate communication along the neural "telephone lines." As their name implies, neurotransmitter molecules transmit messages between neurons at points of contact called synapses. They do so by traveling from the so-called presynaptic neuron across a tiny gap to bind to receptors on the postsynaptic neuron's cell membrane. When enough neurotransmitter molecules trigger the appropriate receptors, the postsynaptic cell's membrane generates an electrical potential that travels down its length to the next neuron in its network. Glutamate, serotonin, norepinephrine and acetylcholine are just a few of the neurotransmitters widely studied for their role in promoting such signaling throughout the central nervous system.

In anesthesia research, however, another neurotransmitter called gamma-aminobutyric acid (GABA) has gained the most attention because of its ability

Overview/Refining a Blunt Instrument

- General anesthetics are powerful nervous system suppressors, but how the drugs produce their broad effects throughout the brain and body is poorly understood.
- Investigation of anesthetics' underlying mechanisms is revealing that individual aspects of the anesthetized state are attributable to different sets of nerve cells, which are themselves distinguished by specific surface proteins that interact with the drugs.
- New compounds designed to target just those proteins, and hence only specific cell types, could be combined to selectively produce the desirable effects of anesthetics—as well as sedatives, sleep aids and memory drugs—with fewer risks and side effects.

to block neural communication. GABA is an inhibitory neurotransmitter: it helps to maintain overall balance in the nervous system by dampening neurons' ability to respond to excitatory messages from other cells. For that reason, GABA is thought to play a central part in the actions of anesthetic drugs.

Most postsynaptic receptors on cells that interact with GABA belong to a class termed ligand-gated ion channels. When GABA (the ligand) binds to the receptor, the receptor changes its conformation, temporarily opening a channel that admits negatively charged ions into the cell. The increased ion concentration generates a negative potential, preventing the cell from being able to produce an excitatory electrical pulse.

The receptor that is believed to be a primary target for anesthetics is the GABA subtype A, or GABA_A, receptor, which is also known to underlie the therapeutic effects of other classes of sedative and hypnotic drugs, most notably benzodiazepines such as Valium. Very low concentrations of benzodiazepines increase GABA_A receptor function, a relation that is easily confirmed because reversal agents that impede benzodiazepine from binding to the GABA_A receptor rapidly blunt the effects of those drugs.

Unfortunately, no such reversal agents exist for general anesthetics that might provide clues to their receptor targets. Nevertheless, studies using slices from many different brain regions and neurons grown in tissue culture have shown that both intravenous and inhaled anesthetics prolong the duration of postsynaptic electric currents generated by GABA_A receptors.

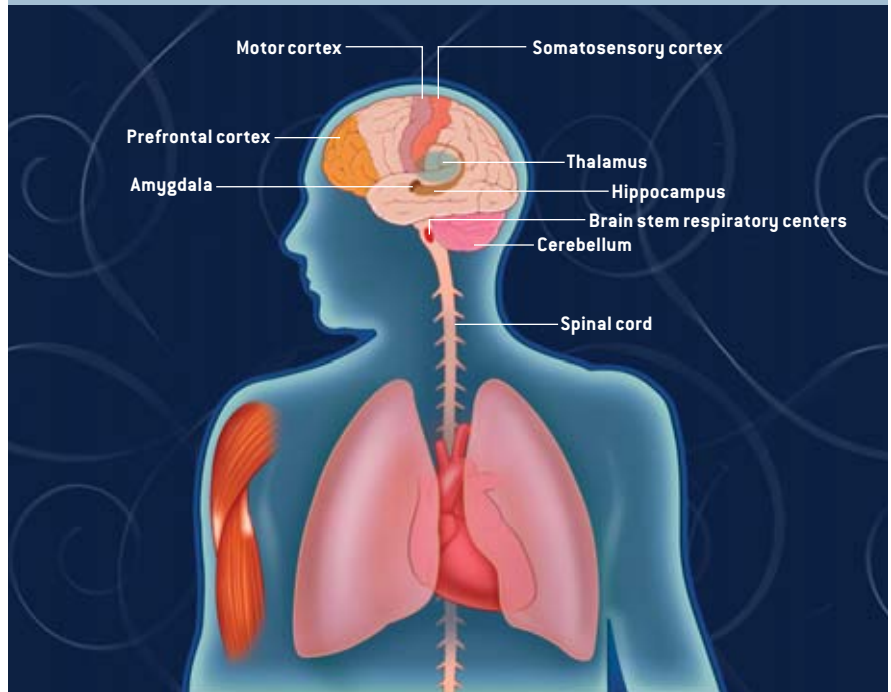
Anesthetics are believed to increase the function of GABA_A receptors by interacting at discrete binding cavities or attaching to specific amino acids in the receptors themselves and prolonging the channel opening, which extends the inhibitory effects of GABA molecules bound to the receptor. At high enough concentrations, anesthetics may even trigger the GABA receptors alone.

The vast majority of neurons contain GABA_A receptors, however, so scientists could not understand how anesthetics

Anesthesia's Broad Impact

Both the desirable and unwanted effects of anesthetic drugs stem from their power to suppress neuronal activity throughout the central nervous system, which encompasses the brain and spinal cord and controls

heart rate and breathing. Ongoing research is attempting to pinpoint the neural structures and regions whose changed functioning produces each of the defining properties of the anesthetized state.



Components of the Anesthetized State

Sedation

Reduced arousability, as evidenced by longer response times, slurred speech and decreased movement. Neuronal activity across brain cortical areas drops.

Unconsciousness (also called hypnosis)

Impaired perception of, and response to, stimuli. Cortical depression is deeper than in sedation. Activity in the thalamus, an area important for integrating brain processes, also falls significantly.

Immobility

Lack of movement in response to stimulation such as shaking or heat. Suppression of spinal cord neuronal activity is the main cause of this temporary paralysis, although the cerebellum, a motor control area, may also contribute.

Amnesia

Lack of recall for the anesthetized period. Many brain structures involved in memory formation, including the hippocampus, amygdala, prefrontal cortex and sensory and motor areas, exhibit anesthetic-induced changes.

Others

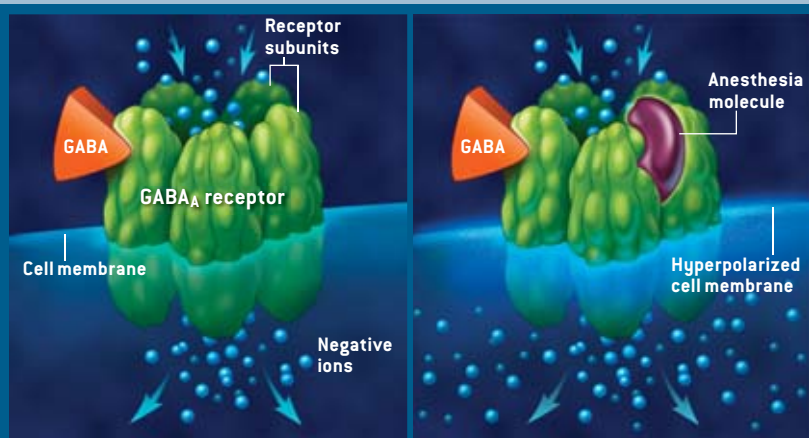
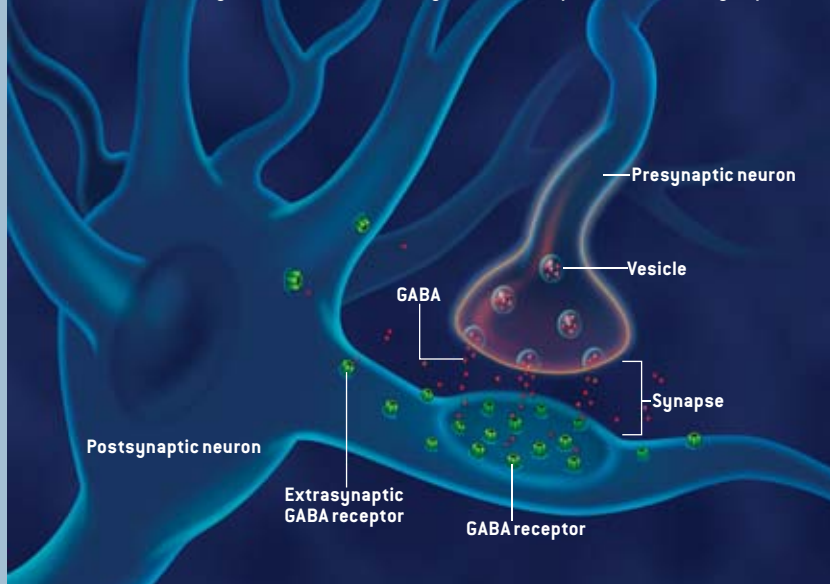
Muscle relaxation and lack of pain (analgesia) are sometimes included in definitions of the anesthetized state and are largely attributed to depression of spinal cord activity.

Jamming Transmission

Anesthetic drugs have been found to dampen neuronal communication, in part, by enhancing the effects of the neurotransmitter GABA, a signaling molecule that inhibits nerve cells from firing. Current research is focused on how the drugs interact with cellular GABA receptors to block neural activity.

Signal to Be Silent

An electrical pulse in the membrane of one neuron provokes release of GABA into the synapse, a juncture with another neuron. The molecules cross a small gap and bind to GABA-specific receptors on the postsynaptic cell. In many brain areas, GABA receptors are also found outside the synapse, along the nerve cell body, and are activated by GABA that spills out of the synapse.



Anesthetics and GABA: Changing Charge

A receptor subtype called GABA_A is a channel into the postsynaptic cell composed of five protein subunits. When GABA binds to it, the receptor opens to admit negatively charged ions, which increases polarization of the cell's membrane and prevents the neuron from generating an electrical pulse (*left*). Anesthetics are thought to act by binding to clefts in the GABA_A receptor and prolonging the channel opening, which causes hyperpolarization of the cell membrane (*right*).

could selectively influence different brain regions until research breakthroughs over the past decade revealed that not all GABA_A receptors are structurally or pharmacologically the same. The GABA_A receptor is a protein complex composed of five subunit parts, which can be mixed and matched in various combinations. At least 19 different GABA_A receptor subunits exist in mammals, and most of those have variant subtypes, so the possible number of combinations is high. The subunits most commonly seen in neurons, however, are the ones designated alpha, beta and gamma. In fact, most GABA_A receptors are composed of two alpha subunits, two betas and one gamma, although sometimes a delta or epsilon subunit replaces the gamma, depending on the brain region. But the key discovery was that the receptor's subunit composition dramatically alters its pharmacological properties: just one subunit difference within a GABA_A receptor's structure can determine whether and how it will respond to a particular anesthetic drug.

Because different GABA_A receptor subtypes predominate in different brain regions, researchers are increasingly able to pinpoint how anesthetics produce specific effects in various parts of the central nervous system by examining how the drugs interact in those regions with their particular target receptors.

Narrowing Down Targets

MY COLLEAGUES and I decided to focus on identifying the receptors that influence the memory-impairing properties of anesthetics, so we concentrated our studies on GABA_A receptors in the hippocampus. Anesthetics are known to cause amnesia at doses considerably lower than those required for unconsciousness or immobility, an effect that is clearly evident to anesthesiologists, for example, because patients rarely remember their own animated conversations as they go under or emerge from anesthesia. Yet, for unknown reasons, some patients experience unexpected recall of events during the surgery itself. Thus, by finding the correct target receptors for the amnesia-inducing effects

of anesthesia, it may become possible to identify patients at risk for intraoperative awareness because they lack those receptors. Alternatively, drug strategies to prevent awareness or at least its recollection could also be developed.

In the course of this work, we discovered to our surprise that even receptors outside the synapse could play a role in anesthetic action. If the synapse serves as a switchboard at the junction between two cells, then receptors at the synapse periphery or scattered along the nerve

Rather than provoking a response at the “switchboard,” the drugs were acting to boost a kind of static or inhibitory buzz in the telephone line itself that interfered with communication.

We found that the injectable anesthetics propofol and etomidate, and even the inhaled anesthetic isoflurane, increased the amplitude of this current by as much as 35-fold at concentrations several times lower than those required to cause immobility. Other investigators, including Stephen G. Brickley, Mark

on other research questions had also indicated that GABA_A receptors containing the alpha-5 subunit are involved in normal hippocampal-dependent memory processes, supporting our theory that the extrasynaptic alpha-5 receptors were responsible for the memory effects of an anesthetic. To test our hypothesis further, we turned to experimenting with genetically modified mice that lacked the alpha-5 subunit and wild-type mice that had the normal receptor. As expected, in behavioral tests the wild-type

Drug strategies to prevent awareness or at least its recollection could be developed.

cell body can be imagined as residing on the telephone line itself. Such extrasynaptic GABA_A receptors are activated by even the very low concentrations of GABA that are naturally present in the extracellular space or that spill over from nearby synapses. As it turns out, high numbers of extrasynaptic receptors are found in certain brain regions, such as the hippocampus and the thalamus (an area involved in consciousness and pain processing), as well as parts of the cortex and the cerebellum.

We serendipitously learned the relevance of extrasynaptic GABA_A receptors as anesthetic targets after struggling unsuccessfully for quite some time to identify postsynaptic receptors that were sensitive to the very low amnesia-inducing concentrations of anesthetics. We had also searched for populations of postsynaptic receptors that were synergistically modulated by midazolam and propofol, two of the most commonly used intravenous neurodepressive drugs, and had not found any of those either. Our work, however, was based on taking electrophysiological recordings of currents generated in hippocampal neurons in tissue culture, and we did notice that amnesia-producing concentrations of anesthetics significantly increased a persistent low-amplitude current generated by extrasynaptic GABA_A receptors.

Farrant and their colleagues at University College London, had already described this steady low current even in the absence of anesthetics. But what surprised our group was the extrasynaptic receptors' remarkable sensitivity to tiny amounts of both inhaled and intravenous anesthetics, whereas the low concentrations of anesthetics caused only negligible changes in postsynaptic currents. Previous studies, such as our own, had apparently focused on the right family of receptor proteins but were looking in the wrong location.

Eventually our experiments determined that the extrasynaptic GABA_A receptors were structurally slightly different from the populations of receptors within the synapse in that they predominantly contained an alpha-5 subunit, which the postsynaptic receptors generally lacked. That single change seemed to account for their sensitivity to even tiny amounts of anesthetics. These results were exciting to us because accumulating evidence from neuroscientists working

mice were sensitive to amnesia-causing doses of etomidate, whereas the alpha-5-deficient mice failed to manifest the drug's effects on memory.

We also established that the loss of alpha-5 GABA_A receptors had no consequences for any of the other anesthesia end points: sedation, immobility, hypnosis and response to a painful stimulus were the same in both groups of mice. These results demonstrated that the memory-impairing effects of etomidate could be dissociated from the drug's other properties based on the pharmacology of specific receptor subunits. They also provided the first animal model for receptor variations that might occur in humans and could explain some cases of resistance to an anesthetic's ability to induce amnesia. Ongoing studies will determine whether other general anesthetic drugs also preferentially target alpha-5 GABA_A receptors to produce amnesia.

At the same time, laboratories in Europe and the U.S. have been employing similar experimental techniques to ex-

THE AUTHOR

BEVERLEY A. ORSER is professor of anesthesiology and physiology at the University of Toronto and a practicing anesthesiologist at the university's hospital, the Sunnybrook Health Sciences Center, where she is also Canada Research Chair in anesthesia. As both a clinician and researcher, Orser focuses on improving patient safety. By studying the molecular mechanisms underlying anesthetic drugs, she hopes to advance development of new agents and related therapies with more finely controlled effects. Orser also consults for the pharmaceutical company Merck, developer of the sleep aid Gaboxadol.

Risk-Management Tools

Before surgery with general anesthesia, a patient will see and be attached to an array of monitoring devices. Many of them are there to guard against side effects of the anesthetics, which depress breathing and heart function, lowering blood pressure and body temperature. An anesthesiologist must constantly calibrate drug delivery to achieve the desired depth of anesthesia without subduing the patient's cardiac and respiratory functions enough to risk death.

Not shown:

Temperature probe

Applied to the skin or inserted into the esophagus or rectum, depending on type and length of surgery

Arterial catheter

Inserted into artery in wrist or groin. Provides beat-to-beat measurement of blood pressure and allows frequent blood sampling

Electrocardiograph electrodes

Applied to chest and limbs to track heart electrical activity and rate

Warming blanket

Provides forced hot air to maintain body temperature. Blood and intravenous fluid warmers are also commonly used

Pulse oximeter probe

Clipped to finger or ear to measure blood oxygen level

Breathing circuit

Delivers oxygen and air mixed with inhaled anesthetics through one tube; removes exhaled gases through the other tube

Blood pressure cuff

Intravenous line

Delivers intravenous anesthetics and fluids

plore the hypnotic and immobilizing effects of anesthetics. Gregg E. Homanics of the University of Pittsburgh School of Medicine, for example, developed a mouse lacking the GABA_A receptor delta subunit, which is known to confer high sensitivity to neurosteroids. His group's investigation found that the delta-deficient mice were, predictably, less sensitive to the steroid-based anesthetic alphaxalone in tests of the drug's power to induce unconsciousness. The mutant mice, however, displayed no differences in their responses to propofol, etomidate and other nonsteroidal anesthetics as compared with wild-type controls. Steroid anesthetics are not commonly used today, but these results also demonstrated the principle that different classes of anesthetics target discrete subpopulations of GABA_A receptors.

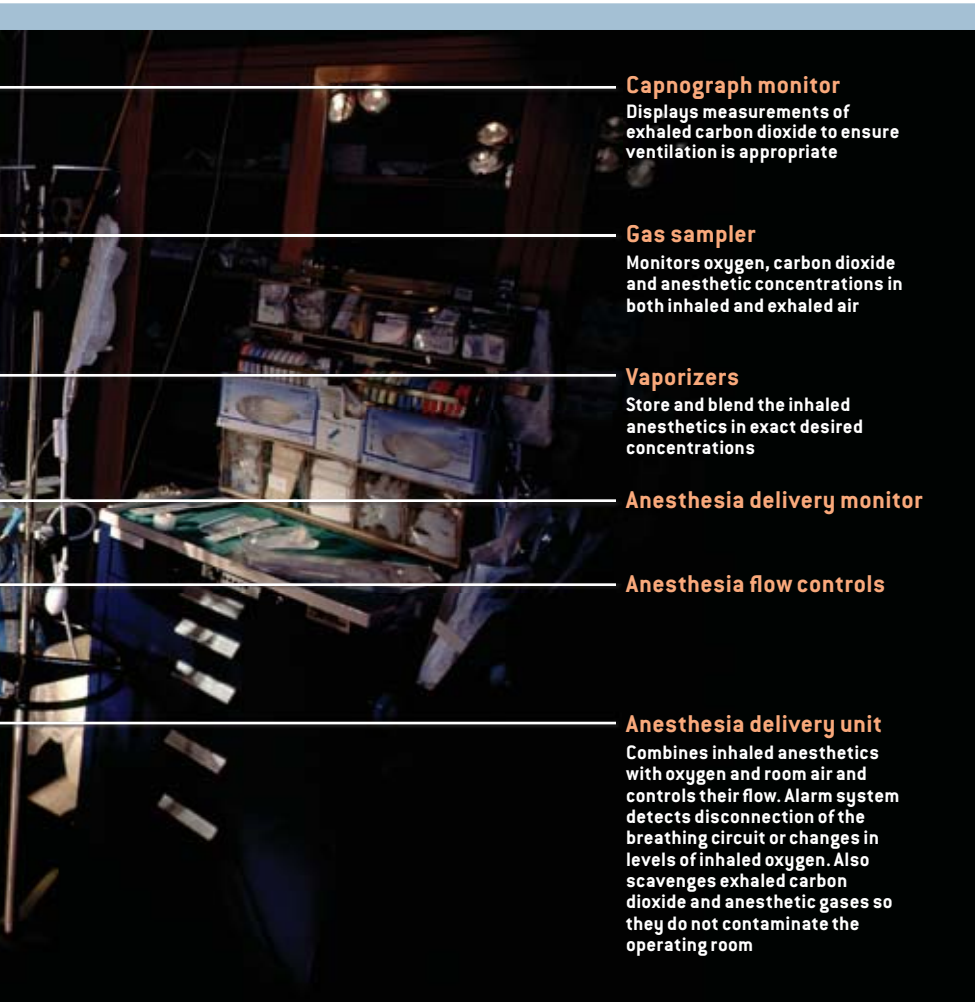
Such experiments have truly over-

turned the old notion that because anesthetics are so chemically different from one another they must produce their multiple effects by some general mechanism. Instead empirical development of anesthetic drugs seems to have stumbled on chemicals that produce similar end points, though each by its own unique mechanisms.

Etomidate, for example, is the only anesthetic in clinical use that is selective for GABA_A receptors containing the beta-2 or beta-3, but not the beta-1, subunit. Indeed, the differences between the beta-subunit variants that respond to etomidate and those that do not involve a single amino acid change at a specific point in the subunit's protein structure. The pharmaceutical company Merck developed transgenic mice with a mutation in that amino acid location within the beta-2 subunit and found

that etomidate was less effective at producing unconsciousness in the animals; the drug's immobilizing properties remained, however. Uwe Rudolph, while at the University of Zurich, also generated transgenic mice with the same mutation in the beta-3 subunit and found that it greatly diminished the effectiveness of both etomidate and propofol in producing unconsciousness and analgesia in the animals. In contrast, he showed that alphaxalone was equally effective in wild-type mice and those carrying the mutation, which indicates that these receptor subunits are probably not important targets of that drug.

Whether the point mutations in beta-2 and beta-3 receptor subunits also influence the drugs' amnesia-inducing properties has not yet been established. And which central nervous system regions in the transgenic mice are affected



Capnograph monitor
Displays measurements of exhaled carbon dioxide to ensure ventilation is appropriate

Gas sampler
Monitors oxygen, carbon dioxide and anesthetic concentrations in both inhaled and exhaled air

Vaporizers
Store and blend the inhaled anesthetics in exact desired concentrations

Anesthesia delivery monitor

Anesthesia flow controls

Anesthesia delivery unit
Combines inhaled anesthetics with oxygen and room air and controls their flow. Alarm system detects disconnection of the breathing circuit or changes in levels of inhaled oxygen. Also scavenges exhaled carbon dioxide and anesthetic gases so they do not contaminate the operating room

by the mutations is another unknown, although some evidence suggests that extrasynaptic GABA_A receptors in the thalamus may be critical. Taken together, however, such studies are confirming the central role of GABA_A receptors in the actions of anesthetics. The next step is to begin translating this knowledge gained from current general anesthetics into drugs that are anything but general.

Tailored Treatment

AS THE WORK of my research group and others has shown, extrasynaptic alpha-5 GABA_A receptors in the hippocampus are vital to the amnesia-inducing effects of etomidate and possibly to other general anesthetics currently in use. These results suggest that drugs that avoid or target that particular receptor could selectively spare or block memory formation as needed.

In fact, such compounds are already in development for other uses. Sedative-hypnotic drugs that do not act on the alpha-5 subunit, and hence should lack the memory-fuzzing effects of benzodiazepine sedatives and certain sleeping pills, are in the preclinical pipeline. And clinical trials of Gaboxadol, the first drug to selectively *target* extrasynaptic GABA_A receptors to enhance their function, are currently under way. Gaboxadol was initially developed as an anticonvulsant but

is now being studied as a sleep-promoting drug. It targets GABA_A receptors containing the delta subunit, primarily found in the thalamus and cerebellum, and therefore may also avoid affecting memory. The memory-blocking potential of similar compounds that do interact with alpha-5 receptors could also prove very useful in the surgical setting, where drugs that cause profound amnesia without depressing respiration, airway reflexes or the cardiovascular system might be highly desirable. In combination with other anesthetics, a potent memory blocker could be used to prevent intraoperative awareness episodes, for example. Alone, such a drug might be helpful in the treatment of patients suffering post-traumatic stress disorder by inhibiting certain distressing memories.

Management of anesthesia's effects on memory is just one example of a new approach to anesthesiology that will become possible with such targeted drugs. In many situations, the broad and profound neurodepression of current anesthetics is unnecessary and undesirable anyway. With a cocktail of compounds, each of which produces only one desirable end point, the future version of anesthesia care could leave a patient conversant but pain-free while having a broken limb repaired or immobile and sedated but aware while having a hip replaced. This polypharmaceutical approach is already widely used for other aspects of surgery-related care, most notably in the treatment of postoperative pain.

Today anesthesia has never been safer, but it is certainly not without risk. A tremendous opportunity now exists for the field to move beyond the ether era and toward a truly modern model of anesthesia care. SA

MORE TO EXPLORE

Anesthesia Safety: Model or Myth? Robert S. Lagasse in *Anesthesiology*, Vol. 97, pages 1609–1617; December 2002.

Molecular and Neuronal Substrates for General Anesthetics. Uwe Rudolph and Bernd Antkowiak in *Nature Reviews Neuroscience*, Vol. 5, pages 709–720; September 2004.

Emerging Molecular Mechanisms of General Anesthetic Action. Hugh C. Hemmings et al. in *Trends in Pharmacological Sciences*, Vol. 6, No. 10, pages 503–510; October 2005.

α5GABA_A Receptors Mediate the Amnestic but Not Sedative-Hypnotic Effects of the General Anesthetic Etomidate. Victor Y. Cheng et al. in *Journal of Neuroscience*, Vol. 26, No. 14, pages 3713–3720; April 5, 2006.

When Fields Collide

The history of particle cosmology, a new branch of physics that has shed light on the origins of the universe, shows that science can sometimes benefit from wrenching changes

By David Kaiser

Particle cosmology, which investigates how the smallest units of matter have determined the shape and fate of the universe, is one of the hottest topics in physics today. In recent years the field has received as much as half a billion dollars in funding from the U.S. Department of Energy, the National Science Foundation and NASA. Scientists have made great strides in understanding the high-energy particle interactions that roiled the universe in the first moments of its history and influenced cosmic evolution in the billions of years afterward.

The dramatic success of particle cosmology is all the more striking given that

this branch of research did not even exist 30 years ago. Before 1975, particle physics and cosmology were treated as separate fields of study (especially in the U.S.), and few scientists considered how discoveries in one specialty could enhance research in the other.

So why did particle cosmology arise? During the mid-1970s, researchers realized that studies of the early universe offered a unique window for investigating high-energy phenomena that cannot be recreated in the laboratory. But a series of changes in the funding and teaching of physics also helped to push cosmological questions to the forefront. The rapid emer-

► FIRST EXPLOSIVE MOMENTS of cosmic history have been reconstructed by particle cosmologists, who study the birth of the universe.

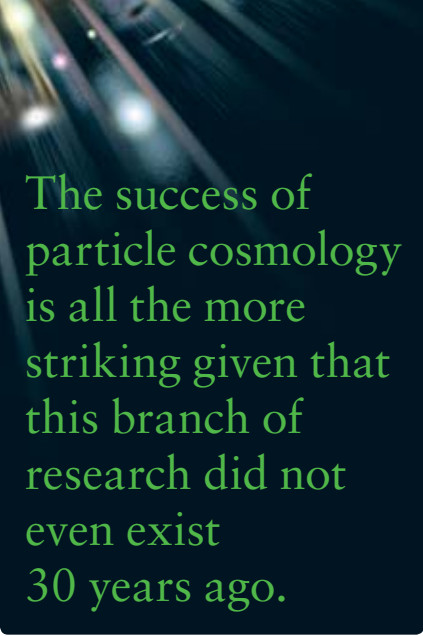


gence of particle cosmology illustrates how government budgets, educational institutions and even the publication of textbooks can radically alter the direction of research. The history of that era also shows that science can reap tremendous benefits when researchers move away from familiar subjects to tackle new challenges.

A good way to tell the story is to focus on the fortunes of two sets of ideas: the Brans-Dicke field, introduced by gravitational specialists, and the Higgs field, puzzled over by particle physicists. Both groups created these concepts in response to a problem that exercised many scientists during the late 1950s and early 1960s: why objects have mass. Although these two theories did not drive the union of particle physics and cosmology, the course of their development demonstrates how the two branches of research converged.

A Tale of Two ϕ 's

MASS SEEMS LIKE such an obvious property of matter that one might not think it requires an explanation. Yet finding descriptions of mass that were compatible with other ideas from modern physics proved no easy feat. Experts on gravitation and cosmology framed the problem in terms of Mach's principle, named for Austrian physicist and philosopher Ernst Mach, a famed critic of Newton and an inspiration to the young Albert Einstein. A good approximation of Mach's principle might be phrased this way: an object's mass—a measure of its resistance to changes in its motion—ultimately derives from that object's gravitational interactions with all the other matter in the universe. Al-



The success of particle cosmology is all the more striking given that this branch of research did not even exist 30 years ago.

though this principle intrigued Einstein and spurred his thinking, his general theory of relativity ultimately departed from it.

To incorporate Mach's principle into gravitational theory, scientists postulated the existence of a new scalar field that interacts with all types of matter. (A scalar field has one value for each point in space and time.) In 1961 Princeton University graduate student Carl Brans and his thesis adviser, Robert H. Dicke, pointed out that in Einstein's general relativity, the strength of gravity is fixed by Newton's constant, G . According to Einstein, G has the same value on Earth as it does in the most distant galaxies and does not change over time. Offering an alternative, Brans and Dicke suggested that Mach's principle could be satisfied if Newton's constant varied over time and space. They introduced a field called ϕ that was inversely proportional to Newton's constant and swapped $1/\phi$ for G throughout Einstein's gravitational equations.

According to the Brans-Dicke theory, matter responds to the curvature of

space and time, as in ordinary general relativity, *and* to variations in the local strength of gravity [see top illustration in box on opposite page]. The ϕ field permeates all of space, and its behavior helps to determine how matter moves through space and time. Any measurement of an object's mass therefore depends on the local value of ϕ . This theory was so compelling that members of Kip Thorne's gravity group at the California Institute of Technology used to joke that they believed in Einstein's general relativity on Mondays, Wednesdays and Fridays and in Brans-Dicke gravity on Tuesdays, Thursdays and Saturdays. (They remained agnostic on Sundays.)

Meanwhile, within the much larger community of particle physicists, the problem of mass arose in a different form. Beginning in the 1950s, theorists found that they could represent the effects of nuclear forces by imposing special classes of symmetries on the equations governing the behavior of subatomic particles. Yet the terms they would ordinarily include in these equations to represent particle masses violated the special symmetries. In particular, this impasse affected the W and Z bosons—the particles that give rise to the weak nuclear force, which is responsible for radioactive decay. If these force-carrying particles were truly massless, as the symmetries seemed to require, then the range of nuclear forces should have been infinite—two protons, for example, should have been able to exert a nuclear force on each other from across the galaxy. Such a long range flagrantly contradicted the observed behavior of nuclear forces, which fall off rapidly for distances larger than the size of atomic nuclei. Only if the force-carrying particles had some mass would the theoretically predicted range come into line with observations.

Many physicists focused on this conundrum, trying to formulate a theory that would represent the symmetry properties of subatomic forces while also incorporating massive particles. In 1961 Jeffrey Goldstone, then at the University of Cambridge, noted that the so-

Overview/A Revolution in Physics

- Until the 1970s researchers considered particle physics and cosmology to be completely separate fields of study.
- Sharp cutbacks in particle physics starting in the late 1960s prompted scientists to expand their horizons and explore topics in gravitation and cosmology.
- By the 1980s researchers had found that studying the early universe offered a new way to explore high-energy phenomena. Since then, the hybrid field of particle cosmology has become one of the most fruitful in physics.

lutions to the equations need not obey the same symmetries that the equations themselves do. As a simple illustration he introduced a scalar field, coincidentally labeled φ , whose potential energy density, $V(\varphi)$, bottoms out at two points: when φ has the values of $-v$ and $+v$ [see bottom illustration in box at right]. Because the energy of the system is lowest at these minima, the field will eventually settle into one of them. The potential energy is exactly the same for both values, but because the field must eventually land at just one value—either $-v$ or $+v$ —the solution to the equations spontaneously breaks their symmetry.

In 1964 Peter W. Higgs of the University of Edinburgh revisited Goldstone's work and found that a theory with spontaneous symmetry breaking would allow for the existence of massive particles. The mass arises from interactions between the φ field and all types of particles, including those that generate the weak nuclear force. The equations governing these interactions, Higgs demonstrated, obey all the requisite symmetries. Before φ settles into one of the minima of its potential energy, the particles skip lightly along, merrily unencumbered. Once φ arrives at either $+v$ or $-v$, however, the newly anchored field exerts a drag on anything coupled to it—the subatomic equivalent of being mired in molasses. In other words, the force-carrying particles (as well as garden-variety matter such as electrons) start to behave as if they have a nonzero mass, and any measurements of their mass depend on the local value of φ .

The Brans-Dicke and Higgs papers were published at about the same time in the same journal, *Physical Review*. Both articles quickly became well known; to this day, both are among the most cited physics articles of all time. Each proposed to explain the origin of mass by introducing a new scalar field that interacted with all types of matter. Given the similarity of the proposals and the quick attention they both received, one might have expected physicists to consider them alongside each other. Yet this pairing rarely happened. Of the 1,083 articles that cited either

Separate Concepts of Mass

The barriers that divided physicists in the early 1960s are illustrated by their parallel attempts to explain why objects have mass. Although cosmologists and particle physicists proposed similar theories, few scientists saw the connection.

From Cosmology: Brans-Dicke Gravity

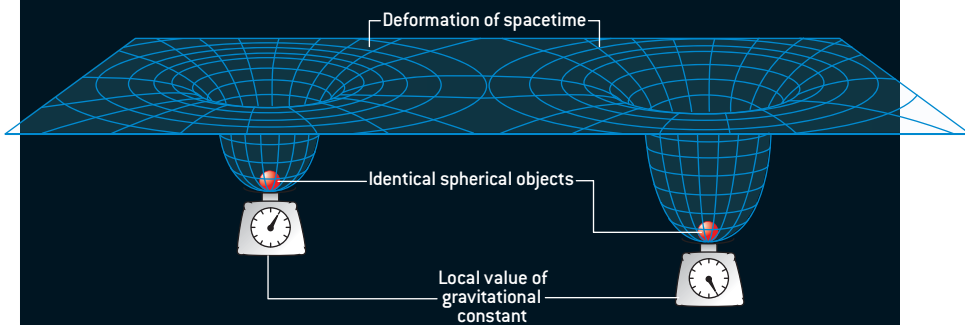


Carl Brans



Robert Dicke

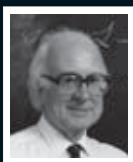
▼ In 1961 Carl Brans and Robert Dicke of Princeton University proposed a field called φ that allows Newton's gravitational constant to vary over time and space. An object at a point in space where the constant is small (left) will be less massive—and warp the local spacetime less—than an identical object at a point where the constant is large (right).



From Particle Physics: The Higgs Field

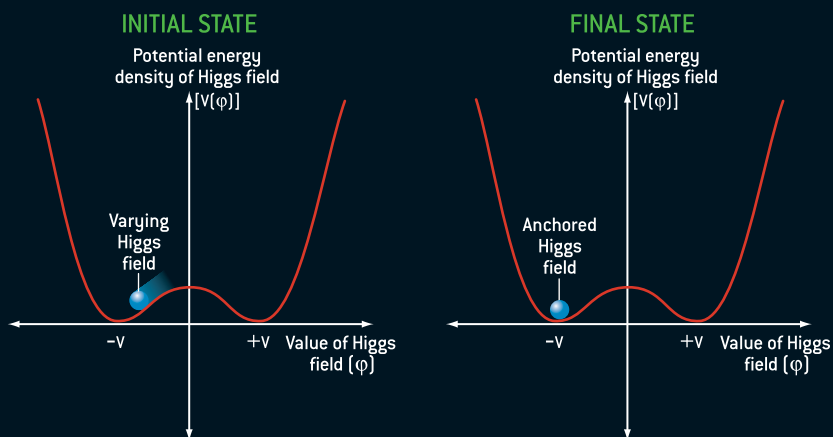


Jeffrey Goldstone



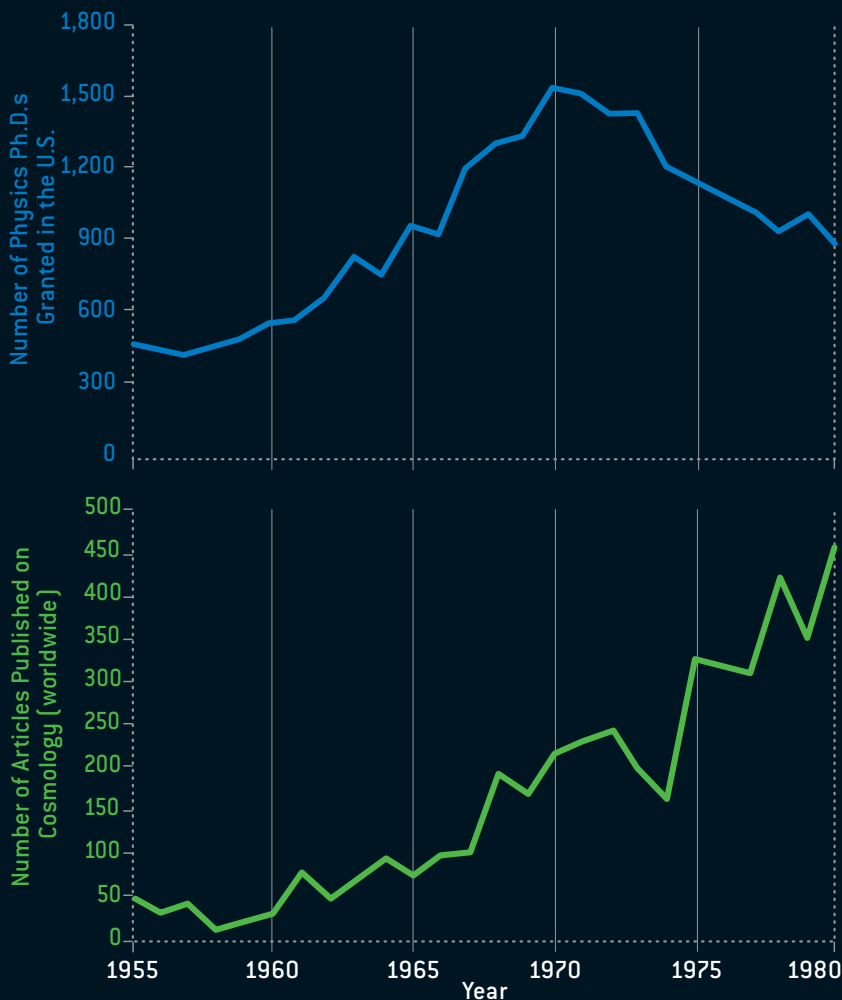
Peter Higgs

▼ In 1961 Jeffrey Goldstone, then at the University of Cambridge, introduced a field—also called φ , coincidentally—whose potential energy density $V(\varphi)$ bottoms out at two points, $-v$ and $+v$. Three years later Peter Higgs of the University of Edinburgh used this field to explain mass. Particles were massless at first when φ varied (left); they acquired mass only after φ settled into one of its minima (right).



A Change of Topic

▼ Government support for physics boomed in the 1950s and 1960s but fell sharply in the late 1960s and 1970s. The number of new Ph.D.s plummeted as well (top). Many particle physicists, who were hit hardest, shifted their attention to cosmology. Research in that field blossomed (bottom).



the Brans-Dicke or Higgs paper between 1961 and 1981, only six—less than 0.6 percent—included both articles in their references. (The earliest instance was in 1972, and the rest came after 1975.) This mutual ignorance highlights the stark boundaries that existed at the time between the particle physicists and the specialists in gravitation and cosmology.

Pushes, Pulls and Pedagogy

CLEARLY, the two communities saw different things in their respective ϕ 's.

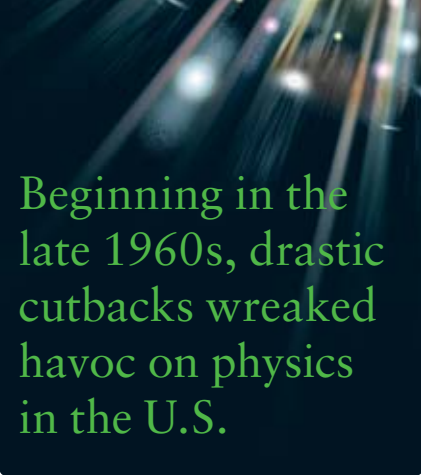
To the experts in gravitation and cosmology, the Brans-Dicke field (ϕ_{BD}) was exciting because it offered an alternative to Einstein's general relativity. To the particle physicists, the Higgs field (ϕ_H) was exciting because it offered hope that their theories might be able to explain the behavior of nuclear forces among massive particles. Before the mid-1970s, nobody suggested that ϕ_{BD} and ϕ_H might be physically similar or even worth examining side by side.

The divide between particle physics and cosmology was especially sharp in

the U.S. when Brans, Dicke, Goldstone and Higgs introduced their respective ϕ 's. The Physics Survey Committee of the National Academy of Sciences, for example, issued a policy report in 1966 that recommended doubling the funding and Ph.D.-level personnel for particle physics over the next few years but called for virtually no expansion in the already small areas of gravitation, cosmology and astrophysics. Furthermore, even though some of the Soviet textbooks published on gravity during that era included speculations about nuclear forces, such mixing of genres was absent from American textbooks.

These research patterns, however, would change radically by the late 1970s. Looking back on the swift rise of particle cosmology, physicists almost always point to two important developments that spurred the merger: the discovery of asymptotic freedom in 1973 and the construction of the first grand unified theories, or GUTs, in 1973 and 1974. Asymptotic freedom refers to an unexpected phenomenon in certain classes of theories governing particle interactions: the strength of the interaction decreases as the energy of the particles goes up, rather than increasing the way most other forces do. For the first time, particle theorists were able to make accurate and reliable calculations of phenomena such as the strong nuclear force—which keeps quarks bound within nuclear particles such as protons and neutrons—as long as they restricted their calculations to very high energy realms, far beyond anything that had been probed experimentally.

The introduction of GUTs likewise directed attention toward very high energies. Particle physicists realized that the strengths of three of the fundamental forces—electromagnetism and the weak and strong nuclear forces—might converge as particle energies increased. Theorists hypothesized that once the energies rose high enough, the three forces would act as a single undifferentiated force. The energy scale at which this grand unification would set in was literally astronomical: about 10^{24} electron volts, or more than one trillion times



Beginning in the late 1960s, drastic cutbacks wreaked havoc on physics in the U.S.

higher than the top energies physicists had been able to probe using particle accelerators. GUT-scale energies could never be achieved in Earth-bound laboratories, but some researchers realized that if the entire universe had begun in a hot big bang, then the average energy of particles in the universe would have been extraordinarily high during early periods in cosmic history.

With the advent of asymptotic freedom and GUTs, particle physicists had an obvious reason to begin studying the early universe: the first moments of the big bang would provide them with “the poor man’s accelerator,” allowing them to observe high-energy interactions that were impossible to re-create on Earth. Scores of scientists, journalists, philosophers and historians have pointed to this development to explain the emergence of particle cosmology.

But is it the whole story? Although the advances in particle theory were certainly important, they are not sufficient to explain the growth of this new subfield. For one thing, the timing is a bit off. Publications on cosmology (worldwide as well as in the U.S.) began a steep rise *before* 1973, and the rate of increase was completely unaffected by the appearance of the papers on asymptotic freedom and GUTs [see box on opposite page]. Moreover, GUTs did not receive much attention, even from particle theorists, until the late 1970s and early 1980s. Three of the earliest review articles on the emerging field of particle cosmology, published between 1978 and 1980, ignored asymptotic freedom and GUTs altogether.

New ideas alone were not enough to pave the way for particle cosmology; governmental and educational changes played major roles as well. Until the mid-1960s, U.S. physicists had benefited from a “cold war bubble,” a period when the federal government lavished funds on education, defense and scientific research. Beginning in the late 1960s, though, drastic cutbacks triggered by anti-Vietnam War protests, a thawing of the cold war and the introduction of the Mansfield Amendment, which heavily restricted Department of

Defense spending on basic research, wreaked havoc on physics in the U.S. Nearly all fields of science and engineering went into decline, but physics fell faster and deeper than any other field. The number of new physics Ph.D.s plummeted, falling nearly as fast from 1970 to 1975 as it had risen during the years after Sputnik.

Federal funding for physics also plunged, dropping by more than one third (in constant dollars) between 1967 and 1976. From the 1950s to the mid-1960s, the number of available jobs had always been greater than the number of physics students looking for work at the placement service meetings held by the American Institute of Physics. But employment prospects quickly turned grim: 989 applicants competed for 253 jobs in 1968, and 1,053 students scrambled for 53 positions in 1971.

Particle physics was hardest-hit by far, with federal spending on the field falling by 50 percent between 1970 and 1974. A swift exodus of talent began: between 1968 and 1970, twice as many U.S. researchers left particle physics as entered the field. The number of new Ph.D.s in particle physics fell by 44 percent between 1969 and 1975—the fastest decline of any branch of physics. At the same time, however, the fortunes of astrophysics and gravitation began to

rise. Spurred in part by a series of breakthroughs during the mid-1960s, including the discovery of quasars, pulsars and the cosmic microwave background radiation, the number of new Ph.D.s in this area grew by 60 percent between 1968 and 1970 and by another 33 percent between 1971 and 1976—even as the total number of physics Ph.D.s fell sharply.

Surveying the wreckage in 1972, the National Academy’s Physics Survey Committee released a new report that highlighted the troubles in particle physics. Many young theorists in that field, the committee noted, were having difficulty switching their research efforts elsewhere because of their “narrow specialization.” The report urged the nation’s physics departments to revamp how particle theorists were trained: “University groups have a responsibility to expose their most brilliant and able students to the opportunities in all subfields of physics.” Changes in university curricula quickly followed, aimed to broaden graduate students’ exposure to other areas of physics—including more emphasis on gravitation and cosmology. Across the country, physics programs began to offer new courses on the subject. After ignoring gravitation and cosmology for decades, American publishers pumped out scores of textbooks on the topic to meet the sudden demand.

Inflating the Ranks

THESE ABRUPT CHANGES left their mark on the way physicists viewed concepts such as the Brans-Dicke and Higgs fields. In 1979, after nearly two decades in which virtually no one had even mentioned the two fields in the same paper, let alone considered them to be physi-

THE AUTHOR

DAVID KAISER is both a physicist and a historian. He received Ph.D.s in theoretical physics and the history of science from Harvard University and is now an associate professor in the Program in Science, Technology and Society at the Massachusetts Institute of Technology and a lecturer in M.I.T.’s physics department. His recent book, *Drawing Theories Apart: The Dispersion of Feynman Diagrams in Postwar Physics* (University of Chicago Press, 2005), traces how Richard Feynman’s idiosyncratic approach to quantum physics entered the mainstream. He is completing a new book on physics during the cold war, looking in particular at changes in the training of graduate students. His current physics research focuses on particle cosmology, working on ways that cosmic inflation might be made compatible with superstring-inspired large extra dimensions.

Making the Connection

By the late 1970s a new generation of physicists, conversant with both particle theory and cosmology, explored possible links between Brans-Dicke gravity and the Higgs field.



◀ ANTHONY ZEE

As an undergraduate, Zee worked with gravitation expert John Wheeler at Princeton University, then pursued a Ph.D. in particle theory. He renewed his interest in cosmology while on sabbatical in Paris in 1974.



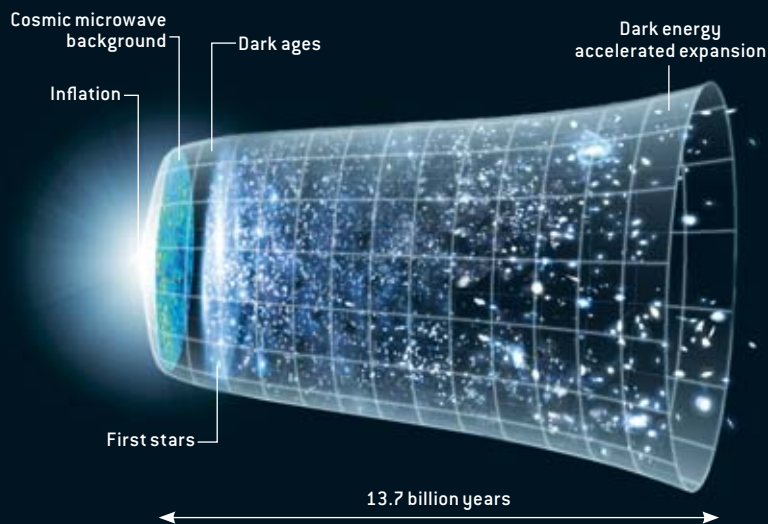
◀ LEE SMOLIN

In the 1970s Smolin studied particle theory and cosmology as a graduate student at Harvard University. He also worked with Stanley Deser of Brandeis University, one of the pioneers of quantum gravity.



◀ ALAN GUTH

Guth earned his Ph.D. in particle physics from the Massachusetts Institute of Technology in 1972. He became interested in cosmology after attending a lecture by Dicke in the late 1970s.



▲ In separate papers published in 1979, Zee and Smolin combined the Brans-Dicke gravitational equations with a Goldstone-Higgs symmetry-breaking potential. In 1981 Guth introduced another field, based on the Higgs, called the inflaton. This field provided the driving force behind a postulated period of superfast expansion—or inflation—during the universe's first moments.

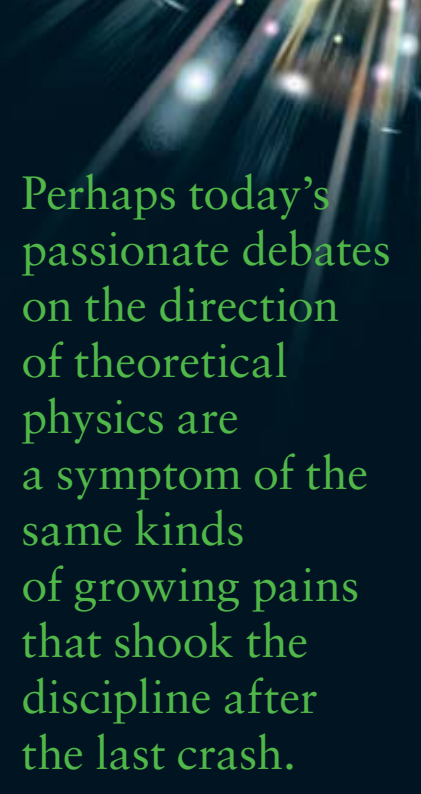
cally similar, two American theorists independently suggested that ϕ_{BD} and ϕ_H might be one and the same field. In separate papers, Anthony Zee, then at the University of Pennsylvania, and Lee Smolin, then at Harvard University, glued the two key pieces of ϕ together by

combining the Brans-Dicke gravitational equations with a Goldstone-Higgs symmetry-breaking potential. (Other theorists, working outside the U.S., had tentatively broached similar ideas between 1974 and 1978, but they received little attention at the time.)

In this model, the local strength of gravity initially varied over space and time, with G proportional to $1/\phi^2$, but its present-day constant value emerged after the ϕ field settled into a minimum of its symmetry-breaking potential, which presumably occurred in the first moments of the big bang. In this way, Zee and Smolin offered an explanation of why the gravitational force is so weak compared with other forces: when the field settled into its final state, $\phi = \pm v$, it anchored ϕ to some large, nonzero value, pushing G (which is inversely proportional to v^2) to a small value.

The career paths of Zee and Smolin illustrate the ways in which physicists focused their attention on cosmology after the collapse of the cold war bubble. Zee had worked with gravitation expert John A. Wheeler as an undergraduate at Princeton in the mid-1960s before pursuing his Ph.D. in particle theory at Harvard. He earned his degree in 1970, at the same time as the biggest declines in that area began. As he later recalled, cosmology had never even been mentioned while he was in graduate school. After postdoctoral work, Zee began teaching at Princeton. He rented an apartment from a French physicist while on sabbatical in Paris in 1974, and in his borrowed quarters he stumbled on a stack of papers by European theorists that tried to use ideas from particle theory to explain various cosmological features (such as why the observable universe contains more matter than antimatter). Although he found the particular ideas in the papers unconvincing, the chance encounter reignited Zee's earlier interest in gravitation. Returning from his sabbatical and back in touch with Wheeler, Zee began to redirect his research interests toward particle cosmology.

Lee Smolin, in contrast, entered graduate school at Harvard in 1975, just as the curricular changes began to take effect. Smolin studied gravitation and cosmology there alongside his course work in particle theory and worked closely with Stanley Deser (based at nearby Brandeis University), who was visiting Harvard's physics department at the time. Deser was one of the few American



Perhaps today's passionate debates on the direction of theoretical physics are a symptom of the same kinds of growing pains that shook the discipline after the last crash.

theorists who had taken an interest in quantum gravity in the 1960s, attempting to formulate a description of gravitation that would be compatible with quantum mechanics. He was also the very first physicist to publish an article that cited both the Brans-Dicke work and the Higgs work (although he treated the two fields rather differently and in separate parts of his 1972 paper). Smolin, who worked on topics in quantum gravity, suggested that φ_{BD} and φ_H might be the same field as he was finishing his dissertation in 1979.

Smolin's experiences marked the new routine for his generation of theorists, trained during the mid- to late 1970s. Physicists such as Paul J. Steinhardt, Michael S. Turner and Edward "Rocky" Kolb studied gravitation as well as particle theory in graduate school. Soon Smolin, Turner, Kolb, Steinhardt and others were training their own graduate students to work in the new hybrid area of particle cosmology. For these young theorists and their growing numbers of students, it was natural to associate φ_{BD} and φ_H . Turner, Kolb and Steinhardt each led research groups that pursued further links between φ_{BD} and φ_H during the 1980s.

Building on his 1979 paper, Zee noted in 1980 that standard cosmological theories, such as the big bang model, remained unable to account for the extraordinary smoothness of the observable universe (at least when viewed on the largest scales). Separately, Dicke concluded that the big bang also could not explain the observed flatness of the universe, whose shape could in principle depart quite far from the minimal curvature that astronomers observed. In 1981 Alan H. Guth—then a postdoctoral fellow at Stanford University and now a professor at the Massachusetts Institute of Technology—introduced inflationary cosmology to address both of these problems. At the heart of Guth's model was another scalar field, modeled on the Higgs. Dubbed the inflaton, this field provided the driving force behind a postulated period of superfast expansion—or inflation—during the universe's first moments.

drei Linde, then at the Lebedev Physical Institute in Moscow, was likewise poised to explore inflationary ideas: having studied in Russia, where particle physics and gravitation had long flourished side by side, Linde was quick to introduce improvements to the theory.

Since then, it has become routine for particle cosmologists to combine the Brans-Dicke, Higgs and inflaton fields, freely adapting the equations to explain a variety of phenomena. This conceptual leap moved from unthinkable to unnoticeable in only a few academic generations. The shift in attitude illustrates the power of pedagogy and the immense influence that institutional changes can have on scientific thought.

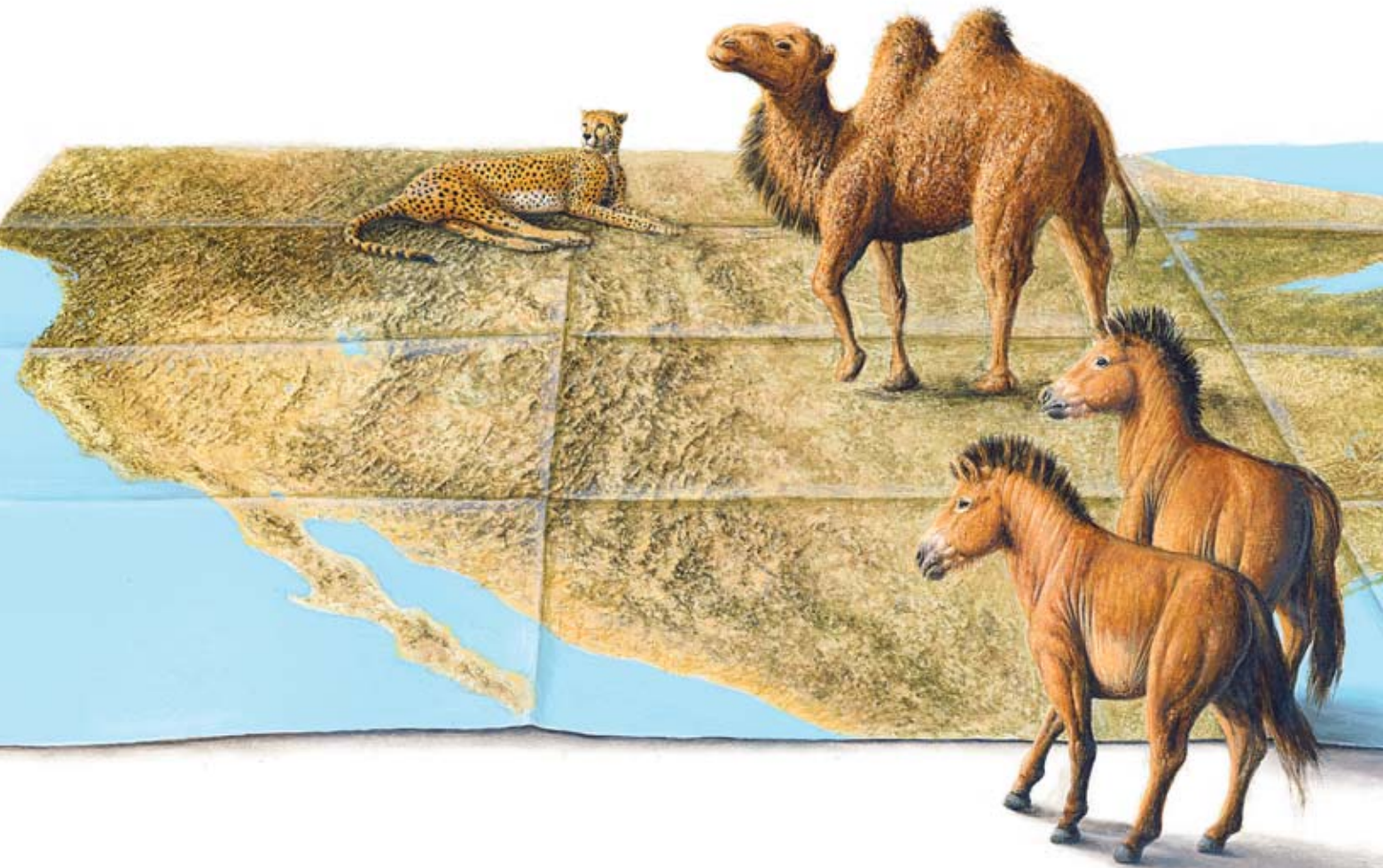
Might history repeat itself? Particle physics was hit hard again in the 1990s (especially with the cancellation of the Superconducting Super Collider, a huge particle accelerator that was under construction in Texas), and funding in the U.S. has continued to slide since then. Perhaps today's passionate debates on the direction of theoretical physics, pitting the advocates of string theory against the proponents of alternative approaches, are a symptom of the same kinds of growing pains that shook the discipline after the last crash.

Physicists are now looking forward to new results from projects scheduled to come online over the next year: the Large Hadron Collider in Switzerland, the Gamma-Ray Large Area Space Telescope and the Planck satellite, which will measure the cosmic microwave background with unprecedented accuracy. With any luck, high-energy physics will emerge just as vibrant and exciting as it did 30 years ago. SA

MORE TO EXPLORE

- Was Einstein Right? Putting General Relativity to the Test.** Second edition. Clifford M. Will. Basic Books, 1993.
- The Inflationary Universe: The Quest for a New Theory of Cosmic Origins.** Alan H. Guth. Addison-Wesley, 1997.
- Einstein's Universe: Gravity at Work and Play.** Anthony Zee. Oxford University Press, 2001.
- Three Roads to Quantum Gravity.** Lee Smolin. Basic Books, 2001.
- Cold War Requisitions, Scientific Manpower, and the Production of American Physicists after World War II.** David Kaiser in *Historical Studies in the Physical and Biological Sciences*, Vol. 33, pages 131–159; 2002.
- Inflationary Cosmology: Exploring the Universe from the Smallest to the Largest Scales.** Alan H. Guth and David Kaiser in *Science*, Vol. 307, pages 884–890; February 11, 2005.

RESTORING AMERICA'S



In the fall of 2004 a dozen conservation biologists gathered on a ranch in New Mexico to ponder a bold plan. The scientists, trained in a variety of disciplines, ranged from the grand old men of the field to those of us earlier in our careers. The idea we were mulling over was the reintroduction of large vertebrates—megafauna—to North America.

Most of these animals, such as mammoths and cheetahs, died out roughly 13,000 years ago, when humans from Eurasia began migrating to the continent. The theory—propounded 40 years ago by Paul Martin of the University of Arizona—is that overhunting by the new arrivals reduced the numbers of large vertebrates so severely that the populations could not recover. Called Pleistocene overkill, the concept was highly controversial at the time, but the general thesis that humans played a significant role is now widely accepted. Martin was present at the meeting in New Mexico, and his ideas on the loss of these animals, the ecological consequences, and what we should do about it formed the foundation of the proposal that emerged, which we dubbed Pleistocene rewilding.

Although the cheetahs, lions and mammoths that once roamed North America are extinct, the same species or close relatives have survived elsewhere, and our discussions focused on introducing these substitutes to North American ecosystems. We believe that these efforts hold the potential to partially restore important ecological processes, such as predation and browsing, to ecosystems where they have been absent for millennia. The substitutes would also bring economic and cultural benefits. Not surprisingly, the published proposal evoked strong reactions. Those reactions are welcome, because debate about the conservation issues that underlie Pleistocene rewilding merit thorough discussion.

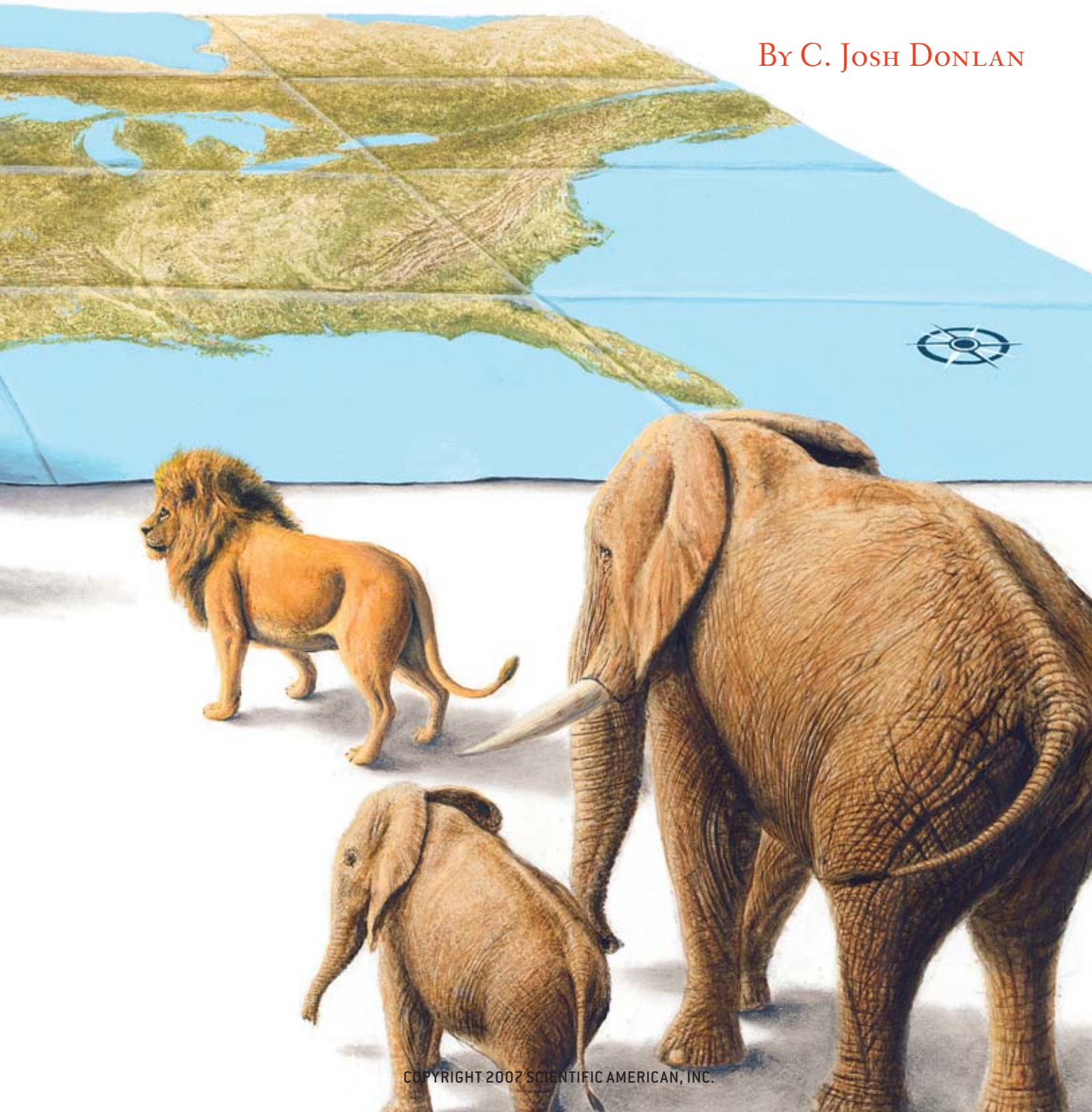
Why Big Animals Are Important

OUR APPROACH concentrates on large animals because they exercise a disproportionate effect on the environment. For tens of millions of years, megafauna dominated the globe, strongly interacting and co-evolving with other species and influencing entire ecosystems. Horses, camels, lions, ele-

BIG, WILD ANIMALS

Pleistocene rewilding—a proposal to bring back animals that disappeared from North America 13,000 years ago—offers an optimistic agenda for 21st-century conservation

BY C. JOSH DONLAN



Questions from Readers

A summary of this article on our Web site invited readers to ask questions about rewilding. Many of the questions helped to shape the article, and answers appear throughout the text. The author replies to a few others here.

phants and other large creatures were everywhere: megafauna were the norm. But starting roughly 50,000 years ago, the overwhelming majority went extinct. Today megafauna inhabit less than 10 percent of the globe.

Over the past decade, ecologist John Terborgh of Duke University has observed directly how critical large animals are to the health of ecosystems and how their loss adversely affects the natural world. When a hydroelectric dam flooded thousands of acres in Venezuela, Terborgh saw the water create dozens of islands—a fragmentation akin to the virtual islands created around the world as humans cut down trees, build shopping malls, and sprawl from urban centers. The islands in Venezuela were too small to support the creatures at the top of the food chain—predators such as jaguars, pumas and eagles. Their disappearance sparked a chain of reactions. Animals such as monkeys, leaf-cutter ants and other herbivores, whose populations were no longer kept in check by predation, thrived and subsequently destroyed vegetation—the ecosystems collapsed, with biodiversity being the ultimate loser.

Similar ecological disasters have occurred on other continents. Degraded ecosystems are not only bad for biodiversity; they are bad for human economies. In Central America, for instance, researchers have shown that intact tropical ecosystems are worth at least \$60,000 a year to a single coffee farm because of the services they provide, such as the pollination of coffee crops.

Where large predators and herbivores still remain, they play pivotal roles. In Alaska, sea otters maintain kelp forest

Won't African elephants and lions freeze to death in the Montana winter? Will these animals be kept indoors in the winter, as they are in zoos? —Jason Raschka

The reader brings up an important point—that many questions, including climate suitability, would have to be answered by sound scientific studies during the process of Pleistocene rewilding. One could imagine scenarios in which animals would be free-living (albeit intensively managed) year-round in expansive reserves in the southwestern U.S. and perhaps other scenarios in locales farther north, where animals would be housed indoors during the coldest months. Even the latter alternative is arguably better than a zoo. Many zoos are struggling to find appropriate space for elephants and are choosing to abandon their elephant programs.



ELEPHANT at the Copenhagen zoo.

Aren't there big-game hunting ranches in Texas? Could we learn anything from them? —Foster

There are many big-game hunting ranches throughout Texas, and some of them hold animals such as cheetahs that we know, from the fossil record, once lived in North America. To my knowledge, no conservation biologists have studied the ranch-held animals, but these ranches might present excellent research opportunities if they are willing to collaborate. Other ranches, however, have animals that are not supported by the fossil record and as such offer no potential for study.

Hasn't this general idea been around for a while? —Kevin N.

Both the concept of rewilding and the term itself have been around for some time. Rewilding is the practice of reintroducing species to places from which they have been extirpated in the past few hundred years. Pleistocene rewilding, in contrast, involves introducing species descended from creatures that went extinct some 13,000 years ago or using similar species as proxies.

ecosystems by keeping herbivores that eat kelp, such as sea urchins, in check. In Africa, elephants are keystone players; as they move through an area, their knocking down trees and trampling create a

habitat in which certain plants and animals can flourish. Lions and other predators control the populations of African herbivores, which in turn influence the distribution of plants and soil nutrients.

In Pleistocene America, large predators and herbivores played similar roles. Today most of that vital influence is absent. For example, the American cheetah (a relative of the African cheetah) dashed across the grasslands in pursuit of pronghorn antelopes for millions of years. These chases shaped the pronghorn's astounding speed and other biological aspects of one of the fastest animals alive. In the absence of the cheetah, the pronghorn appears "overbuilt" for its environment today.

Overview/Rewilding Our Vision

- A group of conservationists has proposed reintroducing to North America large animals that went extinct 13,000 years ago.
- Close relatives of these animals—elephants, camels, lions, cheetahs—survived elsewhere; returning them to America would reestablish key ecological processes that once thrived there, as well as providing a refuge for endangered species from Africa and Asia and creating opportunities for ecotourism.
- The proposal has, understandably, evoked strong reactions, but the bold suggestion has spurred debate and put a positive spin on conservation biology, whose role has been mainly a struggle to slow the loss of biodiversity.

Pleistocene rewilding is not about re-creating exactly some past state. Rather it is about restoring the kinds of species interactions that sustain thriving ecosystems. Giant tortoises, horses, camels, cheetahs, elephants and lions: they were all here, and they helped to shape North American ecosystems. Either the same species or closely related species are available for introduction as proxies, and many are already in captivity in the U.S. In essence, Pleistocene rewilding would help change the underlying premise of conservation biology from limiting extinction to actively restoring natural processes.

At first, our proposal may seem outrageous—lions in Montana? But the plan deserves serious debate for several reasons. First, nowhere on Earth is pristine, at least in terms of being substantially free of human influence. Our demographics, chemicals, economics and politics pervade every part of the planet. Even in our largest national parks, species go extinct without active intervention. And human encroachment shows alarming signs of worsening. Bold actions, rather than business as usual, will be needed to reverse such negative influences. Second, since conservation biology emerged as a discipline more than three decades ago, it has been mainly a business of doom and gloom, a struggle merely to slow the loss of biodiversity. But conservation need not be only reactive. A proactive approach would include restoring natural processes, starting with ones we know are disproportionately important, such as those influenced by megafauna.

Third, land in North America is available for the reintroduction of megafauna. Although the patterns of human land use are always shifting, in some areas, such as parts of the Great Plains and the Southwest, large private and public lands with low or declining human population densities might be used for the project. Fourth, bringing megafauna back to America would also bring tourist and other dollars into nearby communities and enhance the public's appreciation of the natural world. More than 1.5 million people visit San Diego's Wild Animal Park every year to catch a

When the West Was Really Wild

Soon after humans crossed the Bering land bridge into North America some 13,000 years ago, almost 75 percent of the continent's large mammals (those weighing more than 45 kilograms) disappeared (color). One of the goals of Pleistocene rewilding is to restore some of these species or close proxies to the American West. For example, the same species of lion and cheetah that once lived in North America survive today in Africa; the African or Asian elephant could substitute for the extinct mammoth; and Bactrian camels might stand in for the extinct *Camelops*.

Large Mammals of Pleistocene North America

Xenarthra

- Glyptodont (*Glyptotherium floridanum*)
- Harlan's ground sloth (*Paramylodon harlani*)
- Jefferson's ground sloth (*Megalonyx jeffersonii*)
- Shasta ground sloth (*Nothrotheriops shastensis*)



Carnivores (Carnivora)

- Dire wolf (*Canis dirus*)
- Gray wolf (*Canis lupus*)
- Black bear (*Ursus americanus*)
- Brown bear (*Ursus arctos*)
- Giant short-faced bear (*Arctodus simus*)
- Saber-toothed cat (*Smilodon fatalis*)
- American lion (*Panthera leo*)
- Jaguar (*Panthera onca*)
- American cheetah (*Miracinonyx trumani*)
- Mountain lion (*Puma concolor*)



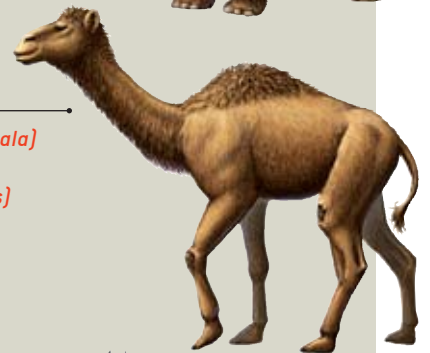
Elephants (Proboscidea)

- American mastodon (*Mammuthus americanum*)
- Columbian mammoth (*Mammuthus columbi*)
- Dwarf mammoth (*Mammuthus exilis*)
- Woolly mammoth (*Mammuthus primigenius*)



Horses (Perissodactyla)

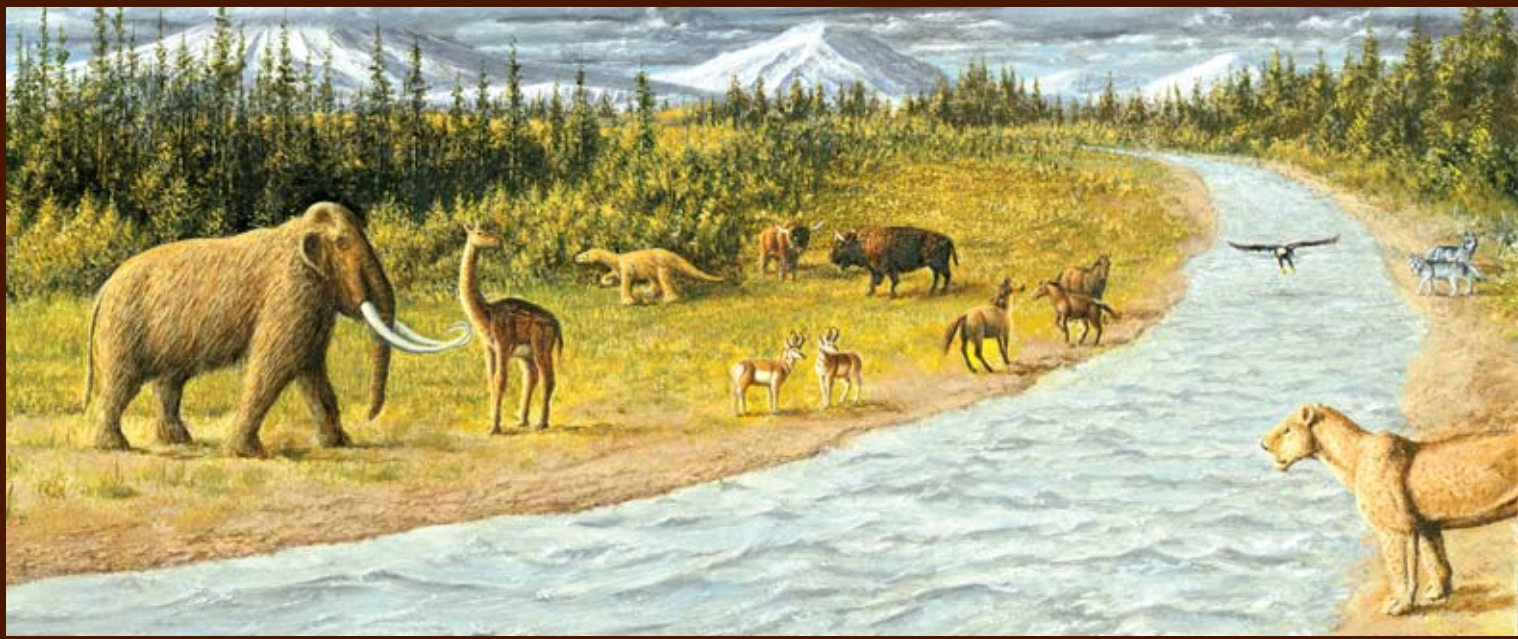
- Mexican horse (*Equus conversidens*)
- Western horse (*Equus occidentalis*)
- Other extinct horses and asses (*Equus spp.*)



Even-Toed Ungulates (Artiodactyla)

- Western camel (*Camelops hesternus*)
- Long-legged llama (*Hemiauchenia macrocephala*)
- Long-nosed peccary (*Mylohyus nasutus*)
- Flat-headed peccary (*Platygonus compressus*)
- Mule deer (*Odocoileus hemionus*)
- White-tailed deer (*Odocoileus virginianus*)
- Mountain deer (*Navahoceros fricki*)
- Woodland caribou (*Rangifer tarandus*)
- Moose (*Alces alces*)
- Wapiti (*Cervus elaphus*)
- Pronghorn (*Antilocapra americana*)
- Harrington's mountain goat (*Oreamnos harringtoni*)
- Mountain goat (*Oreamnos americanus*)
- Bighorn sheep (*Ovis canadensis*)
- Shrub ox (*Euceratherium collinum*)
- Bonnet-headed musk ox (*Bootherium bombifrons*)
- Bison (*Bison bison*)
- Extinct bison (*Bison spp.*)





Past, Present and Future

These artist's conceptions illustrate some of the ways that large animals influence the ecosystem and how Pleistocene rewilding could be beneficial. David Burney of the National Tropical Botanical Garden, who is a member of the Pleistocene rewilding group, worked with artist Larry Felder in conceiving the illustrations.

◀ LATE APRIL DURING THE PLEISTOCENE, 14,000 YEARS AGO

Humans have not yet appeared on the landscape, but glaciers have receded and the climate of what is now the western U.S. is similar to modern conditions, though a little cooler. Dense snow covers the Rocky Mountains. The river, flushed by snowmelt, is fairly high and a milky blue because of the fine particles carried from the young, postglacial soils. The banks are grazed down short but are generally green. A Columbian mammoth, a *Camelops*, a Shasta ground sloth, long-horned bison, horses and pronghorn antelope browse, while an American lion and two wolves watch from the opposite bank and a bald eagle glides over the river.

◀ APRIL 2007

The Rocky Mountains have snow on their upper summits only. The river, muddied by spring rains and by too many cattle on the banks, has widened, and the banks are eroded and broken down by hoofprints and cattle wallows. Grass is short, patches of soil are bare, thorny shrubs predominate. Bison graze in the background, and a wary coyote stands on the opposite shore.

◀ APRIL 2027

Vegetation is more summerlike—global warming conditions are in full swing. The mountains have no snow on the peaks. Drought, low humidity and absence of mountain snow have lowered the level of the river dramatically; the water is clear and dark. Erosion is less apparent; grass looks manicured from diversified grazing and browsing. An Asian elephant, Bactrian camel, bison, antelope and Przewalski's horses graze under the watchful eye of an African lion; a cheetah approaches from the distance. In the background a monorail carries tourists. A high-tech electrified fence with solar panels is out of view.

glimpse of large mammals. Only a handful of U.S. national parks receive that many visitors. Last, the loss of some of the remaining species of megafauna in Africa and Asia within this century seems likely—Pleistocene rewilding could help reverse that.

How It Might Be Done

WE ARE NOT TALKING about backing up a van and kicking some cheetahs out into your backyard. Nor are we talking about doing it tomorrow. We conceive of Pleistocene rewilding as a series of staged, carefully managed ecosystem manipulations. What we are offering here is a vision—not a blueprint—of how this might be accomplished. And by no means are we suggesting that rewilding should be a priority over current conservation programs in North America or Africa. Pleistocene rewilding could proceed alongside such conservation efforts, and it would likely generate conservation dollars from new funding sources, rather than competing for funds with existing conservation efforts.

The long-term vision includes a vast, securely fenced ecological history park, encompassing thousands of square miles, where horses, camels, elephants and large carnivores would roam. As happens now in Africa and regions surrounding some U.S. national parks, the ecological history park would not only attract ecotourists but would also provide jobs related both to park management and to tourism.

To get to that distant point, we would need to start modestly, with relatively small-scale experiments that assess the impacts of megafauna on North American landscapes. These controlled experiments, guided by sound science and by the fossil record, which indicates what animals actually lived here, could occur first on donated or purchased pri-

vate lands and could begin immediately. They will be critical in answering the many questions about the reintroductions and would help lay out the costs and benefits of rewilding.

One of these experiments is already under way. Spurred by our 2004 meeting, biologists recently reintroduced Bolson tortoises to a private ranch in New Mexico. Bolson tortoises, some weighing more than 100 pounds, once grazed parts of the southwestern U.S. before disappearing around 10,000 years ago, victims of human hunting. This endangered tortoise now clings to survival, restricted to a single small area in central Mexico. Thus, the reintroduction not only repatriates the tortoise to the U.S., it increases the species' chance for survival. Similar experiments are also occurring outside North America [see box on page 77].

The reintroduction of wild horses and camels would be a logical part of these early experiments. Horses and camels originated on this continent, and many species were present in the late Pleistocene. Today's feral horses and asses that live in some areas throughout the West are plausible substitutes for extinct American species. Because most of the surviving Eurasian and African species are now critically endangered [see "Endangered Wild Equids," by Patricia D. Moehlan; *SCIENTIFIC AMERICAN*, March 2005], establishing Asian asses and Przewalski's horse in North America might help prevent the extinction of these animals. Bactrian camels, which are critically endangered in the Gobi Desert, could provide a modern proxy for *Camelops*, a late Pleistocene camel. Camels, introduced from captive or domesticated populations, might benefit U.S. ecosystems by browsing on woody plants that today are overtaking arid grasslands in the Southwest, an ecosystem that is increasingly endangered.

LARRY FELDER
THE AUTHOR

C. JOSH DONLAN holds a Ph.D. in ecology and evolutionary biology from Cornell University, where he is a research biologist. Founder and director of Advanced Conservation Strategies, he serves as an adviser to the Galápagos National Park and to Island Conservation and is a senior fellow at the Robert and Patricia Switzer Foundation and Environmental Leadership Program. He was highlighted in the *New York Times Magazine*'s "Big Ideas of 2005" issue and named one of 25 "all-star" innovators for 2005 by *Outside* magazine. He spends much of his time in Tasmania, Australia and Santa Cruz, Calif., trying to halt extinctions on islands.

Another prong of the project would likely be more controversial but could also begin immediately. It would establish small numbers of elephants, cheetahs and lions on private property.

Introducing elephants could prove valuable to nearby human populations by attracting tourists and maintaining grasslands useful to ranchers (elephants could suppress the woody plants that threaten southwestern grasslands). In the late Pleistocene, at least four elephant species lived in North America. Under a scientific framework, captive elephants in the U.S. could be introduced as proxies for these extinct animals. The biggest cost involved would be fencing, which has helped reduce conflict between elephants and humans in Africa.

Many cheetahs are already in captivity in the U.S. The greatest challenge would be to provide them with large, securely fenced areas that have appropriate habitat and prey animals. Offsetting these costs are benefits—restoring what must have been strong interactions with

pronghorn, facilitating ecotourism as an economic alternative for ranchers, many of whom are struggling financially, and helping to save the world's fastest carnivore from extinction.

Lions are increasingly threatened, with populations in Asia and some parts of Africa critically endangered. Bringing back lions, which are the same species that once lived in North America, presents daunting challenges as well as many potential benefits. But private reserves in southern Africa where lions and other large animals have been successfully reintroduced offer a model—and these reserves are smaller than some private ranches in the Southwest.

If these early experiments with large herbivores and predators show promising results, more could be undertaken, moving toward the long-term goal of a huge ecological history park. What we need now are panels of experts who, for each species, could assess, advise and cautiously lead efforts in restoring megafauna to North America.

A real-world example of how the re-introduction of a top predator might work comes from the wolves of Yellowstone National Park [see “Lessons from the Wolf,” by Jim Robbins; *SCIENTIFIC AMERICAN*, June 2004]. The gray wolf became extinct in and around Yellowstone during the 1920s. The loss led to increases in their prey—moose and elk—which in turn reduced the distribution of aspens and other trees they eat. Lack of vegetation destroyed habitat for migratory birds and for beavers. Thus, the disappearance of the wolves propagated a trophic cascade from predators to herbivores to plants to birds and beavers. Scientists have started to document the ecosystem changes as reintroduced wolves regain the ecological role they played in Yellowstone for millennia. An additional insight researchers are learning from putting wolves back into Yellowstone is that they may be helping the park cope with climate change. As winters grow milder, fewer elk die, which means less carrion for scavengers such

What the Critics Say

Since we proposed the rewilding idea in print, in *Nature* in 2005, commentators have pointed out a number of concerns, some legitimate, some not. “We all remember *Jurassic Park*,” wrote Dustin Rubenstein and his colleagues in the journal *Biological Conservation*. “Pleistocene re-wilding of North America is only a slightly less sensational proposal.” We disagree with this assessment because the majority of dinosaurs went extinct 65 million years ago, whereas many of North America’s original megafauna or very close relatives are alive today elsewhere in the world and can be both studied and saved.

Rubenstein and colleagues go on to say, “*Modern day proxies species are wrong ... different genetically from the species that occurred in North America during the Pleistocene.*” True, but not that different. Available evidence indicates that the lions in Africa and Asia today are the same species, albeit of smaller stature, as the lions that prowled North America 13 millennia ago. Recent studies of ancient DNA have elucidated close relationships between extinct elephant and horse species and those alive today. Further, introduction of proxies for now extinct species has proved successful in other experiments. Hundreds of peregrine falcons from Australia, Europe and South America were used, for example, in captive-breeding programs to reintroduce the peregrine falcon to parts of

the U.S. and Canada where DDT had wiped it out. Those birds were certainly different genetically from the ones that once soared over the Midwest, yet they have done well in their new homes.

“One can only imagine ... farmers coping with crop destruction by herds of elephants, or lions and cheetah attacking cattle, or even children,” Rubenstein and his colleagues warn. One

need not merely imagine the challenges of coexisting with large predators and herbivores. Africa and Asia have been struggling with them for centuries, and substantial progress has been made; in our plan the animals will not be unrestrained.

“Global climate change since the Pleistocene extinctions makes the restoration of vanished ecosystems through large-mammal introduction quite unlikely,” wrote Christopher Smith in a letter to *Nature* in the weeks following the initial proposal. Many have expressed concern about the fact that North America’s ecosystems are not the same today as they were 13,000 years ago and that reintroduced animals might therefore be unable to survive in them. Whereas habitats are and will continue to be dynamic on a timescale of thousands of years, very few plants or small mammals went extinct during the late Pleistocene. The major missing component of North American ecosystems today compared with the Pleistocene is the megafauna, which we can infer are critical cogs in the wheels. —C.J.D.



Elsewhere in the World

as coyotes, ravens and bald eagles. Wolves provide carcasses throughout the winter for the scavengers to feed on, bestowing a certain degree of stability.

The Challenges Ahead

AS OUR GROUP on the ranch in New Mexico discussed how Pleistocene rewilding might work, we foresaw many challenges that would have to be addressed and overcome. These include the possibility that introduced animals could bring novel diseases with them or that they might be unusually susceptible to diseases already present in the ecosystem; the fact that habitats have changed over the millennia and that reintroduced animals might not fare well in these altered environments; and the likelihood of unanticipated ecological consequences and unexpected reactions from neighboring human communities. Establishing programs that monitor the interactions among species and their consequences for the well-being of the ecosystem will require patience and expertise. And, of course, it will not be easy to convince the public to accept predation as an important natural process that actually nourishes the land and enables ecosystems to thrive. Other colleagues have raised additional concerns, albeit none that seems fatal [see box on opposite page].

Many people will claim that the concept of Pleistocene rewilding is simply not feasible in the world we live in today. I urge these people to look to Africa for inspiration. The year after the creation of Kruger National Park was announced, the site was hardly the celebrated mainstay of southern African biodiversity it is today. In 1903 zero elephants, nine lions, eight buffalo and very few cheetahs lived within its boundaries. Thanks to the vision and dedication of African conservationists, 7,300 elephants, 2,300 lions, 28,000 buffalo and 250 cheetahs roamed Kruger 100 years later—as did 700,000 tourists, bringing with them tens of millions of dollars.

In the coming century, humanity will decide, by default or design, the extent to which it will tolerate other species and thus how much biodiversity will endure. Pleistocene rewilding is not about trying

In other parts of the world, preliminary efforts have begun to reintroduce species to places from which they have long been absent.

- In April 2006 a team of Canadian and Russian biologists flew 30 wood bison from Canada's Elk Island National Park to Pleistocene Park reserve in the Republic of Sakha, Russia, where the closely related steppe bison vanished 5,000 years ago.
- At the 15,000-acre nature reserve Oostvaardersplassen in the Netherlands, conservationists are restoring horses, roe deer and Heck cattle.
- Beavers are being reintroduced throughout Europe, in some cases in areas where they have been absent for thousands of years.
- In the tropical Pacific, endangered birds from the Marquesas and Tongan islands have been reintroduced to nearby islands where fossils indicate they once lived.
- In the Indian Ocean, scientists from the Mauritian Wildlife Foundation are using giant tortoises from Aldabra Island to replace two extinct species of tortoise on the Mascarene Islands. Scientists have already documented increased seed dispersal for many of the island plants. The tortoises have also brought increased tourism.

—C.J.D.



KONIK HORSES stand in for the extinct tarpan horse at the Oostvaardersplassen reserve in the Netherlands.

to go back to the past; it is about using the past to inform society about how to maintain the functional fabric of nature. The potential scientific, conservation and cultural benefits of restoring megafauna are clear, as are the costs. Although sound science can help mitigate the potential costs, these ideas will make many uneasy. Yet given the apparent dysfunction of North American ecosystems and Earth's overall state, inac-

tion carries risks as well. In the face of tremendous uncertainty, science and society must weigh the costs and benefits of bold, aggressive actions like Pleistocene rewilding against those of business as usual, which has risks, uncertainties and costs that are often unacknowledged. We have a tendency to think that if we maintain the status quo, things will be fine. All the available information suggests the opposite. SA

MORE TO EXPLORE

The Ghosts of Evolution: Nonsensical Fruit, Missing Partners, and Other Ecological Anachronisms. Connie Barlow. Basic Books, 2000.

Twilight of the Mammoths: Ice Age Extinctions and the Rewilding of America. Paul S. Martin. University of California Press, 2005.

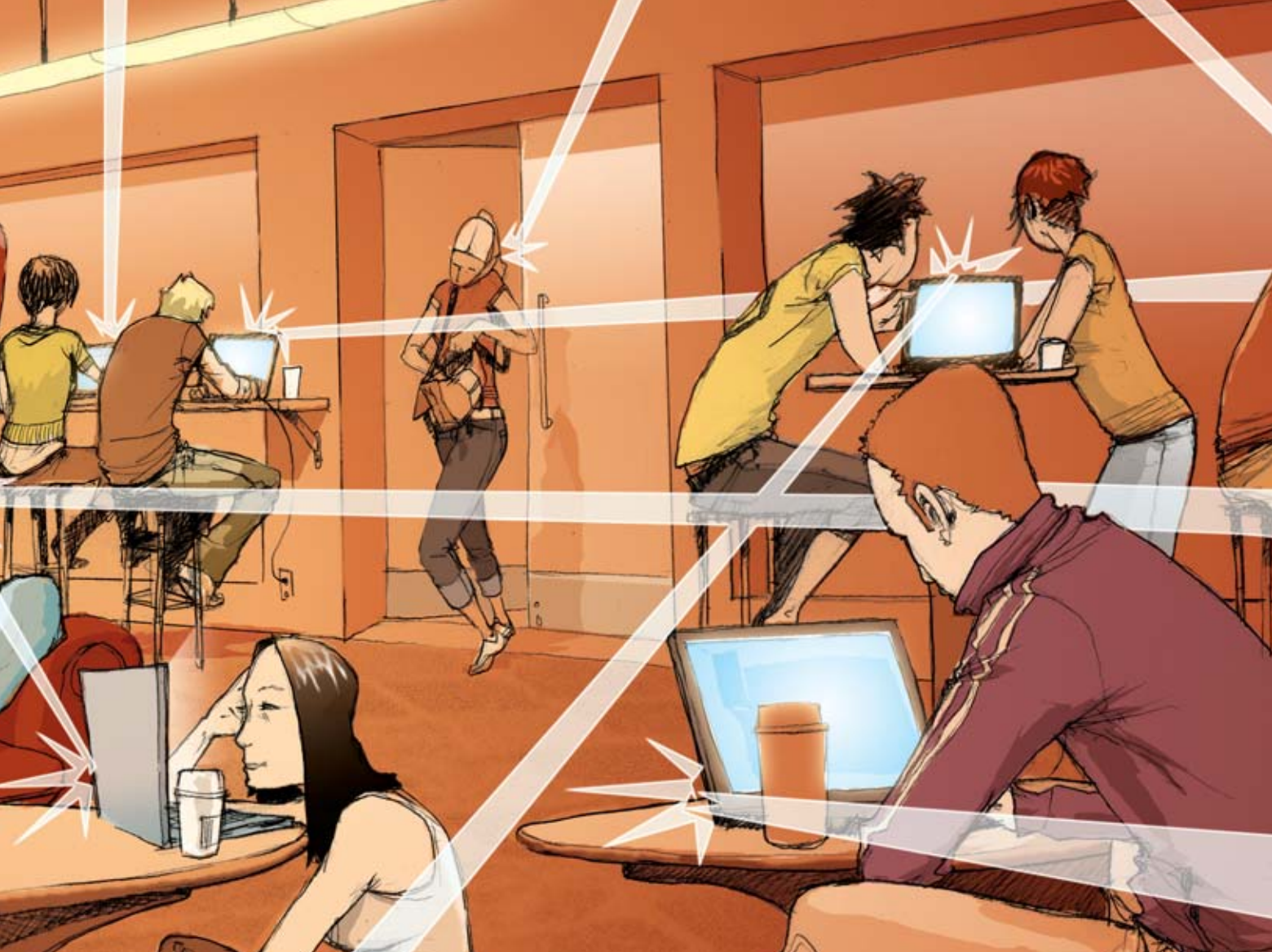
Pleistocene Rewilding: An Optimistic Agenda for Twenty-First Century Conservation. C. J. Donlan, J. Berger, C. E. Bock, J. H. Bock, D. A. Burney, J. A. Estes, D. Foreman, P. S. Martin, G. W. Roemer, F. A. Smith, M. E. Soulé and H. W. Greene in *American Naturalist*, Vol. 168, No. 5, pages 660–681; November 2006.



Breaking Network Logjams

An approach called network coding could dramatically enhance the efficiency and reliability of communications networks. At its core is the strange notion that transmitting evidence about messages can be more useful than conveying the messages themselves

**By Michelle Effros,
Ralf Koetter and
Muriel Médard**



USERS SHARING A NETWORK based on network coding would enjoy many benefits, including seeing fewer delays when downloading videos and accessing Web sites.

The history of modern communications systems has been marked by flashes of startling insight.

Claude E. Shannon, mathematician and engineer, launched one such revolution almost 60 years ago by laying the foundation of a new mathematical theory of communications—now known as information theory. Practical outgrowths of his work, which dealt with the compression and reliable transmission of data, can be seen today in the Internet, in landline and wireless telephone systems, and in storage devices, from hard drives to CDs, DVDs and flash memory sticks.

Shannon tackled communications over phone lines dedicated to individual calls. These days, information increasingly travels over shared networks (such as the Internet), in which multiple users simultaneously communicate through the same medium—be it a cable, an optical fiber or, in a wireless system, air. Shared networks can potentially improve the usefulness and efficiency of communications systems, but they also create competition for communal resources. Many people must vie for access to, say, a server offering downloadable songs or to a wireless hot spot.

The challenge, then, is to find ways to make the sharing go smoothly; parents of toddlers will recognize the problem. Network operators frequently try to solve the challenge by increasing resources, but that strategy is often insufficient. Copper wires, cables or fiber optics, for instance, can now provide high bandwidth for commercial and residential users yet are expensive to lay and difficult to modify and expand. Ultra-wideband and multiple-antenna transmission systems can expand the number of customers served by wireless networks but may still fail to meet ever increasing demand.

Techniques for improving efficiency are therefore needed as well. On the Internet and other shared networks, information currently gets relayed by routers—switches that operate at nodes where signaling pathways, or links, intersect. The routers shunt incoming messages to links heading toward the messages' final destinations. But if one wants efficiency, are routers the best devices for these intersections? Is switching even the right operation to perform?

Until seven years ago, few thought to ask such questions. But then Rudolf Ahlswede of the University of Bielefeld in Germany, along with Ning Cai, Shuo-Yen Robert Li and Raymond W. Yeung, all then at the University of Hong Kong, published groundbreaking work that introduced a new approach to distributing information across shared networks. In this approach, called network coding, routers are replaced by coders, which transmit evidence about messages instead of sending the messages themselves. When receivers collect the evidence, they deduce the original information from the assembled clues.

Although this method may sound counterintuitive, network coding, which is still under study, has the potential to dramatically speed up and improve the reliability of all manner of communications systems and may well spark the next revolution in the field. Investigators are, of course, also exploring additional avenues for improving efficiency; as far as we know, though, those other approaches generally extend existing methods.

Bits Are Not Cars

AHLSWEDE AND HIS COLLEAGUES built their proposal in part on the idea, introduced by Shannon, that transmitting evidence about data can actually be more useful than conveying the data directly. They also realized that a receiver would be able to deduce the original data once enough clues had been gathered but that the receiver would not need to obtain *all* of the evidence emitted. One kind of clue could be replaced by another, and all that was important was receiving some combination of clues that, together, would reveal the original message. (Receivers would be able to make sense of the evidence if they were informed in advance about the rules applied to generate it or if instructions on how to use the evidence were included in the evidence itself.)

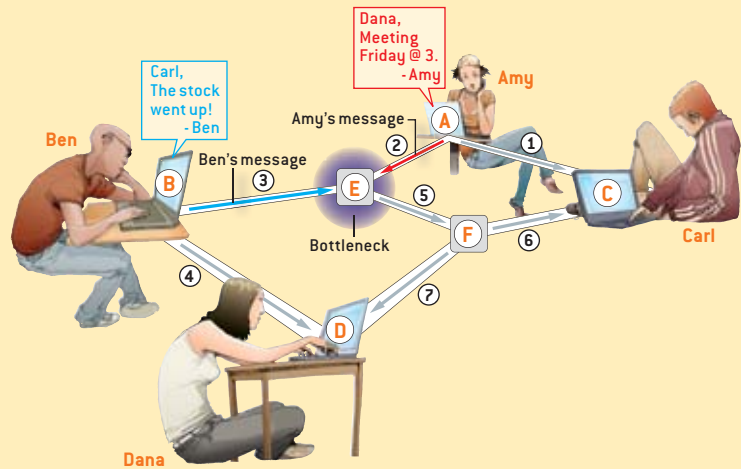
Network coding breaks with the classic view that communications channels are analogous to roads and that bits are like

Overview/Network Coding

- In 2000 investigators proposed a seemingly crazy idea for limiting logjams in communications networks. Called network coding, this potentially revolutionary approach replaces routers, which simply relay messages at intersections, with network coders, which send evidence about the incoming messages instead of the messages themselves.
- Network coding is faring well in experiments, most of which so far focus on sending data across multicast networks, where all receivers need to get the same information simultaneously.
- It promises to make the operations of many networks more efficient (increasing capacity without having to add hardware or bandwidth) as well as faster, more reliable and better resistant to attack.

THE BASIC IDEA

A simple six-node network, in which messages travel along links (numbered) at one bit per second, can illustrate the basic idea of network coding. Nodes are labeled with letters.



WHAT HAPPENS IN A STANDARD NETWORK

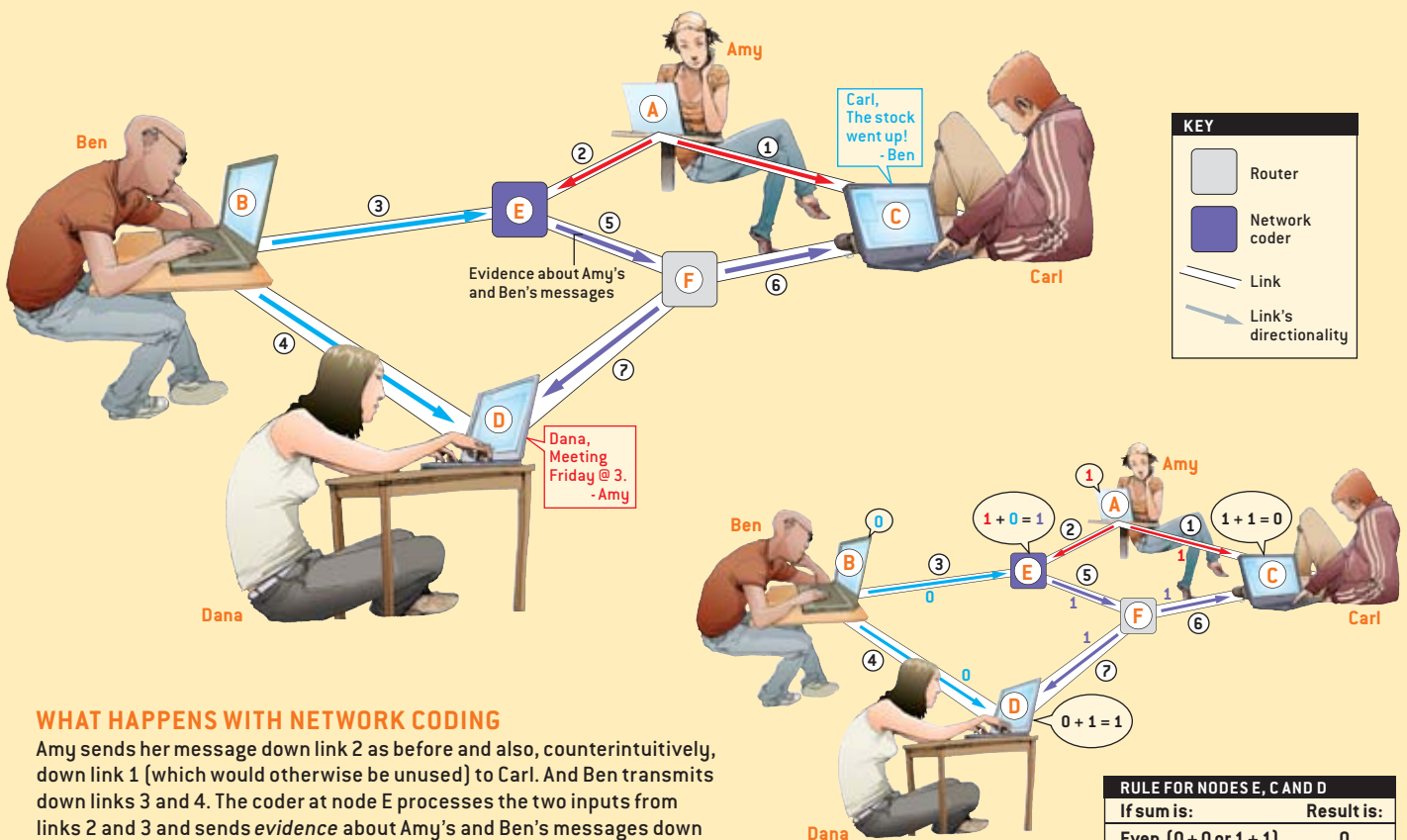
If Amy at node A sends a message to Dana at node D just when Ben at node B transmits to Carl at node C, the messages hit a bottleneck at node E, where a router would have to relay one message at a time down link 5, causing delays.

the cars that travel those roads. But an understanding of the transportation model of communications is useful for grasping how the new scheme works and why it has such promise.

Shannon proved mathematically that every channel has a capacity—an amount of information it can relay during any given time frame—and that communications can proceed reliably as long as the channel’s capacity is not exceeded. In the transportation analogy, a road’s capacity is the number of cars per second it can handle safely. If traffic stays below capacity, a car entering the road at one end can generally be guaranteed to exit at the other end unchanged (barring the rare accident). Engineers have built increasingly complex communications systems based on the transportation model. For example, the phone systems Shannon pondered dedicate a distinct “road” to every conversation; two calls over traditional phone lines never share a single line at the same time and frequency.

Computer networks—and the Internet in particular—are essentially a maze of merging, branching and intersecting roads. Information traveling from one computer to another typically traverses several roads en route to its destination. Bits from a single message are grouped into packets (the car-pools or buses of the information superhighway), each of which is labeled with its intended destination. Routers sit at the intersections of the roads, examine each packet’s header and forward that packet toward its destination.

Ironically, the very transportation model that fueled today’s sophisticated communications systems now stands in



WHAT HAPPENS WITH NETWORK CODING

Amy sends her message down link 2 as before and also, counterintuitively, down link 1 (which would otherwise be unused) to Carl. And Ben transmits down links 3 and 4. The coder at node E processes the two inputs from links 2 and 3 and sends *evidence* about Amy's and Ben's messages down link 5, with no delay. This evidence then gets relayed down links 6 and 7. Carl's computer deciphers Ben's message by examining Amy's message and the evidence received from link 6, and Dana's computer deciphers Amy's message analogously. How can that be? Read on.

RULE FOR NODES E, C AND D	
If sum is:	Result is:
Even (0 + 0 or 1 + 1)	0
Odd (0 + 1 or 1 + 0)	1

HOW DECODING WORKS

Imagine that Amy's message is the number 1 and that Ben's is a 0 and that the coder at node E, as well as Carl's and Dana's computers, processes two incoming data streams according to the code, or rule, in the table. Node E would send a 1 down link 5, and Carl's computer would combine the 1 from link 1 and the 1 from link 6 to produce a 0, thus recovering Ben's message. Meanwhile Dana's computer would combine the 1 from link 7 and the 0 from link 4 to produce a 1—Amy's message. In reality, if nodes C and D were not preprogrammed with the rule, the coder at node E would transmit it to them along with the evidence it sent down link 5.

the way of progress. After all, bits are not cars. When two vehicles converge on the same narrow bridge, they must take turns traversing the bottleneck. When two bits arrive at a bottleneck, however, more options are possible—which is where network coding comes in.

How It Works

THE HYPOTHETICAL SIX-NODE digital network depicted in the box on these two pages can help clarify those options. Recall that in computers, all messages take the form of a string of binary code. Imagine that each link, or road, in this network can carry one bit—be it a 0 or a 1—per second and only in the direction designated by the corresponding arrow. Amy, a network user at node A, hopes to send information at one bit per second to Dana at node D. Meanwhile Ben at node B hopes to send, at exactly the same time and rate, information to Carl at node C. Can both Amy's and Ben's demands be satisfied simultaneously without exceeding any of the links' capacities?

In a router system [see *leftmost illustration*], the outlook seems bleak. Both paths, from Amy to Dana and from Ben to Carl, require traversing link 5. This link becomes the equivalent of a narrow, one-lane bridge. The router at node E, where

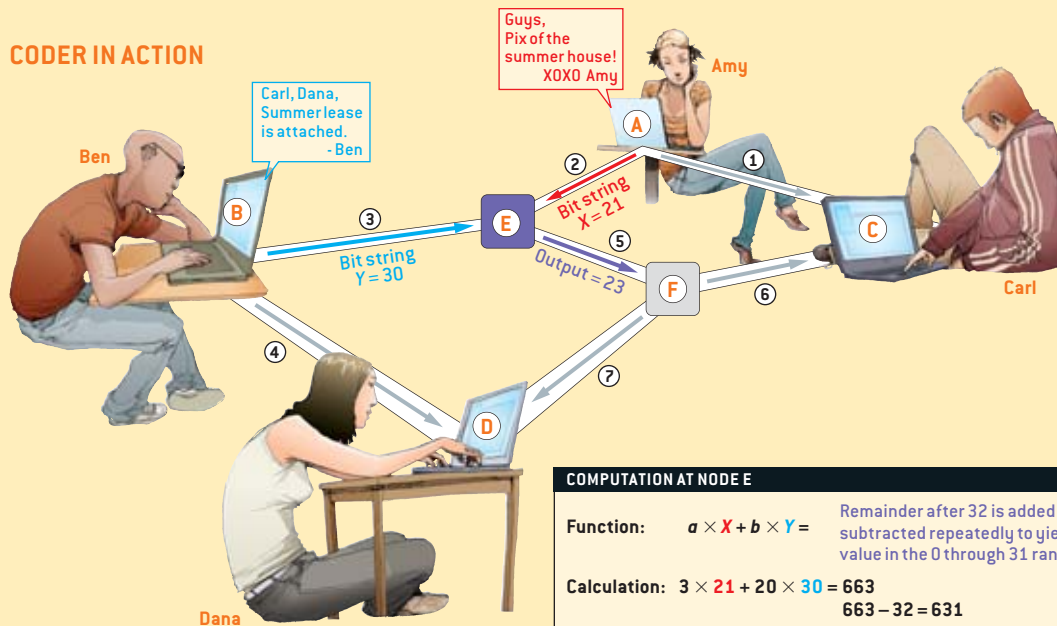
link 5 starts, receives a total of two bits per second (one from link 2 and one from link 3), but because link 5's capacity is one, the router can send only one bit per second along it. In the transportation model, such bottlenecks cause nightmare traffic jams, with more and more bits piling up over time, waiting their turn.

In the new approach [see *illustrations above*], though, the plain router would be replaced by a coder, which would have more options than would be open to a traffic cop. Instead of relaying the actual bit streams collected at the bottleneck, the coder could send quite different information. It could, for example, add up the number of 1s that arrive during any given second and transmit a 0 if that sum is even. If the sum is odd, the device could transmit a 1. So, if link 5 simultaneously re-

A CODE DESIGN EXAMPLE

To gain more of a taste for how coders would process information and to see why designers of multicast networks can easily select the rules that coders will use, consider our six-node network again. This time, though, Amy and Ben want to send information to both Carl and Dana. This is a multicast problem because both receivers hope to gather all the data.

CODER IN ACTION



COMPUTATION AT NODE E

Function: $a \times X + b \times Y =$ Remainder after 32 is added or subtracted repeatedly to yield a value in the 0 through 31 range

Calculation: $3 \times 21 + 20 \times 30 = 663$
 $663 - 32 = 631$
 $631 - 32 = 599$
 \vdots
 $55 - 32 = 23$
OUTPUT = 23

FIVE-BIT BINARY STRINGS		
STRING $m = 5$	INTEGER	VALUE IN EXAMPLE
00000	0	
00001	1	
00010	2	
00011	3	<i>a</i>
00100	4	
00101	5	
00110	6	
00111	7	
01000	8	
01001	9	
01010	10	
01011	11	
01100	12	
01101	13	
01110	14	
01111	15	
10000	16	
10001	17	
10010	18	
10011	19	
10100	20	<i>b</i>
10101	21	<i>X</i>
10110	22	
10111	23	
11000	24	
11001	25	
11010	26	
11011	27	
11100	28	
11101	29	
11110	30	<i>Y</i>
11111	31	

The coding challenge here basically boils down to selecting a coding rule, or mathematical function, for node E, which receives information from links 2 and 3 and transmits the output of the function on link 5 (above). Many options would achieve our communication objective, but let us focus on linear functions, to show how simple the codes can be. In addition, because network coding is typically applied to blocks of bits, we will design the function for link 5 as a random combination of some set number (m) of bits from link 2 and an equal number of bits from link 3.

One approach we can take is to make $m = 5$. There are 32 possible binary strings of length 5, and we represent each one by

an integer from 0 through 31 (table at right). We can also combine the strings at node E using the function $a \times X + b \times Y$, where X and Y are the values between 0 and 31 that correspond to the strings node E receives from links 2 and 3 (say, 21 and 30 at one instant) and where a and b are also binary strings of length m (set arbitrarily here at 3 and 20). We also ensure that the output of this linear combination is itself a member of the finite set 0 through 31. This we do (as shown in computation box) by plugging the values of a , b , X and Y into the formula, then repeatedly adding or subtracting 32 from the initial result until we reach a number in the range 0 through 31. That last number becomes the output.

ceives a 1 and a 0 from links 2 and 3, it carries a 1. If either two 0s or two 1s are received from links 2 and 3, link 5 carries a 0. The result then gets sent by router F down links 6 and 7 to Carl and Dana, respectively.

This approach replaces each pair of bits at node E with a hybrid of the two. Such a bit stream seems ridiculous. Our proposed coder has done the equivalent of combining one phone conversation with another in a way that obscures both. The apparent absurdity of the approach is precisely why it went uninvestigated for so long.

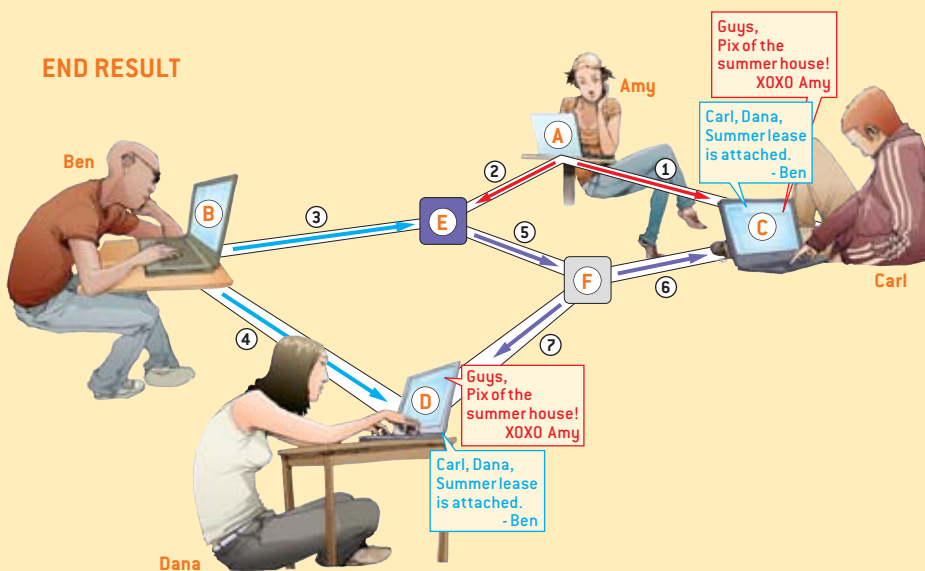
But sometimes apparent madness is true innovation. A hybrid bit stream may describe neither transmission perfectly, yet it can supply evidence about both. Suppose we additionally send Amy's missive to Carl along link 1 and Ben's to Dana along link 4. Sending these two messages uses network re-

sources (links 1 and 4) that the routing system could not usefully employ for meeting Amy's and Ben's demands. Carl's node receives Amy's transmission and knows for each instant (from link 6) whether the number of 1s in the pair of messages issued by Amy and Ben is even or odd. If Carl's node is programmed to also "know" the rule used by the coders at the start of link 5 or if it can infer the rule from the evidence itself, the collected evidence will enable it to decipher the message sent by Ben. And Dana's node will similarly uncover Amy's message.

Clear Benefits

THIS STRATEGY ACCOMPLISHES two goals that were unthinkable given the limitations of the transportation model. First, it enables the bit leaving a node to travel two paths

END RESULT



If Amy sends her message down link 1 to Carl (in addition to sending it down link 2 to node E) and if Ben sends his message down link 4 as well as down link 3 (above), Carl's and Dana's computers will be able to plug the output from E into the same function applied by E and deduce Y and X, respectively. Repeating the same process for each bit string ultimately reveals both messages. [To see exactly how the computers uncover the values of X and Y and acquire the values for a and b, visit www.sciam.com/ontheweb]

simultaneously, something a car cannot do. Second, it allows a pair of bit streams arriving at the head of a bottleneck to combine into a single stream, whereas two cars converging on one narrow bridge cannot become a single entity; one would have to wait for the other to pass before it could proceed across the bridge.

The data-handling approach exemplified by our six-node model (a minor variation on one first given by Ahlswede and his colleagues in 2000) can potentially increase the capacity of a network without requiring the addition of extra conduits because it avoids logjams. Using routing alone, our six-node network could sustain simultaneous transmissions averaging one half of a bit per second. (Because the two competing transmissions would have to share link 5, the effective data rate would be one bit per two seconds, or one half of a bit per second, for each of the competing demands.) With network coding, the same system supports simultaneous transmissions at one bit per second. So, here, network coding doubles capacity.

Sometimes network coding could yield even bigger capacity gains, sometimes none. But the approach would never decrease the capacity of a network because, at worst, it would precisely mimic the actions of router systems. It should also increase reliability and resistance to attacks in relatively sub-

WHY SELECTING CODES IS EASY

At first, selecting codes randomly might seem like a bad idea. Because a and b can each take any value in the set 0 through 31, there are 32×32 , or 1,024, possible functions of the form $a \times X + b \times Y$ that we can use at node E, and not all of those functions will ensure that both receivers can determine the original messages X and Y. For example, if we use a code in which $a = 0$, Dana gets Y on link 4 and $b \times Y$ on link 7, and this information will not suffice for obtaining message X, no matter what the value of b is.

Still, when neither a nor b is 0, both receivers can retrieve the proper messages successfully. As a result, all but 63 of the 1,024 possible functions available to node E are good. Even choosing the coefficients at random yields a good code roughly 94 percent of the time. The probability of success approaches 100 percent very rapidly as we increase m , so all that programmers need do to ensure reliability is select an m that is sufficiently high.

stantial networks, because the interchangeable nature of evidence means that some packets of evidence can be lost without creating problems.

Lessons from Multicast Networks

SO FAR MUCH of the research into implementing network coding has focused on multicast networks—in which all receivers need to get the same information. Internet video games rely on multicast systems to update every player each time one makes a move. Webcasts of videos or live sporting events and new software released electronically to a large group of customers also travel over multicast networks. Today such networks still use routers, and a return to the transportation analogy helps to explain why designing them is usually quite difficult.

Imagine the country's highways teeming with cars. Each router is like a police officer directing traffic at a single intersection. Incoming cars join the queue behind vehicles that arrived before them. The officer reads each car's destination

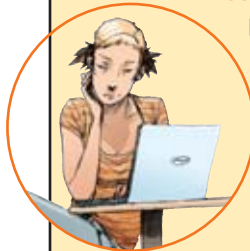
THE AUTHORS

MICHELLE EFFROS, RALF KOETTER and MURIEL MÉDARD are longstanding collaborators and friends. Effros is professor of electrical engineering at the California Institute of Technology [Caltech]. *Technology Review* named her among the top 100 young innovators in 2002. Koetter is professor and head of the Institute for Communications Engineering at the Technical University of Munich. Médard is associate professor of electrical engineering and computer science at the Massachusetts Institute of Technology and associate director of the Laboratory for Information and Decision Systems there.

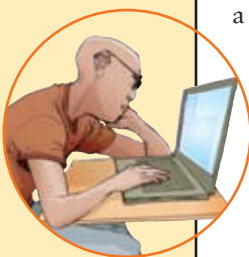
A Brief History of Network Coding

The entries below include some highlights; see “More to Explore” for related references. A fuller bibliography is at www.ifp.uiuc.edu/~koetter/NWC/ —M.É., R.K. and M.M.

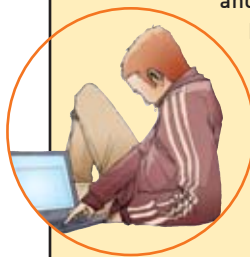
2000: Concept introduced. In a landmark paper, Rudolf Ahlswede, Ning Cai, Shuo-Yen Robert Li and Raymond W. Yeung showed the potential power of network coding in multicast networks, where all receivers get identical information. They proved that good [informative] codes exist, although they did not describe a method for designing them.



2003: Important steps taken toward practical implementation. Li, Yeung and Cai showed network coding for multicast networks can rely on mathematical functions involving only addition and multiplication, which reduces the complexity of designing codes. And two of us [Koetter and Médard] introduced a powerful algebraic framework for analyzing coding approaches and simplifying code design.



2005–2006: Valuable design algorithms published. Sidharth Jaggi, then at Caltech, with Peter Sanders of the University of Karlsruhe in Germany, one of us [Effros] and collaborators and, separately, Tracey Ho of Caltech, with the three of us and others, published low-complexity algorithms for designing the functions used by each node in a multicast network. The first paper gave a systematic approach for designing functions; the second showed that choosing functions randomly and independently for each node should work just as well. (Early versions of both contributions were presented at a conference in 2003.)



2006: Applications for wireless networks explored. At a conference in 2006, Christina Fragouli of the École Polytechnique Fédérale de Lausanne in Switzerland and collaborators demonstrated the potential benefits of network coding for wireless applications and characterized scenarios in which the approach would be particularly helpful.



in turn and directs it on its way. The goal in system design is for each router to direct traffic in a way that not only speeds each subsequent car to its intended destination but also allows the nation’s transportation system as a whole to satisfy as many drivers as possible.

Even a central designer with a complete map of all the nation’s roads in hand would be hard put to determine the best possible strategy for every router to follow. The difficulty increases as the network changes over time: rush hours, road repairs, accidents and sporting events mean the roadways and the demands placed on them change constantly.

Intuition might suggest that designing a system reliant on network coding should be even harder, because there are more options to consider. A node could forward data unchanged, thereby mimicking a router. But it might also mix two or more incoming data streams before sending them on, and how it mixes them might also be open to consideration; further, different nodes might use different algorithms.

Luckily, this logic is flawed. Sometimes adding more options actually simplifies things. Without coding, architects of a multicast system would need to enumerate as many paths as possible from the transmitter to each receiver and then determine how many of those paths the network could support simultaneously. Even for simple networks, finding and testing all combinations of paths would be a dizzying task.

In contrast, a multicast system using network coding would be rather easy to design. The startling truth is that addition and multiplication are the only mathematical functions that coded networks need apply. Also, even if the function, or rule, programmed into each coder in a network is chosen independently of the message and the other coding functions and without any knowledge of the network layout, the system as a whole will, with extremely high probability, operate at peak performance. Even if the system changes over time, as can happen in mobile or reconfigurable networks, the network will continue to perform optimally without requiring redesign. To learn why, see the box on page 82.

Tomorrow’s Networks

THE OPERATION OF NETWORKS, then, will be very different if coders replace routers. The way our messages traverse networks will change: they will not only share “the road” with other transmissions but may become intimately entangled with traffic from a variety of other sources. Some might fear that such entanglement would compromise the security of the messages. More likely, though, traffic traversing networks would become a locally undecipherable algebraic stream. Users on the network would unwittingly collaborate to one another’s mutual advantage, allowing not just higher rates or faster downloads of data but also, in the case of wireless networks, an improvement in energy efficiency. (Because each wireless transmission consumes energy, a node can reduce consumption by mixing together the information intended for several neighbors and sending only a single transmission.)

By changing how networks function, network coding may influence society in ways we cannot yet imagine.

Moreover, delays in downloading videos and lost cell phone calls will be far less common. On the Internet, routers fail or are taken down for maintenance and data packets are dropped all the time. That is why people must sometimes re-request Web pages and why a site sometimes comes up slowly. Reliability will increase with network coding, because it does not require every single piece of evidence to get through.

And network managers will provide such benefits without having to add new communications channels, because better use will be made of existing channels. Network coding will thereby complement other communications technologies, allowing users to get as much as possible out of them.

Sometimes users will know that network coding is operating, because it may modify how some common applications, such as peer-to-peer downloads, function. Today someone seeking to download a file searches for a collaborating user on whose machine the file resides. In a system using network coding, the file would no longer be stored as a whole or in recognizable pieces.

But users would not personally have to figure out how to find the evidence needed to obtain the desired files. A request sent into a network from a user's computer or phone would cause either that individual's computer or a local server to scavenge through the network for pieces of evidence related to a file of interest. The gathered evidence, consisting of algebraically mixed pieces of information relating to the desired file, would help recover that file. Instead of putting together a puzzle whose pieces are recognizable fragments of a whole, the server or an individual's computer would solve a collection of algebraic equations. And, all the while, most people would remain blissfully unaware of these operations—just as most of us are ignorant of the complicated error-correction operations in our cell phones.

The military has recognized the robustness of network coding and is now funding research into its use in mobile ad hoc networks, which can form on the fly. Such networks are valuable in highly changeable environments, such as on the battlefield, where reliable communications are essential and establishing and maintaining an infrastructure of fiber-optic cables or cell towers is difficult. In an ad hoc network, every soldier's radio becomes a node in a communications system, and each node seeks out and establishes connections to neighboring nodes; together these connections establish a network's links. Every node can both send and receive messages and serve as an intermediary to pass along messages intended for other receivers. This technique extends communications capabilities far beyond the transmission range of a single

node. It also allows enormous flexibility, because the network travels with the users, constantly reconfiguring and reestablishing connections as needed.

By changing how networks function, network coding may influence society in ways we cannot yet imagine. In the meantime, though, those of us who are studying it are considering the obstacles to implementation. Transitioning from our router-based system to a network-coded one will actually be one of the more minor hurdles. That conversion can be handled by a gradual change rather than a sudden overhaul; some routers could just be reprogrammed, and others not built to perform coding operations would be replaced little by little.

A bigger challenge will be coping with issues beyond replacing routers with coders. For instance, mixing information is a good strategy when the receiving node will gather enough evidence to recover what it desires from the mixture. This condition is always met in multicast networks but may not be the case in general. Moreover, in some circumstances, such as when multiple multicasts are transmitted, mixing information can make it difficult or impossible for users to extract the proper output. How, then, can nodes decide which information can and cannot be mixed when multiple connections share the same network? In what ways must network coding in wireless networks differ from its use in wired ones? What are the security advantages and implications of network coding? How will people be charged for communications services when one person's data are necessarily mixed with those of other users? In collaborations that span the globe, we and others are pondering how to unravel such knots even as we strive to enhance the capabilities of the communications networks that have become such an integral part of so many lives. SA

SA

For more on network coding, including a fuller explanation of how it can vastly improve transmission rates, visit www.sciam.com/ontheweb

MORE TO EXPLORE

- A Mathematical Theory of Communication.** C. E. Shannon in *Bell System Technical Journal*, Vol. 27, pages 379–423 and 623–656; July and October 1948. Available at <http://cm.bell-labs.com/cm/ms/what/shannonday/paper.html>
- Network Information Flow.** R. Ahlswede, N. Cai, S.-Y. R. Li and R. W. Yeung in *IEEE Transactions on Information Theory*, Vol. 46, No. 4, pages 1204–1216; July 2000.
- Linear Network Coding.** S.-Y. R. Li, R. W. Yeung and N. Cai in *IEEE Transactions on Information Theory*, Vol. 49, No. 2, pages 371–381; February 2003.
- An Algebraic Approach to Network Coding.** R. Koetter and M. Médard in *IEEE/ACM Transactions on Networking*, Vol. 11, No. 5, pages 782–795; October 2003.
- Polynomial Time Algorithms for Multicast Network Code Construction.** S. Jaggi, P. Sanders, P. A. Chou, M. Effros, S. Egner, K. Jain and L.M.G.M. Tolhuizen in *IEEE Transactions on Information Theory*, Vol. 51, No. 6, pages 1973–1982; June 2005.
- A Random Linear Network Coding Approach to Multicast.** T. Ho, M. Médard, R. Koetter, D. R. Karger, M. Effros, J. Shi and B. Leong in *IEEE Transactions on Information Theory*, Vol. 52, No. 10, pages 4413–4430; October 2006.

Seeing Triple

Anticipated for decades, machines are finally displaying real objects in three true dimensions

By
Stuart F. Brown

Inventors have struggled for years to create displays that can conjure vivid three-dimensional images that users can manipulate and interact with. Chemists could exploit such marvels to design new drug molecules. Oil and gas explorers could see exactly where to aim their drills. Surgeons could pass probes or radiation beams through

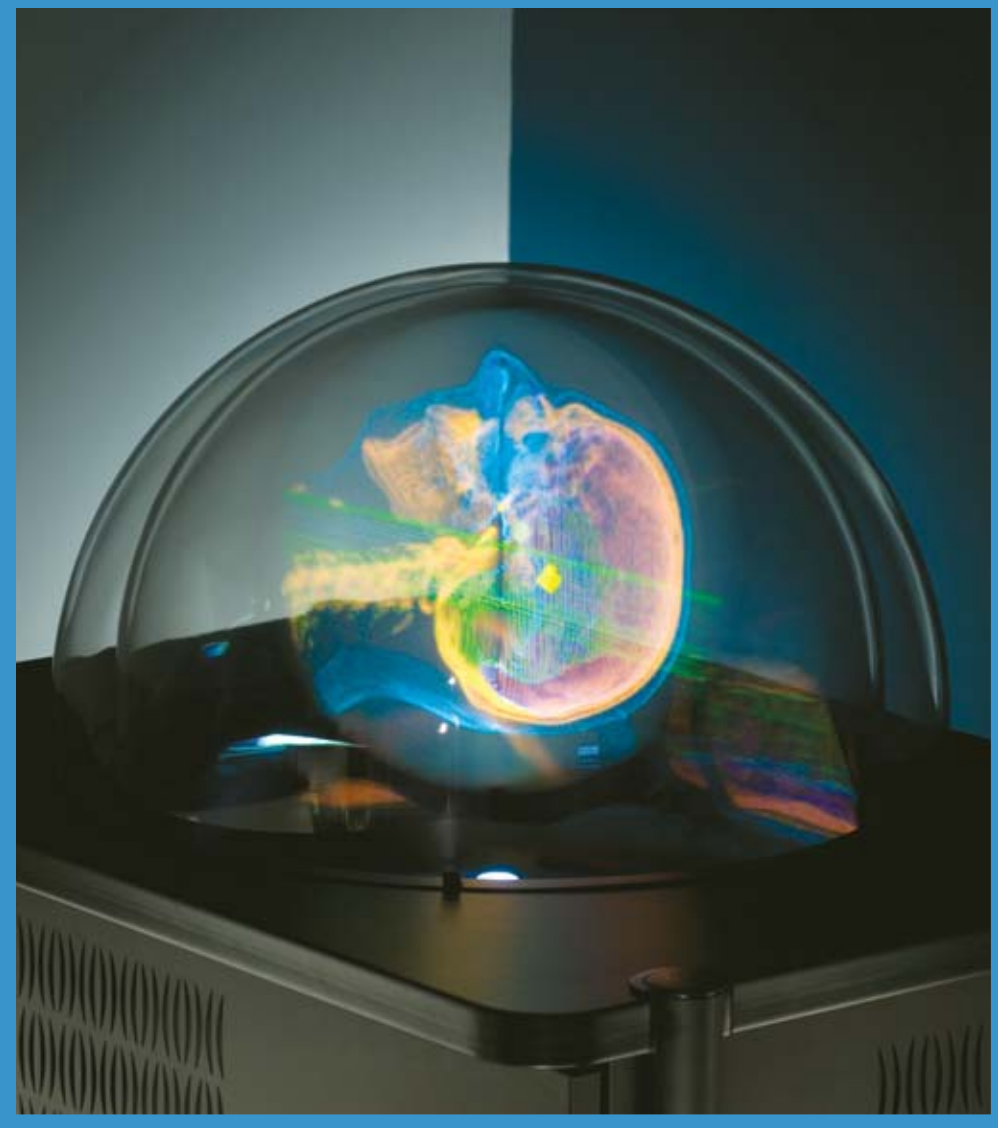
the collated slices of diagnostic data produced by magnetic resonance imaging (MRI) and computed tomography (CT) machines to test procedures before performing an operation. But shortcomings such as a flickering image, a narrow angle of view or the need to wear special glasses have bedeviled the devices.

Two companies have recently mixed their own technology with off-the-shelf components, including the Digital Light Processor (DLP) chip from Texas Instruments, to create interactive systems called 3-D volumetric displays that overcome these limitations. The two firms'

products are just now transitioning from the laboratory to commercial models.

Spinning Algorithms

WAIT A MINUTE. Aren't holograms three-dimensional and viewable without funny glasses? Yes, but they are recorded once as a final image and thus do not allow interactivity. Engineers have also knit together cubes and spinning arrays of light-emitting diodes to give a full-bodied view, but the resolution is coarse, restricted by the connections among the diodes. Other contenders appear to be 3-D but really are not; the Heliodisplay



ANGLE OF ATTACK: Data from a CT scan is projected in 3-D in PerspectaRad, revealing a brain tumor core (yellow diamond, center) and possible paths for radiation treatment (green rays).

data,” says Gregg E. Favalora, Actuality’s chief technology officer. “For instance, we have a patent on how you draw a straight line on a rotating screen—because it is not obvious what dot to pick on that grid as it spins around.”

Perspecta creates a glowing, semi-transparent image. Every volumetric pixel (voxel) that seems to be in a certain point in space really is there, but it becomes visible only when the screen sweeps through that point—hence the fast rotation rate. The screen is made from a plastic that looks like taut tissue paper and is 50 percent reflective and 50 percent transmissive, allowing the imagery to be seen from any angle by observers gathered around it. The unit provides both vertical and horizontal parallax; if a viewer moves his or her head up and down or left and right, background objects that were previously obscured by foreground objects will come into view, as they are perceived in the real world.

A user can also wield a penlike mouse to zoom into or out of the image, rotate and flip it, or change colors. This feature has only recently been made possible by rapid advances in computer graphics. “In our first demo in 2002,” Favalora recalls, “it took 45 minutes of processing time just to perform a click-and-drag. Now we get a video card off the shelf for a few hundred dollars, and it computes the problem just like that.”

It was also a while before Favalora and his colleagues realized that perfecting the core technologies would not be enough to build a business; they had to identify an initial market and develop a ready-to-use system tailored to it. That niche turned out to be radiation therapy

from IO2 Technology in San Francisco projects floating images onto a vertical plane of fine mist suspended above the instrument that seem to possess depth, but the illusion comes from an absence of depth cues, not actual imagery in the third dimension. Users who want to load 3-D medical data into a machine one day and 3-D military scenes the next and then turn, prod or alter them while on view can exploit two inventions worthy of the term “volumetric”: Perspecta and DepthCube.

Perspecta, developed by Actuality Systems in Bedford, Mass., might best be described as a crystal ball for looking inside objects. A transparent polycarbonate dome houses a flat, disk-shaped screen 10 inches in diameter that rotates on a vertical spindle at 900 revolutions per minute. The system takes data generated

by a CT, MRI or PET (positron-emission tomography) scanner and mathematically divides the information into 198 radially disposed segments, like an apple thinly sliced about its core. Held in a frame-buffer memory, the data slices are fed to three DLP chips. DLPs are arrays of hundreds of thousands of tiny mirrors, each of which can be individually tilted by onboard circuitry. They form the heart of projection televisions and new slide projectors, as well as digital movie projectors that may replace film reels in theaters. In Perspecta, each DLP is assigned a color and projects its light through a prism onto the rapidly spinning screen, creating a 3-D apparition.

A lot of mathematical heavy lifting was needed to make Perspecta work. “It took us three or four years to invent the algorithms that let us slice the image

THE AUTHOR

STUART F. BROWN was formerly a staff writer at *Popular Science* and *Fortune*. He calls his now freelance beat “the man-made world,” which includes aerospace, transportation and biotechnology.

for cancer tumors. Doctors need to carefully plan the paths along which they aim radiation beams, trying to maximize the killing effect where the rays converge on a tumor while minimizing damage to nearby healthy tissue. Because oncologists must work with 2-D slices of scanner data, planning the beam paths for a treatment can take several hours. Actuality developed its PerspectaRad system as an add-on to existing radiation therapy equipment manufactured by Philips Medical Systems.

PerspectaRad sports the 3-D display plus software that connects the device to the Philips systems. When a doctor pushes a button, an image of the CT data for, say, a brain tumor appears in 3-D. Another button adds the radiation pathways chosen by a dosimetrist, who plans the treatment. The physician can

see exactly where the beams will strike the tumor, which healthy brain tissues they will pass through, and the dose cloud—the volume of tissue that will be affected by the radiation. This imagery helps doctors adjust the beams to improve treatment or reduce damage. The first PerspectaRad systems cost about \$90,000. Greater production could lower the price to \$65,000, according to Favolora, but the displays are unlikely to reach consumer markets.

Nevertheless, treatment stands to gain. James Chu, head of the medical physics department at Rush University Medical Center in Chicago, recently studied 12 patients with brain tumors for whom treatment plans were developed using both PerspectaRad and conventional methods. The plans were reviewed by doctors who were unaware of

which method had been used. The protocols developed with PerspectaRad turned out better in six cases, equivalent in four cases, and worse in two cases. In one patient, PerspectaRad made it clear how to reduce incidental damage to the optic nerve. Calling the results “interesting,” Chu is planning a larger study that will include patients with tumors in other parts of the body. “When working just with CT data,” he says, “you have to look at individual slices and somehow integrate them all in your head to get a 3-D picture. With Perspecta, you see the 3-D picture directly.”

Chu is also excited by Perspecta’s ability to display moving images of internal body parts. Because internal organs and tissue move as the heart beats and the lungs fill and empty, it is very useful to be able to discern the axis of motion of a tumor. With this information, a doctor could direct a lower-energy radiation beam along the motion axis, rather than a more intense beam across it, lessening collateral damage. Chu says Perspecta could also allow more precise implanting of radioactive “seeds” in the prostate gland to treat cancer there, by allowing a physician to better compensate for motion of tissue that occurs when the needle that delivers the seeds is inserted.

PERSPECTA: CRYSTAL BALL

A transparent dome, translucent screen and optics all rotate at 900 rpm to create a 3-D image. A computer sends graphics data to electronics below the pedestal, which instruct three light-processor chips to focus light from an arc lamp through a projector lens. It reflects a beam up through the spinning shaft and across relay mirrors onto the screen. A second, larger dome (not shown) encases the spinning parts for safety.



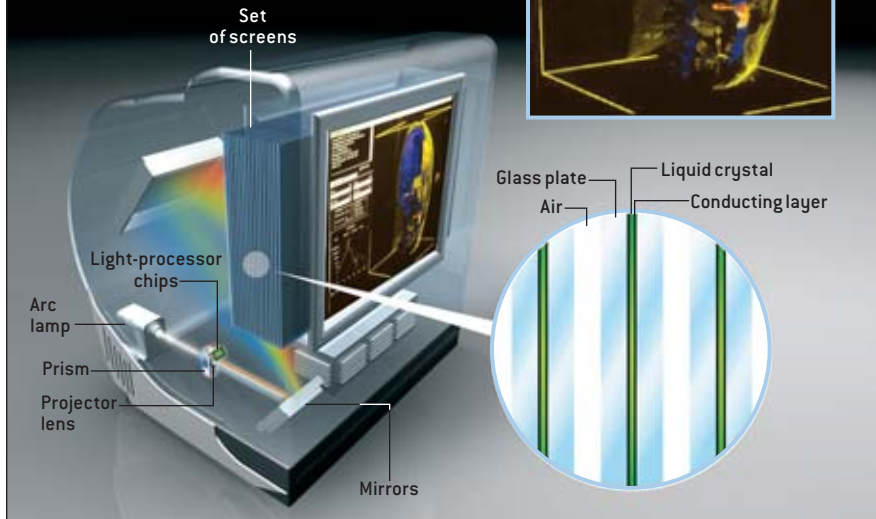
Voxels in Glass Plates

DEPTH CUBE, the other interactive volumetric display, was developed by LightSpace Technologies in Norwalk, Conn. It is a rear-projection monitor shaped somewhat like a computer terminal, with a face that measures 16 by 12 inches. The “screen” is four inches thick and made of 20 vertical, transparent glass plates sandwiched together, yet the imagery looks about 12 inches deep. Collaborators standing randomly in front of the screen will each see objects with the corresponding perspective. Internal structures appear and disappear as a viewer’s angle changes. The system would be handy for, say, a team of product engineers studying how parts drafted on a computer-aided design system would—or would not—fit together.

When LightSpace president Alan Sullivan built his first prototype eight

DEPTHCUBE: THICK-SCREEN IMAGING

An arc lamp shines light through optics and a prism, breaking it into blue, red and green beams that are reflected by light-processor chips through a projector toward relay mirrors. The mirrors direct the beams onto one of 20 screens. Each screen consists of two glass plates that contain conducting layers and a liquid-crystal mixture that scatters light. By illuminating the four-inch-deep set of screens in succession, the system creates a 3-D image that appears 12 inches deep. The sample image of half a human head allows particular structures to be seen [sinuses are bright yellow; cartilage is orange], whereas others can be hidden.



years ago, he managed to coax a trio of Texas Instruments DLP chips into projecting depth-related imagery onto the 20 plates, which are separated by thin air gaps. In DepthCube, each DLP contains 786,432 mirrors that tile an area similar to that of a fingernail.

Sullivan still needed a convenient way to generate depth information, and he was elated when he realized that an affordable, commercial 3-D graphics card could suffice. Graphics cards use a color buffer—a morsel of memory—to assign the appropriate color to every pixel on a two-dimensional screen. But the cards also have a hidden component, called the depth buffer, which describes the depth of every pixel. In a normal application, the depth buffer is largely untapped, because only the frontmost layer of a pixel must be defined to create a 2-D picture. So a place for Sullivan's depth information, he muses, "was in there for

free." The information drives the 20 plates, known as liquid-crystal scattering shutters, which can rapidly change from a transparent state to a scattering state. That trait allows a plate to let pixels pass through to other plates as needed yet also enables it to display a pixel. At any moment, all the plates are blank except for one, but the processors project coordinated image slices 50 times a second onto every plate, creating full depth, height and width.

The prototype DepthCube conveyed three-dimensionality to a viewer nicely but only within the four-inch depth of the screen; items in an image appeared almost like flat scenery elements on a

theater stage, standing in front of and behind one another. That was when Sullivan, who formerly studied ultrahigh-energy lasers at Lawrence Livermore National Laboratory, had a smart attack that won him a patent. It occurred to him that the so-called antialiasing algorithms used to smooth jagged edges in 2-D images could also be applied to smoothing the transitions between the DepthCube's 20 planes. This innovation makes the display's 15.3 million physical voxels look like a whopping 465 million virtual voxels. "We produce 31 subplanes between the physical planes, so the perceived resolution is much higher," Sullivan explains. As a result, to the human brain, the images can appear to be as much as 12 inches deep.

The image data that are fed to the chips can come from nearly any 3-D software that runs the OpenGL application programming interface, a common protocol used by computer-aided design and engineering programs such as Catia or ProEngineer. LightSpace has sold a handful of DepthCubes to research institutions, including the U.S. Air Force Research Laboratory and Hokkaido University in Japan, for about \$50,000 apiece. Sullivan acknowledges that the market is limited at this price but says he can see the path to a product that would cost about \$5,000. "There's nothing in our architecture that's different from what's in a rear-projection TV, except for the liquid-crystal shutters," he says, "and those could be produced quite cheaply in volume."

The products developed by these two young companies are garnering respect from boffins in the 3-D world, and more applications will follow. Optical scientist Steve Hines, owner of HinesLab in Glendale, Calif., states that "both these groups are doing extremely hard things and have pulled them off." The natural places to sell the technology, he adds, will be where the money is: "medicine, the military and the movies." SA

MORE TO EXPLORE

Volumetric 3D Displays and Application Infrastructure. Gregg E. Favalora in *Computer*, Vol. 38, No. 8, pages 37–44; August 2005.

A Method for the Real-Time Construction of a Full Parallax Light Field. K. Tanaka and S. Aoki in *Stereoscopic Displays and Virtual Reality Systems XIII*. Edited by A. J. Woods et al. *Proceedings of the SPIE*, Vol. 6055, Article 605516; January 30, 2006.

By Kaushik Basu

THE TRAVELER'S



Lucy and Pete, returning from a remote Pacific island,

find that the airline has damaged the identical antiques that each had purchased. An airline manager says that he is happy to compensate them but is handicapped by being clueless about the value of these strange objects. Simply asking the travelers for the price is hopeless, he figures, for they will inflate it.

Instead he devises a more complicated scheme. He asks each of them to write down the price of the antique as any dollar integer between 2 and 100 without conferring together. If both write the same number, he will take that to be the true price, and he will pay each of them that amount. But if they write different numbers, he will assume that the lower one is the actual price and that the person writing the higher number is cheating. In that case, he will pay both of them the lower number along with a bonus and a penalty—the person who wrote the lower number will get \$2 more as a reward for honesty and the one who wrote the higher number will get \$2 less as a punishment. For instance, if Lucy writes 46 and Pete writes 100, Lucy will get \$48 and Pete will get \$44.

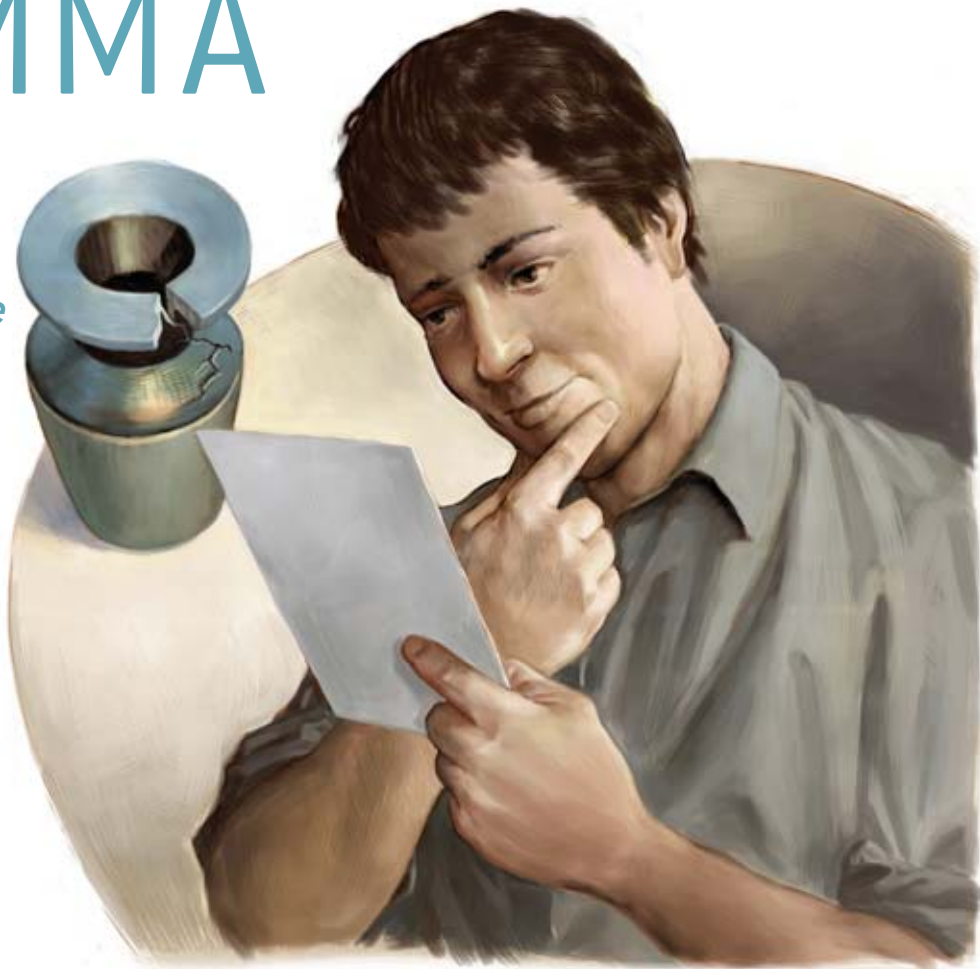
What numbers will Lucy and Pete write? What number would you write?

Scenarios of this kind, in which one or more individuals have choices to make and will be rewarded according to those choices, are known as games by the people who study them (game theorists). I crafted this game, “Traveler’s Dilemma,” in 1994 with several objectives in mind: to contest the narrow view of rational behavior and cognitive processes taken by economists and many political scientists, to challenge the libertarian presumptions of traditional economics and to highlight a logical paradox of rationality.

Traveler’s Dilemma (TD) achieves those goals because the game’s logic dictates that 2 is the best option, yet most people pick 100 or a number close to 100—both those who have not thought through the logic and those who fully understand that they are deviating markedly from the “rational” choice. Furthermore, players reap a greater reward by not adhering to reason in this way. Thus, there is something rational about

DILEMMA

When playing this simple game, people consistently reject the rational choice. In fact, by acting illogically, they end up reaping a larger reward—an outcome that demands a new kind of formal reasoning



choosing not to be rational when playing Traveler's Dilemma.

In the years since I devised the game, TD has taken on a life of its own, with researchers extending it and reporting findings from laboratory experiments. These studies have produced insights into human decision making. Nevertheless, open questions remain about how logic and reasoning can be applied to TD.

Common Sense and Nash

TO SEE WHY 2 is the logical choice, consider a plausible line of thought that Lucy might pursue: her first idea is that she should write the largest possible number, 100, which will earn her \$100 if Pete is similarly greedy. (If the antique actually cost her much less than \$100, she would now be happily thinking about the foolishness of the airline manager's scheme.)

Soon, however, it strikes her that if she wrote 99 instead, she would make a little more money, because in that case she would get \$101. But surely this insight will also occur to Pete, and if both wrote 99, Lucy would get \$99. If Pete wrote 99, then she could do better by writing 98, in which case she would get \$100. Yet the same logic would lead Pete to choose 98 as well. In that case, she could deviate to 97 and earn \$99. And so

on. Continuing with this line of reasoning would take the travelers spiraling down to the smallest permissible number, namely, 2. It may seem highly implausible that Lucy would really go all the way down to 2 in this fashion. That does not matter (and is, in fact, the whole point)—this is where the logic leads us.

Game theorists commonly use this style of analysis, called backward induction. Backward induction predicts that each player will write 2 and that they will end up getting \$2 each (a result that might explain why the airline manager has done so well in his corporate career). Virtually all models used by game theorists predict this outcome for TD—the two players earn \$98 less than they would if they each naively chose 100 without thinking through the advantages of picking a smaller number.

Traveler's Dilemma is related to the more popular Prisoner's Dilemma, in which two suspects who have been arrested for a serious crime are interrogated separately and each has the choice of incriminating the other (in return for leniency by the authorities) or maintaining silence (which will leave the police with inadequate evidence for a case, *if* the other prisoner also stays silent). The story sounds very different from our tale of two travelers with damaged souvenirs, but the mathematics of the rewards for each option in Prisoner's

Dilemma is identical to that of a variant of TD in which each player has the choice of only 2 or 3 instead of every integer from 2 to 100.

Game theorists analyze games without all the trappings of the colorful narratives by studying each one's so-called payoff matrix—a square grid containing all the relevant information about the potential choices and payoffs for each player [see box on opposite page]. Lucy's choice corresponds to a row of the grid and Pete's choice to a column; the two numbers in the selected square specify their rewards.

Despite their names, Prisoner's Dilemma and the two-choice version of Traveler's Dilemma present players with no real dilemma. Each participant sees an unequivocal correct choice, to wit, 2 (or, in the terms of the prisoner story line, incriminate the other person). That choice is called the dominant choice because it is the best thing to do no matter what the other player does. By choosing 2 instead of 3, Lucy will receive \$4 instead of \$3 if Pete chooses 3, and she will receive \$2 instead of nothing if Pete chooses 2.

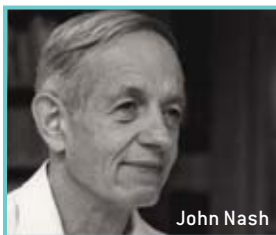
\$101. Because Lucy is better off by this change, the outcome (100, 100) is not a Nash equilibrium.

TD has only one Nash equilibrium—the outcome (2, 2), whereby Lucy and Pete both choose 2. The pervasive use of the Nash equilibrium is the main reason why so many formal analyses predict this outcome for TD.

Game theorists do have other equilibrium concepts—strict equilibrium, the rationalizable solution, perfect equilibrium, the strong equilibrium and more. Each of these concepts leads to the prediction (2, 2) for TD. And therein lies the trouble. Most of us, on introspection, feel that we would play a much larger number and would, on average, make much more than \$2. Our intuition seems to contradict all of game theory.

Implications for Economics

THE GAME AND OUR INTUITIVE prediction of its outcome also contradict economists' ideas. Early economics was firmly tethered to the libertarian presumption that individuals should be left to their own devices because their selfish choices will re-



Game theory predicts that the Nash equilibrium will occur when Traveler's Dilemma is played rationally.

In contrast, the full version of TD has no dominant choice. If Pete chooses 2 or 3, Lucy does best by choosing 2. But if Pete chooses any number from 4 to 100, Lucy would be better off choosing a number larger than 2.

When studying a payoff matrix, game theorists rely most often on the Nash equilibrium, named after John F. Nash, Jr., of Princeton University. (Russell Crowe portrayed Nash in the movie *A Beautiful Mind*.) A Nash equilibrium is an outcome from which no player can do better by deviating unilaterally. Consider the outcome (100, 100) in TD (the first number is Lucy's choice, and the second is Pete's). If Lucy alters her selection to 99, the outcome will be (99, 100), and she will earn

sult in the economy running efficiently. The rise of game-theoretic methods has already done much to cut economics free from this assumption. Yet those methods have long been based on the axiom that people will make selfish rational choices that game theory can predict. TD undermines both the libertarian idea that unrestrained selfishness is good for the economy and the game-theoretic tenet that people will be selfish and rational.

In TD, the "efficient" outcome is for both travelers to choose 100 because that results in the maximum total earnings by the two players. Libertarian selfishness would cause people to move away from 100 to lower numbers with less efficiency in the hope of gaining more individually.

And if people do not play the Nash equilibrium strategy (2), economists' assumptions about rational behavior should be revised. Of course, TD is not the only game to challenge the belief that people always make selfish rational choices [see "The Economics of Fair Play," by Karl Sigmund, Ernst Fehr and Martin A. Nowak; *SCIENTIFIC AMERICAN*, January 2002]. But it makes the more puzzling point that even if players have no concern other than their own profit, it is not rational for them to play the way formal analysis predicts.

TD has other implications for our understanding of real-world situations. The game sheds light on how the arms race acts as a gradual process, taking us in small steps to ever worsening outcomes. Theorists have also tried to extend TD to understand how two competing firms may undercut each other's price to their own detriment (though in this case to

Overview/Sensible Irrationality

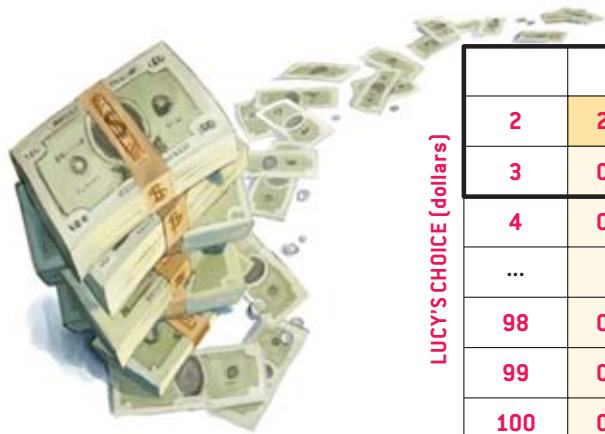
- In the game Traveler's Dilemma, two people separately choose a whole number from 2 to 100, and the player whose number is smaller is rewarded with a larger sum. Game theory insists that rationality should lead players to select 2, but most people pick an integer closer to 100.
- A new kind of reasoning is needed to gain a rigorous understanding of this rational choice not to be rational.
- The results of Traveler's Dilemma contradict economists' assumption that standard game theory can predict how supposedly selfish rational people will behave. They also show how selfishness is not always good economics.

PAYOFF MATRIX OF THE TRAVELER'S DILEMMA

This payoff matrix summarizes everything game theorists need to know about Traveler's Dilemma. Lucy's possible choices are shown in the leftmost column; Pete's run across the top row. The first number in the square at the intersection of the chosen row and column is Lucy's payoff, and the second number is Pete's payoff. For example, if Lucy chooses 98 and Pete 99, then Lucy receives \$100 and Pete receives \$96.

The outcome in which both players choose 2 and each earns \$2 (*gold*) is called the Nash equilibrium: Lucy does worse (earning \$0) if she chooses any other number and Pete still chooses 2. Similarly, Pete does worse if he alone chooses something other than 2.

When restricted to choices of only 2 and 3 (*black outline*), the game becomes equivalent to Prisoner's Dilemma.



		PETE'S CHOICE (dollars)						
		2	3	4	...	98	99	100
LUCY'S CHOICE (dollars)	2	2 2	4 0	4 0	...	4 0	4 0	4 0
	3	0 4	3 3	5 1	...	5 1	5 1	5 1
	4	0 4	1 5	4 4	...	6 2	6 2	6 2

	98	0 4	1 5	2 6	...	98 98	100 96	100 96
	99	0 4	1 5	2 6	...	96 100	99 99	101 97
	100	0 4	1 5	2 6	...	96 100	97 101	100 100

the advantage of the consumers who buy goods from them).

All these considerations lead to two questions: How do people actually play this game? And if most people choose a number much larger than 2, can we explain why game theory fails to predict that? On the former question, we now know a lot; on the latter, little.

How People Actually Behave

OVER THE PAST DECADE researchers have conducted many experiments with TD, yielding several insights. A celebrated lab experiment using real money with economics students as the players was carried out at the University of Virginia by C. Monica Capra, Jacob K. Goeree, Rosario Gomez and Charles A. Holt. The students were paid \$6 for participating and kept whatever additional money they earned in the game. To keep the budget manageable, the choices were valued in cents instead of dollars. The range of choices was made 80 to 200, and the value of the penalty and reward was varied for different runs of the game, going as low as 5 cents and as high as 80 cents. The experimenters wanted to see if varying the magnitude of the penalty and reward would make a difference in how the game was played. Altering the size of the reward and penalty does not change any of the formal analysis: backward induction always leads to the outcome (80, 80), which is the Nash equilibrium in every case.

The experiment confirmed the intuitive expectation that the average player would not play the Nash equilibrium strategy of 80. With a reward of 5 cents, the players' average choice was 180, falling to 120 when the reward rose to 80 cents.

Capra and her colleagues also studied how the players' behavior might alter as a result of playing TD repeatedly.

Would they learn to play the Nash equilibrium, even if that was not their first instinct? Sure enough, when the reward was large the play converged, over time, down toward the Nash outcome of 80. Intriguingly, however, for small rewards the play increased toward the opposite extreme, 200.

The fact that people mostly do not play the Nash equilibrium received further confirmation from a Web-based experiment with no actual payments that was carried out by Ariel Rubinstein of Tel Aviv University and New York University from 2002 to 2004. The game asked players, who were going to attend one of Rubinstein's lectures on game theory and Nash, to choose an integer between 180 and 300, which they were to think of as dollar amounts. The reward/penalty was set at \$5.

Around 2,500 people from seven countries responded, giving a cross-sectional view and sample size infeasible in a laboratory. Fewer than one in seven players chose the scenario's Nash equilibrium, 180. Most (55 percent) chose the maximum number, 300 [see box on next page]. Surprisingly, the data were very similar for different subgroups, such as people from different countries.

The thought processes that produce this pattern of choices remain mysterious, however. In particular, the most popular response (300) is the only strategy in the game that is

THE AUTHOR

KAUSHIK BASU is professor of economics, Carl Marks Professor of International Studies and director of the Center for Analytic Economics at Cornell University. He has written extensively in academic journals on development economics, welfare economics, game theory and industrial organization. He also writes for the popular media, including a monthly column in BBC News online. He is a fellow of the Econometric Society.

“dominated”—which means there is another strategy (299) that never does worse and sometimes does better.

Rubinstein divided the possible choices into four sets of numbers and hypothesized that a different cognitive process lies behind each one: 300 is a spontaneous emotional response. Picking a number between 295 and 299 involves strategic reasoning (some amount of backward induction, for instance). Anything from 181 to 294 is pretty much a random choice. And finally, standard game theory accounts for the choice of 180, but players might have worked that out for themselves or may have had prior knowledge about the game.

A test of Rubinstein’s conjecture for the first three groups would be to see how long each player took to make a decision. Indeed, those who chose 295 to 299 took the longest time on average (96 seconds), whereas both 181 to 294 and 300 took about 70 seconds—a pattern that is consistent with his hypothesis that people who chose 295 to 299 thought more than those who made other choices.

Game theorists have made a number of attempts to explain why a lot of players do not choose the Nash equilibrium in TD experiments. Some analysts have argued that many people are unable to do the necessary deductive reasoning and therefore make irrational choices unwittingly. This explanation must be true in some cases, but it does not account for all the results, such as those obtained in 2002 by Tilman Becker, Michael Carter and Jörg Naeve, all then at the University of Hohenheim in Germany. In their experiment, 51 members of the Game Theory Society, virtually all of whom are professional game theorists, played the original 2-to-100 version of TD. They played against each of their 50 opponents by selecting a strategy and sending it to the researchers. The strategy could be a single number to use in every game or a selection of numbers and how often to use each of them. The game had a real-money reward system: the experimenters would select one player at random to win \$20 multiplied by that player’s average payoff in the game. As it turned out, the winner, who had an average payoff of \$85, earned \$1,700.

Of the 51 players, 45 chose a single number to use in every game (the other six specified more than one number). Among those 45, only three chose the Nash equilibrium (2), 10 chose the dominated strategy (100) and 23 chose numbers ranging from 95 to 99. Presumably game theorists know how to reason deductively, but even they by and large did not follow the rational choice dictated by formal theory.

Superficially, their choices might seem simple to explain: most of the participants accurately judged that their peers would choose numbers mainly in the high 90s, and so choosing a similarly high number would earn the maximum average return. But why did everyone expect everyone else to choose a high number?

Perhaps altruism is hardwired into our psyches alongside selfishness, and our behavior results from a tussle between the two. We know that the airline manager will pay out the largest amount of money if we both choose 100. Many of us do not feel like “letting down” our fellow traveler to try to

earn only an additional dollar, and so we choose 100 even though we fully understand that, rationally, 99 is a better choice for us as individuals.

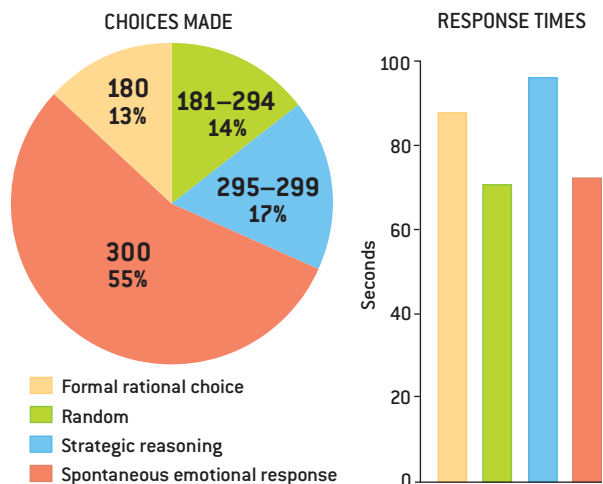
To go further and explain more of the behaviors seen in experiments such as these, some economists have made strong and not too realistic assumptions and then churned out the observed behavior from complicated models. I do not believe that we learn much from this approach. As these models and assumptions become more convoluted to fit the data, they provide less and less insight.

Unsolved Problem

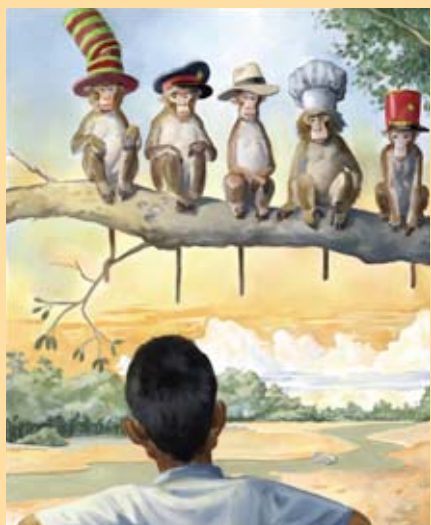
THE CHALLENGE THAT REMAINS, however, is not explaining the real behavior of typical people presented with TD. Thanks in part to the experiments, it seems likely that altruism, socialization and faulty reasoning guide most individuals’ choices. Yet I do not expect that many would select 2 if those three factors were all eliminated from the picture. How can we explain it if indeed most people continue to choose large numbers, perhaps in the 90s, even when they have no dearth of deductive ability, and they suppress their normal altruism and social behavior to play ruthlessly to try to make as much money as possible? Unlike the bulk of modern game theory, which may involve a lot of mathematics but is straightforward once one knows the techniques, this question is a hard one that requires creative thinking.

WHAT WERE THEY THINKING?

Researchers hypothesized that different thought processes lay behind different types of choices that people made playing a version of Traveler’s Dilemma with the options ranging from 180 to 300 (pie chart): a spontaneous emotional response (choosing 300), a strategically reasoned choice (295–299) or a random one (181–294). Players making the formal rational choice (180) might have deduced it or known about it in advance. As expected, people making “spontaneous” or “random” selections took the least time to choose (bar graph).



Game Theory vs. Ordinary Decision Theory: I Know that You Know that I Know . . .



I heard this tale in India. A hat seller, on waking from a nap under a tree, found that a group of monkeys had taken all his hats to the top of the tree. In exasperation he took off his own hat and flung it to the ground. The monkeys, known for their imitative urge, hurled down the hats, which the hat seller promptly collected.

Half a century later his grandson, also a hat seller, set down his wares under the same tree for a nap. On waking, he was dismayed to discover that monkeys had taken all his hats to the treetop. Then he remembered his grandfather's story, so he threw his own hat to the ground. But,

mysteriously, none of the monkeys threw any hats, and only one monkey came down. It took the hat on the ground firmly in hand, walked up to the hat seller, gave him a slap and said, "You think only you have a grandfather?"

This story illustrates an important distinction between ordinary decision theory and game theory. In the latter, what is rational for one player may depend on what is rational for the other player. For Lucy to get her decision right, she must put herself in Pete's shoes and think about what he must be thinking. But he will be thinking about what she is thinking, leading to an infinite regression. Game theorists describe this situation by saying that "rationality is common knowledge among the players." In other words, Lucy and Pete are rational, they each know that the other is rational, they each know that the other knows, and so on.

The assumption that rationality is common knowledge is so pervasive in game theory that it is rarely stated explicitly. Yet it can run us into problems. In some games that are played over time, such as repeated rounds of Prisoner's Dilemma, players can make moves that are incompatible with this assumption.

I believe that the assumption that rationality is common knowledge is the source of the conflict between logic and

intuition and that, in the case of Traveler's Dilemma, the intuition is right and awaiting validation by a better logic. The problem is akin to what happened in early set theory. At that time, mathematicians took for granted the existence of a universal set—a set that contained everything. The universal set seemed extremely natural and obvious, yet ultimately several paradoxes of set theory were traced to the assumption that it existed, which mathematicians now know is flawed. In my opinion, the common knowledge of rationality assumed by game theorists faces a similar demise. —K.B.



Suppose you and I are two of these smart, ruthless players. What might go through our minds? I expect you to play a large number—say, one in the range from 90 to 99. Then I should not play 99, because whichever of those numbers you play, my choosing 98 would be as good or better for me. But if you are working from the same knowledge of ruthless human behavior as I am and following the same logic, you will also scratch 99 as a choice—and by the kind of reasoning that would have made Lucy and Pete choose 2, we quickly eliminate every number from 90 to 99. So it is not possible to make the set of "large numbers that ruthless people might logically choose" a well-defined one, and we have entered the philosophically hard terrain of trying to apply reason to inherently ill-defined premises.

If I were to play this game, I would say to myself: "Forget game-theoretic logic. I will play a large number (perhaps 95), and I know my opponent will play something similar and both of us will ignore the rational argument that the next smaller number would be better than whatever number we choose." What is interesting is that this rejection of formal rationality and logic has a kind of meta-rationality attached

to it. If both players follow this meta-rational course, both will do well. The idea of behavior generated by rationally rejecting rational behavior is a hard one to formalize. But in it lies the step that will have to be taken in the future to solve the paradoxes of rationality that plague game theory and are codified in Traveler's Dilemma. SA

MORE TO EXPLORE

On the Nonexistence of a Rationality Definition for Extensive Games. Kaushik Basu in *International Journal of Game Theory*, Vol. 19, pages 33–44; 1990.

The Traveler's Dilemma: Paradoxes of Rationality in Game Theory. Kaushik Basu in *American Economic Review*, Vol. 84, No. 2, pages 391–395; May 1994.

Anomalous Behavior in a Traveler's Dilemma? C. Monica Capra et al. in *American Economic Review*, Vol. 89, No. 3, pages 678–690; June 1999.

The Logic of Backwards Inductions. G. Priest in *Economics and Philosophy*, Vol. 16, No. 2, pages 267–285; 2000.

Experts Playing the Traveler's Dilemma. Tilman Becker et al. Working Paper 252, Institute for Economics, Hohenheim University, 2005.

Instinctive and Cognitive Reasoning. Ariel Rubinstein. Available at arielrubinstein.tau.ac.il/papers/Response.pdf

WORKING KNOWLEDGE

CHARACTER RECOGNITION

The Write Type

Electronically scan a book to import its content into a word-processing program. Save a snippet handwritten on a personal digital assistant (PDA) screen into a spreadsheet. Decipher a scrawled form or the zip code on an envelope. In all these cases, software translates typed or handwritten characters into digital text that can be edited, e-mailed, stored or used to tell a high-speed machine which direction to route a letter.

That software was originally known as optical character recognition; today the term refers just to recognizing text from a typeset page. Analyzing printed or cursive handwriting is called intelligent character recognition. Regardless of labels, the programs rely on similar algorithms to assess the features of an inkblot [see illustrations]. The programs then compare the blot's features against mathematical models to determine which letter or number it most closely resembles.

Determining characters handwritten with a pen on a PDA is perhaps the easiest task, because the pen or screen can track the stylus's movement. Analyzing type or handwriting on a printed page is tougher "because you must extract a signal from a static image, clean clutter, then discern letters," says David Doermann, co-director of the University of Maryland's Laboratory for Language and Media Processing. Translating unconstrained cursive remains most elusive, and, Doermann says, "it is not being done commercially."

Years ago character-recognition techniques relied on one or two algorithms that compared simple patterns; current software exploits multiple algorithms and then weighs, or votes on, their results to make a final determination. Furthermore, "the old software took several minutes to convert a page; new software takes only a few seconds," says Allan Stratton, a technical director at Nuance Communications in Burlington, Mass., which makes OmniPage, a document conversion program.

Extensive research is under way to devise systems for Arabic, Japanese and Chinese characters, in which strokes can represent whole words, not just letters. Recognition software is also beginning to be included in digital cameras and cell phones. "Just snap a photo of a document," Stratton imagines, "then stream the text as an e-mail or text message." —Mark Fischetti

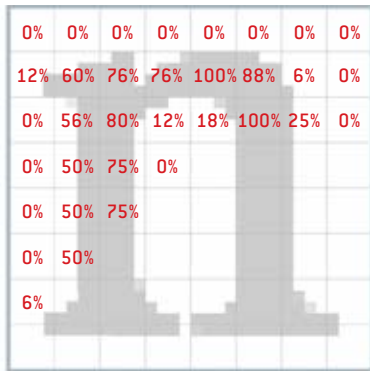
EXISTING TYPE is recognized by software that masks out stray marks, frames a word in a grid and segments it into letters. An algorithm then analyzes the features of each letter (common features are shown at right). In each case, the software compares the feature vector—a set of numbers—against a table of values that indicate a likely character, punctuation mark or clutter (such as a smudge). This "offline" software, as it is called, weighs the results of different feature vectors to determine the final letter. Postprocessing reassembles characters into words in a word-processing program.



DID YOU KNOW...

- LONG RECOGNIZED: European and American patents for optical character recognition have been filed since 1929. The U.S. Armed Forces Security Agency tried the technology in the early 1950s to automate code breaking. In 1965 the U.S. Postal Service began to electronically scan zip codes; today the U.S. Census Bureau processes millions of forms this way. In the early 1990s the Apple Newton PDA and the IBM ThinkPad notebook computer brought handwriting recognition to the masses, although commercial success was marginal.
- NOTATHOMEDEPOT: Customers who sign for a credit-card purchase on a small electronic pad by a cash register may think their signa-

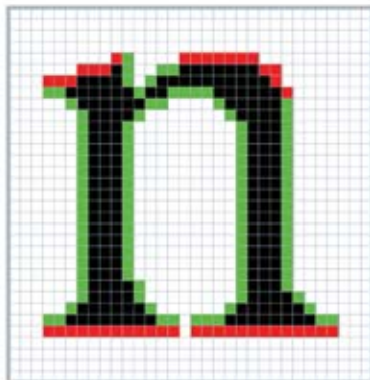
- ture is being analyzed, but it is probably not. Most retailers simply store the whole image "just to save paper," says David Doermann of the University of Maryland. Next time, he proposes, draw something completely different and see if a clerk or machine questions you.
- CLASS: The algorithms that analyze character features feed their findings into classifier programs that compare the results against reference tables. One mainstream classifier is NN (nearest neighbor), which compares the values against every table entry to find the closest match. Another classifier is HMM (hidden Markov model), first applied to speech recognition, which assesses the probability distribution of variables that make up a feature.



FEATURE:
Color

ALGORITHM:
Matrix matching

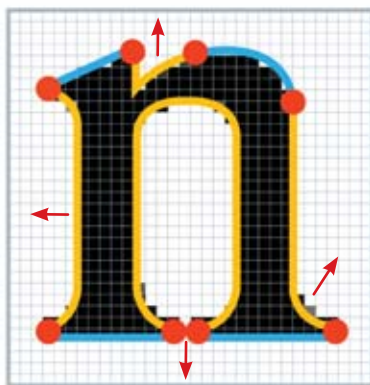
METHOD:
Calculate ratio of black and white pixels within a cell



FEATURE:
Number and length of curves

ALGORITHM:
Contour tracing

METHOD:
Measure curvatures

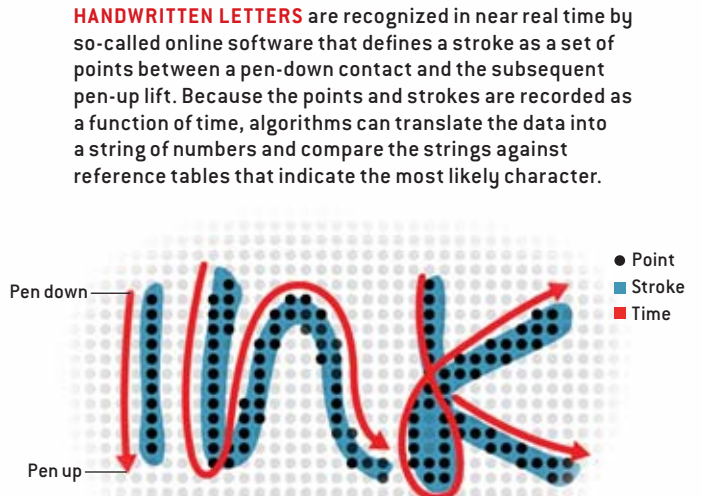


FEATURE:
Direction curves are facing

ALGORITHM:
Contour tracing

METHOD:
Plot end points and inflection points

LETTER features are determined.



PEN or tablet records the vertical and horizontal movement of the stylus.

Topic suggested by reader Morton Nadler. Send ideas to workingknowledge@sciam.com

5W INFOGRAPHICS; SOURCES: NUANCE IMAGING (letter n); VISION OBJECTS (handwritten "ink")

Social Baboons and Oceans' Depths



Baboon Metaphysics

The Evolution of a Social Mind

Dorothy L. Cheney and Robert M. Seyfarth

"Like characters in Edith Wharton's novels about rigidly stratified Old New York, baboons need to constantly navigate the treacherous shoals of fluctuating relationships based on kinship, rank, sexual histories, and opportunistic alliances. In this gem of a book, Cheney and Seyfarth combine in-depth ethological research with ingenious experiments to probe the minds of wild baboons living beside Botswana's Okavango swamp."

—Sarah Blaffer Hrdy, author of *The Woman that Never Evolved*

Cloth \$27.50



The Silent Deep

The Discovery, Ecology, and Conservation of the Deep Sea

Tony Koslow

"I know of no other work that provides such a comprehensive review of the history of deep-sea exploration. *The Silent Deep* will be readily understood and appreciated by the general reader."—Richard Lutz, Rutgers University

Cloth \$35.00



Available in bookstores
The University of Chicago Press
www.press.uchicago.edu

REVIEWS

The Editors Recommend

UNCERTAINTY: EINSTEIN, HEISENBERG, BOHR, AND THE STRUGGLE FOR THE SOUL OF SCIENCE

by David Lindley
Doubleday, 2007 (\$26)



Lindley, an astrophysicist-turned-writer, charts the course of Werner Heisenberg's uncertainty principle. The culmination of Heisenberg's equally perplexing quantum theory, the uncertainty principle posited that in many physical measurements, one can extract one bit of information only at the price of losing another. Heisenberg's mentor Niels Bohr agreed with the basic premises of his startling insights but saw the need to "make sense of the new quantum physics without throwing overboard the hard-won successes of the previous era." The third voice in this argument was Albert Einstein, to whom Heisenberg's ideas were a "monstrous misrepresentation . . . the very idea of a true fact seemed to crumble into an assortment of irreconcilable points of view." Eventually, and reluctantly, Einstein conceded the technical correctness of the system that Heisenberg and Bohr had laid out, but he never fully accepted it. This story has been told before but seldom with such clarity and elegance.

DIRT: THE EROSION OF CIVILIZATIONS

by David R. Montgomery
University of California Press, 2007
(\$24.95)

Montgomery, a geomorphologist at the University of Washington, argues that

good old dirt has always been necessary to sustain civilizations, from ancient times right on through today's digital society. This natural resource is being exhausted at a faster rate than it is being replenished, however. In the past century the effects of long-term soil erosion were masked by bringing new land under cultivation and developing fertilizers, pesticides and crop varieties. But as the population continues to grow and the arable land base continues to shrink, Montgomery believes that society must rethink its relationship with the land. Governments and farmers must rely not only on technological sophistication to protect the soil but on intergenerational land stewardship and conservation. "Civilization's survival depends on treating soil as an investment," he says, "as a valuable inheritance rather than a commodity—as something other than dirt."



FLOWER CONFIDENTIAL: THE GOOD, THE BAD, AND THE BEAUTIFUL

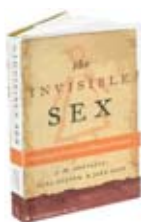
by Amy Stewart
Algonquin Books of Chapel Hill, 2007
(\$23.95)

Rose breeders can spend seven years developing a new variety; one team of scientists has been working for 10 years on a blue rose, using a pigment gene from petunias. Making flowers flawless is a \$40-billion industry. Per capita spending on cut flowers in the U.S. is about \$25 a year; the Swiss have the highest consumption, at \$100 per person annually.

The facts are surprising and intriguing. But it is the way nature writer Stewart packages them that makes *Flower Confidential* the rare nonfiction book that keeps you turning pages. Her own passion for flowers, along with her adventures in the fields, greenhouses, auction houses and laboratories, gives the facts life. An almost perfect book—its only blemish being the off-putting title that seems to promise more muck than this intelligent, evenhanded account rakes.




THE INVISIBLE SEX: UNCOVERING THE TRUE ROLES OF WOMEN IN PREHISTORY
by J. M. Adovasio, Olga Soffer and Jake Page
Smithsonian Books, 2007 [\$26.95]



Funny how no one has to think twice about which sex is invisible. Anthropologists Adovasio and Soffer and science writer Page take the field of archaeology to task for its traditional focus on hard artifacts such as stone tools and weapons, pointing out that archaeologists are not even trained to look for evidence of women's use of more perishable artifacts such as string and netting. They argue for the central importance of the "fiber revolution," which began some 26,000 years ago in Eurasia. In dry caves and other rare places where artifacts do not deteriorate, fiber and wood objects can account for 95 percent of all artifacts recovered. The influence of fiber and its use in nets, baskets and clothing had "profound effects on human destiny—probably more profound than any advance in the technique of making spear points, knives, scrapers and other tools of stone."


If you can't find it here, it doesn't exist.

 **AbeBooks.com**

100 million new, used, rare, and out-of-print books.
Visit our Science Fiction Room at www.abebooks.com/sciencefiction

Bright Horizons™
January 27th – February 3rd, 2008 • W. Caribbean

www.InSightCruises.com/SciAm



Cruise prices vary from \$849 for an Inside Stateroom to \$2,899 for a Full Suite, per person. (Cruise pricing is subject to change. InSightCruises will generally match the cruise pricing offered at the Holland America website at the time of booking.) For those attending our Program, there is a \$1,275 fee. Taxes are \$70 per person. Program subject to change.

Contact Neil Bauman
neil@insightcruises.com
650-787-5665

CST# 2065380-40

CO-PRODUCED BY:

InSight Cruises
EDUCATION THAT TAKES YOU PLACES

SCIENTIFIC AMERICAN



Dumb Cup

RECIPE FOR A STEAMING CUP OF SOMETHING BY STEVE MIRSKY

On a chilly, late March day I was happily sipping a Starbucks half-caf when I caught a glimpse of a friend's cup and narrowly avoided performing a Danny Thomas-style spit take. On the side of the paper cup was printed:

The Way I See It #224 "Darwinism's impact on traditional social values has not been as benign as its advocates would like us to believe. Despite the efforts of its modern defenders to distance themselves from its baleful social consequences, Darwinism's connection with eugenics, abortion and racism is a matter of historical record. And the record is not pretty."—Dr. Jonathan Wells, biologist and author of *The Politically Incorrect Guide to Darwinism and Intelligent Design*

I knew that Starbucks roasted the hell out of their beans, but I didn't realize they published half-baked ideas.

A visit to the Starbucks Web site turned up an explanation: "To get people talking, 'The Way I See It' is a collection of thoughts, opinions and expressions provided by notable figures that now appear on our widely shared cups." Further, the cups are supposed to extend "the coffeehouse culture—a way to promote open, respectful conversation among a wide variety of individuals."

Fair enough, although an open, respectful conversation initiated by a closed, disrespectful assertion is going to be a challenge, especially without any context. (To find some context, see, among myriad sources,

Ernst Mayr's essay "Darwin's Influence on Modern Thought" in the July 2000 *Scientific American* and Michael Shermer's 2006 book *Why Darwin Matters*.)

Nevertheless, I'd like to suggest some other quotes for Starbucks cups in the hopes that they, too, may stimulate piping-hot conversations.

The Way I See It #13.5 billion "On one corner, there was a Starbucks. And across the street from that Starbucks... was a Starbucks. At first I thought the sun was playing tricks on my eyes. But no. There was a Starbucks across from a Starbucks. And that, my friends, is the end of the universe."—Lewis Black, comedian and philosopher

The Way I See It # $e^{i\pi} + 1$ "The ratio of the diameter and circumference is as five fourths to four."—Indiana House Bill No. 246, which would have legally declared the value of π to be 3.2 in 1897

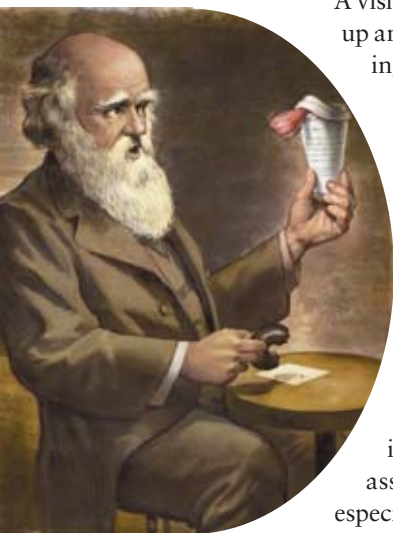
The Way I See It #Venti "Would you drink a Quarter Pounder with Cheese? If you order a venti (20-oz.) Starbucks Caffè Mocha, you might as well be sipping that 500-calorie burger through a straw. And a venti Starbucks Java Chip Frappuccino, with 650 calories and nearly a day's saturated fat, is a McDonald's coffee plus 11 creamers and 29 packets of sugar."—Center for Science in the Public Interest press release, September 5, 2006

The Way I See It #34 Proof "I'll quit coffee. It won't be easy drinking my Baileys straight, but I'll get used to it. It'll still be the best part of waking up."—Megan Mullally as Karen Walker on the TV series *Will & Grace*

The Way I See It #1962 "Mr. Kroger: two C's, two D's and an F. That's a 1.2 grade average. Congratulations, Kroger. You're at the top of the Delta pledge class."—Dean Vernon Wormer

The Way I See It #MVD "This coffee tastes like mud!" "It should—it was ground this morning."—encounter between Ricochet Rabbit and Deputy Droop-a-long in the 1960s, although no doubt in existence in some prior form

The Way I See It #Too "Popular, palatable views of the world and how it came to be do not constitute science or truth. But decent science education requires that we share the truth we find—whether or not we like it."—Lynn Margulis, Distinguished Professor, University of Massachusetts Amherst



How do itches come about, and why does it feel good to scratch them?

—B. ERICSSON, SWEDEN

Mark A. W. Andrews, associate professor of physiology at the Lake Erie College of Osteopathic Medicine, replies:

An itch, also known as pruritus, arises from the irritation of nerve cells associated with the skin. Pruritus serves as an important sensory and self-protective mechanism—as do many other skin sensations—by alerting us to harmful external agents. Constant itching, however, can become unbearable if the underlying condition is not treated.

Pruritus is a dominant symptom of many skin diseases and also occurs in some ailments that affect the entire body. An itching sensation results from the stimulation of pruriceptors—itch-sensing nerve endings—by certain mediators. These stimulating agents include chemicals for immune response (such as histamines) and pain relief (such as opioids); neuropeptides, which include endorphins and other pain-regulating messengers released within the brain; the neurotransmitters acetylcholine and serotonin, which relay impulses between nerve cells; and prostaglandins—lipids that, among other functions, create the sensation of pain in spinal nerve cells. Stimulation by any of these agents is typically related to inflammation, dryness or other damage to the skin, mucous membranes or conjunctiva of the eye.

Itching generally involves activation of the pruriceptors of specialized nerve cells called C-fibers. These C-fibers are identical to those associated with the sensation of pain, but they are functionally distinct and only convey the itch sensation. When stimulated on the skin, the C-fibers carry signals along the nerve to the spinal cord and on to the brain.

Scratching and rubbing interfere with the sensations arising from pruriceptors by stimulating various pain and touch receptors in the same areas. Like many sensory systems in the body, activation of one signal, in this case that of the pain and touch receptors, causes “surround inhibition” of another signal, that coming from the pruriceptor. This lack of pruricep-



tor firing “turns off” the itch sensation for a short period. Although it is helpful in relieving an itch, scratching offers only temporary relief and may cause the skin to become further irritated and possibly even to tear.

Despite approximately a century of research, no single effective antipruritic treatment exists. Several topical and orally administered agents can help suppress itching, however. These include lotions and creams (such as calamine and hydrocortisone), antihistamines, opioid antagonists (such as naltrexone, a drug used to treat narcotic or alcohol dependence), aspirin and ultraviolet light therapy.


How did the sun wind up in the middle of the solar system?

—A. SOMMERS, EAST BRUNSWICK, N.J.

Michael A. Jura, an astrophysicist at the University of California, Los Angeles, explains:

The best model of our solar system’s history states that the planets formed from a spinning disk of dust encircling the sun, leaving a collection of bodies with the sun at its center.

According to this model, the solar system formed from the collapse and flattening of an interstellar cloud. The cloud may initially have measured as much as a light-year across—10 million times wider than the sun. As the cloud compacted and cooled, its own gravity overpowered any forces acting to stabilize the system, causing it to further contract dramatically.

Before this collapse, the original cloud probably began with a fixed mass and a slight, random rotation relative to some central axis. Most of the cloud’s mass helped to form the sun, but some of it flattened out and remained as a disk encircling that newly created star. Observations from elsewhere in the galaxy indicate that this disk most likely gave birth to Earth and the other planets, leaving them in naturally heliocentric orbits. 

For a complete text of these and other answers from scientists in diverse fields, visit www.sciam.com/askexpert