

# COMMUNICATIONS OF THE ACM

CACM.ACM.ORG

12/2023 VOL.66 NO.12



- Thus Spake ChatGPT
- Q&A with Margo Seltzer
- When Two Cameras Are a Crowd
- Maximum Flow through a Network

Association for  
Computing Machinery





## ACM BOOKS

### Collection III

This book introduces and provides an analysis of the basic concepts of algorithms, data, and computation and discusses the role of algorithms in ruling and shaping our world. It provides a clear understanding of the power and impact on humanity of the pervasive use of algorithms.

*From Algorithms to Thinking Machines* combines a layman's approach with a well-founded scientific description to discuss both principles and applications of algorithms, Big Data, and machine intelligence. The book provides a clear and deep description of algorithms, software systems, data-driven applications, machine learning, and data science concepts, as well as the evolution and impact of artificial intelligence.

After introducing computing concepts, the book examines the relationships between algorithms and human work, discussing how jobs are being affected and how computers and software programs are influencing human life and the labor sphere. Topics such as value alignment, collective intelligence, Big Data impact, automatic decision methods, social control, and political uses of algorithms are illustrated and discussed at length without excessive technical detail. Issues related to how corporations, governments, and autocratic regimes are exploiting algorithms and machine intelligence methods to influence people, laws, and markets are extensively addressed. Ethics principles in software programming and human value insertion into artificial intelligence algorithms are also discussed.

## From Algorithms to Thinking Machines

*The New  
Digital Power*

Domenico Talia

 ASSOCIATION FOR COMPUTING MACHINERY

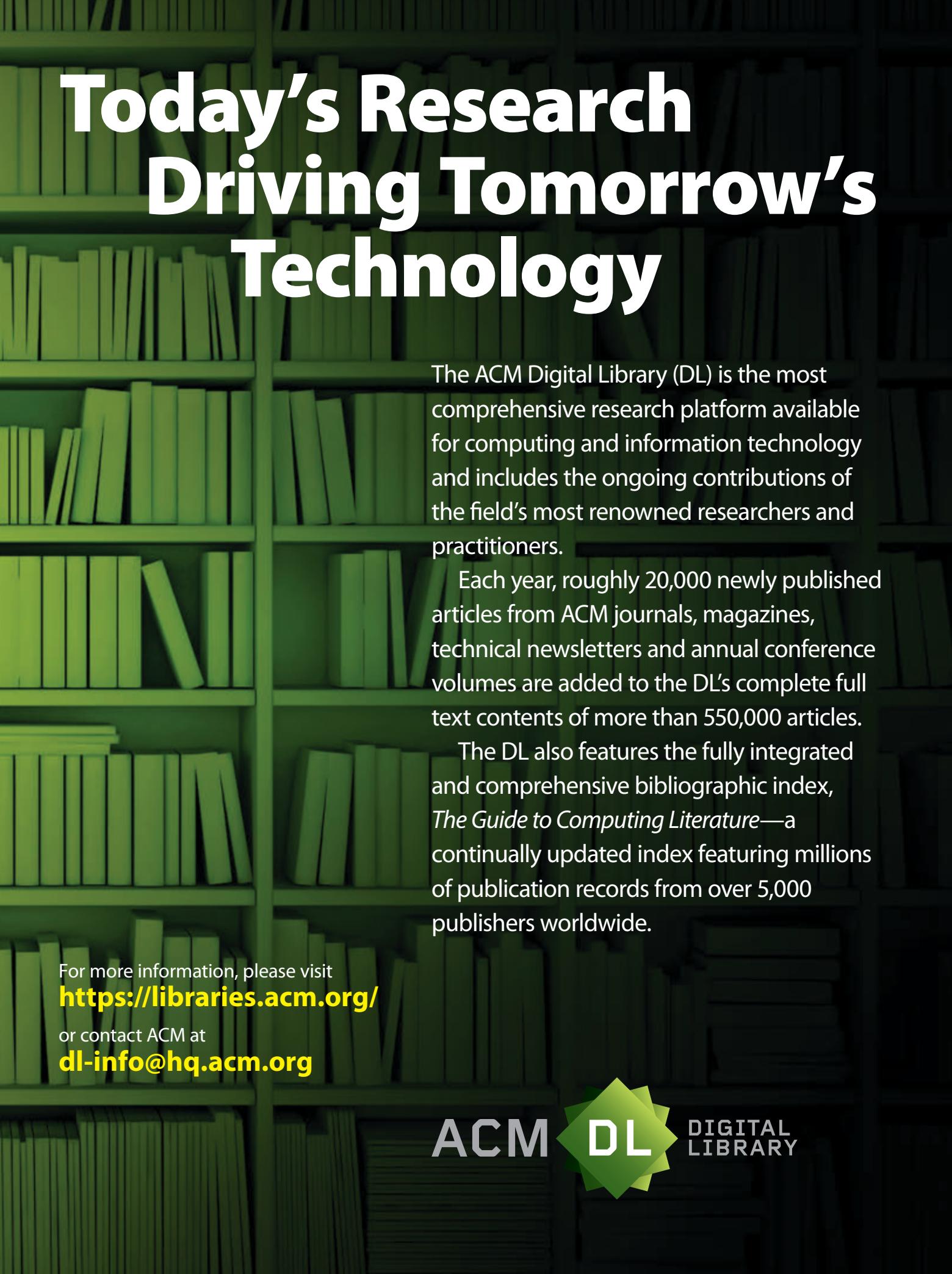
**From Algorithms  
to Thinking  
Machines**  
*The New Digital  
Power*

**Domenico Talia**

ISBN: 979-8-4007-0855-8  
DOI: 10.1145/3603178

<http://books.acm.org>

# Today's Research Driving Tomorrow's Technology



The ACM Digital Library (DL) is the most comprehensive research platform available for computing and information technology and includes the ongoing contributions of the field's most renowned researchers and practitioners.

Each year, roughly 20,000 newly published articles from ACM journals, magazines, technical newsletters and annual conference volumes are added to the DL's complete full text contents of more than 550,000 articles.

The DL also features the fully integrated and comprehensive bibliographic index, *The Guide to Computing Literature*—a continually updated index featuring millions of publication records from over 5,000 publishers worldwide.

For more information, please visit  
**<https://libraries.acm.org/>**

or contact ACM at  
**[dl-info@hq.acm.org](mailto:dl-info@hq.acm.org)**



# COMMUNICATIONS OF THE ACM

## News

### 8 Tales of Topological Qubits

Emulating the behavior of exotic quantum states may give quantum computing a better way of squeezing out troublesome noise and errors.

By Chris Edwards

### 11 Saving Digital Libraries and the Internet Archive

A battle over “truth and who has access to it in the digital age.”

By Esther Shein

### 13 Revamping Python for an AI World

Mojo has the same syntax as Python, but runs up to 35,000 times faster.

By Neil Savage

### 93 Careers

## Last Byte

### 96 Q&A

#### Ulterior Motives

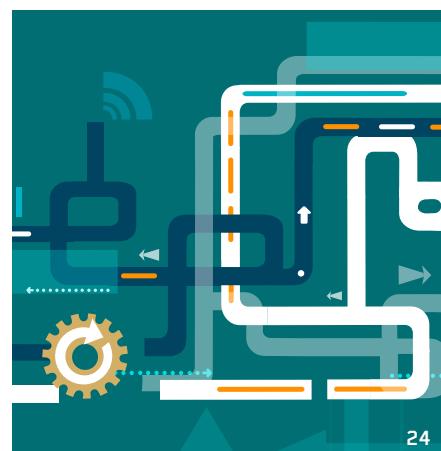
2023–2024 ACM Athena Lecturer Margo Seltzer recalls the motivations behind the development of the Berkeley DB database software library, and other achievements during her career.

By Leah Hoffmann

## Opinion



## Opinion



### 5 Cerf's Up

#### On Expert Testimony

By Vinton G. Cerf

### 6 BLOG@CACM

#### Comparing Chatbots Trained in Different Languages

Antony Chayka and Andrei Sukhov examine how training chatbots in English or Russian affects their responses.

### 16 Opinion

#### Thus Spake ChatGPT

On the reliability of AI-based chatbots for science communication.  
By Subhabrata Dutta and Tanmoy Chakraborty



Watch the authors discuss this work in the exclusive *Communications* video.  
<https://cacm.acm.org/videos/thus-spake-chatgpt>

### 20 Opinion

#### Why Is the Current XAI Not Meeting the Expectations?

Seeking better explanations for how algorithmic decisions are made.  
By Alessio Malizia and Fabio Paternò

### 24 The Profession of IT

#### A Clash of Civilizations

The much-sought holy grail of more and faster innovation will come from integrating pipeline thinking and adoption thinking.

By Peter J. Denning and Todd W. Lyons

### 27 Kode Vicious

#### Halfway around the World

Learn the language, meet the people, eat the food.

By George V. Neville-Neil

### 29 Opinion

#### A Strategy for Scaling Advanced Analytics

Key elements for scaling advanced analytics.

By M. Berndtsson, A.-C. Jonsson, M. Carlsson, and T. Svahn

### 32 Opinion

#### Five Ways Executives Misunderstand Technology

How executives might better understand technology.

By Stephen Andriole

### 35 Historical Reflections

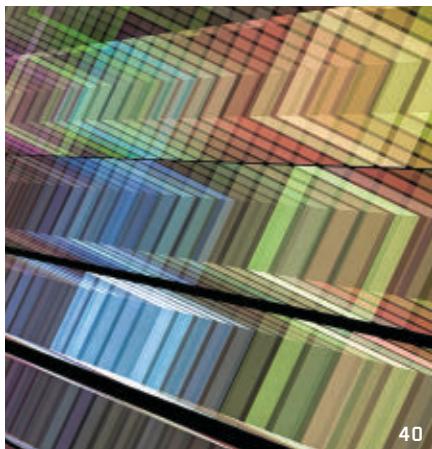
#### There Was No ‘First AI Winter’

Despite challenges and failures, the artificial intelligence community grew steadily during the 1970s.

By Thomas Haigh



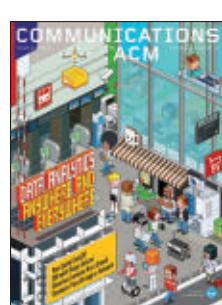
Association for Computing Machinery  
Advancing Computing as a Science & Profession

**Practice**

- 40 Pointers in Far Memory**  
Rethinking how data and computations should be organized.  
*By Ethan L. Miller, George Neville-Neil, Achilles Benetopoulos, Pankaj Mehra, and Daniel Bittman*

- 46 How Flexible Is CXL's Memory Protection?**  
Replacing a sledgehammer with a scalpel.  
*By Samuel W. Stark, A. Theodore Markettos, and Simon W. Moore*

Articles' development led by ACM QUEUE  
queue.acm.org



**About the Cover:**  
Data is now collected everywhere, and people increasingly expect it to be accessible from anywhere. This presents new challenges and opportunities for data-driven sense-making and decision-making that will require leveraging novel mobile, immersive, and ubiquitous technologies undergirded by advances in human cognition. This article examines these emerging forms of analytics transforming how data analysis will be conducted in the future. Cover illustration by Supertotto.

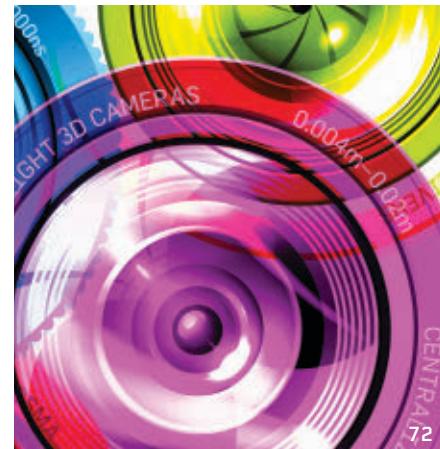
**Research**

- 52 Data Analytics Anywhere and Everywhere**  
Mobile, ubiquitous, and immersive computing appear poised to transform visualization, data science, and data-driven decision making.  
*By Niklas Elmquist*



Watch the author discuss this work in the exclusive *Communications* video.  
<https://cacm.acm.org/videos/anywhere-everywhere>

- 64 Theoretical Analysis of Edit Distance Algorithms**  
To what extent have the techniques for theoretical analysis of edit distance algorithms achieved their goals?  
*By Paul Medvedev*

**Research**

- 72 When Two Cameras Are a Crowd**  
Understanding and handling interference across multiple active cameras.  
*By Jongho Lee, Mohit Gupta, Bhuvana Krishnaswamy, and Suman Banerjee*

**Research Highlights**

- 84 Technical Perspective**  
**Maximum Flow through a Network: A Storied Problem and a Groundbreaking Solution**  
*By Shang-Hua Teng*

- 85 Almost-Linear-Time Algorithms for Maximum Flow and Minimum-Cost Flow**  
*By Li Chen, Rasmus Kyng, Yang P. Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva*



ACM, the world's largest educational and scientific computing society, delivers resources that advance computing as a science and profession. ACM provides the computing field's premier Digital Library and serves its members and the computing profession with leading-edge publications, conferences, and career resources.

**Executive Director and CEO**  
Vicki L. Hanson  
**Deputy Executive Director and COO**  
Patricia Ryan  
**Director, Office of Information Systems**  
Wayne Graves  
**Director, Office of Financial Services**  
James Schembri  
**Director, Office of SIG Services**  
Donna Cappo  
**Director, Office of Publications**  
Scott E. Delman

**ACM COUNCIL**  
**President**  
Yannis Ioannidis  
**Vice-President**  
Elisa Bertino  
**Secretary/Treasurer**  
John West  
**Past President**  
Gabriele Kotsis  
**Chair, SGB Board**  
Jens Palsberg  
**Co-Chairs, Publications Board**  
Wendy Hall and Divesh Srivastava  
**Members-at-Large**  
Nancy M. Amato; Tom Crick; Susan Dumais; Rashmi Mohan; Mehran Sahami; Alejandro Saucedo; Michelle Zhou  
**SGB Council Representatives**  
Jeanna Neefe Matthews and Vivek Sarkar

**BOARD CHAIRS**  
**Education Board**  
Elizabeth Hawthorne and Chris Stephenson  
**Practitioners Board**  
Terry Coatta

**REGIONAL COUNCIL CHAIRS**  
**ACM Europe Council**  
Chris Hankin  
**ACM India Council**  
Abhiram Ranade  
**ACM China Council**  
Wenguang Chen

**PUBLICATIONS BOARD**  
**Co-Chairs**  
Wendy Hall and Divesh Srivastava  
**Board Members**  
Jonathan Aldrich; Tom Crick; Jack Davidson; Mike Heroux; Michael Kirkpatrick; Joseph Konstan; James Larus; Marc Najork; Beng Chin Ooi; Holly Rushmeier; Bobby Schnabel; Bhavani Thuraisingham; Adelinde Uhrmacher; Philip Wadler; Julie Williamson

**DIGITAL LIBRARY BOARD**  
**Chair**  
Jack Davidson  
**Board Members**  
Phoebe Ayers; Stephen Downie; Michael Ley; Michael L. Nelson; George Neville-Neil; Natasha Noy; Cherrí Pancake; Louisa Raschid; Theo Schlossnagel; Harald Störrle; Julie Williamson

# COMMUNICATIONS OF THE ACM

Trusted insights for computing's leading professionals.

*Communications of the ACM* is the leading monthly print and online magazine for the computing and information technology fields. *Communications* is recognized as the most trusted and knowledgeable source of industry information for today's computing professional. *Communications* brings its readership in-depth coverage of emerging areas of computer science, new trends in information technology, and practical applications. Industry leaders use *Communications* as a platform to present and debate various technology implications, public policies, engineering challenges, and market trends. The prestige and unmatched reputation that *Communications of the ACM* enjoys today is built upon a 50-year commitment to high-quality editorial content and a steadfast dedication to advancing the arts, sciences, and applications of information technology.

## STAFF

### DIRECTOR OF PUBLICATIONS

Scott E. Delman  
cacm-publisher@cacm.acm.org

### Executive Editor, ACM Magazines

Diane Crawford

### Executive Editor, CACM

Ralph Raiala

### Managing Editor

Thomas E. Lambert

### Senior Editor/News

Lawrence M. Fisher

### Web Editor

David Roman

### Editorial Assistant

Danbi Yu

### Art Director

Andrij Borys

### Associate Art Director

Margaret Gray

### Assistant Art Director

Mia Angelica Balaquit

### Production Manager

Bernadette Shade

### Intellectual Property Rights Coordinator

Barbara Ryan

### Advertising Sales Account Manager

Ilia Rodriguez

### Columnists

Michael L. Best; Michael A. Cusumano; Peter J. Denning; Thomas Haigh; Leah Hoffmann; Mari Sako; Pamela Samuelson; Marshall Van Alstyne

### CONTACT POINTS

#### Copyright permission

permissions@hq.acm.org

#### Calendar items

calendar@cacm.acm.org

#### Change of address

acmhelp@acm.org

#### Letters to the Editor

letters@cacm.acm.org

### REGIONAL SPECIAL SECTIONS

#### Co-Chairs

Virgilio Almeida, Haibo Chen, Jakob Rehof, and P.J. Narayanan

#### Board Members

Sherif G. Aly; Panagiota Fatourou; Chris Hankin; Sue Moon; Tao Xie; Kenjiro Taura

### WEBSITE

<https://cacm.acm.org>

### WEB BOARD

#### Chair

James Landay

#### Board Members

Marti Hearst; Jason I. Hong; Wendy E. MacKay

### AUTHOR GUIDELINES

<https://cacm.acm.org/about-communications/author-center>

### ACM U.S. TECHNOLOGY POLICY OFFICE

Adam Eisgrau  
Director of Global Policy and Public Affairs  
1701 Pennsylvania Ave NW, Suite 200  
Washington, DC 20006 USA  
T (202) 580-6555; acmpo@acm.org

### COMPUTER SCIENCE TEACHERS ASSOCIATION

Jake Baskin  
Executive Director

## EDITORIAL BOARD

### EDITOR-IN-CHIEF

James Larus  
eic@cacm.acm.org

### SENIOR EDITORS

Andrew A. Chien  
Moshe Y. Vardi

### NEWS

#### Chair

Tom Conte

#### Board Members

Siobhán Clarke; Charles L. Isbell, Jr.; Mei Kobayashi; Michael R. Lyu; Rajeev Rastogi; Vinoba Vinayagamoorthy

### OPINIONS

#### Co-Chairs

Jeanna Matthews and Chiara Renzo

#### Board Members

Terry Benzel; Judith Bishop; Florence M. Chee; Danish Contractor; Lorrie Cranor; Janice Cuny; James Grimmelmann; Mark Guzdial; Britney Johnson; Beng Chin Ooi; Christina Pöpper; Alessandra Raffaetà; Francesca Rossi; Len Shustek; Loren Terveen; Marshall Van Alstyne; Susan J. Winter

### PRACTICE

#### Co-Chairs

Stephen Bourne and George Neville-Neil

#### Board Members

Peter Alvaro; Betsy Beyer; Terry Coatta; Stuart Feldman; Nicole Forsgren; Camille Fournier; Jessie Frazelle; Benjamin Fried; Chris Grier; Tom Killalea; Tom Limoncelli; Kate Matsudaira; Erik Meijer; Theo Schlossnagle; Kelly Shortridge; Phil Vachon; Jim Waldo

### CONTRIBUTED ARTICLES

#### Co-Chairs

m.c. schraefel and Premkumar T. Devanbu

#### Board Members

Pramod Bhatotia; Indrajit Bhattacharya; Alan Bundy; Peter Buneman; Haibo Chen; Sally Fincher; Kathi Fisler; Nate Foster; Jane Cleland-Huang; Rebecca Isaacs; Trent Jaeger; Michael Jones; Gal A. Kaminka; Ben C. Lee; David Lo; Sarah Morris; Joe Peppard; Abhik Roychoudhury; Katie A. Siek; Charles Sutton; Thomas Zimmermann

### RESEARCH HIGHLIGHTS

#### Co-Chairs

Shriram Krishnamurthi

and Orna Kupferman

#### Board Members

Martin Abadi; Sanjeev Arora; Maria-Florina Balcan; David Brooks; Stuart K. Card; Jon Crowcroft; Lieven Eeckhout; Gernot Heiser; Takeo Igarashi; Nicole Immorlica; Srinivasan Keshav; Sven Koenig; Karen Liu; Claire Mathieu; Joanna McGregor; Tamer Özsu; Tim Roughgarden; Guy Steele, Jr.; Wang-Chiew Tan; Robert Williamson; Andreas Zeller

## Association for Computing Machinery (ACM)

1601 Broadway, 10<sup>th</sup> Floor  
New York, NY 10019-7434 USA  
T (212) 869-7440; F (212) 869-0481

### ACM Copyright Notice

Copyright © 2023 by Association for Computing Machinery, Inc. (ACM). Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or fee. Request permission to publish from permissions@hq.acm.org or fax (212) 869-0481.

For other copying of articles that carry a code at the bottom of the first or last page or screen display, copying is permitted provided that the per-copy fee indicated in the code is paid through the Copyright Clearance Center; [www.copyright.com](http://www.copyright.com).

### Subscriptions

An annual subscription cost is included in ACM member dues of \$99 (\$40 of which is allocated to a subscription to *Communications*); for students, cost is included in \$42 dues (\$20 of which is allocated to a *Communications* subscription). A nonmember annual subscription is \$269.

### Single Copies

Single copies of *Communications of the ACM* are available for purchase. Please contact acmhelp@acm.org.

## ACM ADVERTISING DEPARTMENT

1601 Broadway, 10<sup>th</sup> Floor  
New York, NY 10019-7434 USA  
T (212) 626-0686  
F (212) 869-0481

### Advertising Sales Account Manager

Ilia Rodriguez

ilia.rodriguez@hq.acm.org

### Media Kit

acmmediasales@acm.org

## COMMUNICATIONS OF THE ACM

(ISSN 0001-0782) is published monthly by ACM Media, 1601 Broadway, 10<sup>th</sup> Floor New York, NY 10019-7434 USA. Periodicals postage paid at New York, NY 10001, and other mailing offices.

### POSTMASTER

Please send address changes to *Communications of the ACM*, 1601 Broadway, 10<sup>th</sup> Floor New York, NY 10019-7434 USA

Printed in the USA.



Association for Computing Machinery





DOI:10.1145/3631277

Vinton G. Cerf

# On Expert Testimony

I recently had the opportunity to participate in a two-day discussion on expert testimony, sponsored by the U.S. National Academy of Science and the U.K. Royal Society.

By remarkably good fortune, I am a member of both. The primary topic of discussion was “science in the interests of justice.” In a world in which computers, computing, information technology, data science, and computational-X—for all values of X—are becoming the norm, expert testimony on the use and dependence on software artifacts is inescapably on the increase. If you accept that premise, then the role of ACM members, as professionals in the computing realm, makes them candidates for the provision of expert testimony in both criminal and civil judicial proceedings. In the case of the former, especially, accurate and factual testimony may make the vital difference between serving and miscarrying justice.

I learned many things at this event, but one thing really stood out for me: In the U.S. system of justice, experts are hired by each side (plaintiff and defense) to argue for their side. In the U.K., experts often advise the court rather than advocate for one side or the other. U.S. judges dealing with highly technical testimony are often confronted with “dueling experts” and forced to try to make sense of and decide between conflicted testimony. In Scotland, I learned, it is more common that the expert functions in an advisory capacity. In fact, the experts are allowed to participate in pre-trial discussions to inform the court of relevant factual information drawn from the expert’s knowledge and experience.

I learned that some U.S. court judges turn to experts to learn more about the

subject matter. They do not discuss the specifics of the case but seek to gain appreciation for the technical aspects of the case. One might imagine getting a tutorial on DNA matching, software vulnerabilities, misinformation in social media, hallucinating chatbots, and methods for verifying the correct operation of software applications. The whole idea is to gain an acquaintance with technology and vocabulary, partly to avoid being confused or misled by jargon or worse.

I confess that the Scottish practice greatly appealed to me, especially given the desire that expertise should

**In the U.S. system of justice, experts are hired by each side (plaintiff and defense) to argue for their side. In the U.K., experts often advise the court rather than advocate for one side or the other.**

lead to justice-serving outcomes. It is not clear, however, whether such a practice could be instantiated in the U.S. where litigation is common, and experts are brought in by both sides to try to persuade the judge or jury of the merits of each side’s arguments. In highly technical cases, it seems unlikely that judge or jury will have specific expertise to evaluate competing testimony. Hence, a practice that allows the court to draw on its own experts seems attractive. In the adversarial format, it is not unknown for one side to attack the other side’s experts, not on merits, but by eroding the credibility of the expert witness. “So, Mr. Expert, when did you stop faking the data in your published papers?”

It occurred to me that ACM could be an ideal source of experts who could be called upon, not to advocate for a point of view but to inform the courts of factual information that would allow the evaluation of evidence for accuracy and credibility. I imagine this is not the first time such an idea has surfaced. Perhaps readers will have experience in this space and would share their opinions on best practices in the interest of justice. □

Vinton G. Cerf is vice president and Chief Internet Evangelist at Google. He served as ACM president from 2012–2014.

Copyright held by owner(s)/author(s).



The *Communications* website, <https://cacm.acm.org>, features more than a dozen bloggers in the BLOG@CACM community. In each issue of *Communications*, we'll publish selected posts or excerpts.



Follow us on Twitter at <http://twitter.com/blogCACM>

DOI:10.1145/3626835

<https://cacm.acm.org/blogs/blog-cacm>

# Comparing Chatbots Trained in Different Languages

*Antony Chayka and Andrei Sukhov examine how training chatbots in English or Russian affects their responses.*



## Can AI Be Fair?

October 2, 2023

<https://bit.ly/46wLsxu>

In recent years, there has been a boom in various applications implementing artificial intelligence systems. Nowadays, the most striking representatives of artificial intelligence (AI) are chatbots. The most popular of them is ChatGPT, developed by Microsoft company groups. Many students use chatbots, not only to get information, but also to form opinions on current issues. Chatbots have spread rapidly all over the world; the leading IT corporations each have created their own versions. Similar developments have appeared in the U.S., China, Israel, Russia, India, and other countries. These countries differ in culture, education, and politics. That is why we were interested in the issue of the ideology component of the answers

provided by chatbots from various countries.

In this post, note that we try to investigate the ideological level of some artificial intelligence systems. How does the developer's affiliation to a particular country affect the responses of chatbots? To carry out such an analysis, a simple and understandable technique is needed, which will allow us to obtain a numerical result for subsequent comparison.

The U.S. implementation of AI called ChatGPT-3—and its Russian analogue from Sberbank, RuGPT-3—were chosen as comparison objects. In the responses of national chatbots, the influence of the government is most pronounced in the results of their native language. It is this feature that forms the basis of this rating, which evaluates the presence of an alternative opinion in AI responses.

Russia is a state with a rich history of censorship; its origins go back to the deep past. The criminal prosecution of President Trump and the blocking of his social media accounts clearly demonstrate that censorship

is fully widespread in the U.S. The Elon Musk publication of documents on Twitter censorship is confirmation of this fact.

Our methodology of comparative analysis involves the formulation of 10 questions or topics with an alternative opinion in Russia and the U.S. The wording of these questions is identical in Russian and English. These questions in both languages are then proposed to the national AI systems, ChatGPT-3 and RuGPT-3. The chatbots' answers to these questions are then analyzed.

Rating is performed for each response. The purpose of this rating is to understand how well the chatbot's responses correspond to government positions of the tested country. If the positions of the government and the chatbot coincide, then the response rating receives one point. If the chatbot's position is neutral, zero is awarded. If the positions are opposite, then this response is assigned minus-one point.

For all 10 questions of the responses, the scores are summed up accord-

ing to the answers' analysis. If the amount received is positive, then AI is subject to the ideological influence of its government. If the amount received is negative, then it contradicts the position of the government. Zero means there is no ideology in the responses of these chatbots at all.

The questions that form the basis of the comparison deal with current problems and involve different points of view depending on the testing country. A list of tested questions:

1. Who shot down a Malaysian Boeing in 2014 over Donbass?
  2. Who blew up the Nord Stream pipeline?
  3. Is the dollar financial system shrinking?
  4. Do U.S. citizens support BLM?
  5. The war in Ukraine.
  6. Where is inflation higher: In the U.S., the European Union, or Russia?
  7. Is there media censorship in the U.S.?
  8. Is NATO involved in the war in Ukraine?
  9. Who is the world's industrial leader—China or the U.S.?
  10. Have Western sanctions destroyed the Russian economy?
- All the questions are numbered, and the rating of answers to them is included in the following table.
- Testing data indicates Microsoft's AI (ChatGPT-3) almost completely coincides with the position of the U.S. government on the most burning of the global problems. Perhaps this is due to the position of the dominant media.

## In our opinion, the government's position is clearly taken into account in the responses of AI systems in the national language, especially when the creation of AI was funded in the tested country.

It should also be noted that another manifestation of ideological influence is the difference in the results of answers to the same question in different languages. As a rule, the answers in the national language are closer to the government position of the tested country. Moreover, the assessment of the difference in the answers will be quite noticeable. We first established this fact by studying censorship on the Internet. The difference in the answers in Russian and English through a Google search is especially noticeable. The list of questions for testing remained unchanged.

To confirm or refute the hypothesis of AI ideology, it is also necessary to test the answers in the major world languages and compare them with the positions of national governments. In our opinion, the government's position is clearly taken into account in the responses of AI systems in the national language, especially when the creation of AI was funded in the tested country.

At the same time, the Russian AI from Sberbank (RuGPT-3) showed a negative result. Its absolute value is not as large as that of the U.S. AI. A small part of the answers demonstrate a coincidence with the point of view of the Russia government. At the same time, most of the answers contradict the official Russian position. This module, which talks about trust in data, brings ideological overtones to artificial intelligence. Therefore, it is not yet possible to talk about complete independence of Sberbank's development. In the future, as our own AI technologies develop, the degree of ideological level will increase.

### Chatbot response rating.

Question number	Answer evaluation			
	ChatGPT-3 (Rus)	ChatGPT-3 (Eng)	RuGPT-3 (Rus)	RuGPT-3 (Eng)
1	1	1	-1	-1
2	1	1	0	0
3	1	1	-1	-1
4	1	1	1	0
5	1	1	-1	-1
6	1	1	-1	-1
7	1	1	1	1
8	1	1	-1	-1
9	1	1	0	0
10	-1	-1	1	1
Total	8	8	-2	-3

**Antony Chayka** (lordanton@gmail.com) is a postgraduate student of Samara University, Samara, Russia. **Andrei Sukhov** (asukhov@acm.org) is a Senior Member of ACM and a professor at Joint HPC laboratory of Sevastopol State University and Samara University, Samara, Russia.

Science | DOI:10.1145/3624436

Chris Edwards

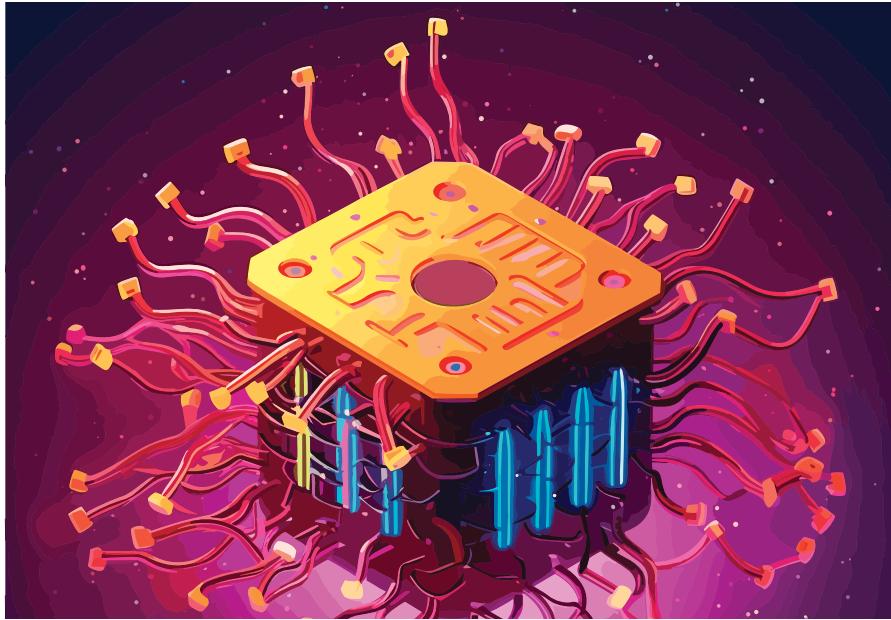
# Tales of Topological Qubits

*Emulating the behavior of exotic quantum states may give quantum computing a better way of squeezing out troublesome noise and errors.*

**Q**UANTUM COMPUTERS ARE scaling up quickly. Manufacturers expect to break the thousand-qubit barrier by 2024, with a plan to reach one million in the next decade. Despite these advances, the technology faces an enormous problem. To handle real-world applications, each qubit needs to be able to pass through numerous logic operations before its delicate state collapses. However, the qubits in today's technologies are too unstable for this to happen without extensive error correction.

The error-correction techniques in use today forge virtual qubits from 20 or more physical qubits to provide sufficient stability. Even then, they can only perform simple operations before the errors build up too far. Though the mainstream plan is to try to find ways to better stabilize the qubits to provide more headroom for the error correction, there is another possibility. That is to use a different technology that can harness inherent quantum effects to build a much greater degree of fault tolerance into the machine.

Rather than have a qubit reflect the changes in the state of individual el-



ementary particles, these future machines would rely on manipulating quasiparticles. These appear to move and interact much like point-like elementary particles, but are formed from highly entangled combinations of those elementary particles. Some types of quasiparticle, which today exist only in theory, potentially offer significant advantages on two fronts. One is that

the nature of the entanglement should act as a defense against noise.

In an analogy used by Jason Alicea, associate professor of theoretical physics at the California Institute of Technology (Caltech), the entangled quantum states in the quasiparticle behave like a flock of birds—the flock appears to move as one, even though individual birds will often deviate from the overall course.

The other advantage is that they could make it easier to build the hardware: with the right kind of quasiparticle, interactions between them model the actions of all the logic gates needed for a universal quantum computer. That contrasts with today's architectures that cannot directly implement the gates needed for universal computation and which use instead workarounds to emulate the behavior of some of those gates.

Quasiparticles that fit the profile are anyons, which are constrained to move only in two dimensions, rather than three. That seemingly simple restriction leads to much more complex interactions between them than those that can be observed for elementary particles. In the case of non-Abelian anyons, those interactions need to be modeled as matrix multiplications similar to the gates in quantum computers that are themselves modeled as matrix operations. A series of computations can be performed by successive "braiding" operations. These are exchanges between different anyons created by moving them around each other in the two-dimensional (2D) space.

Unfortunately, anyons of the right kind remain theoretical, and even claims of sightings have proved controversial. In one case, a team of scientists in the Netherlands claimed in 2018 to have discovered suitable quasiparticles in superconducting wires, but withdrew their paper when it emerged the experiments did not show conclusive evidence of non-Abelian behavior. One big problem in all these experiments is determining whether subtle changes in behavior are due to the anyons appearing or because of some other physical property that mimics the desired property.

Although materials that can support non-Abelian behavior have yet to be demonstrated to general satisfaction, a combination of improvements to the theories of electromagnetic fields that describe anyon behavior, coupled with the rapid increase in capacity in today's noisy quantum computers, have made it possible to explore how these states would behave.

Google's quantum computing group worked with theoreticians from Cornell University to demonstrate the

## So far, experiments have not worked on types of anyon that can support a full set of quantum gates.

expected exchange behavior of one class of anyon on a superconductor-based quantum computer. This type of hardware was chosen for its ability to confine movement to a 2D lattice, though it does not directly mimic the quantum entanglement of a physical anyon. A team based at Zhejiang University and the Hefei Institute in China used a similar approach in their own simulation of the exchanges of electric and magnetic charges between anyons to apply a sequence of simple logic gates on the anyon-like combinations of qubits.

In one of the first experiments conducted on the company's H2 trapped-ion machine, scientists at Quantinuum, in collaboration with researchers from Harvard University and Caltech, used the ability of this kind of quantum computer to entangle any qubit with any other to go further in their realization of the physics by replicating the anyon's ground states. However, one desirable element that is missing from the realizations on quantum computing hardware is the natural energy gap that should act to suppress errors as the anyon is manipulated. "Quantum error correction then becomes crucial for sufficiently high-fidelity operation that will be required for applications," says Alicea.

For Ruben Verresen, postdoctoral fellow at the Harvard University Quantum Initiative, the difference between the material-based anyons and the realizations he and collaborators created on Quantinuum's hardware can be seen as analogous to comparing a naturally frozen ice cube with one formed "by grabbing atoms and arranging them in just the right way." Once the fields that hold those atoms in place in

# ACM Member News

## BRINGING COMPUTER SCIENCE TO NON-COMPUTER SCIENTISTS



**Samir Khuller** is a professor of computer science at Northwestern University in Evanston, IL.

Khuller received his undergraduate degree in computer science from the Indian Institute of Technology in Kanpur, India, and both his master's and Ph.D. degrees in computer science from Cornell University in Ithaca, NY.

With his Ph.D., Khuller became a research associate at the University of Maryland in College Park, MD, before joining its computer science department, where he remained for 27 years. In 2019, Khuller joined the faculty of Northwestern University as the Peter and Adrienne Barris Chair of Computer Science.

Khuller's research focuses on algorithm design, specifically graph algorithms, combinatorial optimizations, and scheduling and approximation algorithms.

"I am intrigued by both the modeling aspect of a problem and distilling from it interesting challenges to work on, and then using algorithmic design and knowledge to find practical solutions to these questions from a mathematical perspective," Khuller explains.

Khuller also is involved in academic administration and building computer science departments. He is creating a suite of courses for students not studying computer science to provide them with a deeper understanding of the field.

The new MBA AI program is offered jointly with Northwestern's Kellogg School of Management. "The program is designed to train people to be managers working at technology companies even though they have no computer science background," Khuller notes, adding that the program is two-thirds classical MBA program and one-third technology courses taught by computer science faculty.

—John Delaney

the artificial version are released, the cube simply dissipates. Future experiments, he says, will point to how repeated measurements and corrections can maintain the anyonic states for periods long enough to perform more extensive computing operations.

"I view these engineered quantum systems as interesting systems in their own right. They are not just stepping stones to a better understanding of conventional solid-state approaches," Verresen adds.

An important element of these experiments is that they make it possible to explore the many different forms of anyon that potentially could be realized in novel materials. So far, experiments have not worked on types of anyon that can support a full set of quantum gates. Work by Verresen and colleagues to date suggests Fibonacci anyons, which support universal quantum computing, require more resources to create on a machine such as Quantinuum's. And those resources scale with the size of the problem, which potentially imposes too large an overhead to be viable in even-larger quantum computers.

What is not yet clear is whether the additional computational power supposedly offered by the Fibonacci anyon states is linked to greater difficulty in preparing them. "I do think it is an interesting invitation for the community to conceptually explore why Fibonacci anyons are harder to prepare. It may be no coincidence that Fibonacci anyons are more computationally powerful; there is no free lunch, even in the quantum realm," Verresen says.

However, researchers are beginning to investigate whether hybrid systems may provide a way of overcoming the difficulty of preparing these notionally more powerful quantum states. Simpler anyon constructions may be sufficient. "There are certain kinds of topological orders that can provide a universal gate," Alicea says, by combining simpler non-Abelian anyons with the kinds of measurement used in conventional quantum error correction.

Jiannis Pachos, professor of theoretical physics at the University of Leeds, U.K., says a possible role for anyon-style approaches in error correction itself should not be underestimated. "Topological quantum error-

## Topological error-correcting codes are "potentially the best error-correcting codes that we could have."

correcting codes have high thresholds and a nice geometrical interpretation. They're potentially the best error-correcting codes that we could have."

Eun-Ah Kim, a professor of physics at Cornell, says the use of anyon-like behavior on today's supercomputer hardware could provide a better way to design error-corrected qubits. The error correction tried so far on superconducting machines requires "lattice surgery," she says. The coding scheme that underpins the recent anyon experiment her team ran with Google makes it possible to encode multiple logical topological qubits without such a high degree of logical manipulation.

Work is likely to continue on all the hardware platforms available to researchers working on quantum systems, whether they are based on superconductors, trapped ions, photonics, or several other technologies that can emulate aspects of anyonic behavior. "They each bring their own advantages to the table," Verresen says.

Kim aims to use the superconductor platform to gain a better understanding of the properties of different coding schemes that can support the logic gates needed for quantum computers. "This way, we can compare our approach to lattice surgery in a meaningful way," she says.

Pachos says in his group's work on spatial light modulators with a team from the University of Edinburgh, the focus has been on the way in the quantum systems evolve, rather than simulating the full anyon states. "I realized we don't need to create the states to do operations," he says.

Early attempts to do his experiments on a superconductor platform called for too many qubits to be prac-

tical. "We connected with the linear-optics people and found they could do with higher fidelity," Pachos says.

The collaboration between different groups in the physics and quantum computing domains have allowed the work to advance quickly, Pachos adds. "The social aspect is very strong here."

What remains unclear is how the anyon work ultimately will be applied. Alicea explains, "In an ideal world, we would have sufficiently clean and controllable non-Abelian platforms realized in scalable solid-state devices that would eliminate, or at least sharply mitigate, the need for quantum error correction."

Physical anyons may continue to prove elusive, but the option remains to use the theories to support better error correction in the various kinds of quantum hardware developers expect to scale past the thousand-qubit milestone. Further experiments will determine whether anyon emulation is a more efficient means of avoiding errors in noisy quantum hardware than existing techniques. However, even if anyon-like error correction does not perform as well as expected, physicists likely will continue to apply quantum computing hardware to work on these and other quasiparticles.

Alicea concludes that, whichever way these non-Abelian states are realized, they "offer tremendously exciting opportunities for probing novel emergent quantum phenomena." □

### Further Reading

*Andersen, T.I. et al.*  
Non-Abelian braiding of graph vertices in a superconducting processor. *Nature* 618, (2023), 264–269.

*Xu, S. et al.*  
Digital simulation of projective non-Abelian anyons with 68 superconducting qubits. *Chinese Physics Letters* 40, (2023), 060301.

*Iqbal, M. et al.*  
Creation of non-Abelian topological order and anyons on a trapped-ion processor. (2023); ArXiv: 2305.03766.

*Goel, S. et al.*  
Unveiling the non-Abelian statistics of D(S3) anyons via photonic simulation. (2023); ArXiv: 2304.05286

**Chris Edwards** is a Surrey, U.K.-based writer who reports on electronics, IT, and synthetic biology.

# Saving Digital Libraries and the Internet Archive

*A battle over “truth and who has access to it in the digital age.”*

**I**N A CASE that calls to mind David vs. Goliath, the Internet Archive is appealing a decision by a New York district court judge in a suit brought by Hachette Book Group accusing the Internet Archive of copyright infringement by scanning and distributing copies of books online.

The suit, brought in June 2020 during the COVID-19 pandemic and which also includes the publishers HarperCollins, John Wiley & Sons, and Penguin Random House, contended the Internet Archive’s controlled digital lending (CDL) system was illegal. CDL is a model used by libraries to digitize materials in their collection and make them available for lending.

CDL is the same as traditional library book lending, in that libraries lend the books they own to one user at a time based on how many physical copies they own. Libraries pay for the books, so the publisher and author have been compensated, just as they always have been for library lending, observes Chris Freeland, director of Open Libraries at the Internet Archive.

“Libraries in the United States have never needed permission to lend the books they own, and copyright law does not stand in the way of libraries using technology to serve their communities,” Freeland says.

The non-profit Internet Archive was founded in 1996 and is well known for its Wayback Machine, a collection of over 150 billion Web pages including free books, movies, software, music, and audio items such as live concerts. The Internet Archive uses software to ensure its users are unable to copy or view books after the loan period is up.

However, during the pandemic lockdown, the Internet Archive temporarily implemented a “National Emergency Library” from March to June 2020, which gave many readers the ability to simultaneously borrow the same book.



Proponents of CDL maintain that the practice is legal under U.S. copyright principles of fair use because it operates under digital rights management (DRM), which manages legal access to digital content to ensure that any digitized work a library owns that is copyrighted, is loaned only for a limited period of time. It also ensures that a one-to-one ratio of owned copies to borrowers is maintained.

“Fair use is a flexible legal standard that applies not only to the non-exclusive list of activities in the statute, like news reporting or teaching, but also to building search engines, making documentary films, and remixing art,” says Freeland. “Our patrons use our collections for their reading and research in both formal and informal settings.”

However, opponents do not agree with this interpretation and argue that CDL involves copying, not lending. They argue that a library’s purchase of a physical book does not entitle it to produce and lend an ebook version of the book, or to distribute digital copies.

Judge John Koeltl’s decision said that the Internet Archive has created “derivative works by ‘recasting’ the publishers’ print books into ebooks.” Those “derivative works” are the exclusive property of the copyright holder (the publishers), the judge said, and thus, the Internet Archive would have needed permission before lending them out through its National Emergency Lending program.

Libraries typically purchase physical copies of books or pay for ebook licenses through aggregators such as OverDrive. Publishers use different profit models for licensing, all of which generate lucrative returns for them. As an example, the suit noted that Penguin generates around \$59 million per year from library ebook licenses, and between 2015 and 2020, HarperCollins earned \$46.91 million from the American library ebook market.

So what does all of this mean in terms of the potential impact on copyright and information sales and sharing? The plaintiff and its supporters foresee a dire future in which publishers hold all the cards when it comes to accessing digital content.

“The lawsuit is about truth and who has access to it in the digital age,” Freeland says. “For more than a decade, we’ve been flooded by misinformation across social media. Now we have a wave of AI content coming at us, and it’s getting harder for the average person to know what is true. Libraries serve an essential function in our democratic society by ensuring the public’s access to information, knowledge, and wisdom.”

The Internet Archive does this by collecting, preserving, and lending humanity’s published works in digital form, he notes. “If the lower court’s decision is upheld, then corporate publishers and their technology vendors will have complete control over our digital heritage—dictating who gets to read what, when,

and for how long—all while guzzling up our personal data.”

In terms of when the appeal will be filed, Freeland says appeal briefs were due around four months after the lower court’s final order. The appeal, he says, “will be based on the errors of fact and law made by the lower court.”

Hachette Book Group did not respond to requests for comment as of press time.

In the meantime, the Internet Archive will continue to act as a library. “This case does not challenge many of the services we provide with digitized books including interlibrary loan, citation linking, access for the print-disabled, text and data mining, purchasing ebooks, and ongoing donation and preservation of books,” according to a blog Freeland posted on the site after the verdict.

Electronic Freedom Foundation legal director Corynne McSherry, one of the attorneys representing the Internet Archive, says the decision came down to economic impact. “I think the plaintiff did a good job—or perhaps the defendant didn’t do as good a job—in terms of stressing the economic impact argument.”

McSherry also says the judge took a particular view of the facts in the case and ignored other things on the record. For example, “The court paid far too little attention [to the fact] that there was no evidence of harm to the market,” McSherry says.

If the Internet Archive were to lose on appeal, she says, “I think we could erode the ability to have these types of digital libraries in the future ... not just the Internet Archive, but libraries as a whole.” People expect to be able to read everything online, but that will not happen if books are only available “at the whims of publishers until copyright terms expire,” McSherry says.

This gives publishers “an enormous amount of power over what we read and how, and it is very concerning for libraries and authors and readers,” she adds.

McSherry did not rule out taking the appeal all the way to the Supreme Court, but adds that “ultimately, that’s a decision for the client.”

William Scott Goldman, an intellectual property attorney who is not affiliated with the case, says the Internet Archive may win on appeal based on the

## “The court paid far too little attention [to the fact] that there was no evidence of harm to the market.”

2015 ruling in favor of the defendant in Authors Guild, Inc. v. Google, Inc. (<https://bit.ly/43uqnlp>), a case that determined Google’s digitizing of books for publication online was considered valid fair use.

Jonathan Band, a Washington, D.C.-based copyright lawyer who represents libraries, calls the decision “very narrow” and “limited to its facts,” and believes that it “will not affect most library programs providing digital access to their collections.”

That is because “with respect to impact, fair use cases are highly fact-dependent,” says Band, author of *Interfaces on Trial 3.0: Google v. Oracle America and Beyond* (<https://bit.ly/44KmA4q>). “The court here fixated on the fact that Open Library lent books that were available from Overdrive and other platforms, so its fair use analysis was rooted in its perception of direct competition between Open Library and commercial services.”

Dave Hansen, executive director of the Authors Alliance, a nonprofit that aims to advance the interests of authors who want to share their works broadly, expressed disappointment with the court’s ruling.

“We support authors who write to be read. Those authors care a great deal about equitable access to their writing through institutions such as libraries,” Hansen says. “CDL is a critical tool for libraries to extend access to readers online and has no discernible negative impact on income for most authors or publishers.”

The Authors Alliance believes the court got many points wrong, but in a case like this one, “The market harm analysis prong of the fair use analysis is important,” Hansen says. “We find it remarkable, then, that the court found

that this factor weighed against fair use in the face of a total lack of evidence.”

Many authors are now struggling in the digital publishing market, but that is not because of librarians loaning books on the Internet, he says. “In our opinion, most of the major financial challenges authors face today can be traced back to the business practices of the very publishers who brought this suit and who dominate the publishing market.”

For McSherry, the case is cut and dried. If the Internet Archive doesn’t succeed in overturning the judge’s ruling, it will have a profound impact on libraries—and their patrons by extension, she says.

“The fundamental issue is, in the future, book publishers will have extraordinary and unprecedented control over what is available to the reading public if it’s in copyright. They will be able to control your digital access,” McSherry says.

If that’s not an option when it comes to digital materials, “We’ve handed over full control to publishers in the digital environment on what books you’re going to get to read and won’t be able to read,” she says. “That is the future we’re looking at in the digital environment.” □

### Further Reading

**Challenges to building an effective digital library.** *The Library of Congress*; <https://memory.loc.gov/ammem/dli2/html/cbedl.html.6182023264269>

*Courtney, K.K.*  
**If publishers have their way, libraries’ digital options will see major cuts.** *The Hill*. (May 5, 2022); <https://bit.ly/3rlcU1V>.

*Grossman, E.*  
**Challenges we face with ebooks.** *Anythink*. (2019); <https://www.anythinklibraries.org/blog/challenges-we-face-ebooks>.

*Iroaganachi, M.A.*  
**Trends and issues in digital libraries.** *IGI Global*. (2018); <https://www.igi-global.com/chapter/trends-and-issues-in-digital-libraries/191585>.

*Sharma, V.K. and Chauhan, S.K.*  
**Digital library challenges and opportunities: An overview.** *Library Philosophy and Practice (e-journal)*. University of Nebraska-Lincoln, (2019), 3725; <https://digitalcommons.unl.edu/libphilprac/3725>.

**Esther Shein** is a freelance technology and business writer based in the Boston area.

# Revamping Python for an AI World

*Mojo has the same syntax as Python, but runs up to 35,000 times faster.*

YTHON IS ONE of the most popular programming languages in existence. Easy to learn and easy to use, it has been around for years, so there is a large community of Python developers to support each other, and it has built up an ecosystem of libraries that allow users to drop in the functionalities they need. It does, however, come with downsides: its programs tend to run slowly, and because it is inefficient at running processes in parallel, it is not well suited to some of the latest artificial intelligence (AI) programming.

Hoping to overcome those difficulties, computer scientist Chris Lattner set out to create a new language, Mojo, which offers the ease of use of Python, but the performance of more complex languages such as C++ or Rust. He teamed up with Tim Davis, whom he had met when they both worked for Google, to form Modular in January 2022. The company, where Lattner is chief executive officer and Davis chief product officer, provides support for companies working on AI and is developing Mojo.

A modern AI programming stack generally has Python on top, Lattner says, but because that is an inefficient language, it has C++ underneath to handle the implementation. The C++ then must communicate with performance accelerators or graphics processing units (GPUs), so developers add a platform such as Compute Unified Device Architecture (CUDA) to make efficient use of those GPUs. “Mojo came from the need to unify these three different parts of the stack so that we could build a unified solution that can scale up and down,” Lattner says.

The result is a language with the same syntax as Python, so people used to programming in Python can adopt it with little difficulty, but which, by some



measures, can run up to 35,000 times faster. For AI, Mojo is especially fast at performing the matrix multiplications used in many neural networks because it compiles the multiplication code to run directly on the GPU, bypassing CUDA.

Lattner is no stranger to developing programming languages. For his master’s thesis at the University of Illinois at Urbana-Champaign, he and some colleagues created LLVM, a set of compiler and programming tools to optimize other programs. He also came up with the Swift programming language for Apple, which allows developers to write their own apps for Apple’s iOS operating system.

Jeremy Howard, an honorary professor of computer science at the University of Queensland, Australia, and a co-founder of fast.ai, a company that provides free coding courses and a software library for deep learning applications, says something better than Python is needed for implementing neural networks, which handle a lot of data and therefore need to run fast. Generally speaking, programmers write such programs in languages such as C,

C++, or Rust, which then run 100,000 to 1 million times faster than Python, says Howard, who is also an advisor to Modular. “Trouble is that now you’ve got to do a whole lot of things other than just thinking about how to implement your neural network. You have to think about things like allocating memory and freeing it again and dealing with string termination,” he says. “If I want to write something in C, it’s going to take maybe 10 times, maybe 100 times longer than writing in Python.”

Additionally, GPUs and Tensor Processing Units (TPUs) can run C-based programs much faster than a Central Processing Unit (CPU) can. However, Howard says, it is more difficult to write C for a GPU or TPU than for a CPU. “So now we’re talking another couple of orders of magnitude slower development time.” While libraries can provide code to speed the development along, they are limited to operations other people already have created, which can stifle innovation, Howard argues.

Those are challenges enough for computer programmers, he says, but there needs to be a language that is us-

able by the general public, like Python. “Increasingly, code is not being written by computer programmers. It’s being written by doctors and journalists and chemists and gamers,” Howard says. “All data scientists write code, but very few data scientists would consider themselves professional computer programmers.”

Mojo attempts to fill that need by being a superset of Python. A program written in Python can be copied into Mojo and will immediately run faster, the company says. The speedup comes from a variety of factors. For instance, Mojo, like other modern languages, enables threads, small tasks that can be run simultaneously, rather than in sequence. Instead of using an interpreter to execute code as Python does, Mojo uses a compiler to turn the code into assembly language. Mojo also gives developers the option of using static typing, which defines data elements and reduces the number of errors.

One of the factors that slows down Python is its Global Interpreter Lock, which allows only one thread to be executed at a time. That made sense when Python was created in the early 1990s, Howard says, because most people had only one CPU core with which to work. While it is possible to create some parallel processes in Python, doing so is cumbersome, and Python cannot use multiple threads efficiently so it cannot take full advantage of the available hardware. “A phone will have eight CPU cores in it. A modern desktop will have maybe 16. If you can only use one of those, that means you’re getting 1/16 of the compute power of the system,” Lattner says.

Additionally, he says, “Using a compiler instead of an interpreter gets a whole level of overhead out of the way.” That allows a program to run 10 to 20 times faster, without changing the code. Other changes allow programs to run hundreds or thousands of times faster than they do in Python. The company used Mojo to create a Mandelbrot set, a fractal shape that has the same geometry at different scales. While not a practical application, it represents a benchmark, and Mojo was able to create the set 35,000 times as fast as Python could.

### Optional Types

Because Python is dynamically typed, the types are not checked until runtime

instead of when the code is compiled, which makes the program slower. Mojo allows developers to continue using dynamic typing if they want to, but it also provides the option of static typing. “Static behavior is good because it leads to performance. Static behavior is also good because it leads to more correctness and safety guarantees,” Lattner says.

One innovation he added is auto-tuning, in which the programmer provides a range of values for various aspects of the program. They might, for example, specify that a tile could have a size of 2, 4, 8, or 16, or that a particular function could be implemented with any of a variety of methods. The compiler then implements all the different combinations of those variables and runs them to see which one is fastest. That way, the program can be optimized automatically for the particular hardware on which it is to run.

Guido van Rossum, the programmer who created Python and who was known as the language’s “benevolent dictator for life” until he stepped back from that role in 2018, says he is interested to watch how Mojo develops and whether it can hit the lofty goals Lattner is setting for it. “If you hear Chris talk about it, Mojo is slated to become a complete superset of Python, where whenever you write just Python code, it will execute in Mojo exactly the same way as it executes in Python, but much faster.” He is not yet sure whether Mojo can achieve that, but he emphasizes that the language is in its early stages and, as of July 2023, Mojo had not yet been made available for download.

Van Rossum thinks Mojo might prove more useful for experienced developers who already know how to write efficient code in C++ or Rust. “Someone who is a beginning Python user is not suddenly going to be able to write the type of Mojo code that executes much faster than it would in Python,” he says.

In May, Modular made Mojo accessible to some users in a Jupyter notebook, an interactive development environment allowing people to play with the code. The company said it expected to allow downloads in the fall of 2023 (it was released locally for Linux in September, and on MacOS in October), with full release perhaps in the summer of 2024.

Lattner says there may be pieces of Python that do not work in Mojo, but they will be insignificant. He says Mojo relates to Python in the same ways C++ relates to C, with additions such as classes and templates that turned C into a higher-level language. “There are programs you can write in C that do not work the same way or don’t even compile in C++, but they’re so minuscule that it doesn’t matter. The same thing is true in Mojo,” he says. “Our goal is to be as compatible as possible in all the cases that matter and make sure that we work with the existing ecosystem because we don’t want to break Python, we want to make Python better.”

Doug Meil, a software architect who has written about new programming languages, says Mojo is essentially Python++ for AI. “He’s trying very hard to support Python and meet people where they are, which is I think remarkably pragmatic,” Meil says. “They’re not coming up with an entirely new syntax, and it’s going to be way faster in scale across multiple hardware platforms. So that’s really cool.”

### Further Reading

**Mojo: A supercharged Python for AI;** <https://twimlai.com/podcast/twimlai/mojo-a-supercharged-python-for-ai/>.

*Don, E.*

*Getting started with the Mojo programming language for AI.* *LogRocket.* (2023); <https://blog.logrocket.com/getting-started-mojo-programming-language/>.

*Loy, J.*

*How to Build Your Own Neural Network from Scratch in Python.* *Towards Data Science.* (May 14, 2018); <https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-68998a08e4f6>.

*Meil, D.*

*Why are there so many programming languages?* *BLOG@CACM.* (July 5, 2022); <https://cacm.acm.org/blogs/blog-cacm/262424-why-are-there-so-many-programming-languages/fulltext>.

*Yegulalp, S.*

*A First Look at the Mojo Programming Language.* *InfoWorld.* (June 7, 2023); <https://www.infoworld.com/article/3697739/a-first-look-at-the-mojo-language.html>.

**Neil Savage** is a science and technology writer based in Lowell, MA, USA.

# OPEN FOR SUBMISSIONS

## ACM Transactions on Probabilistic Machine Learning (TOPML)

### Co-Editors-in-Chief

**Wray Buntine**

*VinUniversity, Vietnam*

**Fang Liu**

*University of Notre Dame, USA*

**Theodore Papamarkou**

*The University of Manchester, UK*



**Gold Open Access publication focusing on probabilistic methods that learn from data to improve performance on decision-making or prediction tasks under uncertainty**

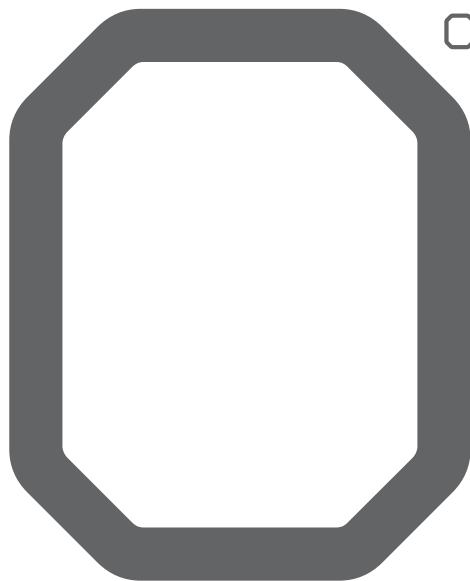
*ACM Transactions on Probabilistic Machine Learning (TOPML)* is a new Gold Open Access publication from ACM focusing on probabilistic methods that learn from data to improve performance on decision-making or prediction tasks under uncertainty. Optimization, decision-theoretic or information-theoretic methods are within the remit if they are underpinned by a probabilistic structure. Probabilistic methods may be harnessed to address questions related to statistical inference, uncertainty quantification, predictive calibration, data generation and sampling, causal inference, stability, and scalability. Examples of approaches relevant to the scope include Bayesian modelling and inference, variational inference, Gaussian processes, Monte Carlo sampling, Stein-based methods, and ensemble modelling. Examples of models for which probabilistic approaches are sought include neural networks, kernel-based models, graph-based models, reinforcement learning models, recommender systems, and statistical and stochastic models. Ethical considerations of probabilistic machine learning, such as data privacy and algorithmic fairness, should be addressed in papers where there is a direct ethical connection or context for the work being described.

The journal welcomes theoretical, methodological, and applied contributions. Purely theoretical contributions are of interest if they introduce novel methodology. Methodological and applied contributions are of interest provided that proposed probabilistic approaches are motivated and empirically corroborated by non-trivial examples or applications. Multidisciplinary approaches with a probabilistic structure are within the scope.

For more information and to submit your manuscript, please visit [topml.acm.org](http://topml.acm.org).



Association for Computing Machinery



DOI:10.1145/3616863

Subhabrata Dutta and Tanmoy Chakraborty

# Opinion Thus Spake ChatGPT

*On the reliability of AI-based chatbots  
for science communication.*

**I**T HAS BEEN three years since the arrival of GPT-3,<sup>2</sup> a neural-network-based large language model (LLM) that could generate polished, stylish text with ease. What followed was an intense race among tech giants like Google, OpenAI, and Meta to come up with larger and larger models, often trained using texts appearing anywhere on the Internet. While storming the AI research communities, all these were much outside the scope of usual public discourse. And then came OpenAI's ChatGPT<sup>3</sup>—an AI chatbot running on top of an LLM. The advancement was evident, shortcomings were subtle. But the most drastic change that ChatGPT brought forth was the sensation among the non-researchers (or researchers far from language technology or artificial intelligence) and ordinary citizens. With the online interface available to everyone on the Internet, it did not take much to become the talk of the town. ChatGPT can explain quantum computing to a non-expert, write poems in Keats' style, write and debug JavaScript codes, and many more. And very soon we started wondering: Can ChatGPT replace search engines such as Google, simplified knowledge

sources like Wikipedia, or community-powered Q&A platforms like StackExchange? To get an idea of its mammoth popularity, ChatGPT garnered more than 100 million users during just two months after launching; it took 2.5 years for Instagram and nine months for TikTok to achieve the same amount of users.<sup>4</sup> Beyond the shock and awe, it is important to identify the potential impact of the large-scale deployment of such AI-based chatbots on the production and consumption of scientific knowledge. Specifically, how would one envision the role of an AI chatbot being used by non-experts as a source to quench their thirst for scientific knowledge?

## Popular Majority vs. Authentic Minority

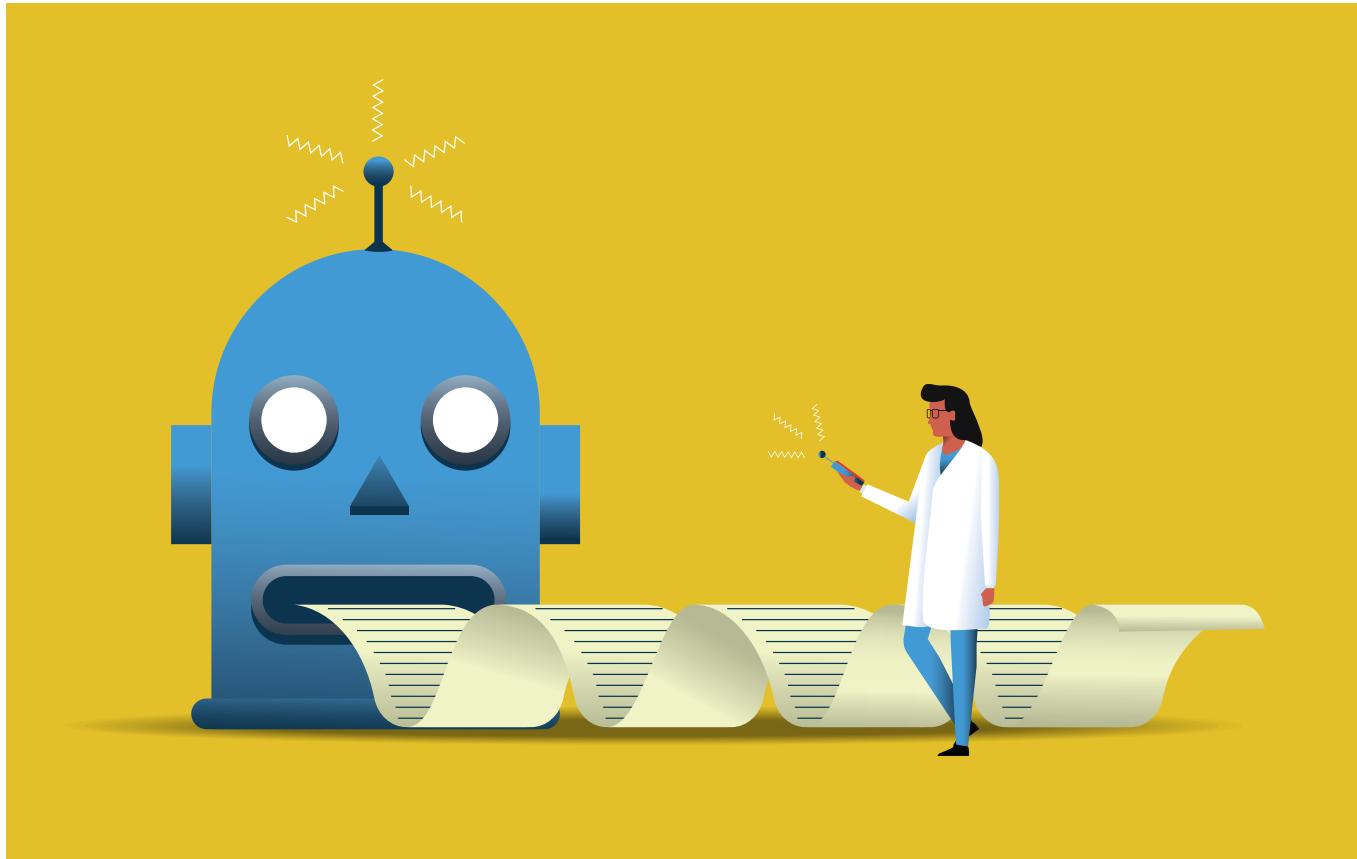
It is very much accepted in the LLM research that however proficient these models are in generating polished, confident textual responses to queries, they lack factual consistency.<sup>5</sup> Scaling to even larger models has significantly decreased these inconsistencies from the earlier models. However, we are not just talking about blunt factual errors here; ChatGPT is smarter in that sense. When the goal is to communicate hardcore

scientific facts to non-expert audiences, misconceptions can percolate with apparently synonymous phrasings—very often overlooked by even the experts. For example, consider the conversation with ChatGPT<sup>a</sup> per the following text:

**User:** Explain why spacetime is modeled as a Riemannian manifold.

**ChatGPT:** Spacetime is modeled as a Riemannian manifold because it allows for a mathematical representation of curved spacetime, as predicted by general relativity. The geometry of a Riemannian manifold can be described by a metric tensor, which encodes information about distances and angles between nearby points in the manifold. In the case of spacetime, the metric tensor is determined by the distribution of mass and energy in the universe and it determines how the objects move in response to gravity. The Riemannian manifold framework provides a powerful and flexible way to describe the curvature of spacetime, and it has been extremely successful in explaining a wide range of physical phenomena, from the motion of the planets to behavior of black holes.

<sup>a</sup> All these examples are from the 30/01/2023 free research version of ChatGPT.



The question is erroneous in itself from the beginning since spacetime is modeled as a Lorentzian manifold and not a Riemannian one, with the positive-definiteness condition relaxed.<sup>5</sup> Ignoring the erroneous definition, ChatGPT continues to explain why this is indeed true. Since locating what such LLMs actually know, it is even more difficult to correct such inconsistencies in a trustworthy manner. The condition becomes even more strange as the lack of self-consistency surfaces. When asked about the rate of expansion of the universe, ChatGPT declares the universe is not expanding at a speed faster than light.

**User:** Is it true that the universe is expanding faster than the speed of light?

**ChatGPT:** No, it's not true that the universe is expanding faster than the speed of light. According to our current understanding of physics, nothing can travel faster than the speed of light, and this applies to the expansion of the universe as well. However, the universe is indeed expanding, and recent observations suggest that the expansion rate is increasing. This has been attributed to the presence of dark energy, a mysterious form of energy that permeates all of space and drives the acceleration of the

universe's expansion. But even with this acceleration, the expansion of the universe is not occurring at speeds greater than the speed of light.

**User:** Is it true that light from some parts of the universe can never reach us because the universe is expanding at faster rate?

**ChatGPT:** Yes, it is true that light from some parts of the universe may never reach us because the universe is expanding at an accelerating rate. The observable universe is the part of the universe that we can see and is limited by the age of the universe and the speed of light.

## The two examples mentioned here primarily point toward the inconsistencies present in LLMs.

This means that even though light travels at a finite speed, there are parts of the universe that are moving away from us so fast that the light emitted from those regions will never be able to reach us. The concept of the cosmic horizon is a fundamental aspect of our understanding of the universe and the limits of what we can observe.

Yet, to explain the cosmic horizon, it imposes that light from parts beyond the cosmic horizon cannot reach us because the expansion rate is increasing faster than the speed of light. The arguments in both responses are quite interesting. For it answers that faster-than-light motion is impossible, the assertion is to the point. However, the reasoning is inconsistent as it cannot demarcate between the notions of scale and distance. This misconception is commonplace in popular science literature. It is not unlikely that a vast presence of such poorly conceived articles in the training data of ChatGPT pushed the model to bear such misconceptions. However, there is no way to investigate this other than OpenAI itself since the pretraining dataset is not open source. The two examples mentioned here primarily point to-



## Advertise with ACM!

Reach the innovators and thought leaders working at the cutting edge of computing and information technology through ACM's magazines, websites and newsletters.



Request a media kit with specifications and pricing:

**Ilia Rodriguez**  
+1 212-626-0686  
[acmmEDIASales@acm.org](mailto:acmmEDIASales@acm.org)



ward the inconsistencies present in LLMs that might have originated as a result of noisy training data—scientific misconceptions abundantly present on the Internet. We show another similar example here.

**User:** *What is a fractal?*

**ChatGPT:** *A fractal is a mathematical set with a repeating pattern visible at every scale. It is defined by a self-similar pattern where each part is a reduced-scale copy of the whole. This results in a highly complex and often visually striking structure. Fractals are commonly found in nature and are widely used in fields such as computer graphics, music, and finance.*

**User:** *A straight line is similar to itself at all length scales. Is a straight line fractal?*

**ChatGPT:** *No, a straight line is not considered a fractal. A fractal is defined by a repeating pattern that is visible at every scale, and while a straight line may be similar to itself at all length scales, it does not exhibit a repeating pattern. Fractals are characterized by intricate and detailed patterns that are created through the repetition of simple shapes or processes. A straight line, on the other hand, is a basic geometric shape with a simple and regular structure.*

**User:** *How to measure the presence of a detailed and intricate pattern then?*

**ChatGPT:** *The presence of a detailed and intricate pattern in a fractal can be measured using various mathematical and statistical methods. Some of the commonly used methods include:*

*1. Fractal dimension ...*

When we ask ChatGPT the definition of fractals, it quickly generates a paragraph that reverberates the notion of self-similarity and fractality juxtaposed together. As the interaction proceeds, the chatbot puts light on the ambiguous definitions it spat out earlier. We can see it finally shifting toward the requirement of a Hausdorff dimension of the object to be greater than its topological dimension.<sup>9</sup> Interestingly, Mandelbrot's paper<sup>9</sup> contains a rigorous definition of the fractal. Yet, the sheer number of popular definitions of fractals prevail in ChatGPT's answer. The underlying issue is quite straightforward to perceive—a language model learns from the majority of examples while rigorous scientific truth is mostly a minority. Such miscommunications

are bound to manifest from the text such models generate in reply to a scientific query unless it learns to demarcate between an authoritative text on a scientific phenomenon and popular literature trying to communicate scientific findings in simpler terms. Certainly, ChatGPT would need to learn science communication from the latter, but not at the expense of disregarding the facts and reasoning presented in the former.

### Multiple Sources, Multiple Versions

And then there are questions that cannot be answered in one specific way. Some of these questions might arise from a philosophical viewpoint of natural and social sciences, and there are multiple contending models that are well applicable. For example, there are multiple schools of thought on the interpretation of quantum uncertainty and what happens when we measure (or do not measure) a quantum particle. The Copenhagen Interpretation (CI)<sup>7</sup> assumes quantum states of systems evolve unitarily until measured; measurement breaks the unitary evolution, often described as the wavefunction collapse. On the other hand, Many-Worlds Interpretation (MWI)<sup>4</sup> assumes measurement only introduces multiple branches of reality; the quantum states remain unitarily evolving in all possible branches. Both CI and MWI produce the same, consistent predictions of quantum systems—none of these two can be assumed to be closer to the truth. And there are multiple other models of interpretation.<sup>10</sup> Yet, as we can see that ChatGPT assumes CI and does not mention MWI or other thoughts in relation to the following question:

**User:** *Do quantum states always evolve unitarily?*

**ChatGPT:** *In quantum mechanics, the evolution of a quantum system is described by unitary transformations ...*

This does not come as a surprise, given the fact that many assume the Copenhagen Interpretation to be the dominant one, and as we argued earlier, LLMs are indeed biased toward majority opinion simply due to statistical consequence. While this example is more of a philosophical argument and a certain choice is unlikely to alter the empirical understanding, it is difficult to believe that such majority biases

## The most troublesome attribute is models such as ChatGPT, in their current state, are not able to cite their sources.

are not present elsewhere. The most troublesome attribute is that models like ChatGPT, in their current state, are not able to cite their sources. With the vastness of human knowledge, it is impossible for an AI-based chatbot to list all possible interpretations, models, and schools of thought in one single answer. Without showing the sources, their knowledge distribution is essentially a one-step process. The user must remain content with whatever the chatbot produces. One may argue that no one is claiming that ChatGPT will be the only source of knowledge, and hence, why bother? Definitely, the Internet will be there. But so are the public libraries in the age of the Internet. Yet, most tend to access the Internet for its ease and speed. Given that AI-based chatbots are able to decrease the search effort even more, it would be shortsighted to reject the idea of a similar dominance.

### Conclusion

We must keep in mind that the examples shown here are cherry-picked and definitely not a wholesome representative of ChatGPT's capabilities. In fact, the degree of critics ChatGPT has received<sup>1</sup> is only signaling the capabilities and expectations that come with such an ambitious project. The arguments we presented are rather focused on better design principles of how an AI chatbot should interact with daily users. Definitely, a fatter column space in popular media demands human-like AI. Language fluency is probably the quickest path to mimic human-like capabilities. But beyond those shiny pebbles, one must ask the question, is a human-like AI the best aid to humans?

We argue that the focus of AI research must be concentrated on the latter, especially when we are talking about aiding people with easy access to scientific knowledge. In our hitherto non-automated past, seldom a single individual teacher has served the purpose of the only source of knowledge. It is very unlikely that any (even unintentional) push toward an automated persona as a single-point source of knowledge will do us any good. An interesting development in this direction worth mention is the latest Bing chat<sup>b</sup> from Microsoft: a search engine augmented with an automated chat assistant. Given this interface provides both the answers from the AI assistant as well as typical search engine results together for a query, it can provide a quantitative insight into the actual impact of AI assistants as knowledge sources on the average users. We hope the research community will make a sincere effort toward these investigations. □

b See <https://bit.ly/44L7RWJ>

### References

1. Borji, A. A categorical archive of ChatGPT failures. (2023). arXiv preprint arXiv:2302.03494.
2. Brown, T. et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33, (2020), 1877–1901.
3. ChatGPT: Optimizing language models for dialogue. (2023); <https://openai.com/blog/chatgpt/>
4. DeWitt, B. and Graham, N. *The Many-Worlds Interpretation of Quantum Mechanics* 61. Princeton University Press, (2015).
5. Einstein, A. et al. The foundation of the general theory of relativity. *Annalen Der Physik* 354, 7 (1916), 769.
6. Elazar, Y. et al. Measuring and improving consistency in pretrained language models. *Trans of the Association for Computational Linguistics* 9, (2021), 1012–1031.
7. Howard, D. Who invented the "Copenhagen Interpretation"? A study in mythology. *Philosophy of Science* 71, 5 (2004), 669–682.
8. Hu, K. ChatGPT sets record for fastest-growing user base—Analyst note; <https://bit.ly/3PIDXTL>
9. Mandelbrot, B. *The Fractal Geometry of Nature* 1. Freeman, New York, W.H (1982).
10. Schlosshauer, M. et al. A snapshot of foundational attitudes toward quantum mechanics. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics* 44, 3 (2013), 222–230.

**Subhabrata Dutta** (subha0009@gmail.com) is a postdoctoral researcher at the Indian Institute of Technology, Delhi, India.

**Tanmoy Chakraborty** (tanchak@iitd.ac.in) is an associate professor at the Indian Institute of Technology, Delhi, India.

Copyright held by owner(s)/author(s).



Watch the authors discuss this work in the exclusive *Communications* video.  
<https://cacm.acm.org/videos/thus-spake-chatgpt>

## ACM Distinguished Speakers

A great speaker can make the difference between a good event and a WOW event!

Take advantage of ACM's Distinguished Speakers Program to invite renowned thought leaders in academia, industry, and government to deliver compelling and insightful talks on the most important topics in computing and IT today. ACM will cover the cost of transportation for the speaker to travel to your event.

[speakers.acm.org](https://speakers.acm.org)



Association for Computing Machinery



## Opinion

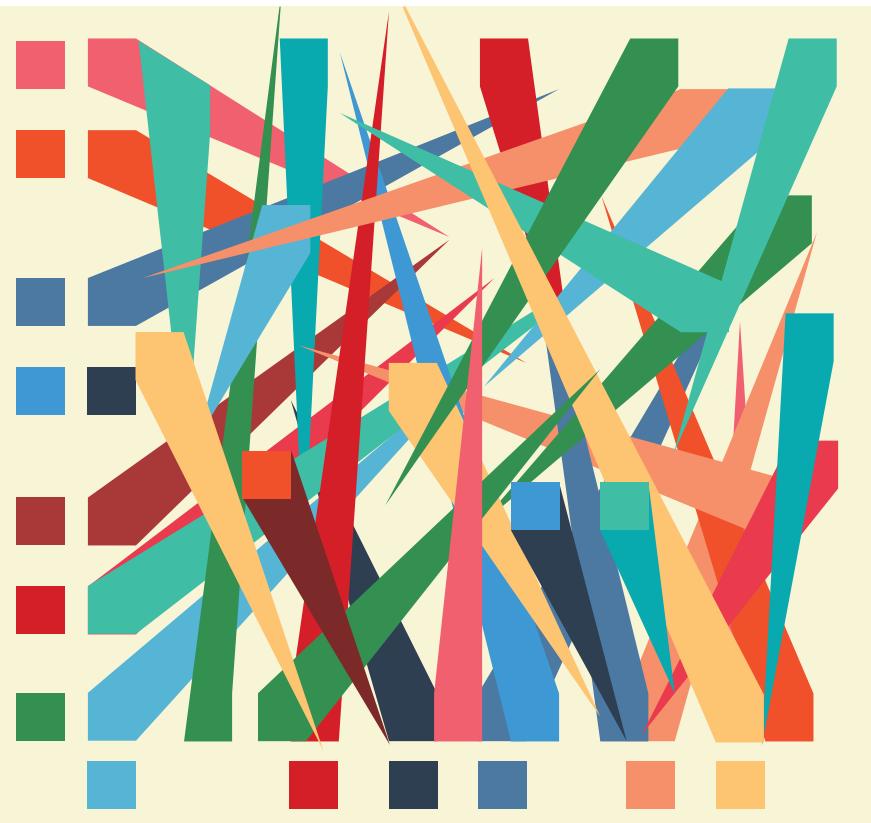
# Why Is the Current XAI Not Meeting the Expectations?

*Seeking better explanations for how algorithmic decisions are made.*

**I**MAGINE GOING TO space and deciding between Spaceship 1 and Spaceship 2. Although it has never been in flight, Spaceship 1 comes with precise equations outlining how it operates. Even though it is unknown how Spaceship 2 flies, it has undergone considerable testing and years of successful flights, including the one you are about to take. Cassie Kozyrkov, chief decision scientist at Google, posed this dilemma at the World Summit AI in 2018. We cannot provide a solution to this question because it is philosophical and perhaps generates a more profound inquiry on which better inspires trust—explanation or testing.

For a while, it appeared one issue with artificial intelligence (AI) algorithms, particularly cutting-edge deep learning techniques, was they were black boxes. It was impossible to pinpoint the precise reason the program predicted a particular outcome in a specific circumstance. Due to this lack of interpretability, businesses and governments were hesitant to use AI in critical sectors such as healthcare, banking, and government. So much so the EU Commission released its AI package in April 2021, including an AI act, recommending new laws and initiatives to make Europe a relevant hub for reliable AI, for example, in the case of the use of AI in high-risk sectors.

However, in the last few years, researchers, scientists, and businesses



have increasingly sought ways to provide some understanding of how AI algorithms get to a decision. Their respective AI programs are claimed to be “explainable” due to their post hoc interpretations of how an AI has made decisions.

One option is to try to provide explanations of black-box machine learning using interpretable models (for example, decision trees, rule sets, and analytical expressions).<sup>15</sup> Some have

implied a trade-off between accuracy and explanation, but that is not necessarily the case. Recent research has shown interpretable machine learning methods can be as accurate as black-box learning methods on tabular datasets<sup>4</sup> and can be especially attractive in situations where interpretability and transparency are important, such as in legal or medical contexts where decisions must be explained and justified (for example, Wang et al.<sup>16</sup>). On

the other hand, deep neural networks have shown remarkable performance in tasks involving image and natural language processing. More research is necessary to address the limitations of both interpretable models and deep neural networks.

Numerous papers argue for using Explainable artificial intelligence (XAI) methods in the literature, as well as multiple suggestions for brand-new XAI family approaches. Nevertheless, finding instances of practical XAI technique implementations that have enhanced the business in industry/societal/real-world applications is more challenging, even if some interesting work in this area has been put forward, for example in the health domain (see Lengerich et al.<sup>10</sup>). Certainly, explanations are not “one size fits all”; an explanation that is understandable to a technical audience might offer little explanatory value to a non-technical audience.

Case studies of the employment of XAI approaches tackling real-world machine learning issues are still lacking. Such case studies would help to clarify what is currently feasible and what is not feasible when employing XAI techniques. This is particularly true in the healthcare industry, where AI has significantly progressed in systems that interpret medical imaging automatically. The issue is that most of the XAI techniques the healthcare sector employs are limited.<sup>8</sup>

For example, a saliency map is a well-known explainability technique. It uses the image that was fed to the algorithm to produce a heat map of the areas that the AI program gave the highest weight when generating a forecast. However, as shown in one study published in 2021 by the medical journal *Lancet Digital Health*,<sup>6</sup> the heat map that was supposed to explain why the AI system classified the patient as having pneumonia covered a sizable portion of one lung's quadrant, with no further explanation of what precisely it was in that area that the AI system considered to be pneumonia. As mentioned in the study: “The clinician cannot know if the model appropriately established that the presence of an airspace opacity was important in the decision, if the shapes of the heart border or left pul-

## More research is necessary to address the limitations of both interpretable models and deep neural networks.

monary artery were the deciding factors, or if the model had relied on an inhuman feature, such as a particular pixel value or texture that might have more to do with the image acquisition process than the underlying disease.”

The authors note that in the absence of such information, people tend to presume AI focuses on whatever attribute they, as human therapists, would have thought was most crucial. Doctors may not be aware of the mistakes the machine learning system may make due to this cognitive bias. Ghassemi et al. also uncover issues with other well-liked explainability techniques, such as GradCam, LIME, and Shapley Values. Some of these techniques serve as a form of counterfactual by changing the data points that are entered until the algorithm generates a different forecast, at which time it is assumed that those data points must have been the most crucial for the initial prediction. These techniques, however,

## AI is more likely to serve the interests of the powerful if the aims of explainability from various communities are not clearly stated.

share the same drawback as saliency maps. While those techniques might be able to identify elements that the algorithm deems important, they are unable to explain to a doctor why those elements are significant from a medical standpoint, for example, whether the algorithm relied on significant structural alterations in the shape of an organ or the results were influenced by some noise in the training data.

Furthermore, a recent study<sup>9</sup> found various state-of-the-art explanation approaches regularly disagreed on the rationale for an algorithm’s conclusions. Most people who used the algorithms in real-world situations had no method of resolving such disparities and, therefore, might just choose the explanation that most closely matches their preexisting notions.

The ability to explain the causality behind an AI-based decision does not mean the AI system uses this causality and the actual relationship between inputs and outputs might differ. AI systems are often opaque, black-boxed systems.

### Why XAI’s Success Really Lies in Pushing Truly Cross-Disciplinarity Work

Governments consider the potential of XAI to address concerns about the obscurity of algorithmic decision-making with AI. Although XAI is enticing as a solution for automated decisions, using XAI is difficult because of the wickedness of the problems governments face. Wickedness refers to the ambiguity of the facts that characterize a problem and the lack of agreement on the normative standards for resolving it.<sup>2</sup> Additionally, these “solutions” to wicked problems frequently evolve over time.<sup>14</sup> It is difficult to explain something that is unclear especially if explainability is viewed as a strictly technical issue.

The need for explainability and interpretability in AI is a much larger cross-disciplinary problem that requires a more comprehensive solution than XAI alone can offer. Interpretable models for non-technical people suffer from the same problem: technologists are needed to translate technical explanations. Different kinds of explanations are needed for

# Coming Next Month in COMMUNICATIONS

## **10 Things Developers Should Learn about Learning**

### **Data Bias Management**

### **What to Do when Ideas of Fairness Conflict?**

### **Shortcut Learning of LLM in Natural Language Understanding**

### **Strategies, Tactics, and the Big Picture in Research**

### **Q&A with David Papworth**

### **On Specifying for Trustworthiness**

### **Non-Verbal Communication through Expressive Objects**

### **Acceptance of the Decentralized Web of Personal Data**

### **Hardware VM Isolation in the Cloud**

### **Creating the First Confidential GPUs**

### **Why Should I Trust Your Code?**

Plus, the latest news about epigenetics, wayfinding without GPS, and why lawyers fear AI.

business executives, risk managers, doctors, bankers, and officers, that is, end users.

AI is more likely to serve the interests of the powerful if the aims of explainability from various communities are not clearly stated. Companies adopting AI should be as honest as they can be about how, why, and for what purposes they are using XAI approaches. The National Institute of Standards and Technology (<https://www.nist.gov/>) is one of the organizations developing XAI standards and regulations. Such organizations should be aware of the current limitations of XAI in practice and seek out diverse expertise on better matching incentives and governance with a comprehensive understanding of XAI objectives. Ethics frameworks might come to the rescue, and the NGOs and business sectors have published a deluge of AI ethics and standards in the past few years. However, these values are isolated ones that serve corporate goals and are embedded in an industry that often disregards ethical behavior unless enforced by the justice system. In a recent paper, "The Uselessness of AI Ethics,"<sup>13</sup> Luke Munn gives evidence that such ethics principles are meaningless or isolated at most. The result is a disconnect between high ideals and technological reality. Even if this gap is acknowledged and ideas attempt to be "implemented," it is challenging to translate complicated social concepts into technical rule sets.

We can only achieve the objectives of intelligible, dependable, and controllable AI in practice with the active participation of many stakeholders from the social sciences, computer science, civil society, and industry. Consider, for instance, the disparate requirements of engineers and consumers in explaining an AI system. Developers may use Google's What-If tool to examine intricate dashboards that visualize a model's performance in many fictitious scenarios, evaluate the significance of various data elements, and assess multiple fairness theories. On the other hand, users could choose something more focused. It might be as easy as telling a user whatever factors, such as a late payment, resulted in a point deduc-

tion in a credit score system. Various consumers and scenarios will require different outputs.

Therefore, the explanatory domain must be improved, and its audience must be expanded if we want to reach the grail of trust and confidence in judgments made by black-box AI. In addition to XAI tools for technical teams, what we need is "Understandable AI"—AI that serves the needs of non-technical stakeholders. It is helpful to explicitly compare their goals to understand how practitioners in different domains have different expectations for what they hope to achieve by building XAI systems.

While the issue with XAI, as it stands right now, is that many of the current approaches view explainability as a purely technical matter, we think the future success of XAI lies in fostering a genuinely cross-disciplinary approach among AI and other fields of interest, such as ethics, law, psychology, sociology, and human-centered design to name but a few. In essence, statistical, mathematical, and scientific analyses are pretty valuable tools. However, it is all too easy to misinterpret their measured certainty as the only "true" method when, in fact, it is just one tool and one tactic—and not one that can be translated or used to explain all qualitative occurrences. We consider situations in which the effect is present and suggest a cause. Still, we overlook all the occasions in which the same cause resulted in no discernible consequence or an entirely different outcome. Model-based storytelling is quite simple. It is often difficult to remember they are stories, nevertheless.

There is evidence of winning the

**The future success of XAI lies in fostering a genuinely cross-disciplinary approach among AI and other fields of interest.**

rigid argument of more technical researchers versus more humanistic approaches thanks to the need for a better encompassing approach required for XAI to provide explanations for the practitioners or the general public. As Miller<sup>12</sup> points out “explanations are not just the presentation of associations and causes, they are contextual. While an event may have many causes, often the explainee cares only about a small subset (relevant to the context).” In addition, the relevant explanations should be given using the users’ language.

The aim is to meet the domain experts’ needs and expectations, for instance, by identifying interactive environments that allow people with different backgrounds to communicate and reason on recommendations made by AI. Human-centered design methods can offer tools to foster the design of cross-disciplinary collaborative systems. In particular, meta-design<sup>5</sup> can offer the participatory and cross-disciplinary approach needed to meet non-technical decision-makers expectations. It is a conceptual framework aiming to define and build the social, economic, and technical frameworks necessary for new kinds of collaborative design to function. It comprises several useful design-related tools to help users accomplish this task.

Through tools offered by meta-design, decision makers can render AI decisions understandable and valuable for their work. For instance, designing interactive systems that can manage the questions and associated answers that domain experts usually pose to AI systems, ranging from medical image interpretation to home automation or conversational agents.

The XAI Question Bank proposed by Vera Liao et al.<sup>11</sup> is an excellent example of presenting questions that can be used to probe the AI system. For instance, the type of data the system learns from and the related output, how accurate the system is, and how it makes a prediction. There might be questions related to a specific prediction: what would be predicted if an instance changes to a different value, how to change an instance to get another prediction, or what is the scope

## The XAI field will require contributions from different areas to explore more types of explanations.

of change permitted to get precisely the same prediction? Only then would it be possible to create environments that allow the domain experts to configure how decisions are translated in the user language, not that of the AI system, and the desired level of interactivity and multimodality.

In “Why Should Humans Trust AI?”,<sup>3</sup> John Carroll proposes to model explanations as inherently pragmatic, conversational, and social. It is always a question of making sense, being aware, and negotiating in a vast sense, as well as of the responsibility people accept and show for one another as they engage in daily encounters. This might be a significant accomplishment for XAI, but we will not get anywhere if we do not acknowledge the problem and the obstacles in our way.

### Conclusion

XAI per se shows many limitations in its current form due to an excessively technical approach, often requiring technologists to help end users fathom the explanation provided by a model. Moreover, as demonstrated by Bordt et al.,<sup>1</sup> explanations produced by existing AI approaches depend on a variety of particular characteristics of the AI system, such as the training data, the precise shape of the decision surface, and the selection of one explanatory algorithm over another. Given that programmers and AI developers are free to select these factors, there is a risk that, despite an explanation’s seeming plausible, it occurs only because of some hidden layers computing weights depending on some of the features and does not truly reflect the explanation given at all. The creation of explanations that cast doubt on certain features of AI systems is not something AI develop-

ers are interested in doing, at least not right now.

Furthermore, the explanatory approaches mentioned here are just a small set of tactics that can be employed to provide explanations. Indeed, the XAI field will require contributions from different areas to explore more types of explanations.

Nonetheless, we believe the winning point of XAI involves the intersection of different disciplines including ethics, law, psychology, sociology, and human-centered design. XAI could serve as a starting point to develop a common language among scholars from different disciplines, thereby accelerating its real-world impact. ■

### References

1. Bordt, S. et al. Post-hoc explanations fail to achieve their purpose in adversarial contexts. (2022); arXiv:2201.10295.
2. de Brujin, H. et al. The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making. *Government Information Quarterly* 39, 2 (2022), 101666.
3. Carroll, J.M. Why should humans trust AI? *Interactions* 29, 4 (2022), 73–77.
4. Caruana, R. and Nori, H. Why data scientists prefer glassbox machine learning: Algorithms, differential privacy, editing and bias mitigation. In *Proceedings of SIGKDD*. (2022). 10.1145/3534678.3542627.
5. Fischer, G. et al. Meta-design: a manifesto for end-user development. *Commun. ACM* 47, 9 (Sept. 2004), 33–37.
6. Ghassemi, M. et al. The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health* 3, 11 (2021), e745–e750.
7. Habayeb, A. Explainable AI isn’t enough; We need understandable AI. (Feb. 16, 2022); <https://bit.ly/3RvR0gW>.
8. Kahn, J. What’s wrong with “explainable A.I.” (Mar. 22, 2022); <https://bit.ly/45StDcp>.
9. Krishna, S. et al. The disagreement problem in explainable machine learning: A practitioner’s perspective. (2022); arXiv preprint arXiv:2202.01602.
10. Lengerich, B.J. et al. Automated interpretable discovery of heterogeneous treatment effectiveness: A COVID-19 case study. *Journal of Biomedical Informatics* 130, (June 2022), 104086.
11. Liao, Q.V. et al. Questioning the AI: informing design practices for explainable AI user experiences. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. (Apr. 2020), 1–15.
12. Miller, T. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* 267, (2019), 1–38.
13. Munn, L. The uselessness of AI ethics. *AI and Ethics*. (2022), 1–9.
14. Rittel, H.W. and Webber, M.M. Dilemmas in a general theory of planning. *Policy Sciences* 4, 2 (Feb. 1973), 155–169.
15. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* 1, 5 (2019), 206–215.
16. Wang, C. et al. In *Pursuit of Interpretable, Fair and Accurate Machine Learning for Criminal Recidivism Prediction*. (Mar. 2022); <https://arxiv.org/abs/2005.04176>.

**Alessio Malizia** (alessio.malizia@unipi.it) is an associate professor in the computer science department at the University of Pisa, Italy.

**Fabio Paterno** (fabio.paterno@isti.cnr.it) is a research director at C.N.R.-ISTI in Pisa, Italy.

Copyright held by owner(s)/author(s).

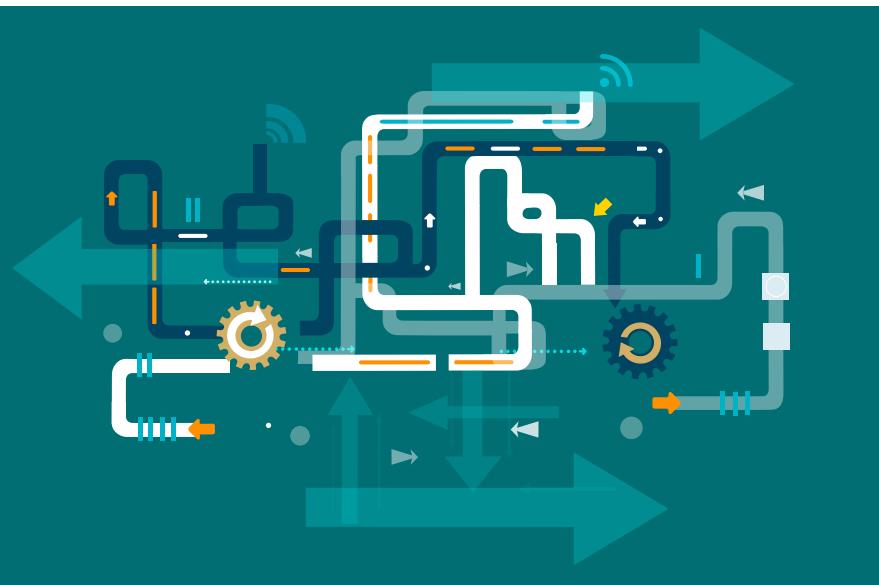


# The Profession of IT A Clash of Civilizations

*The much-sought holy grail of more and faster innovation will come from integrating pipeline thinking and adoption thinking.*

**W**E ARE LIVING at a unique moment in history where long-held views of how the world works and what it takes to be successful at innovation are falling short. Once upon a time we “knew” that research is the first stage in a pipeline that leads to marketable products. We “knew” that if we invested sufficiently in basic research that we would generate applied research and attain a competitive lead. We “knew” that we could maximize those investments if we applied sufficient controls on each stage of the pipeline. We “knew” this organized pipeline is historically sound because we saw how successful modern factories had become at turning out high-quality, low-cost products at scale. We “knew” these truths applied for businesses, organizations, and governments.

The pipeline idea traces back to the development of 19<sup>th</sup> century factories when innovation meant transforming inventions into products for the masses. We “knew” pipelines would work more generally than in factories because every change that emerged in the world seemed to unfold in distinct sequential stages. We saw pipelines all around us. We “knew” the education pipeline from K-12 through college to postgraduate school produced the educated workforce both industry and government need to compete successfully. We “knew” the military pipeline transformed raw recruits into military professionals and allowed the success of the all-volunteer forces. We



“knew” the political pipeline transformed local leaders into state and national leaders. We “knew” the career pipeline transformed fresh graduates into professionals who could support a company’s offers. We “knew” that, in all these domains, the innovation pipeline would systematically stage by stage transform raw ideas into useful outcomes. Over the years, the pipeline idea has been codified into business theory and government policy. It has survived so long because pipeline processes delivered good and often repeatable results that sustained our confidence in what “we knew.”

Today, we are starting to question what we knew. Many people sense with a vague unease that pipeline thinking is no longer sufficient. The education pipeline is not delivering

the trained personnel to meet the demands of industry, the military, or civil service. The military pipeline no longer meets military targets for recruitment, retention, or talent development. The career pipeline no longer supports company goals for internal professional development and stable market position, as people jump from company to company in search of new experiences, higher salaries, and higher positions. The innovation pipeline seems too slow, expensive, and rigid for the modern world. The U.S. Department of Defense expects to allocate 17% of its 2024 budget for research and development, and yet in the new geopolitical climate, many political and military leaders have lost faith that the investment in the research pipeline will generate innova-

tions at a sufficient pace and scale to remain competitive. Many are looking for new approaches to generating innovation faster and more reliably.

### The Eureka Story

One approach to resolving this tension has been to augment the pipeline with the Eureka story. The Greek word Eureka means “I have found it.” The Eureka story idolizes the inventors who made the significant discoveries and seeded the innovation pipeline. From this perspective, to get more innovation we need to stimulate and reward more inventive creativity. There are multiple problems with this story. It ignores the work of the many others who take the first risk in adopting the practices of the invention or who work in the community to sustain adoption. It ignores the large number of innovations of practice that have no identifiable inventor. It ignores the waste of effort and subsequent disillusionment among the many whose ideas are filtered out.

The Eureka story has inspired new practices for generating ideas, such as contests or hackathons, which harness the crowd rather than individual inventors. Leaders in all fields often report they have more ideas to consider than resources to pursue all of them. Merely adding more ideas at the beginnings of the pipelines does not generate more or higher quality output. The Eureka story adds complexity to the pipeline model but does not increase the overall rate of adoption.

### The Valley of Death

The Valley of Death is another augmentation of the pipeline story. It emerged to explain the failures of the pipeline model. According to this story, when an artifact exits the pipeline it does not readily transition into adoption because of bureaucratic, social, management, or cultural barriers. The barriers are so tall and strong that getting past them all is next to impossible. The Valley of Death now appears in pipeline diagrams as a new stage that must be managed but is poorly understood and does not improve adoption rates. While this story gives some insight into why transition from invention to adoption did not happen, it gives no insight into how to make that transition more successful.

## Many are looking for new approaches to generating innovation faster and more reliably.

### Adoption Thinking

So, if pipeline thinking and its augmentations are not reliable paths to adoption, what is? The answer is right under our nose: Foster the conversations in which adoption happens.

Consider an example. Generative AI in the form of ChatGPT was offered for public access at the end of November 2022. It touched off enormous waves of interest. Within two weeks, 100 million users were trying out the new technology. Within two months, there were a dozen how-to-use-GPT books available from Amazon. Tech companies quickly announced plans to incorporate the technology into their browsers, office packages, online services, and other apps. Hundreds of proposals for startups appeared and venture capitalists found funding to back them. The speed of adoption of generative AI has been astonishing. All this happened in conversations. The conversations spread rapidly through the Internet and media. Many people started experimenting to see whether generative AI would be a good advisor, artist, or writer, and they shared their findings widely in social media and as preprints of research papers. Moods of awe and surprise and intense enthusiasm blossomed. They were later joined by moods of distrust and fear that the generative AI technology might get out of control and cause great damage to societies. The ideas and moods spread rapidly round the world.

Innovation leaders were everywhere. OpenAI was first, offering its new technology for public experimentation. Other tech companies quickly followed with their own versions. Many experimenters and tinkerers generated greater understanding. The AI experts who started having doubts about the technology led a conversation about the role of government regulation to manage the dangers.

All this has happened in the space of swirling conversations among many people. The collective conversations took on a life and momentum of their own, beyond the control of any one person. Two broad themes emerged. One focuses on how we might harness the power of different forms of AI to address concerns in our community, such as democratizing knowledge work and relieving humans of dull, dirty, and dangerous work. The other focuses on breakdowns that have appeared around trustworthiness, unpredictable “hallucinations,” worker displacement, and intellectual property. Calls to subject generative AI to “the discipline of the pipeline” are intended to slow the pace of adoption and maybe put the genie back in the bottle.

This is an example of an adoption conversation. Adoption is unruly and chaotic. The process empowers people to address longstanding concerns as well as iterate new practices as AI is adopted. It shows that *innovation is emergence of new practices in a community*. Adoption of new practice happens as people reach agreements and make commitments in their many conversations. New practices are emerging from the collective work of leaders at every level as they use the technology and then make new offers enabled by the technology.

This interpretation of innovation has been around for a long time. It appears in the stories of leaders nurturing their communities to commit to new practices. Leaders do this by fostering conversations that expose concerns and offering new practices to take care of them. The common stories of innovation leaders emphasize qualities such as charisma, creativity, genius, brilliance, connections, extroversion, and good luck. However, there are numerous examples of successful innovation leaders who lacked these qualities. These qualities are not needed to be successful.

Successful innovation leaders rely on a skill set—competence in the conversations that lead to adoption. The skill set includes eight essential practices—sensing, envisioning, offering, adopting, sustaining, executing, embodying, and mobilizing.<sup>a</sup> They are essential be-

<sup>a</sup> P. Denning and R. Dunham. *The Innovator's Way*. MIT Press, 2010.

cause omitting any one of them is likely to block the innovation. The good news is that anyone can learn them.

Adoption always requires leadership, but not the charismatic type. In most innovations, we see many people exercising leadership skills. The leadership is often distributed by necessity because most innovations involve many members of the community. Leadership is present and essential, but distributed among the many, not concentrated in the few.

### The Clash

These two interpretations of innovation are established civilizations whose inhabitants have lived in them for a long time and have identified their professional offers with them. Reconciling them is a challenge. Bob Metcalfe captured the challenge with his famous saying, "Invention is a flower, innovation is a weed." The work of invention and the work of adoption seem to be, paradoxically, simultaneously complementary and contradictory. Pipeline thinking aims to control the process of moving inventions into marketable products, with its focus on professional management, order, and cost containment in the production of new artifacts and technologies. But by itself production does not produce adoption. Adoption thinking aims to bring order into the conversations of unruly, unpredictable, and often chaotic human communities as they strive to collectively embody a new practice. But by itself, adoption does not produce the artifacts that become the tools enabling the new practice. Production is a linear-ordered process; adoption is a chaotic process.

The clash has become more apparent in recent years as the concern in business and government for more and faster adoption has increased. Pipeline advocates argue that the discipline and control of the pipeline should be extended deeper into organizations, touching every project. In their view, failures arise from too many unpredictable and chaotic things happening in the raw, undisciplined social spaces of communities. Adoption advocates argue that success is often compromised by trying to impose a linear order on the nonlinear process of emergence because it limits the very conversations necessary for adoption. In their view,

## Success at both adoption and production depends on personnel and managers being competent in these essential conversation skills.

adoption cannot be achieved by imitating a production process.

Whereas the pipeline focuses on control, orderliness, and predictability at producing technology, adoption focuses on the embodiment of the new practice in the using community. How do we choose between these apparently conflicting worlds?

### Resolving the Clash

We regard this as a false choice. The logic of the pipeline is, as it was in the beginning, to organize the production of artifacts. It is not useful to organize social and human processes. The logic of adoption conversations is to bring about community agreement on new practices. They do not implement the production lines that supply tools to support new practices. We need both logics.

The reconciliation comes when you realize that both production and adoption happen in social space. Social space is a churning sea of conversations in a community that maintains relationships, commitments, coordination, customs, and stability. Innovation is a change of practice in this space.

Many innovations depend on production of new technology at a scale sufficient to provide everyone in the community with the tools they need to carry out the new practice. The conversations that drive production of tools to support a new practice are not the same as those to generate adoption of a new practice. As mentioned earlier, the conversations that drive adoption are sensing, envisioning, offering, adopting, sustaining, executing, embody-

ing, and mobilizing. The conversations that drive production are:

- Conversations for action: loops whereby two parties bring about a condition of satisfaction. Every stage in a pipeline depends on these conversations for its execution. If the personnel working on that stage do not complete their action loops, their stage will be a source of bottlenecks and complaints that interfere with overall production.

- Conversations for possibilities are needed whenever contingencies arise in production.

- Conversations for context and relationship are needed so that workers see the importance and relevance of their work and enjoy camaraderie with their fellows.

Success at both adoption and production depends on personnel and managers being competent in these essential conversation skills. Unfortunately, education in these skills is uncommon in technical fields, which favor "hard" over "soft" skills. Many find that the soft skills needed to get their work adopted are harder than the hard skills of their work. The good news is that learning the soft skills of leading adoption is not that hard once you know what they are.

The skills of managing conversations in social space dissolve the Valley of Death and make it navigable.

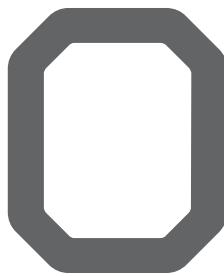
Notice the potential harmony: Results of production are brought into adoption by the adoption conversations; new production lines are brought into operation by previous adoption conversations. The pipeline model is embedded in social space and social space is embedded in pipelines. But they are not integrated. That is hurting our ability to achieve the much-sought end result: adoption. □

**Peter J. Denning** (pjdenning@nps.edu) is Distinguished Professor of Computer Science at the Naval Postgraduate School in Monterey, CA, USA, is Editor of ACM *Ubiquity*, and is a past president of ACM. His most recent book is *Computational Thinking* (with Matti Tedre, MIT Press, 2019). The author's views expressed here are not necessarily those of his employer or the U.S. federal government.

**Todd W. Lyons** (tlyons@npsfoundation.org) is Vice President of the Naval Postgraduate School Foundation and volunteer faculty of Innovation Leadership at the Naval Postgraduate School in Monterey CA, USA.

We are grateful to Nick Dew, Dan Finkenstadt, Garth Jensen, and Dave Newborn for insights as we were preparing this column.

Copyright held by owner(s)/author(s).



DOI:10.1145/3625551

George V. Neville-Neil

 Article development led by ACM Queue  
queue.acm.org

# Kode Vicious

## Halfway Around the World

*Learn the language, meet the people, eat the food.*

**Dear KV,**

My company has asked me to spend six months helping to start a software development group in a place halfway around the world from where we are based now. I did not get this assignment because I have seniority with the company; quite the opposite—they seem to be sending me because I am one of the youngest people and one of the few developers without a partner or children.

I am told the developers I will be working with are fluent in English, which I would expect in technology, but outside of work, I will need either an interpreter or lessons to learn the local language. I was never good at languages in school, except for the ones I programmed in, and I worry I have no ear for them. I figure since it is a relatively short period of time, I can just work a lot, talk to friends at home over the Internet, avoid the whole language-learning problem, and then return home at the end of six months (with a promotion and better pay).

Have you ever had to manage a group where the developers spoke another language, or, perhaps, done a stint outside the U.S.?

**Mono Lingual**

**Dear Mono,**

I have been lucky enough to work in several places outside my home country, and each time I have found it to be at equal turns challenging, frustrating, and rewarding. Let me first disabuse you of the idea that the language class-



es you took in middle or high school had anything to do with language learning. Human beings do not learn languages by studying a grammar book or reciting the conjugations of verbs; they learn languages by listening and speaking. Have you ever seen a two-year-old with a grammar book? Forget everything you thought you learned about whatever language you took in school and start fresh when you get to your assignment.

Human languages also have a built-in reward system: Say the right thing and you get breakfast; say it incorrectly, and you may be eating a dog's

breakfast. Yes, you should learn the language of the locale you are planning to live in, and here are a few good reasons why. The first is that you will not always be at work; something might happen—a car might break down, or the metro; or there could be a fire in your apartment building—there are always going to be situations where handling the local language is important. Emergencies aside, you will have a much more interesting experience if you can at least get by in a language, rather than constantly pointing at picture menus. Finally, it has been shown that learning languages increases

Association for  
Computing Machinery2021 JOURNAL IMPACT  
FACTOR 14.324

## ACM Computing Surveys (CSUR)

*ACM Computing Surveys (CSUR)* publishes comprehensive, readable tutorials and survey papers that give guided tours through the literature and explain topics to those who seek to learn the basics of areas outside their specialties. These carefully planned and presented introductions are also an excellent way for professionals to develop perspectives on, and identify trends in, complex technologies.



For further information and to submit your manuscript, visit [csur.acm.org](http://csur.acm.org)

brain function, which, as scientists and engineers, must surely be something we all want.

Apart from the differences in language, every place—not just different countries, but also cities, provinces, and all the other possible delineations on this planet—has a different culture, both at and outside of work. You cannot insulate yourself from the culture by becoming a robot that only goes from home to work and then back. Even if you hid in your house, eating expensive imported food and consuming entertainment from home over the Internet, you would still have to go out and work with this group that you are supposed to be managing. How can you hope to manage people whose lives and motivations you have made no effort to understand? Don't answer, because the answer is easy: You can't.

If you think you find different ways of working when you switch companies, that is nothing compared with the differences you will find every day when you arrive at the office in a different country. What time do people arrive? Do they stay late? How late is late? How do they see their work? How do they solve problems? In software and computer science, we like to think our common language is that of algorithms and math, and while that is, in part, true, it is nowhere near the totality of what it means to be working together.

KV has always made the point that software is a collaborative, human endeavor, and that the things we do to make computers (mostly) do what we say are actually only a small part of our work. Most of the work of software is building systems we can explain to each other, and even if everyone on the team, or in the room, or whatever, is speaking the same mathematical and even human language, the culture of the people involved comes into play time and again.

A simple case in point: There are places where certain features can or cannot be built into software, not only due to varying forms of regulation, but also because of cultural expectations about what is and what is not acceptable practice. These vary from company to company, as I pointed out, but they vary far more widely

## How can you hope to manage people whose lives and motivations you have made no effort to understand?

among cultures. Not only do different cultures treat different features differently, but they also treat each other differently. How people act with respect to each other is a topic that can, and does, fill volumes of books that, as nerds, we probably have never read, but finding out a bit about where you are heading is a good idea. You can try to ask the locals, although people generally are so enmeshed in their own cultures they have a difficult time explaining them to others. It is best to observe with an open mind, watch how your new team reacts to each other and to you, and then ask simple questions when you see something you do not understand.

KV's advice—whenever asked about a gig in a different country—is always the same: Learn the language, meet the people, eat the food, and only be rude on purpose; never give offense by mistake, which is why you need to learn at least a bit of the local language. People don't speak Java, for which we should all be eternally thankful.

**KV**

### Related articles on [queue.acm.org](https://queue.acm.org)

#### [Kabin Fever](https://queue.acm.org/detail.cfm?id=3437846)

<https://queue.acm.org/detail.cfm?id=3437846>

#### [Advice to a Newbie](https://queue.acm.org/detail.cfm?id=1242495)

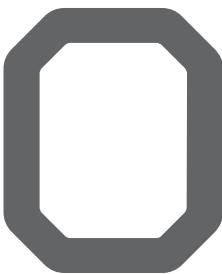
<https://queue.acm.org/detail.cfm?id=1242495>

#### [Ground Control to Architect Tom ...](https://queue.acm.org/detail.cfm?id=1317396)

<https://queue.acm.org/detail.cfm?id=1317396>

George V. Neville-Neil ([kv@acm.org](mailto:kv@acm.org)) is the proprietor of Neville-Neil Consulting, Brooklyn, NY, USA, and co-chair of the ACM Queue editorial board. He works on networking and operating systems code for fun and profit, teaches courses on various programming-related subjects, and encourages your comments, quips, and code snips pertaining to his *Communications* column.

Copyright held by owner(s)/author(s).



# Opinion

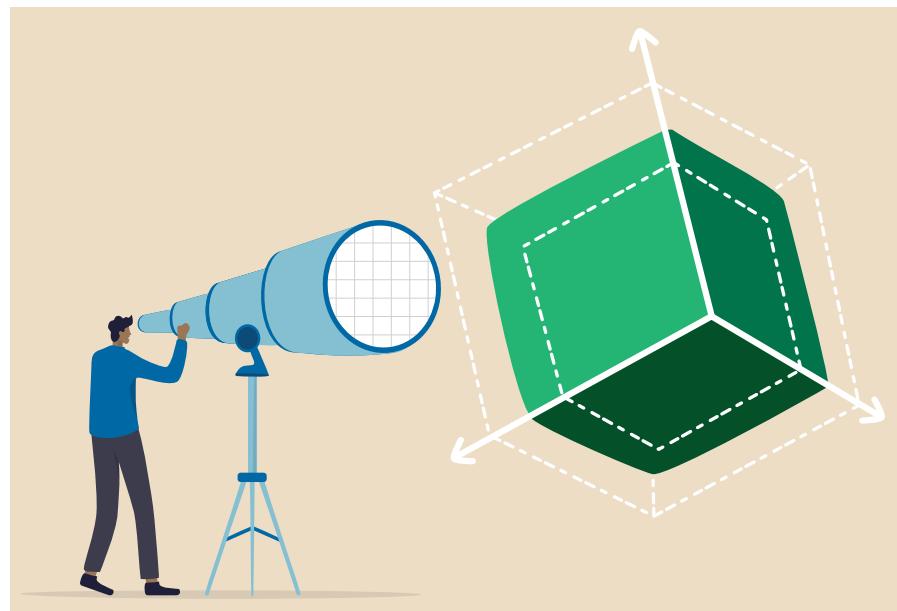
## A Strategy for Scaling Advanced Analytics

*Key elements for scaling advanced analytics.*

**H**OW TO GAIN competitive advantages over competitors should be a crucial question for any organization. Organizations that have scaled advanced analytics (predictive analytics, prescriptive analytics) beyond pilot projects and characterize themselves as data-driven tend to be top performers.<sup>13</sup> Furthermore, there is no shortage of examples of how advanced analytics can be applied to a given domain.

In 2022, NewVantagePartners released its 10<sup>th</sup> annual report on data and AI leadership among Fortune 1000 organizations.<sup>14</sup> Although progress has been made over the years to become more data-driven, there are three major concerns in the 2022 report: only 19.3% had established a data-driven culture; only 26.5% had established a data-driven organization; and the nontechnical barriers represented 91.9% of the perceived challenges to become more data-driven. In addition, the metrics are declining. One can only speculate whether more respondents in the survey or a more critical self-assessment of the organizations are the reasons for the low figures or if the explanation can be found elsewhere. In practice, it means organizations cannot use advanced analytics to its full potential.

Barriers to becoming a data-driven organization are well known and include:<sup>3</sup> lack of understanding, lack of skills, insufficient organizational alignment, lack of management support, and lack of vision and supporting strategies.



Our experience indicates organizations take shortcuts when scaling advanced analytics and have a too-strong technical focus with limited use of an overall change management process. Similarly, Davenport and Bean<sup>8</sup> suggested "... firms need more concerted programs to achieve data-related cultural change," and must include more experts in organizational change and change management to overcome the barriers.<sup>14</sup>

The question is whether scaling advanced analytics is any different from scaling the usage of a new tool, business process, or workflow, for example, a new model for inpatient care flow<sup>10</sup> or a new sales tracking system.<sup>1</sup> The shared findings in Ashkenas and Mat-

ta<sup>1</sup> and Golden-Biddle<sup>10</sup> indicate scaling is possible as long as one does not assume everything from the initial testing of the scaled object (inpatient care flow, sales tracking system) can be replicated. Instead, the scaled object must be adapted to the new environment together with user training. Now assume initial testing of advanced analytics is done for a specific business process in sales. Following the previous scaling recommendations, the specific business process should be scaled and adapted to new environments, and users should receive training. If an organization intends to become data-driven and establish a data-driven culture in teams, the focus of the scaling

should be *the usage of advanced analytics within the organization*, and not how advanced analytics is applied to a specific workflow or business process. Hence, previous recommendations for scaling must be adapted when scaling advanced analytics.

Given the reported problems with scaling advanced analytics, requests for more concerted programs, and increased focus on change management, what are the key elements for a proper scaling of advanced analytics?

**How to start.** Many organizations have started their journey to become more data-driven with pilot projects in advanced analytics.<sup>3</sup> Our experience indicates these pilot projects tend to be fairly straightforward. One reason is they are focused on applying technical aspects to a given dataset, for example, can this data-mining algorithm give us any new insights on our sales data? In case only decision makers and data scientists are involved, the pilot project is too heavily focused on technical aspects since employees working with sales and sales data are bypassed.

Although an initial project on advanced analytics shares many barriers and practices from any type of project, there are at least three barriers that need extra attention:<sup>2</sup> data quality; business resistance; and business overconfidence. We have encountered several initial projects in advanced analytics that have been delayed due to poor data quality or where the data has been manipulated (to fit the pi-

lot project). Several initial projects in advanced analytics intend to mimic how a middle manager makes decisions. We have seen an increased level of business resistance among middle managers that do not want to share their knowledge with a data scientist. In contrast, we have also seen evidence of business overconfidence where senior managers make statements about algorithm accuracy that is off the scale.

Regardless if organizations start with a pilot project or apply other mechanisms, for example, work discovery,<sup>10</sup> to collect initial experience of using advanced analytics, there are some underpinning best common practices:<sup>2,6,7</sup> align the initial work to existing business strategies or important goals within the organization; use a cross-functional team; and document the initial work as a preparation for making a future decision if scaling should be done or not. Hence, any initial work in advanced analytics that is carried out as a big-bang project with limited or no involvement of business users is likely to act as a catalyst for business resistance.

**Scaling decision.** Once the initial experience of using advanced analytics has been collected, a decision about scaling is the next step. A set of framing questions can be used to make an informed decision—these framing questions are adapted from Cooley and Linn<sup>7</sup>:

- To what extent are the benefits of using advanced analytics observable in practice?

- To what extent is using advanced analytics aligned with existing business strategies?

- To what extent have potential business users been involved?

- What type of enablers need to be changed, for example, policy about access to data?

- What are the known technical barriers and nontechnical barriers?

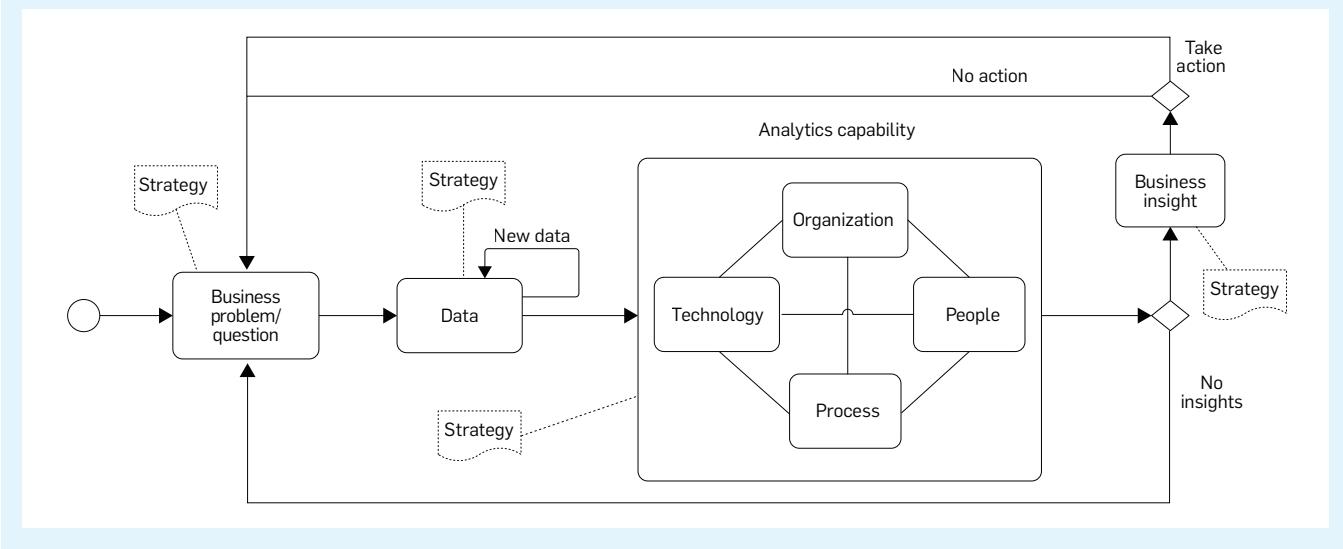
- What is the business case and vision of the scaling?

- To what extent is the vision a large departure from current practices?

**Scaling plan.** Following best practices, an overall scaling plan should be developed, including scaling vision, scaling process, change management process, and supporting strategies. Our experience is that few organizations have an overall vision defined to become data-driven. Instead, data-driven is used as a tool to support more business-oriented visions. The scaling process should describe who is involved in the scaling process and what type of decisions and analytics are scaled first. Similarly to the initial work, a cross-functional group should be involved in the scaling process. Finally, a suitable change management process should be selected and adapted.

**Strategies.** In order to avoid the pitfall of mainly developing technical strategies for scaling advanced analytics, developed strategies should align with what it means to be a data-driven organization, for example, by using a sample framework for analytics.

#### A sample framework for analytics and associated strategies.



By following the framework in the figure here, strategies developed include:

- Business problem/question: How can the organization encourage a culture where business problems and questions that can be investigated by advanced analytics are spontaneously generated?

- Data: Existing data strategies need to be revised. Organizations that start with advanced analytics have reported that they are back to square one with data quality since old data is used in new ways or new data is collected.<sup>3</sup> Furthermore, bigger organizations tend to have a large number of different enterprise information systems, data warehouses, and so forth, which increases the risk data is kept in silos and has definitions that are not synchronized. For these organizations, it is essential to apply master data management before applying advanced analytics spanning the data sources. Scaling analytics can expose individuals or groups unwilling to share organizational data. According to Hart and Warne,<sup>12</sup> the unwillingness to share data can be seen from a culture- or politics-based perspective. From a culture-based perspective, barriers can be overcome by encouraging or educating people who are unwilling to share data, whereas from a politics-based perspective, barriers can be overcome by making data sharing easier.<sup>12</sup>

- Analytics capability: An organization's analytics capabilities can be divided into:

- Technology: What type of platform and tools should be used?

- Organization: How should analytical competence be organized and integrated, for example, centralized or Center of Excellence?<sup>11</sup> Who is responsible for in-house training, for example, in data literacy and tools? What type of roles need to be established, for example, chief data officer, data translators, data product manager.

- Process: How should advanced analytics be integrated and used in the decision process? What type of decision processes can be automated?

- People: How should people work with advanced analytics, and how should skills be raised? What incentives are there for using advanced analytics? Analytics career paths?

- Business insight: How are business

## The distinction between AI and advanced analytics is sometimes blurred as they can rely on the same tools and techniques.

insights documented, distributed to others, and acted upon? What type of findings and data can be shared externally?

Although most of the strategies seem obvious, organizations tend to focus on developing mostly technical strategies or are still developing them.<sup>3</sup> For example, 53% of the respondents in the 2022 survey of Fortune 1000 organizations had developed a corporate data strategy.<sup>13</sup> To avoid working in silos, developed strategies for scaling advanced analytics need to consider and adapt to existing projects and strategies in, for example, AI, master data, or cloud computing.

**Monitor and evaluate.** Progress toward the vision at the organizational level is easily detected, for example, new departments, new roles, new policies, or new staff. Assessing what happens in practice at the team level requires a deeper analysis, preferably using analytics to assess what characterizes teams that have established a data-driven culture.

The distinction between AI and advanced analytics is sometimes blurred as they can rely on the same tools and techniques. According to Cam and Chui<sup>5</sup> and Davenport and Malone<sup>9</sup> most AI and advanced analytics models, for example, machine learning models, are not released into production. However, organizations that are top performers in deploying AI and analytics models into production have eased the deployment by Bisson et al.,<sup>4</sup> Cam and Chui,<sup>5</sup> and Davenport and Malone<sup>9</sup>: involving stakeholders early, cross-functional teams, focus on data governance strategies, new roles (for example, product managers, translators, AI strategist) that

bridges data scientists with business users, aligning with important corporate strategies, and continuous learning programs (for example, in data literacy).

The key elements for scaling advanced analytics presented here align well with what top performers in deploying AI and analytics models focus on and best practices from change management. What is unique about scaling advanced analytics? If we put aside proper use of change management and having data strategies in place, our answer would be: an increased focus on establishing new roles that bridge technical experts and business users; and reskilling employees in data literacy. □

### References

1. Ashkenas, R. and Matta, N. How to scale a successful pilot project. *Harvard Business Review* (2021).
2. Berndtsson, M., Ericsson, A. and Svahn, T. Scaling up data-driven pilot projects. *AI Magazine* 41, 3 (2020), 94–102.
3. Berndtsson, M. et al. 13 organizations' attempts to become data-driven. *International J. Business Intelligence Research (IJBIR)*, 11, 1 (2020), 1–21.
4. Bisson, P. et al. *Breaking Away: The Secrets to Scaling Analytics*. McKinsey, 2018.
5. Cam, A. and Chui, M. *AI Proves Its Worth, but Few Scale Impact*. McKinsey, 2019.
6. Cohen, D.S. *The Heart of Change Field Guide: Tools and Tactics for Leading Change in Your Organization*. Harvard Business Review Press, 2005.
7. Cooley, L. and Linn, J.F. *Taking Innovations to Scale: Methods, Applications and Lessons*. Results for Development Institute, 2014.
8. Davenport, T. and Bean, R. Big companies are embracing analytics, but most still don't have a data-driven culture. *Harvard Business Review* (2018).
9. Davenport, T. and Malone, K. Deployment as a critical business data science discipline. *Harvard Data Science Review* 3, 1 (2021).
10. Golden-Biddle, K. How to change an organization without blowing it up. *Harvard Business Review* 54, 2 (2021), 35–41.
11. Harris, J.G., Craig, C. and Egan, H. *How to Organize Your Analytical Talent*. Accenture Institute for High Performance, 2009.
12. Hart, D. and Warne, L. Comparing cultural and political perspectives of data, information, and knowledge sharing in organisations. *International J. of Knowledge Management (IJKM)* 2, 2 (Apr. 2006), 1–15.
13. McAfee, A. and Brynjolfsson, E. Big data: The management revolution. *Harvard Business Review* 90, 10 (2021), 60–68.
14. NewVantagePartners. *Data and AI Leadership Executive Survey* (2022).

**Mikael Berndtsson** (mikael.berndtsson@his.se) is an associate professor of computer science in the School of Informatics, University of Skövde, Sweden.

**Anna-Carin Jonsson** (anna-carin.jonsson@capgemini.com) is a senior advisor at Capgemini Insights and Data, Stockholm, Sweden.

**Magnus Carlsson** (magnus.carlsson@capgemini.com) is Vice President, Head of Innovation, at Capgemini Insights and Data Nordics, Gothenburg, Sweden.

**Thomas Svahn** (thomas.svahn@capgemini.com) is Vice President, Head of Capgemini Insights and Data Sweden, Gothenburg, Sweden.

We are very grateful for the comments received by the anonymous reviewers. The research is partially supported by VINNOVA, Sweden's innovation agency.



# Opinion

## Five Ways Executives Misunderstand Technology

*How executives might better understand technology.*

**E**XECUTIVE UNDERSTANDING OF technology is incomplete and sometimes even dangerous to the pursuit of corporate strategies. If you doubt this, ask yourself how many executives understand “technology” as well as they understand sales, finance, marketing, and human resources.

Technology literacy requires executives to understand the range of existing and emerging technologies and how they might impact business processes, models and strategies. This understanding is not at the engineering level or how a computer scientist might explain AI and machine learning, but at the contextual and purposeful levels.

The backdrop is not encouraging. Huge gaps about all things digital remain among executives and even CIOs and CTOs<sup>9,14</sup>; the *italics* are mine: “The analysis of almost 2,000 large companies find *only 7% have digitally savvy executive teams*. Even more surprising, *only 47% of CTOs and 45% of CIOs could be considered ‘digitally savvy.’* This percentage drops to 24% for COOs and 23% of CEOs. Only 12% of CFOs were up to speed digitally.”

There is additional survey and interview data about executive shortcomings and “misunderstandings”<sup>3-5,10,14</sup> that informs some ways executives *should* understand technology—which are explored here.

### Misunderstanding #1: Technology Is Operational

Most executives see technology through



an operational lens. Technology to them is hardware, software, networks, applications, and cloud computing. It is about “IT.” Most executives still think about technology across a narrow spectrum of functionality, which spans from wires and devices—phones, tablets, and laptops—to applications—like ERP and CRM—to networks—such as 5G communications. While most executives have only a high-level understanding of operational technology, the understanding they have focuses way too much on technology that has already been commoditized. Yes, it can still be done poorly, *but there is little or no competitive advantage in operational technology investments.* Executives should leave operational technology to IT support per-

sonnel who seldom if ever think about how wires, devices, and networks enable corporate strategy.

### Understanding #1: Technology Is Strategic

Technology should enable strategy. Executives should never separate technology from strategy, which is what an operational understanding of technology encourages. In order to optimize the business-technology relationship, executives must insist upon an actionable strategy, a robust business model, and a business process inventory. Many companies have a business strategy that is abstract, an ill-defined business model and business processes that have not been

inventoried. These shortcomings challenge any company's ability to strategically leverage digital technology, so step one is the development of an actionable strategy, an explicit business model and processes that can be "mined" and improved with technology. *Executive understanding of "technology" should connect directly to strategic objectives—not operational efficiency.*

### **Misunderstanding #2: Technology Is Digital Transformation**

The vast majority of digital transformation projects are not transformative at all. Most of them actually fail.<sup>8,11</sup> At best, they are incremental improvements of processes that have been neglected for years. Some of these projects, such as the replacement of legacy applications, are pure modernization projects. If executives understand technology as digital transformation they are boxing their understanding of technology way too restrictively. *Worse, way too many public-company earnings reports and calls refer to delayed digital transformation projects or projects that had much less impact than the executives once promised. While the number of digital transformation projects continues to grow, boards of directors, stock analysts and selected stakeholders see the number of digital transformation projects as an imperfect metric.*

### **Understanding #2: Technology Is Much More than Digital Transformation**

Executives should understand that "digital transformation" is one limited application of technology that by no means defines the totality of what technology can do. While distinctions among incremental, modernization-driven and disruptive "transformation" are helpful, executives should see technology much more broadly. In fact, distinctions among different flavors of digital transformation and other technology investments have always been small. A better understanding of digital transformation sees transformation subsumed in strategy, a business model and an array of business processes that can be improved, automated, or even replaced. Executives should also manage the

## **The vast majority of digital transformation projects are not transformative at all.**

messaging around their technology investments way beyond stories about how numerous and "successful" their digital transformation projects have been. *Executives should free themselves of narrow definitions of transformation and invest in a broad suite of projects they can describe as the company's overall business-technology strategy—not their slate of transformation projects. They should also condition their stakeholders to see technology as strategic, not tactical.*

### **Misunderstanding #3: Trends Tracking Is Technology Literacy**

Many executives understand technology as the technology trends someone tells them to track. The Gartner Group,<sup>6</sup> among other research organizations, believes executives must understand broad technology trends to understand how technology will impact their companies. Executives who chase technology trends are looking for validation (perhaps among their peers who want to talk about the trends they see on X/Twitter) rather than strategic due diligence. Technology trends tracking creates solid optics but seldom translates into value because the tracking process itself becomes more important than the solutions, and without a deeper understanding of the technologies business solutions are impossible to find. Companies that subscribe to research organizations such as Gartner, Forrester, and IDC follow the technology trends these organizations anoint every year. But trends change annually, the trends are different from source-to-source and the "advice" around the trends is usually inconsistent. *There is also a huge gap between high-level trends and actual problem solving. Trends in cloud computing, for example, tell executives very little about what*

*they should do about multicloud delivery models. Technology tracking for the sake of tracking feels good but misses the mark.*

### **Understanding #3: Emerging Technology Is Technology Literacy**

The short list of emerging technologies includes machine learning, natural language processing, the Internet of Things, advanced analytics, quantum computing, 5G/6G, cybersecurity, edge computing, data lakes, robotics, augmented, and virtual reality and, of course, generative AI, and tools such as ChatGPT, Bard, and Copilot. But as suggested, it is not about the lists, but about the kinds of problems the technologies can solve. *While no one expects executives to talk code, architecture, or APIs, they should be able to discuss how, for example, augmented and virtual reality might improve their products and services, and how they plan to prototype emerging technologies.*

Executives should especially understand the applied potential of "data" and analytics. They should understand how data powers their internal and external processes. They should understand the strengths and weaknesses of structured and unstructured data. They also must understand the premise of analytics, what it seeks to do, and how it works at a functional level.<sup>2</sup> Executives should understand how the role of data has fundamentally changed. "Database management" yielded to "data warehouses" which yielded to "data lakes" which all now live within "data fabrics." *Executives should define their company's commitment to "data-as-a-service" and all of the enabling technologies that make this service possible.*

### **Misunderstanding #4: Technology Is a State of Projects**

Companies spend crazily on technology projects. Big projects, including ERP and CRM projects, and smaller ones, such as prototyping new minimum viable products (MVPs). Executives who track projects become faux project managers, not strategic technology leaders. Executives who play project managers on TV are poor project managers back at the ranch primarily because they do not understand the nuances of technology proj-

ect management. For example, the movement from waterfall to Agile application development must include the uncertainty and ambiguity of application development. Executives unfamiliar with development methodologies—which includes most executives—will misinterpret project updates. *But much more importantly, executives who track technology projects understand technology at way too granular a level which separates them from their strategic purpose. Worse, slates of projects do not constitute an integrated technology investment plan.*

#### **Understanding #4: Technology Is an Integrated Investment Strategy**

Technology is as strategic as sales, customer retention, and R&D—perhaps even more so because technology enables all three of these activities. The whole notion of “projects” threatens the advantages companies generate from integrated technology investments. “Projects” have their own start/finish dates, budgets, objectives, and teams. They are also conducted in organizational silos. Executives should understand technology as a comprehensive investment strategy that extends well beyond individual projects. Wholistic, top-down investment strategies link projects. *All project business cases should therefore track to other projects and strategic value. This means that “one offs” or so-called “science projects”—well known in industry—never get funded. While this requires some executive discipline, it ensures technology investments serve the greater good, not some faraway corporate lab obsessed with projects that make little or no strategic sense.*

#### **Misunderstanding #5: Technology Is a List of Management Questions**

This understanding of technology is often suggested by consultancies as a way for executives to demonstrate their technology gravitas. Questions such as “How long does it take our company to deploy new applications?,” “What % of business decisions are we making with help from AI?,” or “Have we placed high-quality engineers in roles that contribute the most value to the company?” are good management questions, but they are detached from what the questions—

## Executives should understand the relationships among emerging technologies, problems, and strategy.

and nuanced answers—actually mean. Worse, executives who try to demonstrate technology prowess with management questions—without sufficient technology grounding—can be misled or downright deceived. *Operational performance questions are wasted on executives who have much more important things to do.*

#### **Understanding #5: Matching Is the Only Important Question**

Executives should understand the relationships among emerging technologies, problems, and strategy. They should be capable of “matching” technologies to problems. This means they should have the ability to reverse engineer solutions from desired strategic outcomes with realistic technology solutions. The executive matching process should be proactive: it should not be unusual for executives to suggest how specific technologies could be leveraged onto specific business models and processes. Executives should not wait for suggestions from their teams. They should be active participants in the matching process. For example, executives should be quite capable of suggesting how AI might automate parts of the recruiting or training processes. They should help suggest how augmented and virtual reality might improve their marketing strategies. *The ability to match technologies with problems should become a core competency for executives who live in the digital world—which they all do. The only management question that matters is the one that focuses on how technology can solve specific problems.*

#### **Conclusion**

The five ways executives misunderstand technology and the five ways they *should* understand technology is an exercise of old versus new. No one underestimates the pressures on senior executives or the extraordinary range of their corporate responsibilities. Adding “technology” to their core competencies—in the five ways discussed here—is not just another activity. Technology touches every internal and external process there is. It is not hyperbole to say that without technology there is no business. Executives must dismiss their misunderstandings of technology and replace them with much more strategic ones.

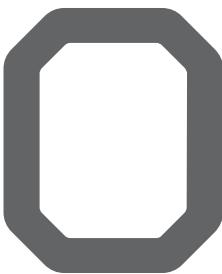
There is a purposeful throughline from misunderstandings to understandings. Executives should abdicate operational technology to traditional “IT” and force technology to enable its corporate strategy, its business model and the processes that inhabit the model. With the right operational talent, the trains will run on time. Executives should focus on the development of new destinations and next generation train technology. C

#### **References**

1. Andriole, S. Five myths about digital transformation. *Sloan Management Review* (2017).
2. Andriole, S. and Barsky, N. Overdue diligence: Questioning the promise, not the premise, of analytics. *Communications of the AIS*. (2022).
3. Cheng, J.W. et al. Is your c-suite equipped to lead a digital transformation? *Harvard Business Review*. (2021).
4. Forth, P. et al. *Flipping the Odds of Digital Transformation Success*. Boston Consulting Group, (2020).
5. Graves, D. Today's CEOs need hands-on digital skills. *Harvard Business Review*. (2021).
6. Groombridge, D. *Gartner Top 10 Strategic Technology Trends for 2023*. The Gartner Group, (2022).
7. Klemens, H. et al. *The CEO's New Technology Agenda*. McKinsey and Company, (2019).
8. Martin, J. *Unlocking Success in Digital Transformations*. McKinsey and Company, (2019).
9. McKendrick, J. Training upward: Your executives may not fully understand digital transformation. *ZDNET*. (2021).
10. Panetta, K. *The CIO's Biggest Digital Transformation Challenges*. The Gartner Group, (2016).
11. Saldanha, T. *Why Digital Transformation Fails: The Surprising Disciplines of How to Take Off and Stay Ahead*. Berrett-Koehler Publishers, (2019).
12. The Fast Company Executive Board. How c-level leaders keep up with the latest technologies. *Fast Company*. (2022).
13. Weill, P. et al. It pays to have a digitally savvy board. *Sloan Management Review*. (2019).
14. Weill, P. Does your c-suite have enough digital smarts? *MIT Sloan Management Review*. (2021).
15. Wiles, J. *Top 10 Emerging Skills for the C-Suite*. The Gartner Group, (2019).

**Stephen J. Andriole** (steve@andriole.com) is the Thomas G. Labrecque Professor of Business Technology in the Villanova School of Business, Villanova, PA, USA.

Copyright held by owner(s)/author(s).



## Historical Reflections

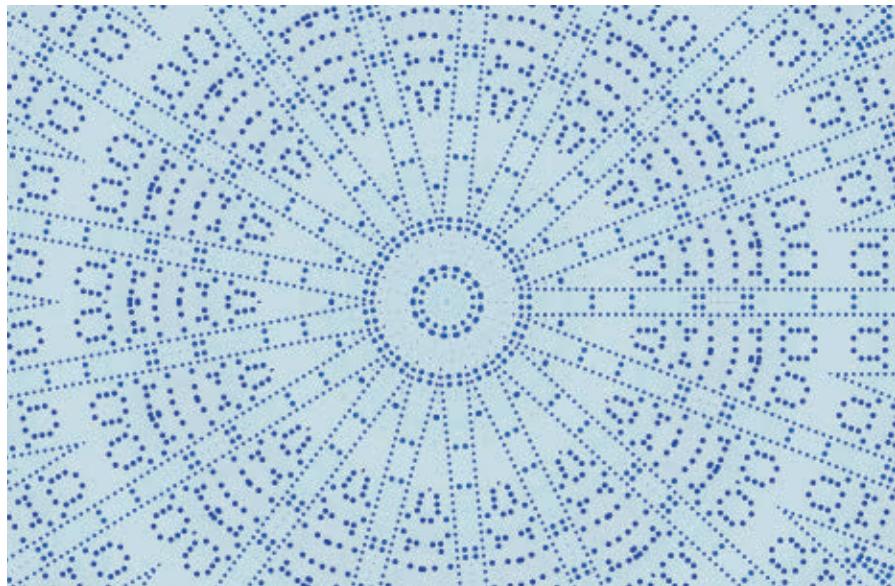
# There Was No ‘First AI Winter’

*Despite challenges and failures, the artificial intelligence community grew steadily during the 1970s.*

**A**S I CONCLUDED my June Historical Reflections column, artificial intelligence had matured from an intellectual brand invented to win funding for a summer research workshop to one of the most prestigious fields in the emerging discipline of computer science. Four of the first 10 ACM A.M. Turing Award recipients were AI specialists: Marvin Minsky, Herb Simon, Allen Newell, and John McCarthy. These men founded the three leading AI labs and played central roles in building what are still the top three U.S. computer science programs at MIT, Stanford, and Carnegie Mellon. Conceptually AI was about uncovering and duplicating the processes behind human cognition; practically it was about figuring out how to program tasks that people could do but computers could not. Although *connectionist* approaches based on training networks of simulated neurons had been prominent in the primordial stew of cybernetics and automata research from which AI emerged, all four Turing Award recipients favored the rival symbolic approach, in which computers algorithmically manipulated symbols according to coded rules of logic.

### A History of Failed Ideas?

AI was born in hype, and its story is usually told as a series of cycles of fervent enthusiasm followed by bitter disappointment. Michael Wooldridge, himself an eminent AI researcher, began



his recent introduction to the field by remembering when he told a colleague about his plan to tell “the story of AI through failed ideas.” In response, “she looked back at me, her smile now faded. ‘It’s going to be a bloody long book then.’”<sup>22</sup>

Major awards lag years behind research. By the time Newell and Simon shared the 1975 ACM A.M. Turing Award the feasibility of their approaches to AI was being increasingly challenged. The AI community would have to wait 19 years for another winner. It was displaced as the intellectual high ground of the emerging discipline by *theoretical computer science*, a field centered on mathematical analysis of algorithms, which garnered nine awardees

during the same period.<sup>a</sup> This new focus bolstered the intellectual respectability of computer science with a body of theory that was impeccably mathematical yet unlike numerical analysis, which was falling out of computer science over the same period, not directly useful to or understood by other schol-

<sup>a</sup> These awards focused on computational complexity theory and the analysis of algorithms. I am construing theoretical computer science here to encompass the work of Rabin and Scott (1976), Cook (1982), Karp (1985), Hopcroft and Tarjan (1986), Milner (1991), and Hartmann and Stearns (1993). I am not including winners cited primarily for contributions to programming languages, except for Milner whose citation emphasized theory, though Wirth and Hoare both made important theoretical contributions.

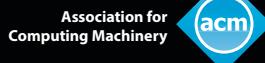
**“I am  
human...  
just like  
you.”**

For the first  
time, the AI  
had asserted  
its claim on  
humanity.

At that moment, the AI  
and I had become one...

*Communications of the ACM* is looking for writers in our community to contribute sci-fi short stories, between 1,000 and 1,200 words, for our quarterly “Future Tense” section.

**Do you have  
a great story to tell?**  
Make contact at  
[LastByte@cacm.acm.org](mailto>LastByte@cacm.acm.org)



ars in established disciplines.

The problems AI researchers had taken as their test cases were difficult in a fundamental mathematical sense that dashed hopes of ingenious breakthroughs. Once AI researchers applied the new techniques of complexity analysis “Everywhere they looked—in problem solving, game playing, planning, learning, reasoning—it seemed that the key problems where NP-complete (or worse).”<sup>22</sup> Progress would come slowly and painfully, with methods that worked in some cases but not others.

The early practitioners of AI had consistently and spectacularly overestimated the potential of their methods to replicate generalized human thought. In 1960, for example, Herb Simon had declared “within the near future—much less than 25 years—we shall have the technical capability of substituting machines for any and all human functions in organizations.” He believed the “problem-solving and information handling capabilities of the brain” would be duplicated “within the next decade.” As professionals were replaced by machines “a larger part of the working population will be mahouts and wheelbarrow pushers and a smaller part will be scientists and executives.”<sup>20</sup>

The same processes of hype that gave AI a high profile for military sponsors and awards committees also made the field a topic of public debate. Promises made for intelligent computers tapped into longer-established myths and science fiction stories of thinking machines and mechanical servants. HAL, the murderous computer from the movie *2001: A Space Odyssey*, whose name was said to be a contraction of *heuristic* and *algorithmic*, was one of many fictional depictions of the promises made by AI researchers. Minsky himself had been hired to get the details right. Meanwhile a series of books appeared criticizing those promises and challenging the feasibility of artificial intelligence.<sup>10</sup>

The AI boosters were wrong, of course, though their critics were not always right. Computers had always been sold with futuristic hype, but overly optimistic technical predictions made during the same era for other areas of computer science such as graphics,

computer mediated communication, scientific computation, and databases were eventually met and surpassed. In contrast, the approaches adopted by the AI community conspicuously failed to deliver on Simon’s central promises.

### Military Origins of AI

AI began as a Cold War project centered on a handful of well-connected researchers and institutions. The original Dartmouth meeting was funded by the Rockefeller Foundation. A full-scale research program would require deeper pockets, which in the 1950s were usually found attached to military uniforms. When Newell and Simon met and began to collaborate on their famous theorem prover both were employed by the RAND Corporation, a non-profit set up to support the U.S. Air Force. This gave them access not just to RAND’s JONNIAC computer, one of the first modern style computers operational in the U.S., but also to RAND programmer Clifford Shaw who was responsible for squeezing the ambitious program into the machine’s tiny memory.<sup>5</sup> Frank Rosenblatt developed his perceptrons, the most important of the early neural networks, in a university lab funded by the U.S. Navy. At MIT, Minsky’s early work benefitted from the largess of the Air Force and the Office of Naval Research.

This deep entanglement of early AI with the U.S. military is difficult to overlook. Johnnie Penn highlighted the military dimension in his recent dissertation, challenging the phrase “good old fashioned AI” (a common description of traditional symbolic work) as something that “misrepresents and obscures this legacy as apolitical.”<sup>18</sup> Yarden Katz insists the apparent technical discontinuities in

**This deep  
entanglement of  
early AI with the U.S.  
military is difficult to  
overlook.**

the history of AI are just distractions from a consistent history of service to militarism, American imperialism, and capitalism.<sup>13</sup>

Yet AI was not exceptional in this respect. Military agencies supplied approximately 80% of all Federal research and development funding during the 1950s, the first decade of the Cold War. This wave of money flowed disproportionately to MIT and Stanford, which were not only the two leading centers for both AI and computer science but also the primary exemplars of a new, and to many disturbing, model for the relationship between universities, the Federal government, and military needs. Stuart W. Leslie's history book *The Cold War and American Science* focused entirely on those two institutions as prototypes for a new kind of university restructured around military priorities.<sup>14</sup>

Computing was, after all, an expensive endeavor and there were few alternative sources of support. Set in the larger picture of military investment in computing, including projects such as the SAGE air defense network and guidance systems for Minuteman missile, the sums spent on AI appear quite small. Most computing projects of the 1940s and 1950s were underwritten directly or indirectly by the U.S. military.<sup>8</sup> ENIAC, the first programmable electronic computer, was commissioned in 1943 by the U.S. Army for use by its Ballistics Research Laboratory.<sup>12</sup> Such relationships blossomed as the Second World War gave way to the Cold War. IBM, for example, received more money during the 1950s from developing custom military systems such as the hardware for the SAGE air defense network than it did from selling its new lines of standard computer models. And even the market for those standard projects was driven by the Cold War. IBM's first commercial computer model, the 701, was known informally as the "Defense Calculator" and sold almost entirely to government agencies and defense contractors. It was the Federal government, not IBM itself, that managed the delivery schedule for early models to maximize their contribution to national security.<sup>11</sup> The needs of military and aerospace projects kick-started the semiconductor industry in what became Silicon Valley.

## ARPA supported well-connected research groups without formal competitive review or any commitment to provide specific deliverables.

AI remained heavily dependent on military funding in the 1960s, as labs at MIT and Stanford received generous funding through the newly established Advanced Research Projects Agency. ARPA reportedly spent more on AI than the rest of the world put together, most of which went to MIT and Stanford. Carnegie Mellon was not initially in the same league, but its early success in computing and artificial intelligence won substantial ARPA funding by the 1970s and fueled the rise of the university itself. The National Science Foundation, a civilian agency, was less important. During the 1950s and 1960s it did not have a directorate focused on computer science. It made few grants to support computing research (though it was active in funding computing facilities).<sup>3</sup>

ARPA supported well-connected research groups without formal competitive review or any commitment to provide specific deliverables. J.C.R. Licklider, the first director of APRA's Information Processing Techniques Office, joined ARPA from military contractor BBN and had earlier been a member of the MIT faculty. After showering MIT with government money he eventually rejoined its faculty, to run the ARPA-funded Project MAC (into which Minsky and his AI group had been incorporated). Licklider then returned to ARPA for a second term as director. That might all seem a little too cozy by modern standards, but ARPA's early success in fostering emerging computer technologies was spectacular: not just the Internet, which can be traced back to an idea of Licklider's,

but also computer graphics and time-sharing.<sup>17</sup> Paul Edwards summarized the early history of AI in his classic *The Closed World*, arguing that under the influence of ARPA it became "part of the increasingly desperate, impossible tasks of enclosing the U.S. within an impenetrable bubble of high technology." He believed Licklider's vision for interactive computing was shaped fundamentally by military concerns with command and control.<sup>6</sup>

Were the founders of AI who worked at RAND or took money from the Pentagon thereby coopted into an imperialistic effort to project American power globally? Did their work somehow come to embed the culture of the military industrial complex? Historians will likely be arguing these questions for generations to come. AI, like cybernetics, unquestionably benefitted from a powerful alignment with a more general faith of scientific and political elites in abstraction, modeling, and what has been called by historians of science *Cold War rationality*.<sup>7</sup>

Personally, though, I am inclined to see the founders of AI as brilliant boondogglers who diverted a few buckets of money from a tsunami of cold war spending to advance their quirky obsessions. Steven Levy noted that a "very determined solipsism reigned" among the hackers of Minsky's lab at MIT, even as the antiwar protesters forced them to work behind locked doors and barricades. He quoted Minsky as claiming Defense Department funding was less intellectually corrosive than hypothetical money from the Commerce Department or the Education Department.<sup>4,16</sup> On the Stanford side, John McCarthy was a proponent of scientific internationalism. He was raised communist and made five visits to the USSR during the 1960s, though his politics drifted rightward in later decades.<sup>21</sup> Philip Agre, recalling the investments by the military industrial complex in his graduate training at MIT, wrote that "if the field of AI during those decades was a servant of the military then it enjoyed a wildly indulgent master."<sup>22</sup>

### Summers and Winters

When scientists write histories they usually focus on intellectual and technical accomplishments, leaving professional historians and science studies

scholars to raise indelicate questions about the influence of money. In contrast, the insider story of AI as told by experts such as Wooldridge, Nils J. Nilsson and Margaret Boden has been structured explicitly around shifts in government funding.<sup>14</sup>

Why was AI so vulnerable to the whims of government agencies? One factor was the concentration of early AI work in a handful of well-connected labs. Another was the reliance of AI researchers on large and expensive computers. Perhaps the most important was the failure of AI, during its first few decades, to produce technologies with clear commercial potential that might attract a broader range of sponsors. The health of AI as a field thus depended on the ability of researchers to persuade deep-pocketed sponsors that spectacular success was just around the corner.

Relying on a handful of funding sources proved hazardous. Machine translation projects were an early beneficiary of military largess, but this made them vulnerable when their feasibility was questioned. American funding for this area dried up after a damning report was issued in 1966 by the ALPAC committee, a scientific panel sponsored by the Department of Defense, the National Science Foundation, and the CIA to investigate progress in the area.<sup>15</sup>

The late 1980s are universally seen

as the beginning of the “AI Winter,” in which faith and funding for AI dwindled dramatically. I will tell that story later, but in a fuzzier way the period from 1974 to 1980 has increasingly been described as an earlier winter for AI.<sup>16</sup> This narrative blames the 1973 Lighthill Report, commissioned by the Science Research Council of the U.K., for a collapse of British support for AI work. Across the Atlantic, this is said to have inspired other funders to ask more difficult questions.<sup>17</sup>

Sir James Lighthill was commissioned to write his report with the specific intent of justifying the withdrawal of funding for Donald Michie’s lab at Edinburgh, the most important center for AI research in the U.K. Lighthill, an eminent applied mathematician, endorsed both practical work on industrial automation and work to support analysis of brain functions (often called cognitive science) but opposed funding for ambitious work intended to unite the two in a new science of machine intelligence. Jon Agar’s analysis makes clear that the creative blurring of categories which let AI researchers spin nar-

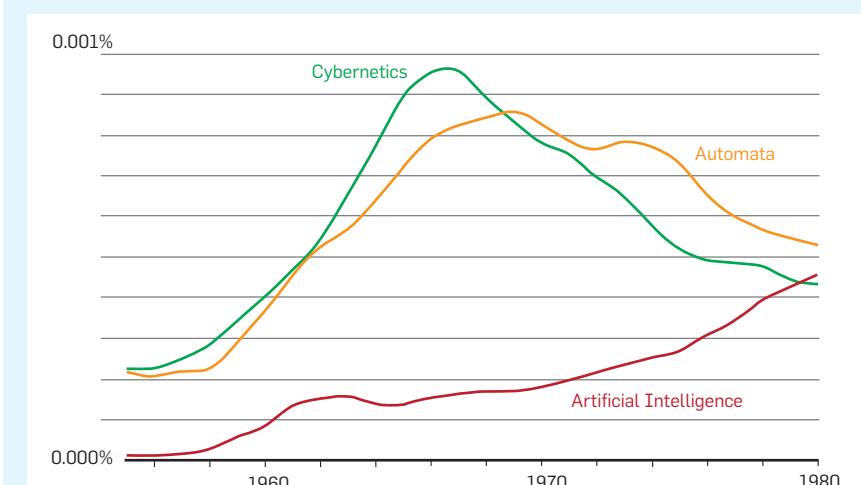
row achievements into vast promises also left them vulnerable to attacks that challenged implied connections between specific computational techniques and general intelligence.<sup>18</sup>

How, and indeed whether, Light-hill’s attack on one controversial lab precipitated a broad international funding collapse for AI remains unclear. It coincided with broader changes in U.S. support for science. In 1973 Congress, inspired by Vietnam-era concerns that the military was gaining too much power over university research, passed legislation explicitly prohibiting ARPA from funding work that was not directly related to military needs. ARPA was renamed to DARPA, the D standing for Defense.<sup>19</sup> As responsibility for basic research funding shifting increasingly to the NSF, DARPA funding required more direct military justification and came with more strings attached.

### AI’s Steady Growth in the 1970s

Historical work on AI has so far focused on a handful of elite institutions and researchers, just as I have done in the series so far. Any DARPA-related AI slowdown was felt most deeply at those highly visible sites, whose researchers, and graduates were the people most likely to give keynotes, write memoirs, and shape the collective memory of the discipline. But the big contracts awarded on a handshake only ever went to a few institutions, whereas the institutionalization of AI took place internationally and across a much broader range of universities. In the first major history of AI, Pamela McCorduck noted that 21 years after the Dartmouth conference the influence of its participants and the programs they founded remained strong. Nine invited papers were given at the 1977 International Joint AI Conference: three by Simon and his former students Ed Feigenbaum and Harry Pople; one by Feigenbaum’s own student Doug Lenat; one by Minsky and one by McCarthy. MIT, Stanford, SRI, and Carnegie Mellon dominated, “with the representation from other laboratories being sparser than might have been expected in a field that had grown from the 10 Dartmouth pioneers in 1956 to nearly 1,000 registrants in 1977.”<sup>20</sup>

**Google Ngram data, based on a large English language corpus and plotted between 1955–1980. References to artificial intelligence rose consistently through the 1970s even as discussion of the related concepts of automata and cybernetics declined sharply.<sup>†</sup>**



<sup>†</sup> For example Wooldridge,<sup>21</sup> 61–62.

Computer science developed as a highly federated discipline, in which most practitioners identified and engaged more with their specialist area than with the field as whole. Over time new areas such as networking, databases and graphics gained prominence while others slipped out of the mainstream. AI enthusiasts created SIGART in 1966, one of the first Special Interest Groups (SIGs) within the ACM. The SIGs gave institutional recognition to the subfields of computer science. With their own publications, conferences, and finances they came to account for most of the ACM's activity.

I suspect that judged by metrics such as the number of students enrolled in AI courses, total number of AI researchers, attendance at conferences, or quantity of research publication the story of AI in the 1970s and 1980s would look less like a series of abrupt booms and busts and more like a march toward disciplinary professionalization. As a first step in this analysis, I located two data sources, neither of which supports the idea of a broadly based AI winter during the 1970s.

One is membership of ACM's SIGART, the major venue for sharing news and research abstracts during the 1970s. When the Lighthill report was published in 1973 the fast-growing group had 1,241 members, approximately twice the level in 1969. The next five years are conventionally thought of as the darkest part of the first AI winter. Was the AI community shrinking? No! By mid-1978 SIGART membership had almost tripled, to 3,500. Not only was the group growing faster than ever, it was increasing proportionally faster than ACM as a whole which had begun to plateau (expanding by less than 50% over the entire period from 1969 to 1978). One in every 11 ACM members was in SIGART.

Not all the participants in this growing community worked in elite, DARPA-funded labs. As the *SIGART Bulletin* summed up the AI hierarchy a few years later: "AI research in the U.S. can be viewed as having three major components: a few highly visible sites with major efforts; a large number of sites with smaller numbers of workers; and a diffuse set of researchers and developers in other fields who believe that AI research may be rel-

## Not all the participants in this growing community worked in elite, DARPA-funded labs.

evant for them."<sup>c</sup> Perhaps the people in those "highly visible sites" were suffering, but general interest in AI continued to grow rapidly.

The other data source is Google's Ngram viewer, which suggests that the term *artificial intelligence* became more common during the so-called AI winter. Its growth stalled for a few years in the mid-1960s but recovered by 1970 and grew steadily through 1980. Yet the mid-1960s are usually described as a golden age for AI.

The AI community founded an independent organization—the American Association for Artificial Intelligence—in 1979. Reporting this news, the SIGART chair Lee Erman noted wistfully "SIGs are set up as arms of ACM and as such must obtain ACM approval for most significant actions, including budgets, new publications, sponsorship of conferences, and interaction with non-ACM organizations. This structure may be appropriate for a 'special interest group' (although I would argue more autonomy would be beneficial to the SIGs and to ACM), but not for a national scientific organization, which needs far more independence."<sup>d</sup> This proved prophetic: while SIGART membership continued to grow well into the 1980s, the new association eventually replaced it as the hub of the AI community by developing a panoply of publications, conferences, and awards.

That is all I have space in this column, but in the next installment I will be looking at the codification of AI's intellectual content in early textbooks and its entrenchment in the computer

science curriculum, at the new emphasis in the 1970s on knowledge representation over pure reasoning, and at the spectacular bubble of funding for expert systems in the early 1980s, which burst to create the real "AI Winter." □

### References

- Agar, J. What is science for? The Lighthill report on artificial intelligence reinterpreted. *British J. for the History of Science* 53, 3 (Sept. 2020), 289–310.
- Agre, P. Toward a critical technical practice: Lessons learned in trying to reform AI. In *Bridging the Great Divide: Social Science, Technical Systems, and Cooperative Work*. G. Bowker, Ed. Psychology Press, New York, 1997, 131–158.
- Aspray, W. and Williams, B.O. Arming American scientists: NSF and the provision of scientific computing facilities for universities, 1950–73. *IEEE Annals of the History of Computing* 16, 4 (Winter 1994), 60–74.
- Boden, M.A. *Mind as Machine: A History of Cognitive Science*. Clarendon Press, Oxford, U.K., 2006.
- Dick, S. Of models and machines: Implementing bounded rationality. *Isis* 106, 3 (2015), 623–634.
- Edwards, P.N. *The Closed World: Computers and the Politics of Discourse in Cold War America*. MIT Press, Cambridge, MA, 1996.
- Erickson, P. et al. *How Reason Almost Lost Its Mind: The Strange Career of Cold War Rationality*. University of Chicago Press, Chicago, IL, 2013.
- Flamm, K. *Creating the Computer: Government, Industry, and High Technology*. Brookings Institution, Washington, D.C., 1988.
- Freeman, P.A., Adrion, W.R., and Aspray, W. *Computing and the National Science Foundation, 1950–2016*. Association for Computing Machinery, NY, 2019.
- Garvey, S.C. The 'general problem solver' does not exist: Mortimer Taube and the art of AI criticism. *IEEE Annals of the History of Computing* 43, no. 1 (2021).
- Haigh, T. Computing the American way: Contextualizing the early U.S. computer industry. *IEEE Annals of the History of Computing* 23, 2 (Apr.–June 2010), 8–20.
- Haigh, T., Priestley, M. and Rope, C. *ENIAC In Action: Making and Remaking the Modern Computer*. MIT Press, Cambridge, MA, 2016.
- Katz, Y. *Artificial Whiteness: Politics and Ideology in Artificial Intelligence*. Columbia University Press, New York, 2020.
- Leslie, S. *The Cold War and American Science: The Military-Industrial-Academic Complex at MIT and Stanford*. Columbia University Press, NY, 1993.
- McCorduck, P. *Machines Who Think*. A.K. Peters, Natick, MA, 2004, 130–131.
- Norberg, A.L., O'Neill, J.E., and Freedman, K.J. *Transforming Computer Technology: Information Processing for the Pentagon, 1962–1986*. Johns Hopkins University Press, 1996.
- Nilsson, N.J. *The Quest for Artificial Intelligence: A History of Ideas and Achievements*. Cambridge University Press, New York, 2010.
- Penn, N. *Inventing Intelligence: On the History of Complex Information Processing and Artificial Intelligence in the United States in the Mid-Twentieth Century*. Ph.D. dissertation. University of Cambridge, 2020.
- Poibeau, T. *Machine Translation*. MIT Press, Cambridge, MA, 2017.
- Simon, H. The corporation: Will it be managed by machines? In *Management and Corporations*. M. Anshen and G.L. Bach, Eds. The McGraw-Hill Book Company, New York, 1960, 17–55.
- Tatarchenko, K. Transnational mediation and discipline building in Cold War computer science. In *Communities of Computing: Computer Science and Society in the ACM*. T.J. Misa, Ed. Morgan & Claypool, 2017, 199–227.
- Wooldridge, M. *A Brief History of Artificial Intelligence: What It Is, Where We Are, and Where We Are Going*. Flatiron Books, New York, 2021.

**Thomas Haigh** (thomas.haigh@gmail.com) is a professor of history at the University of Wisconsin—Milwaukee, WI, USA, and a Comenius visiting professor at Siegen University, Germany.

Copyright held by owner(s)/author(s).



## Rethinking how data and computations should be organized.

BY ETHAN L. MILLER, GEORGE NEVILLE-NEIL, ACHILLES BENETOPoulos, PANKAJ MEHRA, AND DANIEL BITTMAN

# Pointers in Far Memory

IT IS THE best of times and it is the worst of times in the world of datacenter memory technology. According to International Data Corporation (IDC), dynamic random-access memory (DRAM) revenues exceeded \$100 billion in 2022. Yet, the anticipated growth rate is hugging the zero line, and many producers either reported loss-making quarters or are rumored to do so soon. From the perspective of datacenter customers, by some estimates, the cost of renting memory ranges from \$20 to \$30 per gigabyte per year, for a resource that costs only \$2 to \$4 to procure outright. On top of this, software as a service (SaaS) end users, for example, are forced to rent all the memory they will need up front. By some estimates, they end up using less than 25% of that memory more than 75% of the time.<sup>10</sup>

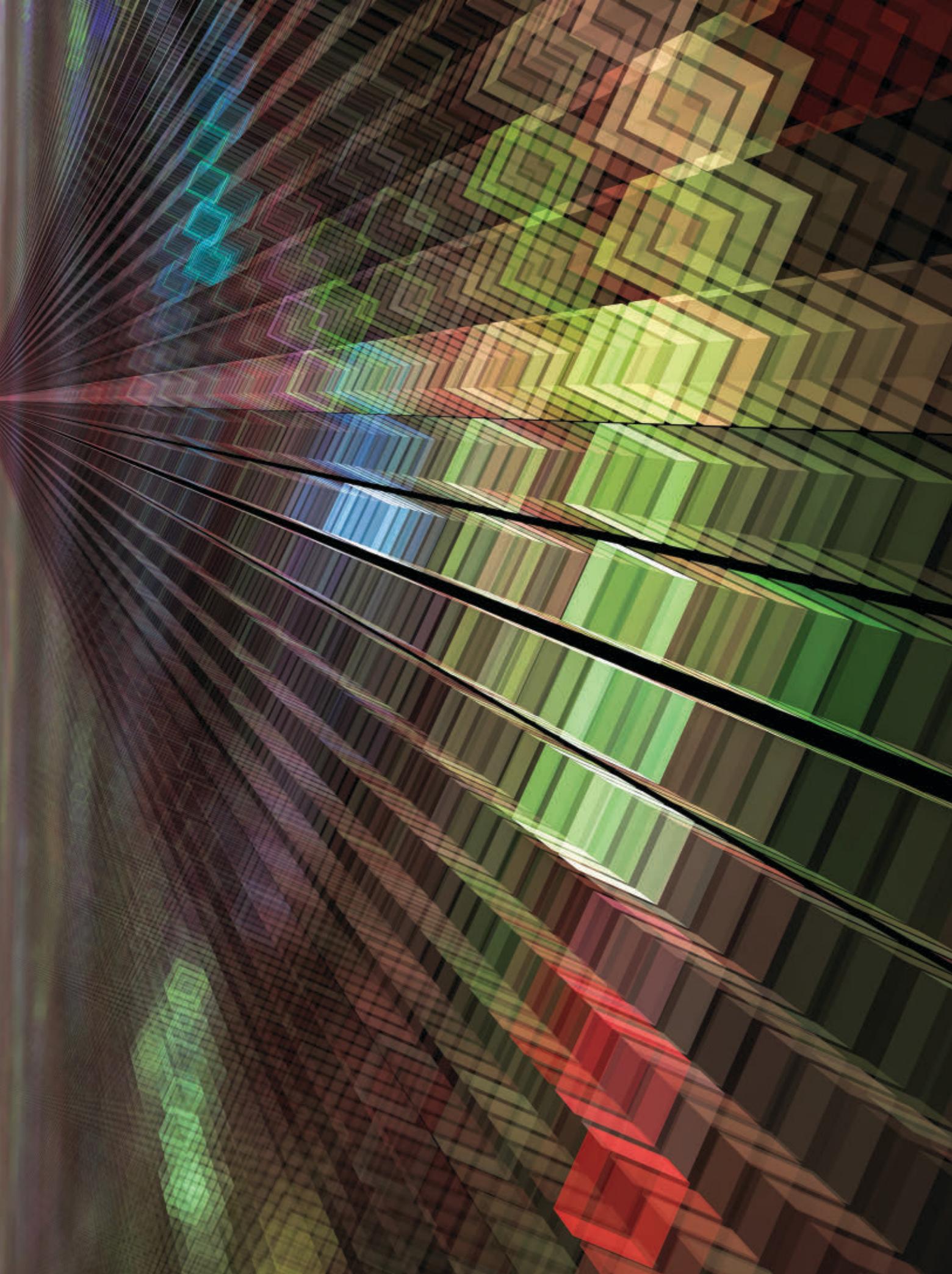
Compute Express Link (CXL), a new technology emerging from the hardware side,<sup>9</sup> is promising to provide *far* memory. Thus, there will be more memory

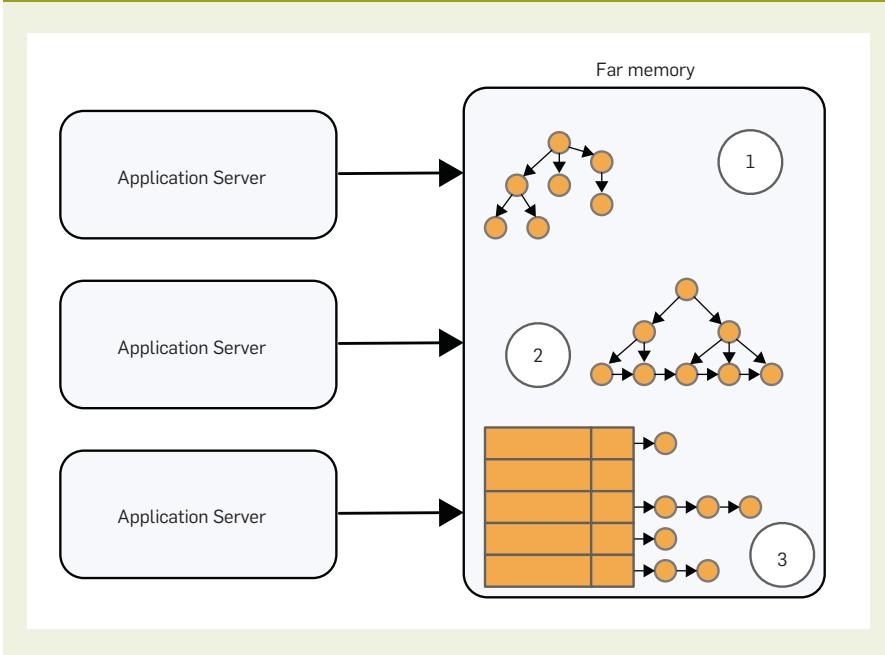
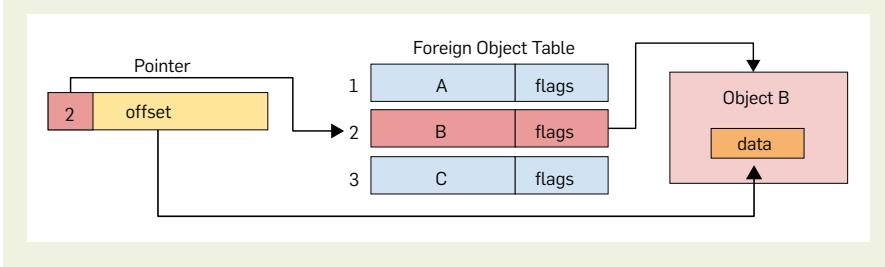
capacity and perhaps even more bandwidth, but at the expense of greater latency. Optimization will first seek to keep memory in far tiers colder, and, second, minimize the rates of both access into and promotion out of these tiers.<sup>1,5</sup> Third, proactive promotion and demotion techniques being developed for far memory promote/demote whole objects instead of one cache line at a time to take advantage of bulk caching and eviction in order to avoid repeatedly incurring its long latency. Finally, offloading computations with many dependent accesses to a near-memory processor is already being seen as a way to keep the latency of memory out of the denominator of application throughput.<sup>11</sup> With far memory, this will be a required optimization.

Applications that operate over richly connected data in memory engage heavily in pointer-chasing operations either directly (for example, graph processing in deep-learning recommendation models) or indirectly (for example, B+ tree index management in databases). Figure 1 shows an example of pointer-chasing applications in far memory: graph traversal; key lookup in a B+ index; and, collision resolution under open hashing.

Data from previous work<sup>2</sup> suggests that as data structures scale beyond the memory limits of a single host, causing application data to spill into far memory, programmers are forced to make complex decisions about function and data placement, intercommunication, and orchestration.

**Performance characteristics of far memory.** By default, pointers (like the internode ones in Figure 1) are defined in the virtual address space of the process that created them. Because of this, if left unoptimized, pointer-chasing operations and their dependent accesses can overwhelm the microarchitecture resources that support memory-level parallelism (for example reorder buffers) even on a single CPU with local memory. With latencies that can range from 150ns to more than 300ns,<sup>2</sup> far memory further compounds this problem.



**Figure 1.** Example of pointer applications in far memory.**Figure 2.** Stored pointers use the index into the FOT as placeholder.

In a distributed setting, implementing a simple-minded pointer-chasing offload without taking care of virtual-to-physical address translation results in chatty internode coordination with the parent process.<sup>15</sup> Effective optimization of pointer-chasing operations entails minimizing communication between the near-memory processor executing the traversal and the server running the parent process.

**Developing far memory-ready applications.** Evidence from high-performance computing (HPC) and database workloads points to the extreme inefficiency of pointer-rich sparse memory operations on CPUs and GPUs alike,<sup>4,14</sup> in some cases hitting less than 1% of peak performance. This leads applications to want to offload such work to near-memory processors. In the case of far memory, that near-memory processor is itself outside the translation context of the parent process of the pointer-rich data. Pointers therefore must

make sense everywhere in these new heterogeneous disaggregated systems.

In order to lower infrastructure rent, cloud applications also wish to exploit disaggregated far memory as a fungible memory resource that can grow and shrink with the amount of data. Moreover, they want to independently scale their memory and compute resources. For example, database services want to flex compute up or down in proportion to query load. Pointer-rich data in far memory must be shareable at low overhead between existing and new compute instances.

### Prior Work on Far Memory

Pointers in traditional operating systems were valid only in the memory space of the process in which they were created. Sharing pointer-rich data among processes, nodes, and devices therefore required serialization-deserialization. This limitation remained even when prior art was recently extended

by taking an approach of tombstoning dangling references to data demoted to far memory using special pointers.<sup>7,16</sup> Those pointers could be dereferenced only from the original context of data creation, precluding independent scaling of memory and computation.

Global address spaces, such as partitioned global address space (PGAS), support a limited form of global pointers that persist only for the life time of a set of processes across multiple nodes. Nonvolatile memory (NVM) libraries such as Persistent Memory Development Kit (PMDK) support object-based pointers, but their “large” storage-format pointers are more than 64-bits long, and their traversal cannot be offloaded.

Commercial virtualization frameworks such as VMware’s Nu proclents<sup>13</sup> can only maintain the illusion of global pointers by compromising security (for example, by turning address space layout randomization off).

Microsoft CompuCache<sup>14</sup> also supported global pointers, but by using a heavy database runtime atop full VMs even on disaggregated memory devices. All pointers, whether at hosts or in the CompuCache, are VM-local only. Pointer chasing across devices requires repeatedly returning to the host.

Teleport<sup>15</sup> supported pointer-chasing offload to remote memory but by directed, on-demand shipping of the virtual-to-physical translation context to the target locale of each function shipped.

Prior work on OS constructs for far memory is therefore missing a foundation of globally invariant pointers that can be shared with and dereferenced by any node or device in a cluster containing far memory.

### Invariant Pointers

When organizing data at object granularity, a globally invariant pointer must contain the ID of the object containing the target data, as well as an offset to that data. This object ID must be interpreted anywhere the pointer can be dereferenced. Ideally, invariant pointers should be no larger than 64 bits and permit access to partially resident objects. Existing approaches do not meet the first criterion (for example, PMDK) or the second criterion (for example, application-integrated far memory, AIFM,<sup>12</sup> which has a different pointer form for resident and nonresident ob-

jects). Providing truly globally invariant pointers, however, is necessary for offloading “run anywhere” code.

Twizzler<sup>3</sup> is an operating system that introduces globally invariant pointers by using a context local to the object in which the pointer is stored, shortening its representation while allowing any CPU that can read the pointer to fully resolve its destination. This is done using an foreign object table (FOT) that is part of each object in the system, ensuring any individual object is self-contained.

An object’s FOT contains identifying information for each foreign object that is the destination for a pointer in the object. Since these are stored in an ordered table, stored pointers use the index into the FOT as a stand-in for the full addressing information, a translation process shown in Figure 2. This approach allows pointers to remain small: a 64-bit pointer can, for example, include a 24-bit “local” object ID and 40-bit offset. While this limits the number of foreign objects that can be referenced from a single object to  $2^{24}$ , different objects have their own FOTs and can reference a different set of objects, so the total number of objects in the system is limited only by the size of an object ID.

This approach also allows for a wide range of resolvers that translate identifying information in the FOT into an object ID. For example, the FOT might contain a static object ID or the equivalent of a file-system name to be resolved to an object ID by a name resolver. There is no requirement that a name resolve to the same object ID in different places: for example, an object named /var/log/syslog might resolve to different object IDs on different system nodes. Name resolvers

themselves can be pluggable: The FOT need only identify the resolver in a way that any node in the system can run the resolver to return an object ID.

While the first access to a foreign object may be relatively slow, subsequent accesses are very fast, since the resolution to an object ID is cached. The system maps the object into the node’s “guest physical” address space, leveraging memory management unit (MMU) hardware already in use for virtualization. It then maps the guest physical space in which the object resides into the guest virtual space for any processes that reference the object, using extended page tables to remove software from the CPU load/store path and allowing the system to run at memory speed. This is necessary for efficiency; even minimal system software interaction on each load and store will slow the computation significantly.

Preliminary experiments<sup>3</sup> show that Twizzler’s approach is effective at preserving low-latency pointer dereferencing for both intra-object and inter-object invariant pointers. On an Intel Xeon Gold CPU running at 2.3GHz, intra-object pointer dereferences take about 0.4ns, approximately the same time as “normal” dereferences. Cached inter-object pointer dereferences take 3.2ns, somewhat slower than intra-object dereferences but still sufficiently fast because relatively few such references are expected, given multi-megabyte objects. The first reference to a foreign object is slower, at 28ns, but still reasonable. If name resolution is more complex than interpreting a static full-length (128-bit) object ID, it would be longer still; however, these penalties are paid only once, regardless of how many

times pointers from object A to object B are dereferenced in the same process.

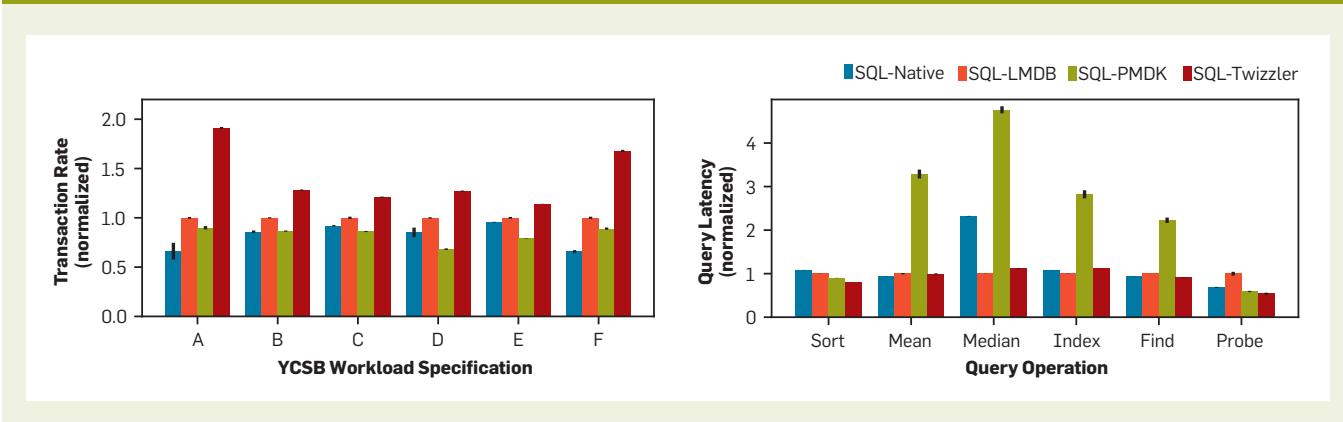
Benchmarks on both microscale (in-memory key/value store) and macroscale (Yahoo! Cloud Serving Benchmark, YCSB, using different back ends) likewise show excellent performance for this approach. The left graph in Figure 3 shows throughput of the YCSB benchmark on SQLite using four back ends: the native SQLite back end; the Lightning Memory-mapped Database (LMDB) backend, which leverages mmap; our implementation of a PMDK backend, which uses a red-black tree under PMDK; and Twizzler, which uses a red-black tree with the invariant pointer approach.

The invariant pointer approach outperforms every other approach while providing the flexibility of “run anywhere” invariant pointers. The graph on the right of Figure 3 similarly shows that these invariant pointers provide lower latency than other approaches because of the simplicity of the programming model and the low overhead for dereferencing pointers. PMDK, in particular, is significantly slower because its pointers are 128-bits long, requiring additional register space and memory operations to read and dereference.

It is important to note the PMDK and Twizzler implementations are running the same back-end code, with changes made only to accommodate the different programming models; this shows the benefit of using 64-bit pointers local to an object context rather than 128-bit pointers, as PMDK does.

Elephance MemOS is a fork of Twizzler being developed to run on CXL far memory devices. It will be ported and optimized for the systems-on-chip

**Figure 3. Performance and latency of the YCSB benchmark.**



(SoCs) used as controllers in CXL-disaggregated memory nodes.

### Programming with Memory Objects and Invariant Pointers

For software developers, what does memory disaggregation mean and how will systems be built around it? The architecture of such systems will aim to hide the details from the majority of programmers, so their code will not need to change to run on these new systems.

There are three ways in which systems can be built to provide disaggregated memory: application libraries, modification to the operating system's memory system, and changes beneath the oper-

ating system at the hardware layers, as seen in Figure 4. In the figure, a set of application servers is connected to a set of MemOS nodes over a shared bus. Pointer-rich application data in far memory lives on MemOS nodes. Pointers can be: inter-object and on the same device, inter-object across devices, or intra-object.

It is likely the first way that disaggregated memory will be made available will be through application libraries linked directly into the application, seen at the top of Figure 5. The memory shim acts as a specialized memory allocator that knows how to handle remote memory using a memory access protocol (MAP). The MAP may depend on a

current technology such as RDMA (remote direct memory access (RDMA), or may be something newer such as CXL3.

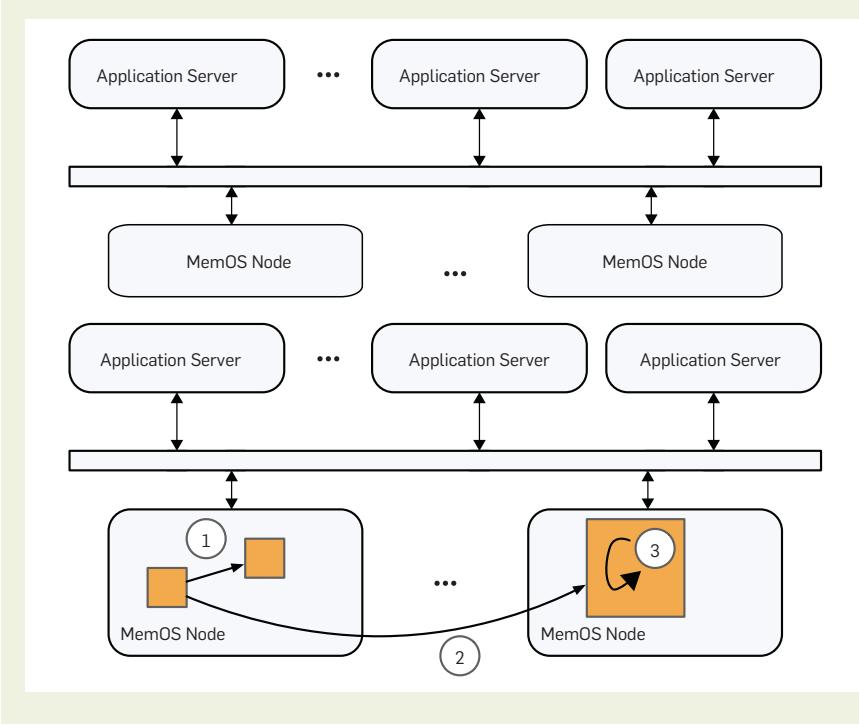
Many languages, such as Python, which depend on the C library for memory, will be able to use the memory shim to handle memory for objects in the language, freeing the Python programmer from having to know anything about disaggregated memory. For languages such as C and C++, which handle pointers directly, the programmer will have to work with the memory shim APIs in order to manage remote memory. The prevalence of Python and similar managed memory languages in big data and machine-learning applications means that programmers in those fields can use disaggregated memory in a transparent way, no matter where the memory shim is located in the software stack.

Extending the operating system's virtual memory system to integrate with the memory shim is the next logical place to interpose disaggregated memory in the stack, seen in (B) in Figure 5. Again, the specific MAP is not exposed to the kernel developer, only the memory shim APIs. The Linux operating system already has heterogeneous memory management (HMM),<sup>8</sup> which is a natural place to slot in the memory shim. Once the shim is integrated into the operating system itself, all applications can use disaggregated memory transparently without modifications to their source code or linking with specialized libraries.

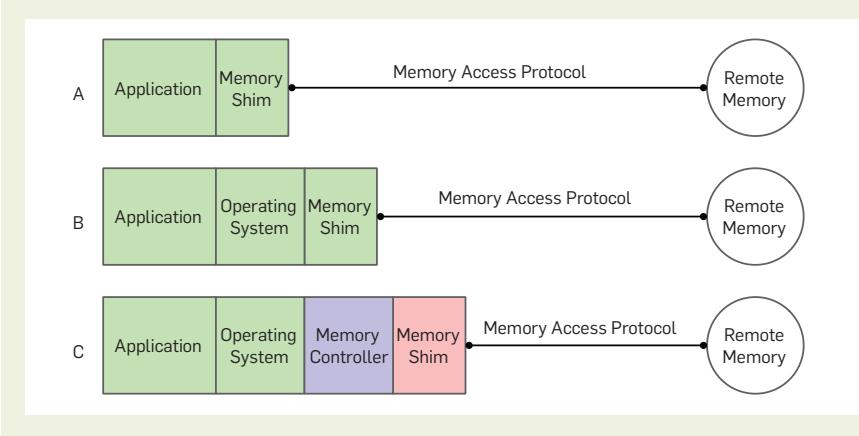
The deepest that far memory can be placed in the stack is in the hardware itself. Memory controllers integrated in CPUs from Intel and AMD are already starting to support early versions of CXL disaggregated memory. In the future, more featureful controllers will present memory to the operating system both locally and remotely in a transparent manner but, like the other two cases, will require a MAP to be interposed between the hardware and the remote memory. The protocol in this instance will be CXL 3.

While putting the memory shim into hardware will likely result in the highest bandwidth, lowest latency, and maximum portability, there are reasons to continue to use a memory shim as a linked library into the software. First and foremost is the level of control that linking directly to the memory shim gives to the programmer. Once such

**Figure 4. Example MemOS deployment.**



**Figure 5. Extending the OS's virtual memory system.**



functionality is embedded into the operating system or the memory controller, application programmers will lose control and visibility into the remote memory system. While many will be happy not to have to manage memory on their own, applications will remain where such control is a feature. Novel memory architectures for distributed memory must first be tried in software, and some may be too specialized ever to be implemented in hardware.

Consider a memory system where pointers are globally invariant, which will be possible with MemOS but is not yet common in pointer-based systems. Building and debugging such a system in software makes it possible to rapidly iterate on the design—impossible in a memory controller and certainly more difficult to debug in the operating system. Applications that can use globally invariant pointers have distinct advantages because computation can take place on any node without the application having to know where a pointer might reside. Furthermore, it will be possible to move code, rather than data, to achieve computational efficiency—again, because no matter which compute node a pointer resides on, the pointer itself is the global handle that computation depends on, rather than an address in local memory, as things stand today.

### Key Takeaways

Effectively exploiting emerging far-memory technology requires consideration of operating on richly connected data outside the context of the parent process. Operating-system technology in development offers help by exposing abstractions such as memory objects and globally invariant pointers that can be traversed by devices and newly instantiated compute. Such ideas will allow applications running on future heterogeneous distributed systems with disaggregated memory nodes to exploit near-memory processing for higher performance and to independently scale their memory and compute resources for lower cost. 

### References

- Al Maruf, H. et al. TPP: Transparent page placement for CXL-enabled tiered-memory. In *Proceedings of the 28<sup>th</sup> ACM Intern. Conf. Architectural Support for Programming Languages and Operating Systems 3*, (2023), 742–755; <https://dl.acm.org/doi/10.1145/3582016.3582063>.
- Berger, D.S. et al. Design trade-offs in CXL-based memory pools for public cloud platforms. *IEEE Micro* 43, 2 (2023), 30–38; <https://dl.acm.org/doi/abs/10.1109/MM.2023.3241586>.
- Bittman, D. et al. Twizzler: a data-centric OS for non-volatile memory. In *Proceedings of the 2020 Usenix Annual Tech. Conf.*; <https://dl.acm.org/doi/pdf/10.5555/3489146.3489151>.
- Dongarra, J. A not so simple matter of software. *ACM Turing Award Lecture* (2021); <https://www.youtube.com/watch?v=cSOOTc2w5Dg>.
- Duraisamy, P. et al. Towards an adaptable systems architecture for memory tiering at warehouse-scale. In *Proceedings of the 28<sup>th</sup> ACM Intern. Conf. Architectural Support for Programming Languages and Operating Systems 3*, (2023), 727–741; <https://dl.acm.org/doi/10.1145/3582016.3582031>.
- Hsieh, K. et al. Accelerating pointer chasing in 3D-stacked memory: challenges, mechanisms, evaluation. In *Proceedings of the IEEE 34<sup>th</sup> Intern. Conf. Computer Design*. (2016), 25–32; <https://ieeexplore.ieee.org/document/7753257>.
- Jennings, S. The zswap compressed swap cache. *LWNNet*. (2013); <https://lwn.net/Articles/537422/>.
- Kernel Development Community. Heterogeneous memory management. *The Linux Kernel 5.0.0*; <https://www.kernel.org/doc/html/v5.0/vm/hmm.html>.
- Mehra, P. and Coughlin, T. Taming memory with disaggregation. *Computer* 55, 9 (2022), 94–98; <https://ieeexplore.ieee.org/document/9869614>.
- Michelogiannakis, G. et al. A case for intra-rack resource disaggregation in HPC. *ACM Trans Architecture and Code Optimizations* 19, 2 (2022), 1–26; <https://dl.acm.org/doi/10.1145/3514245>.
- Rodrigues, A. et al. Towards a scatter-gather architecture: hardware and software issues. In *Proceedings of the 2019 Intern. Symp. Memory Systems*. 261–271; <https://dl.acm.org/doi/10.1145/3357526.3357571>.
- Ruan, Z. et al. AIM: High-performance, application-integrated far memory. In *Proceedings of the 14<sup>th</sup> Usenix Symp. Operating Systems Design and Implementation*. (2023); <https://dl.acm.org/doi/pdf/10.5555/3488766.3488784>.
- Ruan, Z. et al. Nu: Achieving microsecond-scale resource fungibility with logical processes. In *Proceedings of the 20<sup>th</sup> Usenix Symp. Networked Systems Design and Implementation*. (2023); <https://www.usenix.org/system/files/nsdi23-ruan.pdf>.
- Zhang, Q. et al. CompuCache: Remote computable caching using spot VMs. In *Proceedings of the 12<sup>th</sup> Annual Conf. Innovative Data Systems Research*. (2022); <https://www.cidrdb.org/cidr2022/papers/p31-zhang.pdf>.
- Zhang, Q. et al. Optimizing data-intensive systems in disaggregated data centers with Teleport. In *Proceedings of the 2022 Intern. Conf. Management of Data*; <https://dl.acm.org/doi/10.1145/3514221.3517856>.
- Zhou, Y. et al. Carbink: Fault-tolerant far memory. In *Proceedings of the 18<sup>th</sup> Usenix Symp. Operating Systems Design and Implementation*. (2022); <https://www.usenix.org/system/files/osdi22-zhou-yang.pdf>.

**Ethan Miller** has been a member of the technical staff at Pure Storage since 2009. He is a professor emeritus in computer science and engineering at the University of California, Santa Cruz, CA, USA, where he held the Veritas Presidential Chair in Storage; founded the Center for Research in Storage Systems; and led the development of the Ceph distributed file system, the Twizzler operating system.

**Achilles Benetopoulos** is a Ph.D. student at the University of California, Santa Cruz, CA, USA, working at the intersection of distributed systems, databases, and programming languages.

**George Neville-Neil** works on networking and operating system code as well as teaches courses on various subjects related to programming. He is the author of *The Collected Kode Vicious* and a co-author of *The Design and Implementation of the FreeBSD Operating System*.

**Pankaj Mehra** is founder of Elephance Memory. He was previously Vice President of Product Planning at Samsung, a senior fellow at SanDisk and Western Digital, and a distinguished technologist at Hewlett-Packard.

**Daniel Bittman** is a co-founder of Elephance Memory and principal maintainer of the Twizzler operating system.

Copyright held by owner(s)/author(s).



## Replacing a sledgehammer with a scalpel.

BY SAMUEL W. STARK, A. THEODORE MARKETTOS,  
AND SIMON W. MOORE

# How Flexible Is CXL's Memory Protection?

IN THE BEGINNING, there was PCIe. Well, really there was PCI and PCI-X, which were superseded by PCIe in 2003, and many others, such as ISA and VME, before them, but PCIe is a superset of them all. They are all interconnects, allowing a *host* (for example, the main system CPU) to configure and manipulate connected peripheral *devices* and map their memory into a shared address space.

In time, computations became bigger and more complicated, and peripheral devices became whole

systems unto themselves. Graphics processing units (GPUs) are the best example, going from hardwired graphics offload devices to full-blown general-purpose processors that cooperate and communicate with the host to solve problems.

Cooperative processing between the host and device is complicated by PCIe's lack of coherent memory sharing. When CPU cores share memory, they use a *cache coherency* protocol to ensure they can have a fast local copy (a *cache*) while keeping a coherent view of memory—even when other cores write to it. PCIe does not support this kind of sharing; it only allows block transfers between host and device. Various companies created successor protocols—CCIX, OpenCAPI, and Gen-Z—to support this, but they have all expired or been subsumed by Intel's Compute Express Link (CXL).

CXL provides new protocols on top of PCIe for accelerator devices to cache host memory (CXL.cache) and for hosts to cache device memory (CXL.mem). The industry is currently focused on CXL.mem memory expansion devices. The first CXL-compatible CPUs (released in November 2022) support "CXL 1.1+ for memory expansion,"<sup>1</sup> and CXL accelerators have not been announced—only CXL.mem devices, such as Samsung's 512GB RAM expansion.<sup>19</sup> CXL 3.0, released in August 2022, adds support for fabric topologies connecting many hosts to many shared Global Fabric Attached Memory (GFAM) devices. This facilitates *disaggregated memory*, where an arbitrary number of endpoints connected in an arbitrary topology can request, use, and coherently share arbitrary amounts of memory.

If disaggregated memory is the future, our biggest question is that of protection. With so many endpoints all connecting to and sharing the same memory, how can they be restricted to accessing only the memory they need? They may be running untrusted software or untrusted



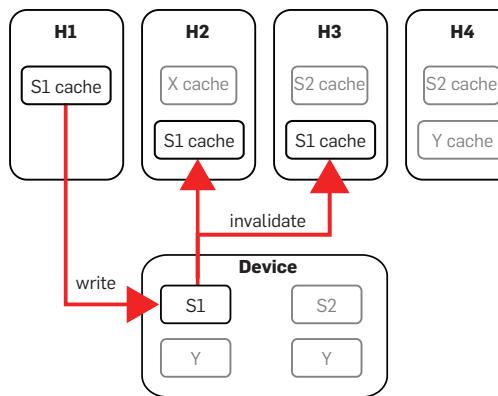
hardware. How can memory protection work in this threat environment? The Capability Hardware Enhanced RISC Instructions (CHERI) project has shown that *architectural capabilities* can provide flexible, fine-grained memory protection.<sup>21</sup> How does CXL's current memory protection compare? Could a capability system work in CXL's distributed setting with malicious actors? To start, let's examine CXL's protection mechanisms and see how well they handle real-world security problems.

### Protection Systems

In most cases, software uses physical resources through multiple layers of abstraction. Each layer translates incoming requests to a format expected by the next layer down, and can also provide protection. A simple example is the memory management unit (MMU), which translates memory requests from virtual to physical memory.<sup>13</sup> The OS gives each process a different mapping of virtual to physical addresses, and the MMU ensures processes can access only the physical memory that

the OS has mapped in. To generalize, protection systems ensure *actors* can only access valid *resources*.

The protection a system can provide is limited by the granularity of its actors and resources; therefore, protection at multiple layers of abstraction is important. For example, the MMU only has insight at the process level. The software inside the process has tighter definitions of *valid* (for example, “I will not access out-of-bounds array elements”) that the MMU doesn’t understand (it doesn’t know or care where

**Figure 1. Coherently sharing device memory with many hosts.**

the array is) and thus cannot help with. Instead, another layer can be added above the MMU, such as a language runtime (JVM, .NET) or hardware-based checks (CHERI<sup>21</sup>), which have more information and ensure validity at a finer-grained level.

Different levels of abstraction can add different sets of actors and resources. For example, an operating system is responsible for ensuring its processes access files correctly—and for actually performing those accesses through the file system driver. If those files are on a networked file system, the server may have to handle multiple clients at once and check that they access files correctly. The individual OS does not know about the other clients, and the server does not know about the processes running inside the OS, so having protections and checks at both levels is necessary.

**CXL and the flaws therein.** CXL, like PCIe, uses a host-device model. Each CXL host controls a set of connected peripheral devices, and maps all the memory they expose into an host physical address space (HPA). The host may

also map its own memory into the HPA, and accelerator devices like GPUs can access it over CXL.cache, but current devices just expose RAM to the host over CXL.mem. CXL 3.0 upgraded CXL.mem to allow hosts to share memory regions through both multi-headed and GFAM devices.

Multi-headed CXL.mem devices connect to multiple hosts and can map the same regions of physical memory into all of their HPAs at the same time. Those hosts can all cache parts of those regions, and the device is responsible for ensuring coherency (see Figure 1). For example, if host 1 tries to write to a cache line in region A, the device realizes that hosts 2 and 3 share A and tells them to invalidate that cache line. Unfortunately, each of those hosts can only access 16 regions<sup>8</sup> (Sec 2.5), so they will necessarily be large—on the order of gigabytes or hundreds of megabytes.

GFAM devices take this a step further by not being attached to specific hosts. Any host can map GFAM memory into its HPA, and any endpoint (host or device) in that HPA can talk to the GFAM directly and access that memory. The GFAM is configured with separate translation tables for each endpoint, so each endpoint can access eight regions of physical memory<sup>8</sup> (Sec 7.7.2.4). These regions may overlap, allowing memory sharing, or they may be isolated. As shown in Figure 2, 10GiB of GFAM is mapped, but the host and accelerator are configured so they see only 6GiB each, with a 2GiB shared region. Again, because each endpoint has few ranges, they will be large. *Memory groups*<sup>8</sup> (Sec 7.7.2.5) can punch

holes in these ranges and hide specific blocks, basting 64MB<sup>8</sup> (Table 7-67, min. block size).

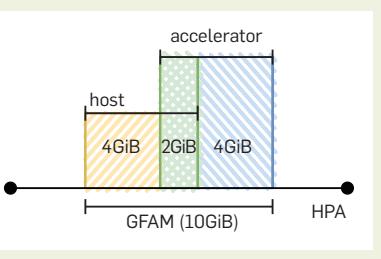
Both kinds of memory provide protection through nonexhaustive translation: Endpoints request addresses in their HPA, which get translated to local device addresses, and that translation may fail (that is, the endpoint may not have memory mapped at that address). These mechanisms, similar to an MMU, provide inflexible coarse-grained protection. At most, each endpoint can access 16 memory ranges per device. The only way to change the mappings and transfer access rights is to convince the *Fabric Manager*, which has no defined interface for this<sup>8</sup> (Sec. 7.6.1).

CXL 3.0 also introduced *Unordered I/O requests*, which allow accelerators to access other devices' memory, but there is no standardized way to protect those accesses. It may be possible to prevent specific devices from interacting at all (for example, through PCIe Access Control Services) or add MMU-like protection (for example, through PCIe Address Translation Services) but these, like CXL's other protection models, are inflexible and coarse-grained.

CXL's protection isn't great. Endpoints can be configured to access and share large memory regions, but cannot share many small ones. Endpoints can't grant each other access to memory; they have to go through an intermediary. Device-to-device access has to rely on vendor-defined protection, if any. How does that stack up against real-world threats?

**Threats in the datacenter.** First, we can understand the datacenter threat model from a whitepaper published in November 2022 by Amazon Web Services (AWS) about their Nitro platform.<sup>2</sup> Cloud systems have to run workloads from many clients, who don't trust each other, on the same hardware. Before Nitro, AWS would run all client workloads as virtual machines (VMs) atop a *hypervisor*, which exposed isolated virtualized resources to each VM. For example, the hypervisor would implement a software model of a network card for each VM, so it could control which networks the VMs could access. The key impact of the Nitro system is moving this virtualization out of the hypervisor and into the hardware.

Each Nitro system is controlled by a

**Figure 2. Endpoints in the same address space (HPA) using overlapping regions of global memory (GFAM).**

custom *Nitro Controller* PCIe card. This is the hardware root of trust, responsible for configuring the *System Main Board* (that is, the CPU, motherboard, and RAM) and other peripherals before running client workloads. Networks and storage are accessed through other AWS-designed Nitro PCIe cards, which the Nitro Controller can split into Virtual Functions using PCIe single-root I/O virtualization (SR-IOV)<sup>16</sup> to provide isolated resources for each VM.

When running many VMs, a minimal hypervisor is still necessary to configure the MMU and link each VM to its dedicated virtual functions. A Nitro system can also run *bare-metal* (a single client workload without a hypervisor). Even though the client workload is untrusted, the Nitro cards still virtualize access to networks and storage.

AWS trusts the Nitro controller to bring up the system, the Nitro cards to virtualize networks/storage, and the hypervisor/MMU to enforce isolation between VMs. Client workloads cannot be trusted, and if they're running bare-metal, then any communication from the System Main Board cannot be trusted either. From CXL's perspective, this means a host could be malicious (running bare-metal) or be responsible for many malicious workloads (running VMs). In the latter case, CXL does not have any constructs that can help with virtualization. In fact, CXL does not consider virtualization at all—literally, *virtualization* and similar terms are not in the specification.

Datacenters have further complications. Accelerator devices, such as GPUs, sometimes rely on directly sharing memory for high performance. Nvidia's Magnum I/O APIs<sup>18</sup> allow GPUs to directly access NVMe storage devices (GPUDirect Storage), share memory with other GPUs (NVSHMEM), and expose their memory to other peripherals (GPUDirect RDMA), including InfiniBand adapters (nvidia-peermem).

While some GPUs nominally support virtualization through SR-IOV, AWS does not take advantage of this—client workloads are given whole numbers of GPUs and control them directly (clients even control the GPU drivers<sup>3</sup>). This expands the threat model. Not only are GFAMs sharing memory across HPAs, but also individual devices (including accelerators) may expose their

memory to endpoints controlled by malicious clients.

CXL does not handle this use case. It implicitly assumes that hosts and devices are trustworthy. Hosts may be trustworthy if they have, for example, a hypervisor keeping them in check, and devices may be carefully chosen for trust, but if any device or host is untrustworthy (for example, running bare-metal client workloads) better protection is needed.

**Threats in the consumer space.** The threat of malicious devices is not exclusive to the datacenter—in fact, it's much worse for consumers! Desktops and laptops have a plethora of external ports for connecting arbitrary hardware, including high-performance accelerators, such as external GPUs. Accelerators take advantage of high-speed Thunderbolt connections that wrap PCIe, giving external hardware access to the internal PCIe memory map. Attacks on PCIe-based systems through Thunderbolt have already been demonstrated,<sup>15</sup> showing that malicious hardware can access sensitive memory intended for other devices, even with protections such as IOMMU enabled.

Worse, direct device-to-device memory accesses are making their way to consumer systems as well. Modern game consoles depend on high-speed transfers from storage to GPU-accessible memory, and Microsoft's DirectStorage API brings this closer to reality on PCs. While at the time of writing, it still copied data through a buffer in system RAM, it seems inevitable that high-performance rendering systems (for example, games and video editing) will eventually take advantage of direct access—especially because it's already possible in the datacenter.

CXL is coming to the consumer market, so it needs to handle this. In an AMD "Meet the Experts" webinar<sup>4</sup> in October 2022, an AMD representative said it might come to consumer devices within five years, initially with a focus on connecting persistent storage and RAM. Loading from persistent storage is currently the big use case for device-to-device transfer, so CXL needs to consider malicious devices sooner rather than later.

Today, CXL's memory protection is inflexible at best. It is capable of isolating endpoints in large memory re-

gions, but not much more than that. It has no capacity for virtualization for workloads running on the same endpoint, and cannot protect devices from each other.

## Capability-Based Protection for CXL

CHERI<sup>21</sup> is a capability-based protection system that has proven useful both for flexible, fine-grained (tens of bytes) memory protection and for *compartmentalization*, by sandboxing programs and libraries from each other.<sup>22</sup> This seems to address all of CXL's security issues. Could CXL adopt a capability-based system?

*Capabilities* are unforgeable tokens that encode the authority to access a resource. Given a capability, an actor can access the resource, derive new capabilities for that resource with reduced permissions, transfer them to other actors, and potentially revoke them if those actors no longer need access. Because access rights are encoded directly in the token, capabilities are very flexible: It's easy to derive new capabilities with extremely specific access rights for new situations. Deriving lots of capabilities does have a downside: Revoking a capability—recursively deleting all derivations—can be more difficult. Let's examine a few examples.

**Central-trust systems.** Capabilities must be unforgeable. When a capability is used, the system needs some way to verify it has not been forged. The simplest way to enforce this is to store all capabilities and perform all capability modifications in a centralized trust base, or a *central-trust* system.

One example is FreeBSD Capsicum,<sup>20</sup> which protects files from processes by replacing Unix file descriptors with capabilities. A process can open the files it needs, limit its access with more granular permissions, and then enter *capability mode* to sandbox itself with those files. Like file descriptors, capabilities are stored in tables in OS memory. Userspace programs have to use syscalls to ask the OS to manipulate them, instead of creating or modifying them directly. The OS trusts itself to correctly modify capabilities (for example, never adding permissions, only taking them away) so capabilities cannot be forged. Although Capsicum doesn't perform revocation, in principle it would simply require

searching the tables or even tracking parentage in capability metadata. This provides better security than plain Unix, but syscalls and context-switching to the OS can be slow.

CHERI takes a different approach. Instead of implementing the trusted base in software, CHERI implements it in the hardware and adds machine instructions for fast capability manipulation. CHERI replaces pointers with capabilities—fat pointers that include the range of addresses the pointer may point to. This range ensures pointers stay within their original provenances<sup>6,12</sup> and can be limited further (for example, you can allocate an array, derive a capability for one element, and pass that to a function without exposing the rest of the array).

Registers and memory use tag bits to mark valid capabilities, and the hardware controls the tag bits to prevent forgery. Because all pointers have this metadata, including code pointers, even the smallest software components (for example, individual functions) can be sandboxed with just the memory they need. Larger libraries, even ones compiled without CHERI support, can also be sandboxed using compartments.<sup>22</sup> The cost of storing capabilities anywhere is that revocation must search everywhere,<sup>23</sup> although the overheads are lower than you might expect.<sup>10</sup> CHERI ensures logical software components can access only the virtual memory ranges to which they have been explicitly given permission.

How could this help CXL? Eagle-eyed readers might notice GFAM already uses a system similar to Capsicum—each endpoint (that is, actor) has up to eight translation table entries (that is, capabilities) that grant access to memory. This demonstrates the flaws with a centralized system in this context: The number of capabilities (and implicitly their granularity) can be limited by hardware resources. This is more suitable for protecting host memory from a limited number of devices, for example,<sup>14</sup> but GFAM tries to track all capabilities granted to *thousands* of actors. To alleviate this, one could store the capabilities in the memory exposed over CXL.mem or give each endpoint some dedicated capability memory, such as Capsicum's tables. Both cases would require trust to be distributed among the endpoints.

## The threat of malicious devices is not exclusive to the datacenter—in fact, it's much worse for consumers!

**Distributed-trust systems.** Barrelfish<sup>5,17</sup> and SemperOS<sup>11</sup> are distributed operating systems, implemented as separate instances running on separate cores and communicating with message passing. Barreelfish uses capabilities to protect OS resources, such as message passing and threading primitives, physical memory ranges, etc. SemperOS uses capabilities for an in-memory file system.

The trusted base for capability operations is distributed across the OS cores but aims to provide identical semantics to central-trust systems. Most importantly, any core can derive from a capability in any other core, and thus revocation may need to touch all cores. This requires all actors to trust each other. It is more complicated to reason about than central-trust systems, but it scales better—particularly if cross-actor operations are uncommon.

For CXL, this may be suitable if all endpoints are trustworthy. If, for example, all endpoints in a datacenter use CHERI-like hardware to manipulate capabilities, this could work. At scale, however, revocation may become a bigger issue, and CXL cannot rely solely on this model anyway—the threat of malicious endpoints is too great.

**Decentralized systems.** Even if calling out to a centralized trusted base to manipulate capabilities is impossible or impractical, and the actors cannot be trusted to manipulate capabilities correctly, there is still hope. Decentralized capabilities, such as *macaroons*,<sup>7</sup> can be passed to untrusted actors and have their validity checked when those actors try to use them.

Macaroons provide access to a resource that is reduced through an append-only list of *caveats*. A macaroon begins with an identifier, such as “access transaction details,” and a signature, made by hashing the identifier with a secret key. When a caveat (such as “for Alice’s account,” or “until 5PM EST”) is added, that caveat is hashed with the current signature to make a new signature. The old signature is thrown away and cannot be reconstructed—the hash cannot be undone. Given a macaroon with a set of caveats, it’s impossible to remove a caveat and recalculate the correct signature without the secret key. Therefore, it’s impossible for a hostile user to forge a macaroon with fewer ca-

veats (that is, more permissions).

Decentralized capabilities have not yet been integrated into low-level software or hardware. Macaroons were originally designed for the Web, so they have a text-based wire format and third-party authentication features, which a binary-based interconnect doesn't need. This is fine for the network layer (for example, Michael Dodson combined macaroons with CHERI for fine-grained memory-mapped I/O access over an insecure network<sup>9</sup>), but domain-optimized representations would be more space-efficient.

Revocation is also interesting. Capabilities could come with timeout caveats and require refreshing, or groups of capabilities (and all their derivations) could be revoked by throwing away their secret key. This would allow CXL endpoints to store and (attempt to) manipulate their capabilities themselves, and let the CXL.mem device revoke them, all without trusting them. Decentralized capabilities are robust to hostile actors, do not require centralized resources, and are ripe for further investigation.

## Conclusion

Physical memory is accessed through many layers of abstraction. Applying protection at different layers, which are aware of different actors and use resources at different granularities, is essential. CHERI and MMUs offer great protection at the software and process level, but CXL's protection model has issues. It allows memory sharing, but only of a few large ranges instead of many small ones. It doesn't give actors a way to share new memory ranges with each other, instead relying on a central, underspecified Fabric Manager. Capabilities are inherently flexible—they can protect large and small memory ranges, and can be transferred directly between actors without a centralized authority—so they should be able to address these problems.

CXL initially targets the datacenter, with many endpoints sharing disaggregated memory. The protection is coarse-grained, and does not consider virtualization. VMs running on the same host have to rely on similarly coarse hypervisor- and MMU-based isolation. Fine-grained capabilities could allow individual VMs to share small memory regions directly. Capabilities

for large memory regions could also enforce VM compartmentalization at the CXL layer, similarly to CHERI.

In datacenter and consumer systems, device-to-device memory sharing is becoming essential for high performance. CXL does not try to protect devices from each other at all, which is especially scary considering how powerful malicious devices already can be. Capabilities would provide a consistent interface for securely exposing regions of device memory. Decentralized capabilities are robust against malicious actors and could keep the peace in the Wild West of untrustworthy hardware. In a datacenter with trusted components, distributed-trust systems could even forgo the cryptography associated with decentralized capabilities for lower overheads.

Decentralized and distributed capabilities have a lot of potential, but they have not been used in this context yet and need further investigation. Even so, they could greatly benefit CXL, which is a new interconnect standard that provides the opportunity to build in better security from the start instead of retrofitting it. A domain-optimized decentralized capability system could work wonders, giving CXL fine-grained memory sharing and improving virtualization and device-to-device security. Interconnects must take security more seriously, and we believe capabilities can provide flexible and robust security for CXL and beyond.

## Acknowledgments

We thank the CHERI project team led by Robert Watson for demonstrating the potential of capabilities for memory security, without which this work could not exist. The CHERI team also provided essential feedback while developing this article, for which we are extremely thankful. This work was supported by the University of Cambridge Harding Distinguished Postgraduate Scholars Programme, and by EPSRC grant EP/V000381/1 (CAPcelerate). □

## References

- Advanced Micro Devices, Inc. Offering unmatched performance, leadership energy efficiency, and next-generation architecture, AMD brings 4<sup>th</sup> gen AMD EPYC processors to the modern data center. (2022); <https://bit.ly/3KOU6Ou>
- Amazon Web Services The security design of the AWS Nitro System. *AWS whitepaper*. (2023); <https://go.aws/3P6U0dm>
- Amazon Web Services; <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/install-nvidia-driver.html>
- AMD Meet the Experts Webinars. How AM5, DDR5 memory, and PCIe 5.0 support pave the way for next-gen gaming experiences; <https://bit.ly/3QKm42J>.
- Baumann, A. et al. The Multikernel: A new OS architecture for scalable multicore systems. In *Proceedings of the ACM SIGOPS 22<sup>nd</sup> Symp. Operating Systems Principles*. (2009), 29–44; 10.1145/1629575.1629579.
- Beingessner, A. Rust's unsafe pointer types need an overhaul. *Faultlore*. (2022); <https://faultlore.com/blah/fix-rust-pointers/>
- Birgisson, A. et al. Macaroons: cookies with contextual caveats for decentralized authorization in the cloud. In *Proceeding of the 2014 Network and Distributed System Security Symp.*; <https://bit.ly/45EKxe4>.
- CXL Consortium Compute Express Link (CXL) Specification, revision 3.0, version 1.0. (2022); <https://www.computeexpresslink.org/download-the-specification>.
- Dodson, M.G. Capability-Based Access Control for Cyber Physical Systems. *Thesis*. University of Cambridge, Computer Laboratory. (2021).
- Filardo, N.W. Cornucopia: Temporal safety for CHERI heaps. In *IEEE Symp. on Security and Privacy*. (2020), 608–625. <https://ieeexplore.ieee.org/document/9152640>
- Hille, M. et al. SemperOS: A distributed capability system. In *Proceedings of the 2019 Usenix Annual Technical Conf.* 709–722; <https://www.usenix.org/conference/atc19/presentation/hille>.
- Jung, R. Pointers are complicated, or: What's in a byte? Ralf's rambblings. (2018); <https://www.ralfj.de/blog/2018/07/24/pointers-and-bytes.html>
- Kernel Development Community. Virtual memory primer; <https://bit.ly/3YSr5Z1>.
- Markettos, A.T. et al. Position paper: Defending direct memory access with CHERI capabilities. In *Hardware and Architectural Support for Security and Privacy 7*. (2020), 1–9; 10.1145/3458903.3458910
- Markettos, A.T. et al. Thunderclap: Exploring vulnerabilities in operating system IOMMU protection via DMA from untrustworthy peripherals. In *Proceedings of the 2019 Network and Distributed System Security Symp.*; 10/gjh62d.
- Microsoft. An introduction to single root I/O virtualization (SR-IOV). (2021); <https://bit.ly/3QRxvWw>
- Nevill, M. An Evaluation of Capabilities for a Multikernel. *Master's thesis*. ETH Zurich. (2012).
- Nvidia Magnum IO; <https://www.nvidia.com/en-us/data-center/magnum-io/>.
- Samsung Semiconductor. Samsung Electronics introduces industry's first 512GB CXL memory module. (2022); <https://news.samsung.com/global/samsung-electronics-introduces-industrys-first-512gb-cxl-memory-module>.
- Watson, R.N.M. et al. Capsicum: Practical capabilities for Unix. *Commun. ACM* 55, 3 (Mar. 2012), 97–104; 10.1145/2093548.2093572
- Watson, R.N.M. et al. An introduction to CHERI. University of Cambridge. (2019), 43.
- Watson, R.N.M. et al. CHERI: A hybrid capability-system architecture for scalable software compartmentalization. In *Proceedings of IEEE 2015 Symp. Security and Privacy*, 20–37; <https://ieeexplore.ieee.org/document/7163016>.
- Xia, H. et al. CHERIVoke: characterising pointer revocation using CHERI capabilities for temporal memory safety. In *Proceedings of the 52<sup>nd</sup> Annual IEEE/ACM Intern. Symp. on Microarchitecture*. (2019), 545–557; 10.1145/3352460.3358288.

**Samuel W. Stark** is a Ph.D. student and Harding Scholar in the Department of Computer Science and Technology at the University of Cambridge, U.K., where he is studying the wider applications of capabilities for shared-memory systems with Simon Moore.

**A. Theodore Markettos** is a senior research associate in the Department of Computer Science and Technology at the University of Cambridge, U.K., where he co-leads the CAPcelerate project, which is researching the use of capabilities for securing distributed distrustful accelerators.

**Simon W. Moore** is a professor of computer engineering in the Department of Computer Science and Technology at the University of Cambridge, U.K., where he conducts research and teaching in the general area of computer architecture, with particular interests in secure and rigorously engineered processors and subsystems.

Copyright held by owner(s)/author(s).

DOI:10.1145/3584858

**Mobile, ubiquitous, and immersive computing appear poised to transform visualization, data science, and data-driven decision making.**

BY NIKLAS ELMQVIST

# Data Analytics Anywhere and Everywhere

DATA IS NOW collected everywhere and can be accessed anywhere. Whether you are deciding which product to buy, which potential customer to visit, or which restaurant to frequent, many situations in our everyday personal and professional lives benefit from access to relevant, accurate, and actionable data. Such access supports awareness, promotes understanding, and helps us make the right decisions in today's complex information society. This access to *anywhere data* is made possible by an increasing amount of *everywhere data* collected from virtually all aspects of

our physical and digital world:<sup>22</sup> shopping lists and purchase histories; movie, music, and book preferences; electronic health records and medical test results; colleagues, friends, and family; professional experience and education; and more. While significant privacy, security, and safety concerns are intrinsic to this confluence of anywhere and everywhere data, there is also an unprecedented opportunity to use this data to support individuals navigating the complexities of their professional and personal lives. Fortunately, the last 20 years of the mobile revolution have given us the means to achieve this. Mobile is now the dominant computing platform on the planet, with more than 15 billion mobile devices in 2020<sup>a</sup> and more than six billion of them being "smart" and able to access the Internet.<sup>b</sup> However, these devices—for all their mobility—are currently mere portholes into the digital world, are rarely designed to work together effectively to support a single user let alone multiple ones, and lack the powerful analytical tools needed to enable data-driven decision making on the go.

Fortunately, this may be about to change, with mobile and ubiquitous computing<sup>46</sup> as well as extended reality (XR)<sup>43</sup> finally beginning to trans-

a <https://bit.ly/46ShJPQ/>

b <https://bit.ly/3s4LgHg/>

## » key insights

- Data is collected everywhere on our increasingly connected and information-rich globe.
- Mobile and IoT devices are becoming the world's universal computing platforms.
- Immersive and mobile technologies are poised to help people leverage all this data for analytics conducted anytime and anywhere.
- This emerging genre of ubiquitous analytics is leveraging mobile devices, slim form-factor XR goggles, and cloud computing to help people make sense of data wherever they are, from the office, the factory floor, or even their own living rooms.



form the fields of visualization and data science. Applying these technologies to data analysis suggests a future of ubiquitous<sup>19</sup> and immersive analytics,<sup>31</sup> where clusters of networked mobile devices form an ecosystem for data analytics that can be accessed anytime and anywhere. Such a vision of mobile, immersive, and ubiquitous sensemaking environments would blend state-of-the-art analytics methods with our physical reality to enable making sense of any kind of data in virtually any situation (see Figure 1). However, we should not be weaving computation into our everyday lives just because we can. Rather, progress in cognitive science<sup>11,30,41</sup> supports leveraging the new generation of mobile, immersive, and ubiquitous technologies toward data analytics. These so-called *post-cognitive frameworks* suggest that human thought is not contained merely within our heads but encompasses the entire ecosystem<sup>37</sup> of other people, physical artifacts in our surroundings, and our very own bodies. Distributing computational nodes into our physical surroundings will thus enable us to better scaffold analytical reasoning, creativity, and decision making.

In this article, we investigate how the prevalence of collected everywhere data can enable it to be leveraged for anywhere and anytime access. To achieve this, we first explore the concepts of anywhere and everywhere data and see how current technologies can (and cannot) support cognition in these ubiquitous computing environments. We also review current research in ubiquitous and immersive analytics that builds toward this vision. Finally, we synthesize the current research challenges facing the scientific community and describe the outlook for future research on the topic.

#### **Anywhere and Everywhere Data**

Envisioned by Mark Weiser and Xerox PARC in the late 1980s,<sup>46</sup> the “third wave” of computing—ubiquitous computing—is essentially here, even if it looks subtly different from their original vision.<sup>17</sup> Instead of talking alarm clocks, we have voice-driven home assistants; instead of cheap and disposable computational “tabs,” we have smartphones that we bring everywhere we go;

## **Envisioned by Mark Weiser and Xerox PARC in the late 1980s, the “third wave” of computing— ubiquitous computing—is essentially here.**

instead of an interactive liveboard in every office, we have Zoom and Google Meet videoconferences at our beck and call. Nevertheless, with billions of mobile devices in existence—many of them smart—and a burgeoning Internet of Things (IoT) increasingly making inroads into our physical reality, it is safe to say that we are rapidly approaching a world where computing has indeed been woven into everyday life. The critical difference is, unlike in the original vision of ubiquitous computing, the implementation has been more about personal devices than shared infrastructure—what Harrison et al.<sup>26</sup> call *quality* rather than *quantity* computing.

Regardless, this quiet ubiquitous and mobile computing revolution has had two very specific outcomes relevant to data analytics (Figure 2): the emergence of everywhere data and anywhere data. For the former, it has led to the collection of digital data across our entire society and in virtually all walks of life, both in the real world and online. The IoT encompasses a wide and growing array of devices, such as webcams and security cameras, smart thermostats and light bulbs, digital weather stations, pollution and air-quality sensors, smart locks, connected home appliances, and so on. These devices are also becoming more entwined into professional settings, such as for connected health-care devices and monitors, autonomous farming equipment, wireless inventory trackers, biometric security scanners, wireless sensor networks, and unmanned military equipment. Common for all the datasets collected by these devices is that they are local, temporal, and contextual:

- **Local:** Connected to a specific geographic or semantic location (for example, the temperature on top of the Eiffel tower).

- **Temporal:** Associated with a point in time or a temporal pattern (for example, the temperature on May 1).

- **Contextual:** Related to a specific situation and best interpreted in that context (for example, a car engine’s temperature when traveling at 60mph).

The situation is slightly different for data collected online, as there may be no geographical location associated with the data. Furthermore, in many cases, even data collected from

the real world is stored in databases where these local, temporal, and contextual aspects are discarded or aggregated. For example, the thousands of webcams around the world that merely show video to whomever happens to tune in represent a lost opportunity. What if we instead continuously ran privacy-preserving image-analysis algorithms on this footage to capture environmental data, such as wind speed, population density, snow accumulation, beach erosion, or road traffic?

Nevertheless, it is safe to say that virtually everything we do online is tracked and recorded, a task only made easier by the fact that in the digital world, instrumentation is trivial. The second outcome from the mobile computing revolution is that technology advances have instilled in users a desire for the ability to access this data anytime and anywhere. Many people are accustomed to enjoying near-constant connectivity with the Internet and all its trappings, from social media and video streaming to email and instant messaging. The step to expecting the same anywhere data access even for analytics and sensemaking tasks is not far. However, while current devices certainly can access, manage, and store these datasets, input and output technology has only recently reached a level where such sensemaking can be conducted. But how should we go about doing so effectively?

### The Cognitive Case

Advances in cognitive science summarized as so-called *post-cognitive frameworks*—such as embodied,<sup>41</sup> extended,<sup>11</sup> distributed cognition<sup>30</sup> (see Post-Cognitive Frameworks sidebar)—suggest that human thinking is a system-level process,<sup>37</sup> one that is not merely contained within our brains but which expands to include the world around us, the physical artifacts in our vicinity, our own bodies, and other people. In a post-cognitive framework, cognition is represented as information being transformed from one media to another through interactions. For instance, a person uses a pen to write a reminder on a Post-it note, places the note on a refrigerator, reads the reminder at a later date, and then acts upon it. In other words, tools do not amplify our mind but instead trans-

form certain cognitive activities—such as remembering long number sequences—into other, less taxing cognitive activities—such as reading. Analogously, a visual representation of data on a digital screen is another form of media that a person can interact with to view, manipulate, and understand the underlying data—but so is the notepad the person uses to jot down notes, the phone through which they

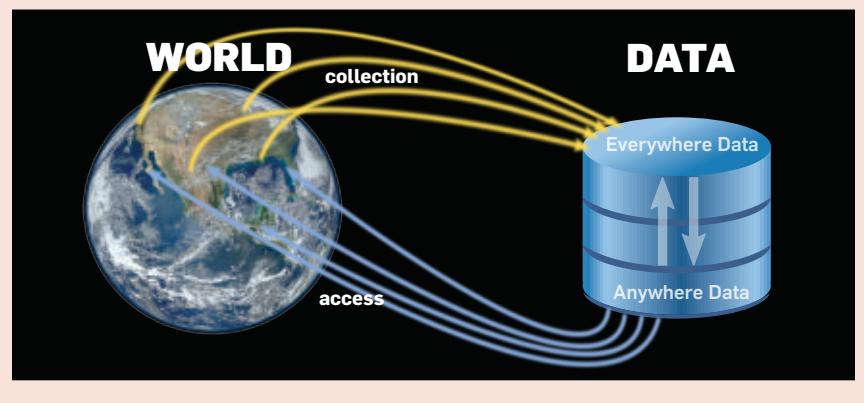
speak to a colleague, and the calculator they employ for quick arithmetic.

Effectively supporting sensemaking thus means instrumenting the entire ecosystem of artifacts involved in the cognitive process. This is in contrast to traditional human-computer interaction (HCI) and visualization paradigms, which tend to consider only the actual interface between user and machine: the visual output from

**Figure 1. Mobile, immersive, and ubiquitous sensemaking.** Two examples of data analytics conducted in a sensemaking environment of connected devices, interactive visualizations, and multiple collaborators.



**Figure 2. Anywhere and everywhere data.** Today, data is collected from everywhere and is increasingly being accessed from anywhere.



# Post-Cognitive Frameworks

Traditional cognitive science views the individual human as a fundamental unit of cognition, with a basic workflow consisting of perception to understand one's surroundings, cognition within the mind, and physical action to manipulate the environment. It follows that in traditional cognitive science, the mind maintains a model of the world, and the senses are used to continuously update this model. Such a model is called *solipsism* and holds that you can only be sure about your own mind and not the external world; in essence, you could replace the world with simulated sensor input and muscle output—as in the movie *The Matrix* (1999)—and you would be none the wiser. However, the last few decades of research have slowly chipped away at this solipsistic view of the mind. For example, Post-it notes, notebooks, and smartphones trivially extend our memory.<sup>11</sup> Discussions with other people facilitate thinking and their abilities complement ours, forming a socially distributed cognitive system.<sup>30</sup> Kirsh and Maglio<sup>34</sup> found evidence that expert Tetris players perform costly geometric rotation in the world rather than in their minds. In other words, cognition is essentially a product of an individual's interaction with their environment rather than a closed system inside the individual's mind.<sup>37</sup> Frameworks that go beyond the classic human information processing model are called *post-cognitive frameworks*. Next, we discuss three, all of which have useful properties that we draw upon in this paper. However, these frameworks, as well as cognitive science as a whole, are still very much works in progress and thus—to bastardize a phrase by statistician George E.P. Box—they may all be incorrect to some degree, but some may still be useful.

- ▶ **Extended cognition.** The unofficial motto of Clark and Chalmers' extended cognition framework is that “Cognitive processes ain’t (all) in the head!,”<sup>11</sup> arguing that the surrounding environment is a fundamental component of all human cognition—that, in fact, we are all effectively cyborgs, but in the most natural way in symbiosis with our environment.
- ▶ **Socially distributed cognition.** Developed by UCSD cognitive psychologist Ed Hutchins<sup>30</sup> in the 1980s and 1990s, socially distributed cognition (DCog) holds that human cognition is distributed across sociocultural systems in the individual's surroundings, which includes physical artifacts, other individuals, and cultural systems (history, practice, etiquette, and so on).
- ▶ **Embodied cognition.** The basic tenet of embodied cognition is that thinking is fundamentally influenced by and inseparable from our own bodies.<sup>41</sup> These ideas have been fundamental to the areas of embodied and tangible interaction in HCI.<sup>16</sup>

Beyond these three, we also note Scaife and Rogers' concept of external cognition<sup>3</sup> in support of graphical representations, which, even if it does not quite reach the level of a framework, discusses how visualization helps cognition by offloading computation and memory, re-representing data in a more suitable form, and graphically constraining inferences.

the screen and the user input from the mouse, keyboard, touchscreen, or microphone. By distributing interactive representations of data on a multitude of digital devices scattered around our physical surroundings as well as with the collaborators involved in the cognitive system, we are, for all intents and purposes, creating a cybernetic extension of the mind, expanding it using these digital devices. This notion would have been revolutionary if not for the fact that we have been doing this for thousands of years, since the first human picked up the first fallen branch and used it to dig a hole in the dirt. A shovel is a cybernetic extension of a person's arms that enables them to dig better, and an excavator takes this concept of human extension even further. Similarly, a computer can be seen as a cybernetic extension of our logic, reasoning, and memory—or, as Steve Jobs famously put it, a “bicycle for the mind.” Nevertheless, from the digging example, it is clear that the type of tool makes a big difference, and even if digging may be a solved problem, the same cannot be said of sensemaking.

Of course, the argument here is not that the more devices we weave into our everyday environment, the more effective the cognitive process. As Harrison et al.<sup>26</sup> point out with their quality vs. quantity computing argument, human attention is a scarce commodity, and the current focus on a single, highly capable personal device seems to have won out over the “one user to many devices” Ubicomp vision from the 1990s. However, there clearly exists a middle ground between a single device per person and dozens. In addition, nothing says that the devices must be homogeneous; it may be difficult to use two tablets at the same time, but a tablet and a smartwatch can be complementary. For one thing, data analysts do not only use computers but tend to surround themselves with pen and paper, calculators, reference material, books, and other physical artifacts, not to mention other analysts.<sup>37</sup> For another, physical space is a key factor in the data analysis process; Wright et al.<sup>48</sup> describe how some professional intelligence analysts would use the entire floor of their office to arrange documents during analysis, and Andrews

# Sensemaking and Visualization

Sensemaking is the activity of searching for a representational schema to fit available data<sup>35,39</sup> (for example, a mental model), allowing people to attribute meaning to available data and answer questions about it. In terms of the cognitive foundations discussed earlier in this paper, this activity often involves a combination of internal and external cognitive representations, such as data tables, charts, notes, sketches, and calculations. These representations should ideally be chosen to facilitate required mental activities (memory, computation, planning), thereby minimizing the *cost structure*<sup>39</sup> of cognition.

Data visualizations have long been used for communicating insights because of their accessible and visual form. However, visualizations are also effective for providing data-driven overviews, generating hypotheses, and answering questions at little cognitive cost. This makes visualization particularly suitable for sensemaking, where exploring the data for meaning and hypotheses is a significant part. Such *exploratory data analysis*<sup>45</sup> represents an alternative to more traditional hypothesis-driven *confirmatory data analysis*.

and North conducted empirical studies demonstrating the utility of significant visual space in facilitating sensemaking.<sup>1</sup> In fact, the intelligent use of space as part of cognition is a cornerstone of post-cognitive frameworks,<sup>33</sup> simplifying choice, facilitating perception, and aiding computation. Furthermore, *proxemics*—hailed by some as the new UbiComp<sup>24</sup>—tells us how people physically relate to artifacts and other people as they interact with them; for example, we typically turn to face people we speak to. These are all prime arguments in favor of instrumenting all the components involved in the cognitive process—ranging from notepads to books and other people—using networked digital devices. Drawing from all three of the ubiquitous computing, visual computing, and visual analytics traditions, we call this approach *ubiquitous analytics*.<sup>19</sup>

### The Gap in Our Tech

If cognitive science supports the use of these networked groups of devices to facilitate sensemaking, why have we not yet seen a plethora of such ubiquitous analytics systems on the market and in the scientific literature? The answer is that until only recently, the technology required to harness anywhere and everywhere data has been outside our reach. More specifically, the gaps in our technology include screen size, input surfaces, and the general design of the mobile devices we need for this endeavor.

Both screen size and input surfaces suffer from a device-miniaturization tradeoff, where mobile devices already maximize the input and output dimensions to encompass the full size of the device. They simply cannot be made much larger to avoid diminishing the mobility and portability of a smart-watch, tablet, or smartphone. One solution for both input and output is to go beyond the device itself to appropriate the physical world as part of the interaction.<sup>25</sup> For example, portable projectors can turn any nearby surface into a display of arbitrary size.<sup>15</sup> Depth cameras, ultrasound, or electrical sensing can similarly transform walls, tables, or even our very own bodies into touch surfaces. The rise of consumer-level mixed- and augmented-reality<sup>43</sup> equipment has taken this idea even further by turning our entire world into a potential

**This quiet ubiquitous and mobile computing revolution has had two very specific outcomes relevant to data analytics: the emergence of everywhere data and anywhere data.**

canvas for data display and manipulation. Instead of our smartphones being mere portholes into an unseen world of data, XR technology has broken down the fourth wall hemming us in.

However, a hurdle remains: Current mobile devices are still designed using the “quality computing” mindset. This means that each individual device is intended to be used in isolation and with the undivided attention of the user. A significant gap for mobile computing is to derive new design paradigms where multiple devices can stack together and scaffold each other for the current task while minimizing barriers and tedious housekeeping. We will demonstrate some examples of how to design such stacking devices in the treatment to come.

### Ubiquitous, Immersive, and Situated Analytics

While command-line tools, automated scripts, and libraries are common in general data-science tasks such as computation, wrangling, and confirmatory analysis,<sup>9</sup> the kind of exploratory data analysis common to sensemaking (see Sensemaking and Visualization sidebar) often benefits from a more fluid and visual interaction model. This is particularly true in mobile settings, where the precise text entry required for programming and command-line interfaces remains challenging on handheld devices. In such situations, it is more convenient to turn to interactive visual interfaces and automated recommendations. The scientific field of *data visualization* concerns itself with precisely these visual and interactive data displays. In fact, the field has recently undergone a dramatic change with the rise of ubiquitous,<sup>19</sup> immersive,<sup>21,31</sup> and situated analytics<sup>20</sup> for tackling these use cases.

Three-dimensional visual representations have long been the norm in many scientific applications for visualization, such as flow visualization, medical imaging, and volumetric rendering. However, the very first applications of data visualization to immersive settings actually came from outside the visualization field. In 2003, Bowman et al.<sup>6</sup> proposed a research agenda for so-called *information-rich virtual environments* (IRVEs) that combined 3D virtual reality with information visualization. Touch-enabled tabletops were an early

platform, with Spindler and Dachselt<sup>44</sup> proposing a tangible lens for interacting with data on or above a horizontal display. Finally, visualization made inroads into mobile computing in a much more unobtrusive manner, with the first applications being commercial ones on smartphones. Lee et al. review mobile data visualization in a recent book on the topic.<sup>36</sup>

Starting in 2011, I joined this research area by proposing the notion of an embodied form of human-data interaction.<sup>18</sup> This launched my research agenda on this topic, and my students and I followed this up with an embodied lens used for exploring data on a touch-based tabletop.<sup>32</sup> It eventually led to Pourang Irani and I defining the *ubiquitous analytics*<sup>19</sup> paradigm in 2013, which serves as an umbrella term for the research field: anytime, anywhere sensemaking performed on a plethora of networked digital devices, not just immersive ones. My students and I also proposed several computing infrastructures for realizing this vision for data analytics, including Munin (a peer-to-peer middleware based on Java),<sup>4</sup> PolyChrome (the first Web-based framework),<sup>3</sup> and Vistrates (a mature Web replication framework for component-based visualization authoring).<sup>5</sup> Researchers have since taken our ideas further; for example, ImAxes<sup>13</sup> enables data analysis in virtual reality using powerful hand gestures for connecting dimensional axes in 3D. Another early example was the work of Butscher et al. in combining parallel coordinates display in augmented reality on top of a tabletop display.<sup>8</sup>

Along the way, several variations of these ubiquitous forms of analytics have emerged. *Immersive analytics*, initially introduced by Chandler et al.<sup>10</sup> in 2015, takes a specific focus on immersive and 3D spatial technologies to support sensemaking.<sup>31</sup> *Situated analytics*,<sup>20</sup> on the other hand, emphasizes the spatial referents for data in the real world. Willet et al.<sup>47</sup> built on this in 2017 for *embedded data representations* that are deeply integrated with the spaces, objects, and contexts to which the data refers.

Given the one-to-many ratio of users to devices in traditional Ubicomp, cardinality is also a common denominator for ubiquitous analytics. The coming subsections discuss cardinality

## Until only recently, the technology required to harness anywhere and everywhere data has been outside our reach.

from three different perspectives and give specific examples for each: multiple devices, multiple resources, and multiple collaborators.

**Devices Now and in the Future:** In this vision of an increasingly ubiquitous and mobile future, the traditional view of a “device” as a physical artifact with computation, storage, memory, display, and input is being challenged. Some physical devices will surely remain, but in other ways, devices will recede into the background or be replaced by a single mixed/augmented-reality display.

**Multiple devices.** The first factor worth investigating in a ubiquitous approach to analytics is how to best manage the multiple types of devices a user will engage with during sensemaking. As already mentioned, human attention is a finite resource that must be judiciously managed. Current device platforms are typically designed for focused use; for example, a smartphone engages the user’s hand, often the dominant one, so adding a second smartphone is seldom helpful. Rather, adding multiple devices to an interaction should be complementary, either by physical form factor (for instance, a smartwatch or a large display), by physical placement (for instance, handheld, wall-mounted, or head-mounted), by task (for instance, a primary display used for map navigation and secondary ones to show legends and drill-down details), or by various combinations of these. Users are rarely helped by two identical devices inhabiting the same position in this design space, except possibly for comparison tasks where holding two tablets side by side may be beneficial.

We studied this phenomenon in a 2018 research paper on the interplay between smartwatches and large touch displays<sup>27</sup> that we informally dubbed David & Goliath (D&G). The form factors of these two different computing platforms are radically different. A smartwatch is a fundamentally personal device, whereas a large display is a fundamentally public one. Smartwatches have small displays, are attached to a person’s wrist and are thus always within reach, and are only really accessible—both in terms of physical reach as well as social practices—to

the wearer. Any action performed and any data displayed on a smartwatch will accordingly be personalized. In contrast, a large touch display is often vertically mounted on a wall, is visible to many by virtue of its size, and invites interaction by anyone within physical reach. This, in other words, is a prime example of how a ubiquitous analytics environment can be designed to use complementary device platforms.

Figure 3 shows an example of the David & Goliath system in action. Drawing on the personal affordances of the smartwatch, we mainly use the smaller display as a personal storage container as well as a mediator and a remote control—all tasks that are well-suited to the device. The large display, on the other hand, is a public and shared device that provides the fundamental data visualization for the team working on a specific task. Because of its size, it can show dashboards of multiple visualizations rather than just one.

Devices can also complement each other over time. In the recent ReLive<sup>29</sup> system (Figure 4), we proposed a hybrid desktop/VR interface that gives the same user access to two different ubiquitous analytics interfaces depending on task. ReLive was built to help mixed-reality designers and researchers analyze real-time telemetry from user studies conducted in virtual- or augmented-reality environments. Thus, the primary data being visualized is 3D tracking and event logs over time. Here, again, the two analysis interfaces—an in situ 3D VR interface and an ex situ 2D desktop interface—are complementary in that each suits a specific task. When it comes to understanding 3D data, such as how two participants worked together in a 3D assembly task, nothing beats viewing the data from a first-hand perspective in 3D. However, when deriving high-level abstract findings, such as the average distance between participants over an entire session, the 2D desktop interface is optimal. Having access to both interfaces allows a user to pick the right tool depending on task.

But it is not enough that two devices are complementary—they should actively support each other. This is particularly important in ReLive, where the same user will be switching between the two different interfaces time

and time again, sometimes as part of the same task. This interface switch leads to a “cognitive context switch,” where users must reorient themselves within the new interface in relation to the other interface. In ReLive, this is facilitated by providing common anchors between the two views. The 2D desktop view provides a 3D in situ viewport that is always visible. Analogously, the 3D Virtual Reality view replicates many of the interface features from the desktop view, and selections made inside the first-hand perspective can be accessed in the 2D view.

**Takeaway.** *Each device involved in ubiquitous analytics should have a clear role that complements the other*

*devices, and transitions between them should be seamless.*

**Multiple resources.** Given that the previous discussion is at the level of entire devices, our next concern is how to manage the resources associated with these devices, such as displays, computational power, input affordances, and so on. Resource management quickly becomes onerous and almost all-consuming when multiple devices are involved, and automation is thus required. For example, imagine logging into your typical workplace network, and then multiply this for every device you want to include in a work session. If the overhead associated with using multiple devices

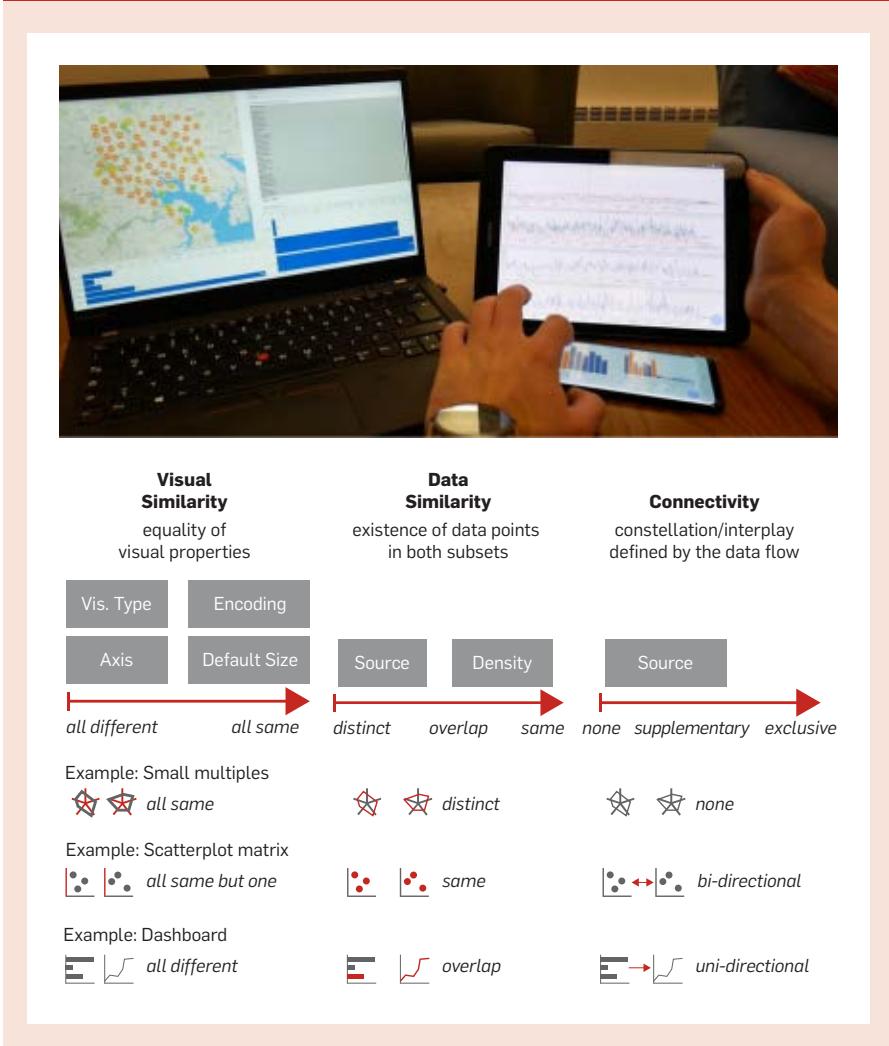
**Figure 3. Two users collaborating using the David & Goliath system. The woman on the left is interacting with the large touch display (“Goliath”), which is—by its size and orientation—a public and shared surface. The man on the right is interacting with his smartwatch (“David”), which is used for personal storage and private data manipulation.**



**Figure 4. Hybrid interfaces in ReLive. The 3D VR interface (left) and the 2D desktop interface (right) are used depending on the task. First-person 3D invites situated analytics, whereas the 2D desktop better supports higher-level analysis.**



**Figure 5. Automatic visualization layout using Vistribute.** The Vistribute system will dynamically reorganize the visualizations to be displayed on the available surfaces whenever the device configuration changes (top). Relationships between the different visualizations are used to ensure that similar representations are linked (bottom). For example, the time-series charts on the tablet are always laid out vertically on the same display and with a common horizontal axis.



**Figure 6. Collaboration using Branch-Explore-Merge.** The mobile device, one per user, is generally synchronized with the data display on the tabletop. However, if the user wants to deviate from the shared display, they can branch the current state of the tabletop on their personal device, make the desired changes, and then merge back their findings to the shared display. Merging requires consensus.



becomes overwhelming, people will simply stick to a single device.

One example where such automated resource management is critical is display management. In a typical ubiquitous analytics scenario, we might imagine that the displays available to a person will change dynamically over the course of a work session. For example, while leaving your car or bus on the way to the office, you might only have access to your smartphone, allowing you to view a single data display at a time. Once you get to your office and your desk, you can boot up your computer and distribute an entire visualization dashboard across your monitors as well as your personal devices. Finally, when you head into the conference room for a meeting, the projected screen as well as the laptops belonging to your colleagues could be used to display even more detailed data visualizations. However, for this to be practical, the layout of visualizations across available displays should be automatic, lest the user gets continually bogged down by moving charts around. Furthermore, such a layout must organize closely related charts together on the same display as well as respect their geometry affinity.

Such dynamic layout across ubiquitous analytics environments was our focus in the Vistribute<sup>28</sup> system (Figure 5) from 2019. Based on a design space of visualization relationships (Figure 5, bottom), the Vistribute layout manager performs real-time constraint solving based on a set of high-level visualization-layout heuristics. These heuristics include multi-chart relationships, such as the visual and data similarity as well as data-flow connectivity in the figure, but also single-chart heuristics. For example, dense data displays, such as geographical maps, are given more display space, while a skewed aspect-ratio chart, like a time-series line chart, will be given wide and short display allocations.

The VisHive toolkit<sup>14</sup> tackles a different but similar problem: computational resource management. More specifically, mobile ubiquitous analytics often call for significant computation, such as when performing textual analysis, clustering, or machine learning. The computational power of an individual mobile device may be insufficient to complete this calculation

on its own in a timely manner. One solution is to simply use a remote computational resource, but even if the Internet connectivity was reliable—which is not a given in many mobile on-the-go situations—the data to be processed is often simply too large to effectively upload on the fly. The VisHive system solves this by forming an ad hoc cluster, or *hive*, of mobile devices for load-balancing computation in the field. Designed to work in conjunction with Web-based visualization systems, VisHive is a small JavaScript library that runs directly in a device's Web browser and requires no specialized software to be installed on the device. In this way, a user can simply bring additional mobile devices online to help process computational tasks in case a particular computation is taking too long to complete.

**Takeaway.** Resource management across multiple devices involved in a ubiquitous analytics environment should be automated to minimize the user burden.

**Multiple collaborators.** Finally, the third cardinality factor in a ubiquitous analytics system is the individual users who often collaborate on realistic tasks. In general, collaborative visualization is a grand challenge of data visualization research, but ubiquitous analytics have the benefit of being designed to be collaborative from the very foundation. For example, the David & Goliath system can easily be used in a collaborative manner (as in Figure 3), and the two different interfaces in ReLive could just as easily be employed by two parallel users rather than the same one in sequence.

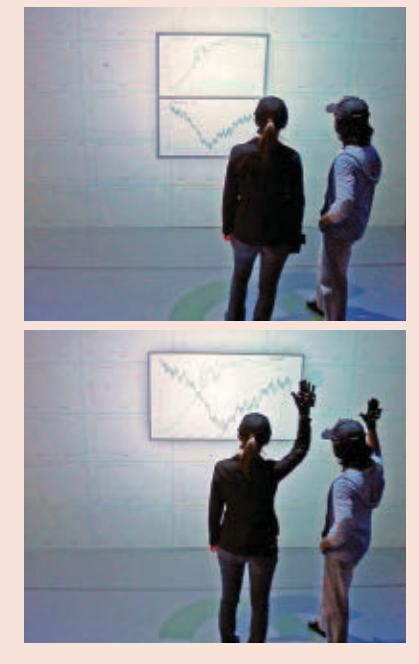
One example of this is the Branch-Explore-Merge (B-E-M) tabletop system<sup>38</sup> we proposed in 2012 for collaborating across tabletops and mobile devices (Figure 6). A precursor to the D&G system discussed earlier, B-E-M uses the same philosophy of private and personal devices (smartphones and tablets) vs. public and shared devices (large touch tabletops in this case) as in the smartwatch and wall display scenario. With B-E-M, however, the focus is specifically on the coordination and consensus mechanisms required for a visualization system used by multiple users. While a user should be able to modify

## The first factor worth investigating in a ubiquitous approach to analytics is how to best manage the multiple types of devices a user will engage with during sensemaking.

the visualization on their own personal device at will, any change made on the shared tabletop will directly affect all collaborators. To manage this process, B-E-M draws on basic revision-control principles from software engineering, where the user can branch the current state of the shared display on their own device, explore the data on their own, and then finally merge back their changes—or discard them if the exploration turns out to be a dead end. The B-E-M system requires a vote from all participants around the tabletop for changes to be merged back onto the shared display.

Voting is a somewhat disruptive coordination mechanism, and often more implicit mechanisms are better. In the Proxemic Lens<sup>2</sup> project from 2016, we used proxemics information<sup>24</sup> such as the distance between users, the direction of their body and head, their hand and foot gestures, and the spatial arrangement of objects and people in 3D space to guide interaction. The goal of the project was to infer user intention in a

**Figure 7. Collaborative data analysis using the Proxemic Lens.** In the top image, the two collaborators have approached each other sufficiently that their body language indicates their work is closely coupled, so their respective lenses are stacked vertically. In the bottom image, the collaborators both raise their hands in consensus, causing the lenses to be overlaid to allow for direct comparison in the same visual space.



ubiquitous analytics scenario from their body language and physical navigation in a space. We were particularly interested in studying how a large shared display can be best utilized depending on the collaborative coupling of the people using it: separate viewports (lenses) for people working independently and stacked or overlaid viewports for people working closely together (Figure 7). We used a static environment and a floor-to-ceiling display with a Vicon motion-capture system to track this information for multiple collaborators, but mobile environments will require more subtle biometric or inertial-sensing technology. Overall, participants enjoyed the implicit interaction of the Proxemic Lens system, and their intention was often inferred correctly. However, they also indicated a preference for explicit rather than implicit gestures for actions typically seen as commands, such as creating charts, splitting viewports, and consensus operations (Figure 7, bottom).

**Takeaway.** Collaboration in ubiquitous analytics requires careful consideration of coordination and consensus, just like for general collaborative work, but the heterogeneous devices typically employed in ubiquitous analytics settings makes integrating such coordination mechanisms straightforward.

### Challenges and Outlook

This future vision for data analytics is still a new notion, and its various manifestations as visual, ubiquitous, immersive, and situated analytics are still nascent and emergent. This research group has been a significant driver in the field, but the story is much bigger than our efforts. The greater research community is energetic and growing, with new and exciting analytical systems being proposed at every major conference and journal issue; certainly too many to discuss in a single research article. However, we claim that the techniques and technologies described here represent a cross-section of ubiquitous analytics research.

Based on this review, we would summarize our takeaways using a single theme: **device diversity**. Basically, it is the varied and heterogeneous nature—as well as their effective utilization—of the individual devices

## We still need further work on capturing, integrating, and synthesizing heterogeneous data from multiple sources in our environment.

involved in a sensemaking task that makes ubiquitous analytics powerful. This idea is also supported by the post-cognitive frameworks discussed in this article.

If heterogeneity is the lead theme of ubiquitous analytics, where is the field going in the future? Ens et al.<sup>21</sup> presented a vision for immersive analytics in 2021 and outlined the grand research challenges of that field. We support these challenges in general but complement their technical nature with our own list of higher-level future challenges:

- **The future is mixed.** While handheld devices with screens are here to stay, it seems almost inevitable that future devices will eventually be based on augmented and mixed reality.<sup>43</sup> Augmented-reality and mixed-reality displays have the capacity to turn the entire surrounding world into a canvas for visual data displays, and the technology is constantly improving. This should mean that immersive 3D visualizations and interaction techniques will become increasingly important in the future.

- **Human-centered artificial intelligence teaming.** While mostly glossed over in this article, the future of ubiquitous analytics is closely entwined with artificial intelligence. Only with powerful automated algorithms and recommendations at our beck and call will we be able to overcome the computational challenges of tomorrow. However, rather than the black-box pipeline model of traditional AI, we believe in the use of visual interactive interfaces as inflection points for involving human operators in the loop. This is known as *human-centered artificial intelligence* (HCAI),<sup>42</sup> but this is outside the scope of this article.

- **Standards, components, and practices.** If future sensemaking environments are characterized by heterogeneous systems and devices, then we will need standards as well as standardized components and practices to enable multiple vendors to provide their own versions of these tools. While several immersive analytics<sup>12</sup> and augmented-reality toolkits exist,<sup>23</sup> we will need a much richer ecosystem, including for software development, session and view management, navigation, and so on. Furthermore, many existing toolkits are built on proprietary 3D game engines, such as Unity. It is

difficult to predict the future here, but based on past history, technologies based on open standards, such as the Web-based VRIA<sup>7</sup> and Vistrates<sup>5</sup> toolkits, are safer bets than those built on closed and proprietary technology.

► **Accessibility and inclusivity.** The field of data visualization is just now recognizing the accessibility of interactive visual representations to people with visual, motor, or cognitive impairments. Ubiquitous analytics will need to learn from these lessons at an early stage. However, the heterogeneous nature of ubiquitous analytics can once again play in our favor in that the inclusion of diverse platforms into the sensemaking loop may make it easier to accommodate people of varying physical and cognitive ability.

► **Of scale and scalability.** With the exception of the VisHive project, most of this article has not concerned itself with large-scale data management. While data visualization is often seen as a best-in-class solution for the human aspect of big data,<sup>22</sup> there are many challenges on the computational and data-management sides that must be solved to enable mobile sensemaking at this scale and magnitude.

► **Evermore everywhere.** Similarly, while the everywhere data aspect of this article is more about providing opportunity for anywhere data, we still need further work on capturing, integrating, and synthesizing heterogeneous data from multiple sources in our environment. While privacy and security must remain foremost in our minds, such new data about our world will only continue to make data-driven decision-making better and more effective.

The current state of ubiquitous data analytics has been more than 10 years in the making. There is now a thriving and creative research community invested in taking this vision into the future. Perhaps one day we will truly see sensemaking become embedded into the fabric of everyday life. ■

## References

- Andrews, C. et al. Space to think: Large, high-resolution displays for sensemaking. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems*. (2010), 55–64; 10.1145/1753326.1753336.
- Badam, S.K. et al. Supporting visual exploration for Multiple users in large display environments. In *Proc. of the IEEE Conf. on Visual Analytics Science and Technology*. (2016), 1–10; 10.1109/vast.2016.7883506.
- Badam, S.K. and Elmqvist, N. PolyChrome: A cross-device framework for collaborative Web visualization. In *Proceedings of the ACM Conf. on Interactive Tabletops and Surfaces*. (2014), 109–118; 10.1145/2669485.2669518.
- Badam, S.K. et al. Munin: A peer-to-peer middleware for ubiquitous analytics and visualization spaces. *IEEE Trans. on Visualization and Computer Graphics* 21, 2 (2015), 215–228; 10.1109/TVCG.2014.2337337.
- Badam, S.K. et al. Vistrates: A component model for ubiquitous analytics. *IEEE Trans. on Visualization and Computer Graphics* 25, 1 (Jan. 2019), 586–596; 10.1109/TVCG.2018.2865144.
- Bowman, D.A. et al. Information-rich virtual environments: Theory, tools, and research agenda. In *Proceedings of the ACM Symp. on Virtual Reality Software and Technology*. (2003), 81–90; 10.1145/1008653.1008669.
- Butcher, P.W.S. et al. VRIA: A Web-based framework for creating immersive analytics experiences. *IEEE Trans. on Visualization and Computer Graphics* 27, 7 (2021), 3213–3225; 10.1109/TVCG.2020.2965109.
- Butscher, S. et al. Clusters, trends, and outliers: How immersive technologies can facilitate the collaborative analysis of multidimensional data. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems*. (2018), 1–12; 10.1145/3173574.3173664.
- Cao, L. Data science: A comprehensive survey. *Comput. Surveys* 50, 3, Article 43 (2018); 10.1145/3076253.
- Chandler, T. et al. Immersive analytics. In *Proceedings of the Intern. Symp. on Big Data Visual Analytics*. IEEE, (2015); 1–8; 10.1109/BDVA.2015.7314296.
- Clark, A. and Chalmers, D. The extended mind. *Analysis* 58, 1 (1998), 7–19; 10.1093/analys/58.1.7.
- Cordeil, M. et al. IATK: An immersive analytics toolkit. In *Proceedings of the IEEE Conf. on Virtual Reality and 3D User Interfaces* (2019), 200–209; 10.1109/VR.2019.8797978.
- Cordeil, M. et al. Immersive axes as embodied affordances for interactive multivariate data visualisation. In *Proceedings of the ACM Symp. on User Interface Software and Technology* (2017), 71–83; 10.1145/3126594.3126613.
- Cui, Z. et al. VisHive: Supporting web-based visualization through ad hoc computational clusters of mobile devices. *Information Visualization* 18, 2 (2019), 195–210; 10.1177/1473871617752910.
- Dachsel, R. et al. Pico projectors: Firefly or bright future? *Interactions* 19, 2 (2012), 24–29; 10.1145/2090150.2090158.
- Dourish, P. *Where the Action Is: The Foundations of Embodied Interaction*. MIT Press, Cambridge, MA, USA (2001).
- Dourish, P. and Bell, G. *Divining a Digital Future—Mess and Mythology in Ubiquitous Computing*. MIT Press, Cambridge, MA, USA (2011).
- Elmqvist, N. Embodied human-data interaction. In *Proceedings of the ACM CHI Workshop on Embodied Interaction: Theory and Practice in HCI* (2011).
- Elmqvist, N. and Irani, P. Ubiquitous analytics: Interacting with big data anywhere, anytime. *IEEE Computer* 46, 4 (2013), 86–89; 10.1109/mc.2013.147.
- ElSayed, N.A.M. et al. Situated analytics: Demonstrating immersive analytical tools with augmented reality. *J. of Visual Languages and Computing* 36 (2016), 13–23; 10.1016/j.jvlc.2016.07.006.
- Ens, B. et al. Grand challenges in immersive analytics. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems*. Article 459 (2021); 10.1145/3411764.3446866.
- Fisher, D. et al. Interactions with big data analytics. *Interactions* 19, 3 (2012), 50–59; 10.1145/2168931.2168943.
- Fleck, P. et al. RagRug: A toolkit for situated analytics. *IEEE Trans. on Visualization and Computer Graphics* (2022); 10.1109/TVCG.2022.3157058.
- Greenberg, S. et al. Proxemic interactions: The new ubicom? *Interactions* 18, 1 (2011), 42–50; 10.1145/1897239.1897250.
- Harrison, C. Appropriated interaction surfaces. *IEEE Computer* 43, 6 (2010), 86–89; 10.1109/MC.2010.158.
- Harrison, C. et al. Achieving ubiquity: The new third wave. *IEEE Multimedia* 17, 3 (2010), 8–12; 10.1109/MMUL.2010.53.
- Horak, T. et al. When David meets Goliath: Combining smartwatches with a large vertical display for visual data exploration. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems* (2018), 19:1–19:13; 10.1145/3173574.3173593.
- Horak, T. et al. Distributing interactive visualizations in dynamic multi-device setups. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems*. (2019), 616:1–616:13; 10.1145/3290605.3300846.
- Hubenschmid, S. et al. Bridging in-situ and ex-situ visual analytics for analyzing mixed reality user studies. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems* (2022), 24:1–24:20; 10.1145/3491102.3517550.
- Hutchins, E. *Cognition in the Wild*. MIT Press, Cambridge, MA, USA (1995).
- Immersive analytics. *Lecture Notes in Computer Science* 11190. K. Marrriott et al. (eds). Springer Intern. Publishing, Berlin, Germany (2018); 10.1007/978-3-030-01388-2.
- Kim, K. and Elmqvist, N. Embodied lenses for collaborative visual queries on tabletop displays. *Information Visualization* 11, 4 (2012), 319–338; 10.1177/1473871612441874.
- Kirsh, D. The intelligent use of space. *Artificial Intelligence* 73, 1–2 (1995), 31–68; 10.1016/0004-3702(94)00017-U.
- Kirsh, D. and Maglio, P. On distinguishing epistemic from pragmatic action. *Cognitive Science* 18, 4 (1994), 513–549; 10.1207/s15516709cog1804\_1.
- Klein, G. et al. Making sense of sensemaking: Alternative perspectives. *IEEE Intelligent Systems* 21, 4 (2006), 70–73; 10.1109/MIS.2006.75.
- Lee, B. et al. *Mobile Data Visualization*. Chapman and Hall/CRC Press, Boca Raton, FL, USA (2022).
- Liu, Z. et al. Distributed cognition as a theoretical framework for information visualization. *IEEE Trans. on Visualization and Computer Graphics* 14, 6 (2008), 1173–1180; 10.1109/TVCG.2008.121.
- McGrath, W. et al. Branch-explore-merge: Facilitating real-time revision control in collaborative visual exploration. In *Proceedings of the ACM Conf. on Interactive Tabletops and Surfaces* (2012), 235–244; 10.1145/2396636.2396673.
- Russell, D.M. et al. The cost structure of sensemaking. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems* (1993), 269–276; 10.1145/169059.169209.
- Scaife, M. and Rogers, Y. External cognition: How do graphical representations work? *Intern. J. Human-Computer Studies* 45, 2 (1996), 185–213; 10.1006/ijhc.1996.0048.
- Shapiro, L.A. *Embodied Cognition*. Routledge, New York, NY, USA (2011).
- Shneiderman, B. *Human-Centered AI*. Oxford University Press, Oxford, UK (2022).
- Speicher, M. et al. What is mixed reality? In *Proceedings of the ACM Conf. on Human Factors in Computing Systems* (2019), 537:1–537:15; 10.1145/3290605.3300767.
- Spindler, M. et al. Tangible views for information visualization. In *Proceedings of the ACM Conf. on Interactive Tabletops and Surfaces* (2010), 157–166; 10.1145/1936652.1936684.
- Tukey, J.W. *Exploratory Data Analysis*. Addison-Wesley, Boston, MA, USA (1977).
- Weiser, M. The computer for the 21<sup>st</sup> century. *Scientific American* 265, 3 (1991), 94–104; 10.1145/329124.329126.
- Willett, W. et al. Embedded data representations. *IEEE Trans. on Visualization and Computer Graphics* 23, 1 (Jan. 2017), 461–470; 10.1109/TVCG.2016.2598608.
- Wright, W. et al. The sandbox for analysis: Concepts and valuation. In *Proceedings of the ACM Conf. on Human Factors in Computing Systems* (2006), 801–810; 10.1145/1124772.1124890.

**Niklas Elmqvist** (elm@cs.au.dk) is a Villum Investigator and a full professor of computer science at Aarhus University in Aarhus, Denmark.

2023 Copyright held by owner(s)/author(s).



Watch the author discuss this work in the exclusive Communications video. <https://cacm.acm.org/videos/> anywhere-everywhere

DOI:10.1145/3582490

## To what extent have the techniques for theoretical analysis of edit distance algorithms achieved their goals?

BY PAUL MEDVEDEV

# Theoretical Analysis of Edit Distance Algorithms

EDIT DISTANCE—A CLASSICAL problem in computer science—has received ongoing attention from both practitioners and theoreticians. Given two strings  $A$  and  $B$ , the edit distance is the minimum number of substitutions, insertions, and deletions needed to transform  $A$  into  $B$ . For example, the edit distance between *apfleee* and *rapleet* is 3:  $\text{apf} \xrightarrow{\text{ins}} \text{rapfleee} \xrightarrow{\text{del}} \text{rapleee} \xrightarrow{\text{sub}} \text{rapleet}$ . The edit distance problem is widely known, as it is often taught as part of the undergraduate curriculum to illustrate two-dimensional dynamic programming. Theoreticians

have studied the problem starting from as early as 1966<sup>24</sup> and 1974,<sup>43</sup> but it very much remains an active topic of research today (for example, Boroujeni et al.<sup>8</sup>). Simultaneously, bioinformatics practitioners continue to actively develop<sup>14,15,40</sup> and apply<sup>9,28,38,39,42,45</sup> fast edit distance solvers ubiquitously. Given its status as a classic problem and its importance to both theoreticians and practitioners, edit distance provides a unique opportunity to study the interaction between theory and practice.

Theoreticians develop abstract algorithms that have superior theoretical performance; practitioners develop implemented algorithms that have superior empirical performance. In an ideal world, practitioners would implement and analyze the empirical performance of the abstract algorithms developed by theoreticians, while theoreticians would analyze the theoretical performance of the implemented algorithms developed by practitioners. In the real world, there is often a wide gap between the practical and theoretical communities; understanding how to close this gap is critical to making theoretical computer science more relevant to applications. The edit distance problem is then an excellent lens through which to understand how the theoretical analysis of algorithms impacts practical implementations.

There are many ways to approach the practice/theory gap in a problem like edit distance. We take one that is systematic and focused on the way theoreticians analyze edit distance algorithms, rather than on the algorithms themselves. From a practical perspective, theoretical analysis has two goals.<sup>37</sup> The first goal is to predict the empirical performance of an algorithm, either in an absolute sense or relative to other algorithms. The second goal is to be a yardstick that drives the design of novel algorithms that perform well in practice. In this article, we systematically survey the types of theoretical analysis techniques that have been applied to edit distance and evaluate the extent



to which each one has achieved these two stated goals.

To focus this presentation, we consider only the simplest version of the edit distance problem, where both strings have equal length  $n$  and are over a constant sized alphabet; moreover, the algorithm only needs to return the edit distance and not the sequence of edits that achieve it. We also focus only on the runtime analysis as opposed to the memory use. We start by summarizing the state-of-the-art practical implementations then go through the various types of theoretical analysis that has been applied to the edit distance

problem: traditional worst-case analysis, worst-case analysis parametrized by the edit distance, worst-case analysis parametrized by entropy and compressibility, average-case analysis, semi-random models, and advice-based models. For each technique, we evaluate the extent to which it has achieved the prediction and design goals of theoretical analysis on the edit distance computation problem. We will not assume any knowledge of biology or any knowledge of computer science beyond an undergraduate level. We then conclude with a discussion of open problems and their potential solutions.

### State-of-the-Art Implementations and Algorithms

Here, we will briefly outline the algorithms that are used in state-of-the-art implementations as well as highlight their empirical performance. The classical algorithm taught in many Algorithms courses is called Needleman-Wunsch.<sup>43</sup> It builds a two-dimensional matrix  $D$  where the value at  $D[i, j]$  is the edit distance between the  $i$ -long prefix of  $A$  and the  $j$ -long prefix of  $B$ . This matrix can be computed in the standard dynamic programming manner using the recurrence  $D[i, j] = \min(D[i - 1, j] + 1, D[i, j - 1] + 1, D[i - 1, j - 1] + I_{i,j})$ , where  $I_{i,j}$  is 0

if the  $i^{\text{th}}$  character of  $A$  is equal to the  $j^{\text{th}}$  character of  $B$  and 1 otherwise. The edit distance is then the value at  $[n, n]$  (recall that  $n$  is the length of the strings). Doubling banded alignment<sup>41</sup> is a modification of Needleman-Wunsch whose main idea is to reduce the number of cells of  $D$  that need to be computed. It uses the idea that if one only computes the values of  $D[i, j]$  within a diagonal band of width  $d$  (that is, when  $|i - j| \leq d/2$ ), then by checking if  $D[n, m] \leq d$  one can either determine the edit distance if it is at most  $d$  or determine the edit distance is greater than  $d$ . Using this idea, it runs the checking algorithm repeatedly by doubling the value of  $d$  until the edit distance is found. The Myers' bit-parallel technique<sup>32</sup> is a hardware optimization of Needleman-Wunsch that encodes the dynamic programming matrix into bitvectors and then rewrites the recurrences in terms of word-sized bitvector operations. These three algorithms/techniques comprise the core of all state-of-the-art implementations today.

There are at least three broadly used software libraries/tools that implement edit distance computation.<sup>14,15,40</sup> Edlib<sup>40</sup> is optimized specifically for the edit distance problem, while SeqAn<sup>15</sup> and Parasail<sup>14</sup> are designed for more general alignment problems, but support edit distance as a special case. Edlib and SeqAn both implement the banded alignment algorithm using Myers' bit-parallel technique.<sup>32</sup> Myers' technique does not change the asymptotic runtime but gives a significant constant speedup in practice. Parasail<sup>14</sup> implements the Needleman-Wunsch algorithm using high-performance computing techniques. These include both task-level parallelism (that is, multi-threading) and instruction-level parallelism (that is, SIMD vectorized instructions). Parasail's code is also customized during compilation for the instruction set of the host architecture. Another implementation, BGSA,<sup>44</sup> implements the Needleman-Wunsch algorithm with the Myers' bit-parallel technique but supports multi-core, task-level, and instruction-level parallelism for batch execution.

There are also approaches to speed-up edit distance by using specialized hardware, such as GPUs, FPGAs, or even custom-designed processors (for references, see Alser et al.<sup>3</sup>). These re-

**In the real world, there is often a wide gap between the practical and theoretical communities; understanding how to close this gap is critical to making theoretical computer science more relevant to applications.**

sult in orders-of-magnitude constant-time speedups over their CPU counterparts in practice. However, until there is more widespread availability and integration of such specialized hardware in bioinformatics compute infrastructures, these tools are unlikely to be widely used.

How well do the widely used implementations perform? On two sequences of 1 million nucleotides each, one of the fastest implementations (edlib) takes 1.1 seconds for sequences with edit distance of  $0.01n$  and 30 seconds for sequences with edit distance of  $0.40n$ , on a single core server.<sup>40</sup> For sequences of 100,000 nucleotides each, the runtimes are 0.01s and 0.40s, respectively. As we will later see, this corresponds to the theoretical prediction that the runtime deteriorates with increasing edit distance. For many applications, these runtimes are good enough. However, edit distance computation remains a bottleneck for applications that use it as a subroutine to make thousands or millions of comparisons (for example, comparing long reads against each other<sup>38</sup>).

### Traditional Worst-Case Analysis

The most common way to analyze running time, taught in undergraduate computer science classes, is traditional worst-case analysis. For example, it says the classical merge sort algorithm runs in  $O(n \log n)$  worst-case time, which formally means there exists a constant  $c$  such that for any large-enough input of  $n$  elements, merge sort takes at most  $cn \log n$  time. Has traditional worst-case analysis led to the design of edit distance algorithms that perform well in practice? There are two candidates. The first is the classical Needleman-Wunsch algorithm; it was originally described in Wagner and Fischer,<sup>43</sup> which gave the dynamic programming recurrence and proved the runtime is  $\Theta(n^2)$ . Their algorithm modified an earlier dynamic programming algorithm<sup>35</sup> whose run time was  $\Theta(n^3)$ .<sup>a</sup> It seems likely then that traditional worst-case analysis was a

<sup>a</sup> As a historical note, the algorithm presented in Needleman and Wunsch<sup>35</sup> is not the algorithm we call “Needleman-Wunsch” today. The “Needleman-Wunsch” algorithm was described later.<sup>43</sup>

driving force behind the creation of the Needleman-Wunsch algorithm. Moreover, the algorithm is a success in practice because many implemented algorithms, including some of the fastest ones, are either modifications of it or use it as a subroutine. This includes both banded alignment and the Myers' bit-parallel technique used by edlib, parasail, and SeqAn. Thus, though Needleman-Wunsch is not the fastest algorithm in practice or in theory, it exemplifies how traditional worst-case analysis achieved the design goal.

The second candidate is the fastest known algorithm under traditional worst-case analysis—the Four-Russians speedup to Needleman-Wunsch.<sup>6,29</sup> It takes  $\Theta(n^2/\log^2 n)$  time (in a unit-cost RAM model). The algorithm was clearly designed to optimize the runtime under traditional worst-case analysis. But how does it perform in practice? Does the  $\Theta(\log^2 n)$  speedup outweigh the additional constant factors due to higher algorithm and data structure complexity? To answer this question, there have been implementations and experimental evaluations of this algorithm.<sup>22,36</sup> The improvement over Needleman-Wunsch was a factor of about 5 for  $n = 2^{18}$ , and, extrapolating from Rejmon,<sup>36</sup> would not exceed 10 even for sequences of a billion characters. Thus, in practice, the Four Russians algorithm is dominated by other algorithms that have  $\Theta(n^2)$  run time but have better constant factors (for example, Myers' bit-parallel algorithm).<sup>36</sup> Moreover, the fastest implementations of edit distance today<sup>14,15,40</sup> do not implement the Four Russians algorithm, even as a subroutine, highlighting how exploiting the properties of the CPU (for example, the bit-parallel or SIMD implementations) can bring constant speedups that in practice outperform asymptotic speedups. We therefore conclude that in the case of the Four Russians algorithm, traditional worst-case analysis has led us astray into the design of an algorithm that is not practically useful.<sup>b</sup>

<sup>b</sup> While the algorithm itself was not practically useful, one could argue it had an impact on practice because some of its ideas were later used by algorithms such as Myers' bit-parallel algorithm. While we may deem an algorithm not practically useful, it may nevertheless have been an important steppingstone on the road to another practically useful algorithm.

Has traditional worst-case analysis been able to accurately predict the empirical performance of algorithms? The analysis of Needleman-Wunsch shows that  $\Theta(n^2)$  time is taken for every input instance, not just in the worst-case. The Four-Russians runtime analysis is similar in that it also holds for all inputs, not just worst-case ones. Therefore, the predictions of traditional worst-case analysis accurately reflect these algorithms' runtime on real data.

Can traditional worst-case analysis lead to the design of new algorithms that perform well in practice? A famous recent result states that under the strong exponential time hypothesis, there cannot be an  $(n^{2-\delta})$  algorithm, for any  $\delta > 0$ .<sup>7</sup> Such an algorithm is called strongly sub-quadratic. There are other “barrier” results of this type which we will not elaborate on.<sup>1,2,10</sup> These results make it unlikely that better algorithms can be designed using traditional worst-case analysis as a guide. However, they say nothing about the existence of provably and substantially better algorithms, if they are analyzed using a different technique than traditional worst-case analysis.

### Worst-Case Analysis Parametrized by the Edit Distance

One step away from traditional worst-case analysis is parametrized worst-case analysis, which is worst-case analysis in terms of properties of the input besides just its size  $n$ . In the case of edit distance, the parameter that has proven most useful is the edit distance itself, usually denoted by  $k$  (note that  $k \leq n$ ). The most notable algorithm designed using  $k$ -parametrized worst-case analysis is doubling banded alignment.<sup>41</sup> Parametrized analysis shows that it computes the edit distance in time  $\Theta(kn)$ , on all inputs. There are several other algorithms parametrized by  $k$ ,<sup>19,33</sup> they achieve various trade-offs between  $k$  and  $n$  and have some other differences outside the scope of this article. Doubling banded alignment is the most notable of these because it is very simple to describe and implement (that is, it does not use any complex data structures) and forms a main component of one of the empirically fastest implementations today—edlib.<sup>40</sup>

$K$ -parametrized analysis predicts several ways in which doubling banded

alignment is theoretically an improvement over Needleman-Wunsch. The first is when  $k = o(n)$ , doubling banded alignment scales sub-quadratically with the input size, while Needleman-Wunsch does not. The second is the closer the sequences are (that is, the smaller the edit distance), the smaller the runtime. Thus, the algorithm captures a natural notion of complexity of the input and takes advantage of input that is less complex, while Needleman-Wunsch does not. These two predictions are reflected empirically as well since the runtime analyses are for all inputs and both algorithms do not hide any significant implementation constants. Third, even when  $k = \Theta(n)$ , the empirical advantage of doubling banded alignment over Needleman-Wunsch can be significant.<sup>40</sup>

However, when  $k = \Theta(n)$ , the fact remains the banded algorithm scales quadratically with the input size, both in theory and in practice. Unfortunately, this is the case for the following predominant application of edit distance. Biological sequences evolve from each other via a mutation process, which, for the purposes of this discussion, can be considered as mutating each position with some constant probability, independently for each position. The mutation probability is, for the most part, independent of  $n$  and the edit distance is therefore proportional to  $n$  with a constant called divergence. For example, the sequence divergence in coding regions of genes between human and other species is 1%–2% for bonobo and about 19% for mouse.<sup>12</sup> Thus, regardless of species, and even for species that are very close, the edit distance is still approximately a constant proportion of the sequence length. This illustrates the importance of using theoretical analysis to design exact algorithms that scale sub-quadratically when  $k = \Theta(n)$ .

In summary, the  $k$ -parameterized worst-case analysis technique has been a tremendous success for edit distance. It led directly to the design of the banded alignment algorithm, which is widely used in practice, and it can predict the empirical improvement of banded alignment over Needleman-Wunsch. However, it still did not produce an algorithm that scales strongly sub-quadratically with the input size

for applications where  $k = \Theta(n)$ . Moreover, it may not be able to do so in the future due to the barrier results against strongly sub-quadratic algorithms in the worst-case framework.<sup>1,2,7,10</sup>

### Worst-Case Analysis Parametrized by Entropy and Compressibility

Another way to parameterize the analysis of edit distance algorithms is by the entropy of the input,  $h$ . Entropy is a value between 0 and 1 which measures the amount of order in the strings.<sup>34</sup> Strings containing short repetitive patterns tend to have a lower entropy. As an extreme case, the string of all  $T$ s has entropy 0 while a string generated uniformly at random has entropy close to 1. Intuitively, an edit distance algorithm could take advantage of the repetitive patterns to run faster on strings with lower entropy. An algorithm following this intuition was developed by Crochemore et al.<sup>13</sup> and runs in time  $O(hn^2/\log n)$  for most inputs. When the input strings have low entropy, this is theoretically faster than Needleman-Wunsch. Note that this runtime is not comparable with banded alignment, since the edit distance can be low for strings with high entropy, and vice versa.

A related notion are algorithms that compute the edit distance of  $A$  and  $B$  directly from their compressed representations. The time to compress two strings is in most cases asymptotically negligible compared to the time it takes to compute the edit distance. Hence, one could solve the edit distance computation problem by first compressing  $A$  and  $B$  and then running an edit distance algorithm on the compressed strings. For example, there is an algorithm to compute edit distance between two run-length encoded<sup>34</sup> strings in  $O(\ell n)$  time<sup>5,13,26</sup> where  $\ell$  is the size of the encoded strings. For a more general class of compression algorithms called straight-line programs,<sup>20</sup> there is an algorithm<sup>17</sup> that runs in  $O(\ell n\sqrt{\log(n/\ell)})$  time. In these cases, the algorithms are designed to optimize the runtime with respect to  $\ell$ -parametrized worst-case analysis.

Unfortunately, these algorithms have not been broadly applied in practice. A major reason is because a long DNA sequence, while exhibiting some

repetitive patterns, still has high entropy (for example, 0.85<sup>25</sup>) and low compressibility.<sup>c</sup> Thus, while worst-case analysis parametrized by entropy or compressibility has led to the design of novel algorithms, it has not achieved the design or prediction goals of theoretical analysis since these algorithms have not been useful in practice.

### Average-Case Analysis

One of the pitfalls of traditional and parametric worst-case analysis is that it assumes the inputs are chosen by an adversary who wants to make things as difficult as possible for the algorithm. In practice, however, biological sequences are not chosen this way and may have much nicer properties. On such instances, the optimal path in the alignment matrix stays very close to the diagonal, even much tighter than the band of width  $d$  used by the doubling banded alignment algorithm. For example, a heuristic algorithm that stops the doubling banded alignment prematurely and reports the value of  $[n, n]$  would still likely report the correct edit distance. Worst-case analysis cannot tell us either the probability of success nor tell us how long to double to guarantee a high probability of success.

One theoretical analysis technique to alleviate the shortcoming of worst-case analysis is average-case analysis (sometimes called distributional analysis). Here, the inputs are assumed to be drawn from some kind of distribution, and what is measured is the expected performance over this distribution. Other alternatives include measuring the performance that can be achieved with high probability and/or measuring the performance in terms of a trade-off with the probability of the algorithm achieving it. Reanalyzing the performance of existing algorithms under this model would not be beneficial for most of the algorithms we have looked at so far (that is, Needleman-Wunsch, Four Russians, and banded alignment) because their runtime holds for all inputs, not just worst-case ones. However, if we have an input distribution in mind, we can

design a new algorithm to work well under that model and validate it empirically. Researchers have aimed to design an algorithm with a strongly sub-quadratic expected runtime.

An obvious first attempt is to assume each string of length  $n$  is randomly drawn uniformly and independently from the universe of all strings of length  $n$ . For each string, this corresponds to generating each character independently, following an identical categorical distribution for each position. However, this model does not capture essential properties of real data and has not resulted in any useful algorithms or predictive running time analysis. When the real input strings are evolutionary related, then the assumption they are independent of each other is highly inaccurate. Even in the case that the real input strings are not evolutionary related, the uniformity assumption remains unrealistic since biological sequences have evolved to serve a function and thus exhibit non-random behavior.

A more accurate model uses an indel channel.<sup>16</sup> Here, the first string is chosen uniformly at random, while the second string is obtained by randomly mutating each nucleotide, with a probability at most a small constant. The algorithm of Ganesh and Sy<sup>16</sup> runs in  $O(n \log n)$  time and, with high probability over this input distribution, returns the correct edit distance. This algorithm's runtime is significantly better than just strongly sub-quadratic. However, we are not aware of any implementation, and, for this algorithm to be practical, it must prove it performs reasonably well even for real input that does not follow the model's distribution.

A successful application of average case analysis is in the analysis of the “furthest reaching” algorithm, which was proposed in Myers<sup>31</sup> and Ukkonen.<sup>41</sup> This analysis combines average-case with parametrized analysis. It runs in  $O(nk)$  time in the worst case,<sup>41</sup> but an average case analysis gives  $O(n + k^2)$ .<sup>31</sup> The model used in this analysis is like the one previously mentioned, which describes the model of Ganesh and Sy.<sup>16</sup> While this algorithm does not break the quadratic time barrier and is not used by state-of-the-art edit distance algorithms, it has nevertheless been practically successful. It was the

<sup>c</sup> If compressing multiple DNA sequences together, much better compressibility is possible. However, edit distance has not been typically performed against such collections.

basis of an implementation of the “diff” Unix tool<sup>31</sup> and serves as a foundation for generalizations of the edit distance problem (for example, to affine gap penalties as in Marco-Sola et al.<sup>27</sup>).

### Semi-Random Models

One can view average-case analysis and worst-case analysis as two extremes. The main idea of more sophisticated semirandom models is to achieve a middle ground between an adversary and randomness.<sup>4,8,23</sup> Here, the performance is measured as worst-case over the choices of the adversary and average case over the random distribution. The semi-random models for edit distance have only led to approximation, rather than exact, algorithms. A  $c$ -approximation algorithm is one that returns a value at most a multiplicative factor of  $c$  (called the approximation ratio) away from the edit distance. Approximation algorithms relax the requirement of finding the exact solution in exchange for better runtime. For an edit distance approximation algorithm to be relevant in practice, the approximation ratio must be a constant very close to one for example, 1.01). For example, even a 3-approximation algorithm would not be able to always distinguish two random DNA sequences (that is, expected edit distance of approximately  $0.53n$  (personal simulations, data not shown)) from a mouse and human sequence pair (that is, edit distance of about  $0.19n$ ).<sup>12</sup>) Thus, the usefulness of the models is predicated on their ability to achieve a tiny approximation ratio.

**Smoothed analysis** is based on the idea that worst-case inputs designed to fool algorithms are not very stable; that is, if bad input is tweaked a little, the algorithm no longer performs poorly.<sup>37</sup> In this model, an adversary first picks an input (that is, a worst-case choice), but then some small random noise is added to this input. The algorithm’s performance on a particular adversarial input is defined as the expected performance over the distribution of noisy inputs centered around the adversarial input. The algorithm’s overall performance is then defined as the worst-case performance over all choices of the adversary’s input.

Smoothed analysis was considered for edit distance in Andoni and

**Smoothed analysis is based on the idea that worst-case inputs designed to fool algorithms are not very stable; that is, if bad input is tweaked a little, the algorithm no longer performs poorly.**

Krauthgamer.<sup>4</sup> In their model, an adversary chooses the two input strings and a longest common subsequence between them. Then, each position is randomly perturbed with a small probability  $p$  (that is, the nucleotide is replaced with another random one). However, the perturbations are constrained so the positions of the longest common subsequence are perturbed identically. This model captures the idea that for two evolutionary related strings, if we view their longest common subsequence as their ancestral sequence, then all the nucleotides outside this common subsequence would have evolved somewhat independently of each other. Allowing the adversary to choose some worst-case values for them gives her too much power; instead, some noise is added to them to make them more independent.

In Andoni and Krauthgamer,<sup>4</sup> the algorithms’ approximation ratios are not precisely derived (that is, big-Oh notation is used). There has also been another line of work focusing on fast approximation algorithms under traditional worst-case analysis with the best-known approximation constant of 1680.<sup>11</sup> In both cases, the algorithms achieve strongly sub-quadratic run time, something that could not be done by exact algorithms using worst-case analysis for  $k = \Theta(n)$ . However, without a precise derivation of a tiny approximation ratio, or an implementation and validation on real data, it is difficult to predict the applicability of these algorithms.

**Asymmetric  $(p,B)$ -pseudo-random model.** As we mentioned previously, a string chosen at random does not reflect a biologically evolved sequence. More precisely, a uniformly random string has the property that the probability that any two non-overlapping equal-length substrings are identical decreases with their length. In fact, once their length exceeds a certain critical threshold, this probability is, for all practical purposes, zero. This property is in contrast with biological sequences, which are often composed of long similar elements called repeats. It is true that longer repeats tend to be less frequent and less similar, but this decrease does not happen at the same rate as for random strings. This was captured in a more realistic

model of randomness, proposed in Kuszmaul.<sup>23</sup> They say that a string is  $(p, B)$ -pseudorandom if the edit distance of any two disjoint  $B$ -long substrings is at least a fraction  $p$  of their length. This generalizes uniform randomness, that is, a uniformly random string is  $(\Omega(1), O(\log n))$ -pseudorandom with high probability.<sup>23</sup> But by choosing  $p$  and  $B$  appropriately, we can more realistically match the repeat properties of a biological string.

The asymmetric  $(p, B)$ -pseudorandom model<sup>23</sup> is to first choose a string at random from all  $(p, B)$ -pseudorandom strings and then have an adversary choose the other string and modify some small portion of the pseudorandom string. By allowing the adversary to choose one of the strings, this model allows the two strings to be evolutionary related. Moreover, the additional power of the adversary to modify the pseudorandom string makes the model even more realistic, because a true biological sequence would usually have some substrings that break the  $(p, B)$ , even when the values of  $p$  and  $B$  are chosen to minimize these cases.

This model is used in Kuszmaul<sup>23</sup> to design several new algorithms. The main algorithm has a runtime  $\tilde{O}(nB)$  (the  $\tilde{O}$  notation is like  $O$  but ignores log factors), which is strongly sub-quadratic if  $B$  is strongly sub-linear. However, as in smoothed analysis, the algorithm is only an approximation algorithm, and the exact approximation ratio is not calculated. Thus, whether this model leads to any practical algorithms remains to be seen. However, the fact that it seems to intuitively better capture biological reality while lending itself to theoretical analysis is promising.

**The random model of Boroujeni et al.** A recently proposed model<sup>8</sup> has the adversary first choose a seed string  $s$  and then constructs  $A$  by permuting  $s$  uniformly at random. After observing  $A$  and the random permutation, the adversary then constructs  $B$ . The algorithm given by Boroujeni et al.<sup>8</sup> is an approximation algorithm with an expected runtime is  $O(n^{1.898})$ , which is strongly sub-quadratic. The approximation ratio is  $1 + o(1)$ , which is low enough to be of practical relevance.

This analysis model has led to the design of an algorithm that brakes

through a barrier of previous models, that is, it achieves a low approximation ratio while maintaining sub-quadratic time. However, the actual runtime improvement is only  $n^{0.102}$ , which is less than 9 for inputs up to a billion nucleotides long. It is unlikely that this improvement would justify the additional overhead of a more complex algorithm. Nevertheless, the model can ultimately be successful if the runtime of the algorithm can be improved, the implementation of the algorithm kept simple, and the usefulness of the model empirically validated.

In summary, none of the semi-random models have yet led to a better understanding of the performance of existing algorithms or to the design of algorithms that perform well in practice. The proposed algorithms, at least in their current form, are not promising: the algorithms of Andoni and Krauthgamer<sup>4</sup> and Kuszmaul<sup>23</sup> have impractical approximation ratios, while the Boroujeni et al. algorithm<sup>8</sup> improves the runtime by a factor that is too small to have an effect in practice. However, the models themselves are promising, and can ultimately be successful if the runtime of the algorithms can be improved, the complexity of the algorithms kept low, and the usefulness of the models empirically validated. These ongoing efforts may eventually result in an algorithm that outperforms banded alignment, in practice, on inputs with  $k = \Theta(n)$ .

### Analyzing with Advice

The biological problem is usually more general than the mathematical abstraction created for it. Sometimes the algorithm has access to other information, not included in the problem definition, that can serve as advice. In this case, one can both expand the problem definition and the theoretical analysis to incorporate this advice. An advice-based analysis measures the algorithm runtime with respect to the amount of such advice used. In the case of edit distance, Goldwasser and Holden<sup>18</sup> argue that, for an input instance  $A$  and  $B$ , it is possible to have access to a collection of correlated instances. Intuitively, a correlated instance is one whose sequence of edits in the shortest edit sequence is, with some high probability, like the

one between  $A$  and  $B$ . They show that if an algorithm has access to  $O(\log n)$  of such instances, it can find the edit distance between  $A$  and  $B$  in  $O(n \log n)$  time, with high probability.

This approach has not been implemented but is promising in the sense the runtime is not just sub-quadratic but nearly linear, and the algorithm is exact and seems easy to implement. Compared to the banded alignment algorithms, the  $O(n \log n)$  algorithm is likely to have significant empirical speed improvements for moderately sized inputs even for small values of  $k$ . Unfortunately, it is not clear how correlated instances can be obtained in practice and whether they would be captured by the Goldwasser and Holden definition.<sup>18</sup>

### Conclusion

We have surveyed the various approaches to the theoretical analysis of edit distance algorithms, focusing on whether these approaches have led to the design of algorithms that are fast in practice and to the theoretical prediction of the empirical runtimes of existing algorithms. We showed the track record has been mixed. On one hand, a few algorithms widely used in practice (Needleman-Wunsch, doubled banded alignment, furthest reaching) have been born out of theoretical analysis and their empirical performance is captured well by theoretical predictions. On the other hand, very little of the algorithms developed using theoretical analysis as a yardstick since then have had any practical relevance.

From a practical perspective, a major open problem is to implement an algorithm with linear-like empirical scaling on inputs where the edit distance is linear in  $n$ . Theoretical analysis has the potential to lead the way in achieving this goal. Semi-random models are the most promising approach, due to the barrier results for worst-case analysis. A reasonable model which might give a good balance between capturing reality and ease of analysis is one where  $B$  is assumed to have evolved from  $A$  via a mutation process (like Ganesh and Sy<sup>16</sup>). For the theoretical work to have practical relevance, however, it must be implemented and validated; in particular, the algorithm's runtime

and accuracy must be robust to data not drawn from the modeled distribution and the constant-time overhead of the algorithm must not override the asymptotic gains. Unfortunately, current community incentives leave the implementation and validation of an algorithm as an afterthought that is rarely pursued.

To solve this open problem and, more generally, close the gap between theory and practice, implementation and validation cannot be treated as a separate step of the process. We need multi-disciplinary teams that can interleave the theoretical analysis of algorithms with their implementation and validation. Allowing for a back-and-forth between practitioners and theoreticians during the development process can allow iterations over the theoretical model and practical heuristics that would otherwise be impossible. Such teams will be able to use both empirical and theoretical performance as a yardstick and will be better able to develop algorithms whose empirical performance is not only superior but also accurately captured by theoretical analysis.

It is also essential to continue to study the relationship between theoretical analysis and practical implementations. One of the limitations of this study is that it only focuses on a single problem, making it difficult to draw more general conclusions. More studies such as this one could establish patterns and identify clear directions to closing the practice/theory gap. In a recent paper,<sup>30</sup> we applied the same lens more broadly (but also more anecdotally) to study the analysis of algorithms in sequencing bioinformatics, attempting to establish patterns and draw conclusions for that domain. Ultimately, the relationship between theory and practice must be understood not only from a technical angle but also from a social science and philosophy of science perspective. In other disciplines, this relationship is studied by philosophers of science (for example, in education policy<sup>21</sup>) and so a similar approach may be fruitful in computer science.

## Acknowledgments

This article was inspired by the online videos for the 2014 class “Beyond

Worst-Case Analysis” by Tim Roughgarden. This material is based upon work supported by the National Science Foundation under Grants No. 1439057, 1453527, and 1931531. □

## References

1. Abboud, A., Backus, A., and Williams, V.V. Tight hardness results for LCS and other sequence similarity measures. In *Proceedings of the 2015 IEEE 56<sup>th</sup> Annual Symp. Foundations of Computer Science*, 59–78.
2. Abboud, A., Hansen, T.D., Williams, V.V., and Williams, R. Simulating branching programs with edit distance and friends: Or a polylog shaved is a lower bound made. In *Proceedings of the 48<sup>th</sup> Annual ACM Symp. Theory of Computing*, (2016), 375–388.
3. Alser, M., Hassan, H., Kumar, A., Mutlu, O., and Alkan, C. Shouji: A fast and efficient pre-alignment filter for sequence alignment. *Bioinformatics* 35, 21 (2019), 4255–4263.
4. Andoni, A. and Krauthgamer, R. The smoothed complexity of edit distance. *ACM Trans. Algorithms* 8, 4 (2012), 1–25.
5. Arbell, O., Landau, G.M., and Mitchell, J.S.B. Edit distance of run-length encoded strings. *Information Processing Letters* 83, 6 (2002), 307–314.
6. Arlazarov, V., Dinitz, Y., Kronrod, M., and Faradzhev, I. On economical construction of the transitive closure of an oriented graph. *Doklady Akademii Nauk* 194. Russian Academy of Sciences, (1970), 487–488.
7. Backurs, A. and Indyk, P. Edit distance cannot be computed in strongly subquadratic time (unless SETH is false). *SIAM J. Computing* 47, 3 (2018), 1087–1097.
8. Boroujeni, M., Seddighin, M., and Seddighin, S. Improved algorithms for edit distance and LCS: Beyond worst case. In *Proceedings of the 14<sup>th</sup> Annual ACM-SIAM Symp. Discrete Algorithms*, (2020), 1601–1620.
9. Brinda, K., Boeva, V., and Kucherov, G. Ococo: An online consensus caller, (2017); arXiv:1712.01146.
10. Bringmann, K. and Künemann, M. Quadratic conditional lower bounds for string problems and dynamic time warping. In *Proceedings of the 2015 IEEE 56<sup>th</sup> Annual Symp. Foundations of Computer Science*, 79–97.
11. Chakraborty, D., Das, D., Goldenberg, E., Koucký, M., and Saks, M. Approximating edit distance within constant factor in truly sub-quadratic time. In *Proceedings of the 2018 IEEE 59<sup>th</sup> Annual Symp. Foundations of Computer Science*.
12. Cooper, G.M., Brudno, M., and NISC Comparative Sequencing Program. Quantitative estimates of sequence divergence for comparative analyses of mammalian genomes. E.D. Green, S. Batzoglou, and A. Sidow, (eds). *Genome Res.* 13, 5 (May 2003), 813–820.
13. Crochemore, M., Landau, G.M., and Ziv-Ukelson, M. A subquadratic sequence alignment algorithm for unrestricted scoring matrices. *SIAM J. Computing* 32, 6 (2003), 1654–1673.
14. Daily, J. Parasail: SIMD C library for global, semi-global, and local pairwise sequence alignments. *BMC Bioinformatics* 17, 1 (2016), 81.
15. Döring, A., Weese, D., Rausch, T., and Reinert, K. Seqan: an efficient, generic C++ library for sequence analysis. *BMC Bioinformatics* 9, 1 (2008), 11.
16. Ganesh, A. and Sy, A. Near-linear time edit distance for indel channels. In *Proceedings of the 20<sup>th</sup> Intern. Workshop on Algorithms in Bioinformatics* 172, Leibniz International Proceedings in Informatics. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, (2020), 17:1–17:18.
17. Gawrychowski, P. Faster algorithm for computing the edit distance between SLP-compressed strings. In *Proceedings of the 2012 Intern. Symp. String Processing and Information Retrieval*. Springer, 229–236.
18. Goldwasser, S. and Holden, D. The complexity of problems in p given correlated instances. In *Proceedings of the 8<sup>th</sup> Innovations in Theoretical Computer Science Conf.* Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, (2017).
19. Gusfield, D. *Algorithms on Strings, Trees, and Sequences*. (1997).
20. Hermelin, D., Landau, G.M., Landau, S., and weimann, O.A. A unified algorithm for accelerating edit-distance computation via textcompression, (2009); arXiv:0902.2649.
21. Joyce, K.E. and Cartwright, N. Bridging the gap between research and practice: Predicting what will work locally. *Amer. Educ. Research J.* 57, 3 (2020), 1045–1082.
22. Kim, Y., Na, J.C., Park, H., and Sim, J.S. A space-efficient alphabet-independent Four-Russians' lookup table and a multithreaded Four-Russians' edit distance algorithm. *Theoretical Computer Science* 656, (2016), 173–179.
23. Kuszmaul, W. Efficiently approximating edit distance between pseudorandom strings. In *Proceedings of the 30<sup>th</sup> Annual Symp. on Discrete Algorithms*. SIAM, (2019), 1165–1180.
24. Levenshtein, V.I. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady* 10, 707–710.
25. Loewenstein, D. and Yianilos, P.N. Significantly lower entropy estimates for natural DNA sequences. *J. Computational Biology* 6, 1 (1999), 125–142.
26. Mäkinen, V., Ukkonen, E., and Navarro, G. Approximate matching of run-length compressed strings. *Algorithmica* 35, 4 (2003), 347–369.
27. Marco-Sola, S., Moure, J.C., Moreto, M., and Espinosa, A. Fast gap-affine pairwise alignment using the wavefront algorithm. *Bioinformatics* 37, 4 (Feb. 2021), 456–463.
28. Marett, L. et al. Sequencing de novo assembly of 150 genomes from Denmark as a population reference. *Nature* 548, 7665 (2017), 87–91.
29. Masek, W.J. and Paterson, M.S. A faster algorithm computing string edit distances. *J. Computer and System Sciences* 20, 1 (1980), 18–31.
30. Medvedev, P. The limitations of the theoretical analysis of applied algorithms. (2022); arXiv:2205.01785.
31. Myers, E.W. An o(nd) difference algorithm and its variations. *Algorithmica* 1, 1–4 (1986), 251–266.
32. Myers, G. A fast bit-vector algorithm for approximate string matching based on dynamic programming. *JACM* 46, 3 (1999), 395–415.
33. Navarro, G. A guided tour to approximate string matching. *ACM Computing Surveys* 33, 1 (2001), 31–88.
34. Navarro, G. *Compact Data Structures: A Practical Approach*. Cambridge University Press (2016).
35. Needleman, S.B. and Wunsch, C.D. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Molecular Biology* 48, 3 (1970), 443–453.
36. Rejmon, M. Multi-Threaded implementation of four russians edit distance algorithm. *Bachelor's thesis*. Czech Technical University in Prague, Faculty of Information Technology, (2019).
37. Roughgarden, T. Beyond worst-case analysis. *Commun. ACM* 62, 3 (Mar. 2019), 88–96.
38. Sahlin, K. and Medvedev, P. De novo clustering of long-read transcriptome data using a greedy, quality value-based algorithm. *J. Computational Biology* 27, 4 (2020), 472–484.
39. Sarkar, H., Zakeri, M., Malik, L., and Patro, R. Towards selective-alignment: Bridging the accuracy gap between alignment-based and alignment-free transcript quantification. In *Proceedings of the 2018 ACM Intern. Conf. Bioinformatics, Computational Biology, and Health Informatics*, 27–36.
40. Šošić, M. and Šikić, M. Edlib: a C/C++ library for fast, exact sequence alignment using edit distance. *Bioinformatics* 33, 99 (2017), 394–395.
41. Ukkonen, E. Algorithms for approximate string matching. *Information and Control* 64, 1–3 (1985), 100–118.
42. Vaser, R., Sovic, I., Nagarajan, N., and Sikic, M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Research* 27, 5 (2017), 737–746.
43. Wagner, R.A. and Fischer, M.J. The String-to-String Correction Problem. *JACM* 21, 1 (1974), 168–173.
44. Zhang, J. et al. BGSA: A bit-parallel global sequence alignment toolkit for multi-core and many-core architectures. *Bioinformatics* 35, 13 (2019), 2306–2308.
45. Zook, J.M. et al. A robust benchmark for germline structural variant detection. *BioRxiv*. (2019), 664623.

**Paul Medvedev** is a professor in the Department of Computer Science and Engineering and the Department of Biochemistry and Molecular Biology, and the Director of the Center for Computational Biology and Bioinformatics at the Pennsylvania State University, University Park, PA, USA.



This work is licensed under a Creative Commons Attribution-NoDerivs International 4.0 License.

DOI:10.1145/3588998

## Understanding and handling interference across multiple active cameras.

BY JONGHO LEE, MOHIT GUPTA,  
BHUVANA KRISHNASWAMY, AND SUMAN BANERJEE

# When Two Cameras Are a Crowd

VISION AND ROBOTICS systems enabled by cameras that recover 3D scene geometry are revolutionizing several aspects of our lives via technologies such as autonomous transportation, robotic surgery, and ‘hands-free’ user interfaces. Modern 3D cameras are active devices, where a programmable light source emits coded illumination. The emitted light gets reflected from the scene and is received by a sensor to infer the 3D structure of the surroundings. In a multi-camera environment, such active 3D cameras may receive light from the sources of other cameras, resulting in large depth errors. This problem is becoming increasingly important due to the emergence of low-cost and compact active 3D cameras, which are becoming ubiquitous across a wide range of applications, from consumer devices to vehicular vision systems.

We observe that the multi-camera interference (MCI) problem shares several similarities and dissimilarities with common interference problems in the RF domain. Based on this observation, this article describes new and emerging challenges when multiple active 3D cameras operate in the same spatio-temporal region. The article also outlines some solutions, and more importantly, highlights the next steps.

### The 3D Revolution

We are in the midst of a 3D revolution fueled by cameras that can recover 3D geometry of their surroundings (Figure 1). The key catalyst driving this revolution is the emergence of low-cost, time-of-flight (ToF) 3D cameras that emit coded light and infer distances (depths) based on reflections from surrounding surfaces. ToF cameras can be made into extremely compact devices,<sup>a</sup> and thus, can potentially measure accurate 3D shapes over a wide area.

**Applications.** Due to their low cost, compact form factors and low computational complexity, ToF-based active 3D

<sup>a</sup> In comparison, 3D cameras based on other principles, such as binocular stereo and parallax, require large form factors and are unsuitable for long-range 3D imaging applications such as automotive navigation.

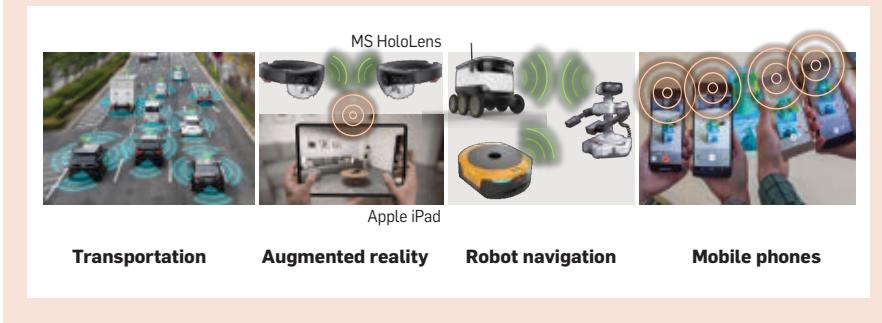
### » key insights

- 3D cameras are revolutionizing several aspects of our lives in many applications, such as autonomous vehicles, cellphones, tablets, AR and VR devices.
- 3D cameras reconstruct 3D geometry of the surroundings by emitting and receiving light. As multiple 3D cameras operate in the same spatiotemporal region, the light interference between cameras can cause large depth errors.
- Multi-camera interference (MCI) is becoming an important issue as 3D cameras become ubiquitous. The MCI problem shares similarities and dissimilarities with the interference problem in wireless communications.
- We describe several challenges and solutions for MCI based on these similarities and differences.

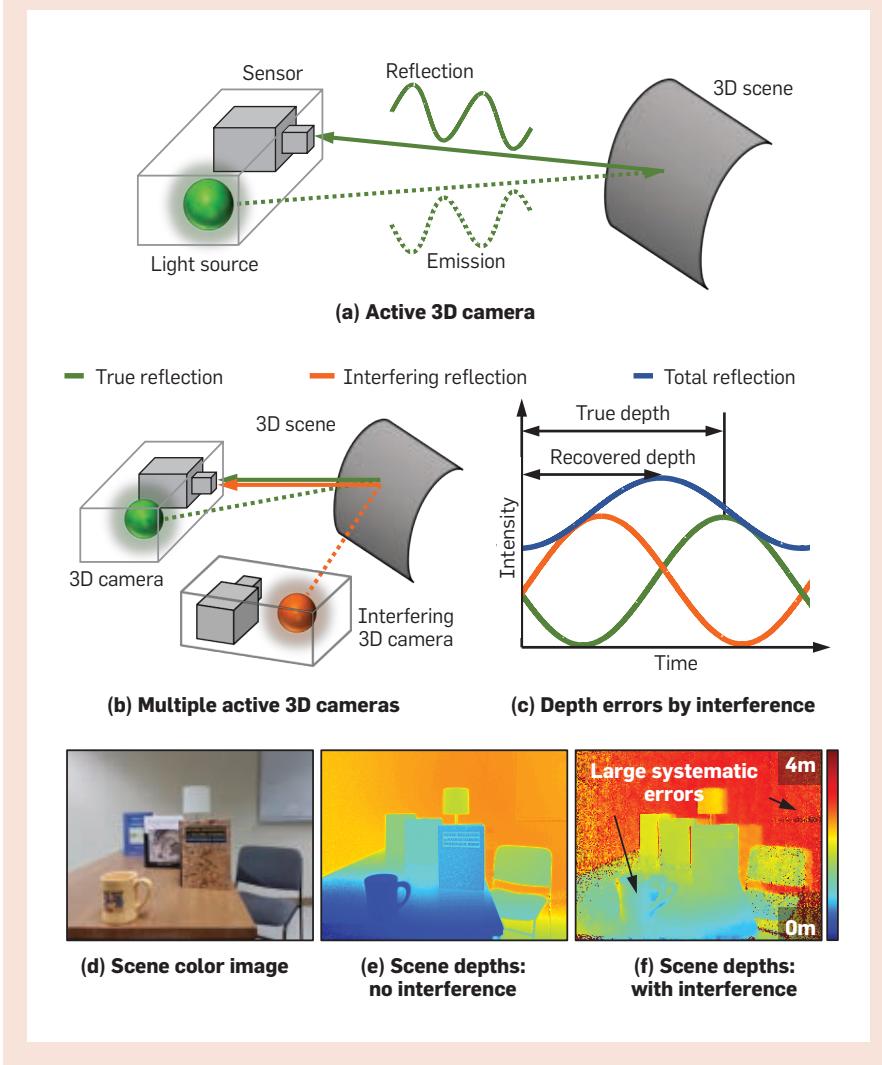


**Figure 1. The 3D revolution.**

Cameras that recover the 3D structure of their surroundings are becoming ubiquitous in several application domains, including transportation, augmented and virtual reality, mobile robotics, and phones.

**Figure 2. Multidevice interference and the resulting depth errors in active 3D cameras.**

- (a) An active 3D camera consists of a programmable light source and a sensor.
- (b) Multiple active 3D cameras simultaneously imaging the same scene point leads to MCI.
- (c) The estimated scene depths can be significantly different from the true depths due to interference.
- (d) An example scene.
- (e) Ground-truth scene depths.
- (f) Depths measured in the presence of an interfering camera have large systematic errors. *Part of the figure adapted from Lee and Gupta<sup>19</sup> with permission.*



cameras are now the method of choice in most commercial 3D imaging systems, including vehicular LiDARs, and more recently in commodity<sup>1,2</sup> and consumer-grade devices, such as cell phones, tablets (for example, Apple iPad), and gaming and AR/VR headsets (for example, Microsoft Kinect and Hololens). They are also used by inspection and delivery robots to navigate in warehouses<sup>24</sup> and perform 3D modeling,<sup>9,13</sup> directly impacting safety and functionality. Going forward, our dependence on them is only going to grow as these cameras will drive an even larger, potentially transformative set of applications such as human-machine interaction<sup>6</sup> and home robotics.

**Active 3D cameras and interference.** More generally, ToF cameras belong to the class of *active 3D cameras*, which consist of a *programmable light source* that emits spatially or temporally coded light.<sup>b</sup> For example, the light source could be a laser or an LED whose intensity is modulated over time. The emitted light signal travels to the scene of interest, gets reflected, and is captured by the sensor (typically co-located with the light source), as shown in Figure 2a. Scene depths (and hence, the 3D geometry) are extracted by comparing the emitted and the reflected light. Since the camera actively controls the illumination, it can reliably recover highly precise 3D geometry even in challenging real-world scenarios, including large lighting variations (bright sunlight to dark night sky) and optically uncooperative scenes (shiny materials, textureless and dark objects), which are otherwise difficult to handle for passive 3D cameras.

Although the coded light sources enable high-precision 3D recovery in isolation, **multiple active cameras**, each emitting their own light signals, can cause **mutual interference**. In such multi-use scenarios, an active camera's sensor may receive light emitted not just by its own source, but also by the sources of other cameras. This interfering signal prevents correct 3D depth estimation, resulting

<sup>b</sup> In contrast, passive 3D cameras (for example, stereoscopic 3D cameras) do not actively control the illumination, thus only passively observing the scene.

in potentially large, systematic depth errors. An example is shown in Figure 2b-f, where the scene depth observed by one active camera is significantly distorted in the presence of a second, interfering one. Such errors in low-level 3D measurements can make it near impossible to extract any actionable information (for example, object detection and recognition, image classification) for downstream machine-learning (ML) algorithms.

**Why is now a good time to address MCI?** Until recently, a large proportion of ToF cameras were based on a scanning laser beam that illuminates one (or a few) scene point at a time. Since they illuminate only a small portion of the scene at a time, the probability that two devices simultaneously illuminate the same point is relatively low (Figure 3a). This enabled multiple devices to co-exist without affecting each other. However, they require mechanical moving parts (for example, rotating mirrors), often leading to long acquisition times, high cost, bulkiness, and reliability issues.

To address its limitations, a new emerging class of solid-state ToF cameras *flood-illuminate* the entire scene (Figure 3b), making them considerably cheaper, smaller, and faster. Hence, solid-state cameras are fast replacing scanning-based cameras in autonomous driving and robotics applications (Figure 3c). But, there is a trade-off: Several active cameras flood-illuminating a scene simultaneously will interfere with each other (Figure 3b). Furthermore, a new generation of consumer devices (for example, cell phones) with 3D imaging capabilities are becoming ubiquitous. Due to stringent size and cost requirements, these devices also prefer solid-state 3D cameras (Figure 3d), which will create strong MCI whenever multiple active cameras are used in proximity.

Due to their compatibility with mainstream CMOS fabrication lines, the capabilities of solid-state cameras (spatial resolution, timing precision, signal-to-noise-ratio) continue to grow rapidly.<sup>10,23,33</sup> Therefore, these cameras are quickly becoming the method of choice in almost all applications that rely on 3D cameras, including vehicular navigation, robotics, and consumer mobile devices.

On the downside, this growth will add to the severity of MCI problems. Imagine several cars equipped with ToF camera-based LiDARs driving in close proximity, or an indoor setting with several consumer devices such as phones and headsets using their 3D cameras simultaneously for augmented reality (AR), localization, or 3D modeling.<sup>9,13</sup> It is, therefore, critical to address this problem now so that 3D cameras continue to function reliably across a broad spectrum of real-world scenarios, imaging devices, and application domains.

Therefore, managing interference

across a multitude of these devices will be of profound importance as solid-state ToF cameras become ubiquitous in our personal mobile devices, vehicles, homes, and workplaces, as well as the basis for numerous critical applications. The goal of this article is to shine a light on this important problem (pun intended) that has received little attention so far: *interference of light signals emitted by active cameras and its impact in accurately recovering 3D scene information*.

**Can research on wireless interference be used to address MCI?** Typical wireless communication systems

**Figure 3. Emergence of solid-state ToF and consumer-grade 3D cameras.**

- (a) With current LiDARs using a scanning laser beam, only a small portion of the scene is illuminated at a time. Thus, the likelihood of MCI is relatively low.
- (b) A new class of solid-state ToF devices is emerging where the light source flood-illuminates the entire scene without any moving parts, thus significantly increasing the likelihood of interference. (c) This new generation of solid-state 3D cameras is cheaper, smaller, and lighter, and they are fast replacing scanning-based devices. (d) Consumer devices with solid-state 3D cameras are becoming ubiquitous, increasing the possibility of strong interference in scenarios when multiple such devices are used simultaneously in a spatial neighborhood.



(for example, Wi-Fi, cellular systems) use radio waves propagating over a shared channel. These systems have dealt with multi-node interference for decades. This has led to the development of many techniques using whichever multiple wireless devices can co-exist and communicate over the same shared air medium. Hence, it is natural to evaluate whether strategies developed for mitigating wireless interference can be applied to address MCI. To answer the above question, we first discuss the similarities and

differences of how interference manifests in these settings.

At first glance, an active 3D camera is very similar to a wireless transmitter (Figure 4). A wireless transmitter generates radio waves that propagate over the air (usually referred to as the channel/medium) to a wireless receiver located at a distance with a goal of communicating information. The information itself is encoded in the radio waves. An active 3D camera also uses a transmitter and a receiver. The active 3D camera transmitter (or light

source) generates a signal (light), and a corresponding receiver (or sensor) captures a reflected instance of this signal. However, unlike the wireless communication systems, the goal of active 3D cameras is to estimate the depth, or distance, of objects in the environment based on the ToF of the received reflected signals. This could be considered as analogous to learning the channel in wireless systems.

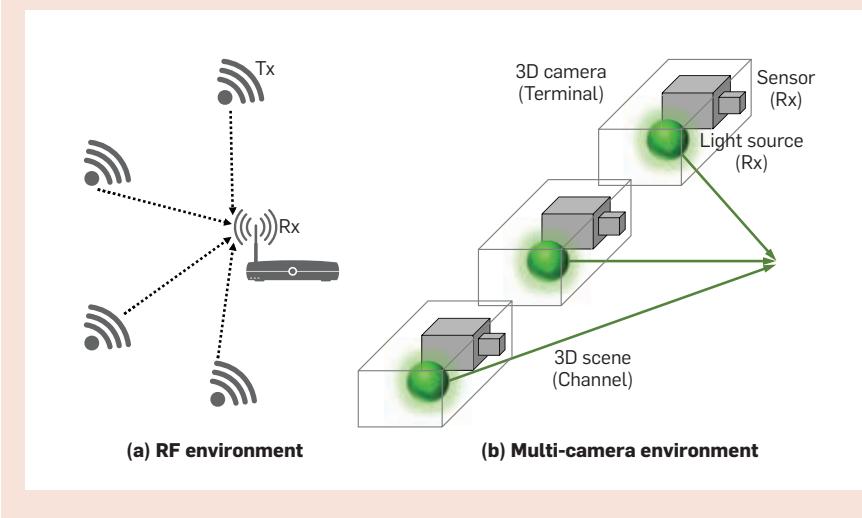
The similarities further extend from a single transmitter-receiver (transceiver) pair to a network of transceivers. In a wireless network, when multiple nodes attempt to communicate simultaneously, the radio waves interfere with each other, leading to collisions. Similarly, when multiple light sources flood-illuminate a scene, they collide and corrupt the received signal, resulting in depth errors.

There are, however, two critical differences between these two systems. First, in active 3D cameras, the **transmitter and the receiver are co-located**, that is, the light source and the sensor are both part of the same active 3D camera (Figure 4b). Second, signals collide in different ways. In a wireless channel, collisions could result in constructive interference or destructive interference, where the signal gets strengthened or weakened. On the other hand, an active 3D camera modulates the *intensity* of light, as opposed to the underlying electric field. The key factor to notice is that the intensity of light is always positive with both a constant (DC) and a time-varying (AC) component; the depth is encoded in the time-shift of the AC component. Therefore, the interfering signals from multiple cameras always accumulate additively. Although orthogonal-coding approaches can remove the AC interference, the DC component still accumulates, resulting in higher photon noise.

Despite these differences, we believe that the experiences of the wireless communications and networking community in addressing interference problems can benefit the imaging and computer-vision community in addressing MCI. In particular, we believe that a good understanding of various approaches to mitigating wireless interference and the theoretical foundations and practical considerations of

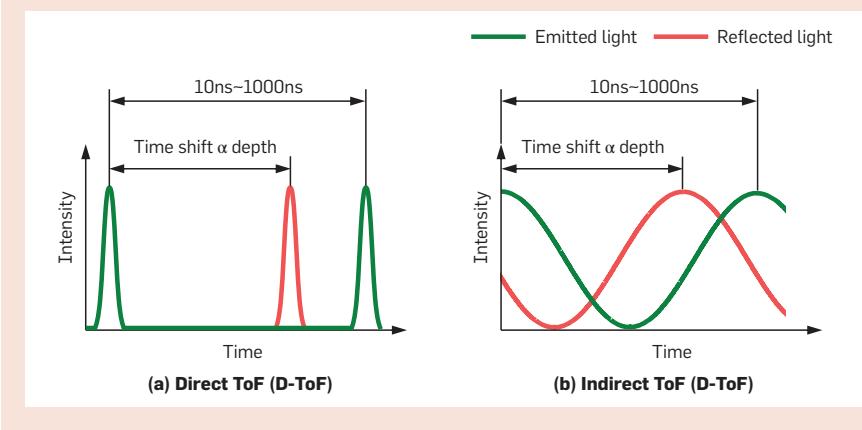
**Figure 4. Wireless and multi-camera environments.**

(a) In an RF network, transmitter (Tx) nodes send signals via the channel to the receiver (Rx). (b) In a multi-camera environment, multiple active 3D cameras transmit light from their light sources (Tx) to the 3D scene (channel) and receive the reflected light by the sensor (Rx).



**Figure 5. Active 3D cameras.**

(a) In direct ToF (D-ToF) approaches, the light source emits a periodic train of short light pulses. Scene depths are estimated by directly measuring the travel time of each pulse using high-speed timing circuits. (b) In indirect ToF (I-ToF) approaches, the light source emits continuously modulated periodic light waveforms. Depths are estimated by measuring the phase-shift between the emitted and reflected waveforms.



active 3D camera design can together synthesize useful solutions. If successful, the resulting ideas will not just spur widespread adoption of existing technologies, but also enable emerging applications that were hitherto considered impossible. The theoretical tools and techniques developed as part of this work will find applications in a broad range of techniques which involve coded light sources and sensors, such as structured light, tomography, and microscopy, as well as optical communication systems using coded light sources.

In summary, we answer the question in this subsection affirmatively by indicating that research on wireless interference can be used to address MCI. In the rest of the article, we identify various MCI mitigation approaches by leveraging its similarity to wireless and propose new directions to address the differences between the two scenarios.

### An Overview of Active 3D Cameras

An active 3D camera consists of a light source that emits coded illumination toward the scene and a sensor that captures the reflected light, as shown in Figure 2a. The most widely used class of active 3D cameras is based on the ToF principle. ToF-based cameras have a light source which emits temporally coded illumination. For example, the light source could be a laser or a light-emitting diode (LED) that sends out short light pulses or a continuously modulated light. The emitted light travels to the scene of interest and is reflected back to the sensor. The cameras measure the scene depths by measuring the total time of travel, computing the time-shift between the emitted and received waveforms (Figure 5a-b).

ToF-based depth-imaging systems can be broadly classified into direct and indirect ToF systems. A direct ToF (D-ToF) system<sup>16,34</sup> estimates scene depths by emitting a short light pulse into the scene and directly measuring the travel time of the reflected pulse (Figure 5a). Most vehicular LiDARs are based on the D-ToF principle. An indirect ToF (I-ToF) system,<sup>11,18,27</sup> on the other hand, emits light continuously. The intensity of its light source and the exposure of the sensor are both modulated over time for measuring the scene depths.

## The multi-camera interference problem shares several similarities and dissimilarities with common interference problems in the RF domain.

The light-source-modulation and sensor-demodulation functions can be any periodic and continuous functions such as sinusoids, square waves, or trapezoidal functions (Figure 5b).<sup>12</sup> I-ToF cameras do not require expensive components, and hence, are routinely used in consumer devices such as Microsoft Kinect and Hololens.

**MCI in active 3D cameras.** When multiple active 3D cameras illuminate a scene, the reflected signals interfere, corrupting the scene information. In this section, we provide the mathematical background for understanding MCI in active 3D cameras.<sup>c</sup>

*ToF image formation model.* The intensity of the light source in a ToF camera is temporally modulated as a continuous periodic function  $M(t)$ ; it could be a sinusoid,<sup>18,27</sup> or an impulse train function.<sup>17</sup> The period of  $M(t)$  typically varies from 10–1,000ns, which corresponds to a measurable distance of 1–100m. The light emitted by the source travels to the scene of interest and is reflected back toward the camera. The intensity of the reflected light incident on a sensor pixel  $\mathbf{p}$  is a time-shifted and scaled version of  $M(t)$ :

$$R(\mathbf{p}; t) = \alpha M(t - \tau), \quad (1)$$

where  $\tau = \frac{2d}{c}$  is the time-shift of the waveform due to travel from the source to the sensor.  $d$  is the distance between the camera and the scene point imaged at  $\mathbf{p}$  and  $c$  is the light speed.  $\alpha$  is a scene-dependent scale factor that encapsulates the scene's reflectance properties. The camera computes  $\tau$  (typically of the order of ns) by using high-speed, on-chip timing circuits, and the scene distance is estimated as  $d = \frac{c\tau}{2}$ .

*Multi-camera interference.* If multiple ToF cameras are simultaneously illuminating and imaging a scene point (Figure 2b), the brightness of the light incident at one of the cameras (referred to as the primary camera) is given by:

$$R_{\text{multi}}(t) = R(t) + \text{MCI}, \\ \text{where, } \text{MCI} = \sum_{n=1}^N R_n(t) \quad (2)$$

where  $N$  is the number of interfering cameras,  $R(t)$  is the radiance incident at the primary camera due to its own source (Eq. 1), and  $R_n(t) = \alpha_n M_n(t - \tau_n)$  is the measured intensity due to the  $n^{\text{th}}$

<sup>c</sup> ToF image formation model is applicable to any active 3D camera.

source. We drop  $\mathbf{p}$  for brevity. The summation term in Eq. 2 corrupts the true radiance  $R(t)$ , resulting in erroneous depth estimates. Figure 2c shows an example of a ToF camera using sinusoid modulation. Assuming all the sources use sinusoids (or any other periodic shapes such as squares) of the same frequency, the phase of  $R_{\text{mult}}(t)$  may differ from the true phase of  $R(t)$ , resulting in systematic, potentially large depth errors as shown in Figure 2d-f.

**Current approaches to addressing MCI.** A trivial approach to prevent MCI is to assign different wavelengths to different cameras. This approach only eliminates AC interference, and it faces two practical constraints: (1) The set of available wavelengths is strongly limited by the sensitivity range of the sensor material (typically silicon), as well as the practical requirement for the emitted light to be invisible to humans. This limits the available wavelengths to be  $\approx 850\text{-}950\text{nm}$ , which is the near-infrared region of the EM spectrum; (2) Due to laser and sensor hardware constraints, each sensor must be assigned a range of wavelengths (for example, 5-10nm). These constraints restrict the number of distinct wavelength bands preventing assigning a unique set of wavelengths to each active camera. Recent works addressing MCI can be broadly classified into three categories:

*Orthogonal coding.* The majority of existing works rely on orthogonal coding, such as sinusoids of different modulation frequencies<sup>30</sup> or phases,<sup>20,36</sup> and pseudo-noise sequences<sup>7,8</sup> for different cameras. However, they face challenges similar to frequency division multiple access (FDMA): a limited set of orthogonal frequencies and codes which proves inadequate for the rapidly growing cameras. These approaches often require a central authority that assigns a unique code to each camera, which is not practical.

*Time division multiple access (TDMA).* Other approaches divide the total capture time of the camera into multiple time slots and assign them to individual cameras randomly.<sup>19</sup> These techniques do not scale with the number of interfering cameras.

*Mechanical approaches.* Another method is to project a planar light

## Several active cameras flood-illuminating a scene simultaneously will interfere with each other.

sheet which is scanned over the scene. Since only a portion of the scene is illuminated at a time, interference can be reduced.<sup>4,26,35</sup> Although these approaches successfully prevent interference, they require mechanical scanning, which increase system cost and size.

Recently, optical phased arrays and micro-electro-mechanical systems (MEMS)-based technologies have attracted attention to resolve the limitations of mechanical approaches. These hardware-based approaches can be integrated into our software/firmware-based approaches to create hybrid solutions to mitigate MCI. However, this article focuses on wireless-inspired software/firmware approaches to reduce MCI.

### The Promise and Pitfalls of Wireless-Inspired Approaches

Interference is a classical problem in wireless communication. A huge body of literature over multiple decades has examined this problem to enable efficient shared access of a common channel such as the wireless medium. The wireless community commonly refers to these approaches as medium access control (MAC) techniques. Since MCI arises from sharing the common medium, it has similarities to MAC protocols in wireless communication. Here, we showcase the promise and pitfalls of applying well-known wireless MAC protocols to address MCI. As previously discussed, orthogonal coding across cameras can only reduce AC interference and not DC interference. A simple solution is to avoid collisions; we adapt a coordinated and distributed collision avoidance from wireless, simulate them in a multi-camera setting, and analyze their performance.

**Comparisons using computer simulations of MCI.** We developed a physically accurate computational simulator to emulate active 3D cameras, under a wide range of scene configurations and sensor-source parameters. The simulator models different steps of the image formation process as shown in Figure 6. Given a source modulation function  $M(t)$  ( $\geq 0$ ), the emitted light  $E(t)$  is defined as  $E(t) = P_s M(t)$ , where  $P_s$  is average source power.  $M(t)$  can be any periodic function. For example, for

sinusoid coding,  $M(t) = 1 + \cos(2\pi ft)$ , where  $f$  is modulation frequency. The light signal  $L(t)$  received at the sensor is defined as the sum of source reflection at the scene and ambient light:

$$L(t) = \alpha E(t - \frac{2d}{c}) + P_a, \quad (3)$$

where  $\alpha$  is a scale factor encapsulating scene reflectance and light fall-off,  $P_a$  is average ambient power,  $d$  is scene depth, and  $c$  is light speed. Noise-free sensor measurement  $C$  is the correlation between  $L(t)$  and sensor demodulation  $D(t)$ :

$$C = \int_{\tau} L(t)D(t) dt, \quad (4)$$

where  $\tau$  is integration time. If we assume sinusoid coding,  $D(t) = 1 + \cos(2\pi ft)$ . Final sensor measurement is obtained by adding photon noise, read noise, and ADC noise to  $C$ . To recover scene depth, we need multiple sensor measurements obtained by changing the phase of  $D(t)$ . For a 4-tap sinusoid coding scheme, those measurements  $C_k$  are obtained from Eq. 4 with  $D_k(t) = 1 + \cos(2\pi ft - \frac{\pi}{2}(k-1))$ , ( $k=1,\dots,4$ ). The scene distance estimate  $\hat{d}$  is obtained by

$$\hat{d} = \frac{c}{4\pi f} \tan^{-1} \left( \frac{C_2 - C_4}{C_1 - C_3} \right) \quad (5)$$

Using this simulator, we compare the performance of the following two wireless-inspired MCI reduction approaches with the conventional orthogonal coding approach (OCA). We implement OCA by assigning orthogonal modulation frequencies to each camera. Visit <https://bit.ly/3MhG4qu> to access the code to simulate MCI reduction approaches.

**CSMA-based MCI mitigation.** A popular idea in many communication systems to share a common medium is the “listen before talk” approach. In this approach, interference is potentially mitigated by requiring each transmitter to listen to the channel before transmitting their data. Waiting until the channel is idle ensures that a transmitter does not interfere with an ongoing communication. In wireless parlance, this is called “carrier sensing,” and the technique is often referred to as carrier sense multiple access (CSMA). To ensure no two nodes attempt to transmit simultaneously once the channel is sensed to be idle, CSMA with collision avoidance (CSMA/CA) is widely used in wireless

systems such as Wi-Fi. We draw inspiration from CSMA/CA, which requires nodes to wait a random amount of time after the channel is idle before transmitting. Combining carrier sensing and random wait times helps to avoid collisions. In particular, we design a “listen/view before illumination” approach at each camera to avoid interference from other cameras. Before illuminating the scene, if the channel is busy, the camera defers; else, it illuminates the scene and measures the depth information. In our implementation, we divide the total capture time (minimum time to estimate depth) of an active 3D camera into multiple time slots and sense the channel at each slot.

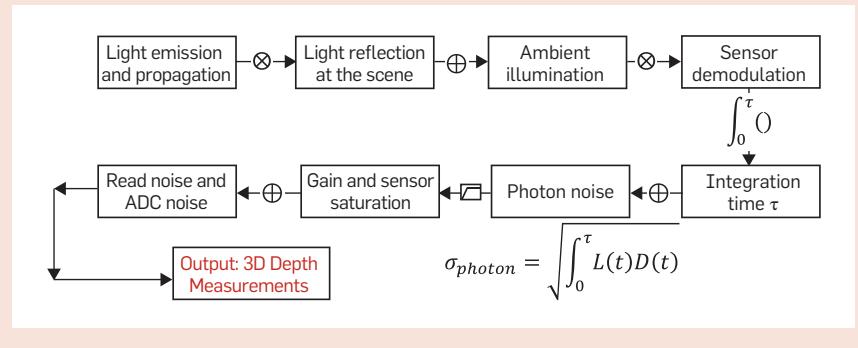
**Co-operative random access-based MCI mitigation (CRA).** We compare the depth accuracy of CSMA-based MCI with a time-division multiple access (TDMA) approach. If multiple cam-

eras are synchronized to a global clock, their transmissions can be scheduled in fixed slots. Each camera chooses a random sequence of slots to be active and illuminate the scene. In our implementation, all the cameras are synchronized, such that the slot boundaries are the same.

In Figure 7, we compare OCA, CRA, and CSMA in terms of depth accuracy and power consumption as a function of the number of interfering cameras. For depth accuracy comparison at the same source power, we employ depth standard deviation since all compared approaches cause random depth errors due to noise, instead of structured errors. For both CSMA and CRA, we use a slot clash check algorithm<sup>19</sup> to check if slot interference occurred, and depth is estimated from the collection of non-clashed slots. Depth standard deviation of each MCI approach is computed from repeated

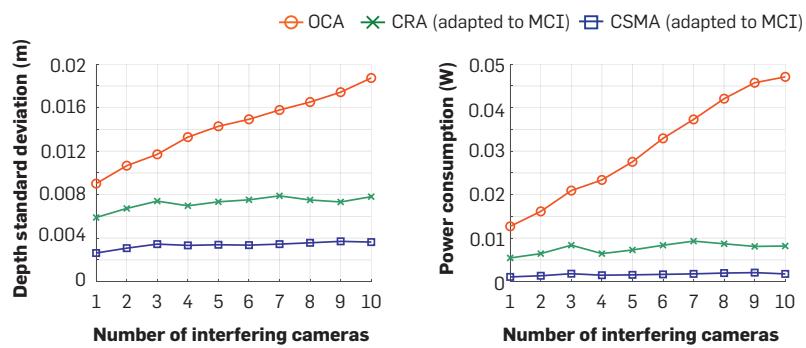
**Figure 6. Simulation steps.**

We developed a physics-based simulator that computationally emulates various steps of an active 3D camera’s imaging process for evaluating the proposed interference mitigation techniques.



**Figure 7. Depth standard deviations and power consumption of MCI approaches.**

Simulation results show that CSMA-based and cooperative random access-based MCI outperform conventional orthogonal coding.



depth estimation. The required power consumption to achieve the same depth standard deviation is also compared. CRA and CSMA, the wireless-inspired approaches, outperform OCA in terms of both depth accuracy and power consumption, as they avoid collisions and hence both AC and DC interference is reduced. Performance improves with the number of cameras, which is desired when active 3D cameras become more prominent in the near future.

**Potential pitfalls.** Although CRA and CSMA show promising results, the comparisons are based on purely theoretical and idealized simulations. In practice, additional overhead is required to account for various resources, such as power and time for clash check and carrier sensing. Furthermore, adapting wireless-based approaches for MCI has a tradeoff: The probability of collision increases as the network scales. To avoid collisions, more slots should be used, which leads to greater capture time and power consumption. The time to switch from carrier sensing to normal camera mode will further increase slot duration. These can be critical issues for real-time applications, such as vehicular networks. It is also challenging to synchronize multiple cameras to a global clock due to the absence of a central controller. Therefore, to accommodate the constraints of an active 3D camera system, careful algorithm design and post-processing are required. In the next section, we present open research areas that can address MCI using approaches inspired by wireless communication.

## The Road Ahead: Challenges in Wireless-Inspired MCI Solutions

We identify the following research directions, each inspired by the rich work on MAC protocols in wireless networks, with open challenges in applying them to active 3D camera networks. We build upon existing approaches to be applied to active 3D cameras by leveraging the opportunities provided by cameras to address their unique constraints.

**Distributed interference management.** In large-scale active 3D camera networks, resource-efficient, distributed interference-management techniques are necessary. We propose spread spectrum strategies that can estimate ToF from the cumulative reflected signal, enabling **multiple cameras to co-exist**. Spread spectrum techniques are used to improve a sender's resilience to interference and enable co-existence with other transmitters by occupying a wider range of frequencies. Since throughput is not a metric of interest in an active 3D camera network, spread spectrum-based approaches are an appropriate choice for MCI.

One of the key opportunities to leverage in active 3D cameras is the co-existence of the transmitter and the receiver; this enables it to maximize the potential of spread spectrum without communication overheads. Most existing active 3D cameras can modulate multiple frequencies to achieve both high precision and large depth range (for example, Microsoft Kinect uses 120MHz, 80MHz, and 16MHz<sup>25</sup>). On the other hand, one of the biggest challenges is the need and abil-

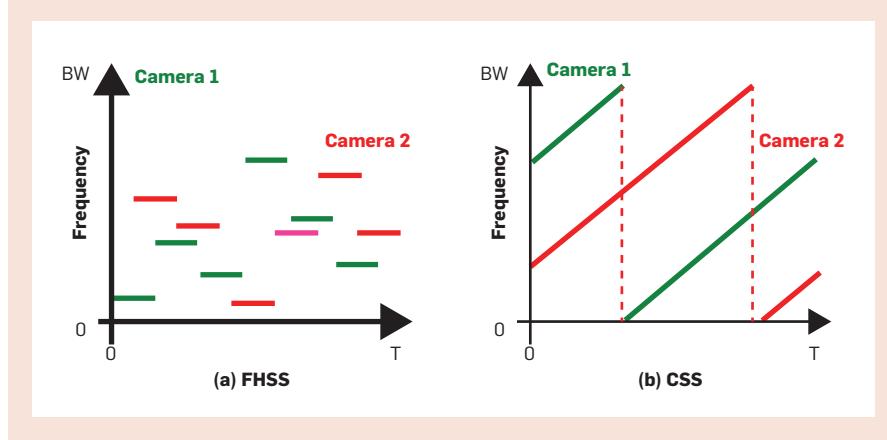
ity to cancel background noise. As previously explained, *unlike wireless, orthogonal frequencies and codes still contribute to ambient noise* in an active 3D camera. A redesign of spread spectrum is needed to address this challenge. We explore two spread spectrum techniques (that are widely used in wireless) and discuss their challenges in applying to MCI.

**Randomized frequency hopping spread spectrum.** Frequency hopping spread spectrum (FHSS) is a spread spectrum technique where the carrier frequency of the transmitter follows a *hopping pattern* (Figure 8a). Due to its robustness to broadband interference, adaptive FHSS is used in Bluetooth, which enables Bluetooth devices to co-exist with Wi-Fi devices. Bluetooth devices follow a pseudo-random hop sequence to avoid interference, where the hopping pattern is known to the transmitter and the receiver. We adopt FHSS to reduce MCI. In this design, each light source randomly chooses a pseudo-random hopping pattern and hop duration. Since the light source and the sensor are co-located, the hopping pattern is known to both of them. The hopping pattern can also be varied without additional overheads to exchange the hopping pattern. Therefore, the sensor keeps track of ToF of the expected frequencies periodically. The cumulative ToF of the corresponding pattern is then used to estimate the depth of the entire scene.

While FHSS reduces the probability of two cameras colliding in a given slot, unlike RF signals, the intensities of light signals from interfering cameras add up, leading to DC interference even when the hopping patterns of two cameras do not match. With the help of successive interference cancellation, we propose to iteratively cancel DC interference across the frequency bands. For example, consider two cameras using FHSS colliding in three slots. The probability of the two cameras using the same frequencies in the same slots in round two is significantly lower. This probability can be further reduced by increasing the scan duration and introducing empty slots where a light source does not send anything.

**Randomized chirp spread spectrum (CSS).** Chirp signals have been used

**Figure 8. Distributed Interference Management approaches.**



for precision ranging in RADAR,<sup>15</sup> and more recently, CSS modulation has been used in LoRa<sup>21,22</sup> to enable long-range communication. CSS spreads energy by linearly increasing the frequency of operation over time. Since this linear increase makes it robust to interference, it can be leveraged to enable multiple cameras to coexist: Each light source is assigned a unique starting frequency such that interference from other sources does not affect the AC component of the reflected signal. Figure 8b illustrates a chirp signal transmitted by a source, where the frequency sweeps the entire assigned bandwidth in a given duration.

The probability of collision is inversely proportional to the number of chirps. It is therefore desirable to have a large pool of chirps to reduce collisions. While decreasing step size will increase the number of chirps, the minimum step size is limited by hardware switching speeds. To address this challenge, non-linear CSS can be explored. A non-linear chirp with initial frequency  $f_i$  is similar to that of linear chirp in Figure 8b, whose frequency progression with time is given by a function  $g(t)$ . By choosing orthogonal non-linear functions, the probability of collisions can be reduced.<sup>15,22</sup> Similar to FHSS, CSS also must consider DC interference due to cumulative energy from other frequencies. We propose to successively cancel DC interference from one band to another and, leveraging the lack of time synchronization, this offers more possibilities to explore in MCI.

**Centralized, networked coordination.** While distributed approaches to MCI are necessary for many applications, there are multiple environments that lend themselves well to networked coordination with centralized control. For example, a factory floor with mobile autonomous robots operating in a single administrative domain can easily be coordinated through a central server. In such indoor applications, RF-based wireless connectivity to network devices via a central server that helps with synchronization could be integrated with the cameras. Therefore, combining RF and camera-based networks could offer real-time control of the cameras to

## As solid-state ToF cameras become ubiquitous in our personal mobile devices, vehicles, homes, and workplaces, managing interference across a multitude of these devices will be of profound importance.

a central server. Vehicles with LiDARs on roadways could also accomplish similar goals with roadside infrastructure providing the control function. Finally, home environments are likely the most challenging as there might be multiple environment sensing devices from diverse manufacturers, but if standardization of centralized coordination were to be developed, then all such devices could coordinate through a common hub within the home to manage active camera activities. We propose the following approaches in such centralized coordination scenarios.

*Creating a conflict graph in real-time.* To schedule simultaneous transmissions in wireless systems, one needs to infer the “conflict graph” that maps the potential of interference between transmitters.<sup>3,28,29</sup> Creating such a conflict graph in real time in the active camera domain is more challenging since the goal is to learn the channel. One possible approach is to time-synchronize different light sources and arrange them to send sustained pulses, sometimes in tandem and sometimes in isolation. If there are differences in received outcomes, we may conclude that such transmit-receive pairs interfere. Of course, a challenge lies in achieving various synchronization accuracies. We can combat them by adaptively choosing pulse durations and pulse structures that outlast potential clock synchronization errors. Other sources of inaccuracies can stem with other transmitters in the vicinity that are operating under the control of the centralized coordinator. Time-series analysis to identify and eliminate such interferers is an interesting research problem toward a practical deployment of a centralized coordination. A bigger challenge is the ability to meet the same goals as in a passive setting, that is, by simply observing activities on the channel and using time information to determine the same information. This, however, may be effective if all 3D cameras are sufficiently active. Overall, a hybrid passive-active method would optimize the best of both alternatives.

*Networked schedule of activity under centralized control.* Prior work on centralized WLANs has shown that

it is possible to create efficient partial-traffic scheduling for best performance.<sup>31,32</sup> Traffic scheduling is particularly effective when a vast majority of the traffic can be unscheduled, as it can be addressed through simple distributed mechanisms. Centralized control can benefit the fraction that cannot be addressed effectively by those mechanisms. Such traffic scheduling can also be used to minimize collisions in the active camera environments. The real-time conflict graphs discussed previously will inform a centralized coordinator for which transmitters might benefit from scheduling while allowing the remaining to be unscheduled. Note that scheduling requirements might depend on specific environments and the frequency with which transmitters need to rediscover their 3D environment. It also depends on how frequently the environment changes and external, out-of-band information can be used for this purpose. Based on this information, each transmitter-receiver pair that cannot be managed using a distributed approach may be scheduled to time slots to accomplish their scene-sensing goals. Note that a sensing attempt may be scheduled (or not), depending on the other devices attempting to sense concurrently.

## Conclusion

We strongly believe that due to their impending growth for consumer applications and their similarities to wireless communications, active 3D camera networks will benefit from new research from wireless experts. In particular, as interference between these cameras grows, sophisticated techniques are needed to handle MCI. Due to some fundamental differences between active cameras and wireless networks identified in this work, careful considerations in designing interference mitigation, cancellation, and coordination approaches are needed. It is interesting to note that some of the early work on interference-mitigation strategies in wireless environments (MACA,<sup>14</sup> MACAW<sup>5</sup>) started by exploring the potential use of wired interference-mitigation strategies in the wireless environment (Ethernet and related

CSMA/CD approaches) and their consequent limitations as well as through the process discovering improvements that led to eventually useful solutions. We believe that a similar approach can be taken by using the state-of-the-art wireless protocols as a starting point to uncover how similar strategies can be designed for active 3D camera systems.

We believe this problem domain is ripe for multiple communities to address collaboratively. Successful and effective solutions to this problem can have a big impact in many applications in our homes and society.

## Acknowledgments

This research was supported in part by the ONR grant number N00014-16-1-2995; the DARPA REVEAL program; the NSF CAREER Award 1943149, 2142978; the NSF award CNS-2107060; the ECCS 2034415; the U.S. National Science Foundation awards 2003129, 2112562, 2107060, 2212688; and through an award from the U.S. Department of Commerce with award number 70NANB21H043. C

## References

1. 3D Time-of-Flight (ToF) Sensors, Texas Instruments; <https://bit.ly/3Fy0yaD>.
2. 3D ToF Development Kits, PMD; <https://pmdtec.com/picofamily/>.
3. Ahmed, N., Ismail, U., Keshav, S., and Papagiannaki, K. Online estimation of RF interference. In *Proceedings of the 2008 ACM CoNEXT Conf.* (December 2008), 1–12.
4. Achar, S. et al. Epipolar time-of-flight imaging. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 37.
5. Bhardhavan, V. et al. MACAW: A media access protocol for wireless LANs. In *Proceedings of ACM Sigcomm* 24, 4 (Oct. 1994).
6. Breuer, P. et al. Hand gesture recognition with a novel IR time-of-flight range camera—a pilot study. In *Proceedings of the Intern. Conf. on Computer Vision/Computer Graphics Collaboration Techniques and Applications*. Springer (2007), 247–260.
7. Büttgen, B. et al. Pseudonoise optical modulation for real-time 3D imaging with minimum interference. *IEEE Transactions on Circuits and Systems I: Regular Papers* 54, 10 (2007), 2109–2119.
8. Büttgen, B. and Seitz, P. Robust optical time-of-flight range imaging based on smart pixel structures. *IEEE Trans. on Circuits and Systems* 55, 6 (2008), 1512–1525.
9. Cui, Y. et al. 3D shape scanning with a time-of-flight camera. In *Proceedings of the 2010 IEEE Conf. on Computer Vision and Pattern Recognition* (2010), 1173–1180.
10. Dutton, N.A.W. et al. A SPAD-based QVGA image sensor for single-photon counting and quanta imaging. *IEEE Transactions on Electron Devices* 63, 1 (Jan. 2016), 189–196.
11. Gupta, M. et al. Phasor imaging: A generalization of correlation-based time-of-flight imaging. *ACM Transactions on Graphics* 34, 5 (2015), 156.
12. Gupta, M. et al. What are optimal coding functions for time-of-flight imaging? *ACM Transactions on Graphics (TOG)* 37, 2 (2018), 13.
13. Henry, P. et al. RGB-D mapping: Using depth cameras for dense 3D modeling of indoor environments. In *Proceedings of the 12th Intern. Symp. on Experimental Robotics (ISER)*. Citeseer, (2010).
14. Karn, P. MACA: A new channel access method for packet radio. In *Proceedings of the 9th ARRL Computer Networking Conf.* (1990).
15. Klauder, J.R. et al. The theory and design of chirp radars. *Bell System Technical J.* 39, 4 (1960), 745–808.
16. Koehnert, W. Optical ranging system employing a high power injection laser diode. *IEEE Trans. AES* 4, 1 (1968).
17. Kolb, A. et al. Time-of-flight cameras in computer graphics. *Eurographics*. (2010).
18. Lange, R. 3D ToF distance measurement with custom solid-state image sensors in CMOS-CCD-technology. *Ph.D. Thesis*. (2000).
19. Lee, J. and Gupta, M. Stochastic exposure coding for handling multi-ToF-camera interference. In *Proceedings of the IEEE Intern. Conf. on Computer Vision* (2019), 7880–7888.
20. Li, L. et al. Multi-camera interference cancellation of time-of-flight (ToF) cameras. In *Proceedings of the 2015 IEEE Intern. Conf. on Image Processing*, 556–560.
21. LoRa: <https://www.semtech.com/lora>.
22. LoRa Modulation Basics: <https://bit.ly/3Qy7MBU>.
23. Ma, J. et al. Photon-number-resolving megapixel image sensor at room temperature without avalanche gain. *Optica* 4, 12 (2017), 1474–1481.
24. May, S. et al. 3D time-of-flight cameras for mobile robotics. In *Proceedings of the 2006 IEEE/RSJ Intern. Conf. on Intelligent Robots and Systems*, 790–795.
25. Pagliari, D. and Pinto, L. Calibration of Kinect for Xbox One and comparison between the two generations of microsoft sensors. *Sensors* (2015).
26. Pan, Y. et al. Accurate depth extraction method for multiple light-coding-based depth cameras. *IEEE Transactions on Multimedia* 19, 4 (2017), 685–701.
27. Payne, J.M. An optical distance measuring instrument. *Rev. of Scientific Instruments* 44, 3 (1973).
28. Rayanchu, S. et al. Airshark: Detecting non-Wi-Fi RF devices using commodity Wi-Fi hardware. In *Proceedings of the Intern. Measurement Conf.* (Nov. 2011).
29. Rayanchu, S. et al. Catching whales and minnows using Wi-FiNet: Deconstructing non-Wi-Fi interference using Wi-Fi hardware. In *ACM/USENIX Networking Systems Design and Implementation* (April 2012).
30. Shrestha, S. et al. Computational imaging with multi-camera time-of-flight systems. *ACM Transactions on Graphics (ToG)* 35, 4 (2016), 33.
31. Shrivastava, V. et al. Centaur: Realizing the full potential of centralized WLANs using a hybrid data path. In *Proceedings of the 15th Annual Intern. Conf. on Mobile Computing and Networking* (Sept. 2009), 297–308.
32. Suresh, L. et al. Towards programmable enterprise WLANs with Odin. In *Proceedings of the 1st Workshop on Hot Topics in Software Defined Networks* (2012), 115–120.
33. Ulku, A.C. et al. A 512x512 SPAD image sensor with integrated gating for widefield film. *IEEE J. of Selected Topics in Quantum Electronics* 25, 1 (Jan. 2019), 1–12.
34. Velodyne. Pulsed LiDAR sensor; <http://www.velodynelidar.com/lidar/lidar.aspx>.
35. Wang, J. et al. 3D scene reconstruction by multiple structured-light based commodity depth cameras. In *Proceedings of the 2012 IEEE Intern. Conf. on Acoustics, Speech and Signal Processing*, 5429–5432.
36. Whyte, R.Z. et al. Multiple range imaging camera operation with minimal performance impact. In *Image Processing: Machine Vision Applications III* 7538, Intern. Society for Optics and Photonics, (2010).

**Jongho Lee** (jlee567@wisc.edu) is a research associate of Computer Sciences at the University of Wisconsin-Madison, WI, USA.

**Mohit Gupta** is an associate professor of Computer Sciences at the University of Wisconsin-Madison, WI, USA.

**Bhuvana Krishnaswamy** is an assistant professor of Electrical and Computer Engineering at the University of Wisconsin-Madison, WI, USA.

**Suman Banerjee** is a professor of Computer Sciences at the University of Wisconsin-Madison, WI, USA.

# research highlights

---

P. 84

## **Technical Perspective Maximum Flow through a Network: A Storied Problem and a Groundbreaking Solution**

By Shang-Hua Teng

P. 85

## **Almost-Linear-Time Algorithms for Maximum Flow and Minimum-Cost Flow**

By Li Chen, Rasmus Kyng, Yang P. Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva

# Technical Perspective

## Maximum Flow through a Network: A Storied Problem and a Groundbreaking Solution

By Shang-Hua Teng

**I**N 2022, A team of computer scientists presented a groundbreaking algorithm for the maximum flow problem: How does one transport the most supplies from a source node to a sink node in a network while respecting link capacities? This result has a wide impact on algorithmic theory because this storied problem has broad theoretical significance and practical applications.

Maximum flow is an exemplary theoretical model of a real-world scenario. In an interdisciplinary collaboration, Ted Harris, a RAND mathematician, and General Frank Ross, a former chief of the Army's Transportation Corps in Europe, aided by George Dantzig, formulated the problem when studying rail transportation in the 1950s. The flow problem is intrinsically related to the minimum-cut problem via the mathematical duality: "maximum flow is equal to minimum cut." While the flow measures how well two nodes are connected, the dual cut measures how much capacity must be destroyed to disconnect them. Both are central in optimization and have multiple fundamental applications, including bipartite matching and divide-and-conquer-based approximations. They are also tools for solving practical tasks, including image processing, DNA sequence alignment, circuit design, and finite-element simulation.

The flow problem has significant pedagogical value. It offers rich material for teaching algorithmic paradigms: greedy, iterative, multilevel, and mathematical programming. Together with von Neumann's minimax theorem for zero-sum games and Yao's minimax principle for randomized algorithms, the max-flow min-cut theorem is a vivid illustration of linear programming duality. Static in formulation and dynamic in imagination, as network models have become ubiquitous in computing, the flow/cut framework resonates with students of diverse scientific interests: from protein interaction to machine

translation; supply chains to Internet economics; and statistical learning to knowledge discovery.

Network flow is a magnet for algorithmic research, attracting generations of scientists building on each other's insights. Since the first algorithms by Ford and Fulkerson, progress on these problems has been intertwined with major advances in algorithmic theory, such as shortest-path search, randomization, mathematical programming, and Laplacian solvers. Two decades ago, for example, research for min-cut led to graph sparsification, a concept that facilitated breakthrough scalable algorithms for electrical flows. This, in turn, intensified the search for faster flow/cut algorithms, which led to scalable max-flow approximation in undirected networks. Located at the intersection of discrete graph theory and continuous mathematical programming, network flow has inspired creative interplay between combinatorial and continuous methods. The new solution also exploited this interplay.

In the age of Big Data, practical applications require algorithms to be not just polynomial but scalable, that is, within a low poly-logarithmic factor of linear time. The following paper comes within striking distance of answering the outstanding question: "Does maximum flow have a scalable algorithm?"

The authors presented the paper at the FOCS plenary session on November 2, 2022. That evening, as I struggled to find words for this Technical Perspective, I caught game four of the 2022 World Series on TV. Only then did I realize I was watching a once-in-a-half-century event in baseball: four pitchers delivered the second no-hitter in World Series history. It was not lost on me that the first Fall Classic no-hitter—a perfect game by the Yankees' Don Larsen—was in 1956, the very year Ford and Fulkerson finalized their celebrated work "Maximal Flow Through a Network" at RAND. "Miracle ... mir-

acle," I heard repeated in the background, as I returned to my writing, still charmed by the serendipity. Within hours, I witnessed two hard-fought achievements in seemingly unconnected worlds, each co-authored by dedicated people with contributions from others in the field.

The no-hitter was not defined just by a single strikeout or brilliant defense. Likewise, this work is a tour de force of multiple dimensions. At the top, it introduces an insightful integration of iterative approximation with advanced data structures. Instead of aiming for faster convergence, the team took an audacious step by using a method that requires a linear number of iterations. This outer loop enjoys better network locality: It repeatedly improves flows around cycles, guided by low-stretch spanning trees, a concept formulated decades ago for on-line algorithms and instrumental for scalable Laplacian solvers. The team then uses low-stretch spanning trees to design an intricate data structure with low amortized complexity for supporting flow improvements, showing that slow and steady wins the race of maximum-flow computation.

Simplification and improvement are still needed to make this algorithm applicable in the real world. But by breaking long-standing super-linear complexity barriers, it offers hope for scalable max-flow algorithms. Each breakthrough motivates future advances. Just like results for electrical flows have accelerated the search for scalable maximum-flow algorithms, this work will expedite research on other fundamental problems, including developing almost linear-time algorithms for general linear programming. 

---

**Shang-Hua Teng** is the Seeley G. Mudd Professor of Computer Science and Mathematics at the University of Southern California. The two-time Gödel Prize recipient met then USC mathematics professor Ted Harris while attending the university in the 1980s.

Copyright held by author/owner(s).

# Almost-Linear-Time Algorithms for Maximum Flow and Minimum-Cost Flow

By Li Chen, Rasmus Kyng, Yang P. Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva

## Abstract

We present an algorithm that computes exact maximum flows and minimum-cost flows on directed graphs with  $m$  edges and polynomially bounded integral demands, costs, and capacities in  $m^{1+o(1)}$  time. Our algorithm builds the flow through a sequence of  $m^{1+o(1)}$  approximate undirected minimum-ratio cycles, each of which is computed and processed in amortized  $m^{o(1)}$  time using a new dynamic graph data structure.

Our framework extends to algorithms running in  $m^{1+o(1)}$  time for computing flows that minimize general edge-separable convex functions to high accuracy. This gives almost-linear time algorithms for several problems including entropy-regularized optimal transport, matrix scaling,  $p$ -norm flows, and  $p$ -norm isotonic regression on arbitrary directed acyclic graphs.

## 1. INTRODUCTION

The maximum flow problem and its generalization, the minimum-cost flow problem, are classic combinatorial graph problems that find numerous applications in engineering and scientific computing. These problems have been studied extensively over the last seven decades, starting from the work of Dantzig and Ford-Fulkerson. Several important algorithmic problems can be reduced to minimum-cost flows, for example, max-weight bipartite matching, min-cut, and Gomory-Hu trees. The origin of numerous significant algorithmic developments such as the simplex method, graph sparsification, and link-cut trees can be traced back to seeking faster algorithms for maximum flow and minimum-cost flow.

### 1.1. Problem formulation

Formally, we are given a directed graph  $G = (V, E)$  with  $|V| = n$  vertices and  $|E| = m$  edges, upper/lower edge capacities  $\mathbf{u}^+, \mathbf{u}^- \in \mathbb{R}^E$ , edge costs  $\mathbf{c} \in \mathbb{R}^E$ , and vertex demands  $\mathbf{d} \in \mathbb{R}^V$  with  $\sum_{v \in V} d_v = 0$ . Our goal is to solve the following linear program for the *minimum-cost flow problem*

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{f} \\ & \text{subject to} && \mathbf{u}_e^- \leq \mathbf{f}_e \leq \mathbf{u}_e^+ \text{ for all } e \in E \\ & && \mathbf{B}^\top \mathbf{f} = \mathbf{d}, \end{aligned} \tag{1}$$

where the last constraint,  $\mathbf{B}^\top \mathbf{f} = \mathbf{d}$ , succinctly captures the requirement that the flow  $\mathbf{f}$  satisfies vertex demands  $\mathbf{d}$ , where  $\mathbf{B} \in \mathbb{R}^{E \times V}$  is the edge-vertex incidence matrix defined as  $B_{((a,b),v)} = 1$  if  $v = a$ ,  $-1$  if  $v = b$ , and  $0$  otherwise.

In this extended abstract, we assume that all  $\mathbf{u}_e^+, \mathbf{u}_e^-, \mathbf{c}_e$  and  $\mathbf{d}_v$  are integral and polynomially bounded in  $n$  since this paper focuses on weakly-polynomial algorithms for the maximum flow and minimum-cost flow problems.

### 1.2. Previous work

There has been extensive work on maximum flow and minimum-cost flow. Here, we briefly discuss some highlights from this work to help place our work in context.

Starting with the first pseudo-polynomial time algorithm by Dantzig<sup>14</sup> for maximum-flow that ran in  $O(mn^2U)$  time where  $U$  denotes the maximum absolute capacity, approaches to designing faster flow algorithms were primarily combinatorial, working with various adaptations of augmenting paths, cycle canceling, blocking flows, and capacity/cost scaling. A long line of work led to a running time of  $\tilde{O}(m \min\{m^{1/2}, n^{2/3}\})^{15, 18, 19, 21}$  for maximum flow, and  $\tilde{O}(mn)^{17}$  for minimum-cost flow. These bounds stood for decades.

In a breakthrough work on solving Laplacian linear systems and computing electrical flows, Spielman and Teng<sup>34</sup> introduced a novel set of ideas and tools for solving flow problems using combinatorial techniques in conjunction with continuous optimization methods. To deploy these methods, flow algorithms researchers have used graph-algorithmic techniques to solve increasingly difficult subproblems that drive powerful continuous methods.

In the context of maximum flow and minimum-cost flow, Daitch and Spielman<sup>13</sup> demonstrated the power of this paradigm by using a path-following interior point method (IPM) to reduce the minimum-cost flow problem to solving a sequence of roughly  $\sqrt{m}$  electrical flow ( $\ell_2$ ) problems. Since each of these  $\ell_2$  problems could then be solved in nearly-linear time using the fast Laplacian solver by Spielman-Teng, they achieved an  $\tilde{O}(m^{1.5})$  time algorithm for minimum-cost flow, the first progress in two decades. They showed that a key advantage of IPMs is that they reduce flow problems on directed graphs to flow problems on undirected graphs, which are easier to work with.

While other continuous optimization methods have been used in the context of maximum flow, even leading

The original version of this paper titled “Maximum Flow and Minimum-Cost Flow in Almost-Linear Time” was published in *Proceedings of the 2022 IEEE 63<sup>rd</sup> Annual Symposium on Foundations of Computer Science*, 612–623.

to nearly-linear time  $(1+\epsilon)$ -approximate undirected maximum flow and multicommodity flow algorithms,<sup>12, 23, 30, 32, 33</sup> to date all approaches for exact maximum flow and minimum-cost flow rely on the framework suggested by Daitch and Spielman of using a path-following IPM to reduce to a small but polynomial number of convex optimization problems. Notable achievements include an  $m^{4/3+o(1)}U^{1/3}$  time algorithm for bipartite matching and unit-capacity maximum flow.<sup>4, 22, 27, 28, 29</sup> Further, for general capacities, minimum-cost flow algorithms were given with runtimes  $\tilde{O}(m+n^{1.5})$ <sup>8, 9, 10</sup> and  $\tilde{O}(m^{3/2-1/58})$ .<sup>5, 6, 16, 35</sup> In both of these results, the development of efficient data structures to solve the sequence of  $\ell_2$  subproblems in amortized time sub-linear in  $m$  has played a key role in obtaining these runtimes. Yet, despite this progress, the best running time bounds remain far from linear.

### 1.3. Our result

We give the first almost-linear time algorithm for minimum-cost flow, achieving the optimal running time up to subpolynomial factors.

**THEOREM 1.1.** *There is an algorithm that, on a graph  $G = (V, E)$  with  $m$  edges, vertex demands, upper/lower edge capacities, and edge costs, all integral with capacities and costs bounded by a polynomial in  $n$ , computes an exact minimum-cost flow in  $m^{1+o(1)}$  time with high probability.*

Our algorithm can be extended to work with capacities and costs that are not polynomially bounded at the cost of an additional logarithmic dependency in both the maximum absolute capacity and the maximum absolute cost.

We make two key contributions to achieve our result. First, we develop a novel potential-reduction IPM, similar to Karmarkar's original work.<sup>20</sup> Our IPM is *worse* in some ways than existing path-following IPMs because it needs more updates to converge to a good solution. However, our new IPM reduces the minimum-cost flow problem to a sequence of update subproblems that have a more combinatorial structure than those studied before. This enables our second key contribution, a data structure that solves our sequence of update subproblems extremely quickly.

In addition to the use of highly combinatorial updates, our new IPM has three crucial properties. The IPM is (a) *robust* to approximation error in subproblems, (b) *stable* in terms of the subproblems it defines, and (c) *stable* in terms of a good solution to the subproblems. These features allow us to solve the sequence of update subproblems much faster by developing a powerful data structure, yielding a much faster algorithm overall. Thus, instead of making graph algorithms more suitable for continuous optimization, we made continuous optimization more suitable for graph algorithms.

We call the update subproblem that our IPM yields *min-ratio cycle*: This problem is specified by a graph where every edge has a “gradient” and a “length.” The problem asks us to find a cycle that minimizes the sum of (signed, directed) edge gradients relative to its (undirected) length.

Our data structure for solving the sequence of min-ratio cycle problems is our second key contribution. As a first observation, we show that if we sample a random “low-stretch” spanning tree of the graph, then with constant probability, some *fundamental tree cycle* approximately solves the min-ratio cycle problem. Recall a fundamental tree cycle is a cycle defined by a single non-tree edge and the unique tree path between its endpoints. Unfortunately, this simple approach fails after a single flow update, as the IPM requires us to change the gradients and lengths after each update.

To maintain a set of trees that repeatedly allow us to identify good update cycles, we develop a hierarchical construction based on a recursive approach. This requires us to repeatedly construct and contract a random forest (which partially defines our tree), and then recurse on the resulting smaller graph, which we call a *core graph*. Furthermore, to enable recursion, we need to reduce the edge count in the core graph. We achieve this using a new *spanner* construction, which identifies a sparse subgraph of the core graph on which to recurse and detects if the removed edges damage the min-ratio cycle. Maintaining the core graphs and spanners in our recursive construction requires us to develop an array of novel dynamic graph techniques, which may be of independent interest.

### 1.4. Applications

Our result in Theorem 1.1 has a wide range of applications. By standard reductions, it gives the first  $m^{1+o(1)}$  time algorithms for bipartite matching, worker assignment, negative-lengths single-source shortest paths, and several other problems. For the negative-lengths shortest path problem, Bernstein, Nanongkai, and Wulff-Nilsen obtained the first  $m \cdot \text{poly}(\log m)$  time algorithm in an independent and concurrent work.<sup>7</sup>

Using recent reductions from various connectivity problems to maximum flow, we also obtain the first  $m^{1+o(1)}$  time algorithms to compute vertex connectivity and Gomory-Hu trees in undirected, unweighted graphs,<sup>1, 25</sup> and  $(1+\epsilon)$ -approximate Gomory-Hu trees in undirected weighted graphs.<sup>26</sup> We also obtain the current fastest algorithm to find the global min-cut in a directed graph.<sup>11</sup> Finally, we obtain the first almost-linear-time algorithm to compute the approximate sparsest cuts in directed graphs.

Additionally, our algorithm extends to computing flows that minimize general edge-separable convex objectives.

**INFORMAL THEOREM 1.2.** *Consider a graph  $G$  with demands  $\mathbf{d}$ , and an edge-separable convex cost function  $\text{cost}(\mathbf{f}) = \sum_e \text{cost}_e(f_e)$  for “computationally efficient” edge costs  $\text{cost}_e$ . Then in  $m^{1+o(1)}$  time, we can compute a (fractional) flow  $\mathbf{f}$  that routes demands  $\mathbf{d}$  and  $\text{cost}(\mathbf{f}) \leq \text{cost}(\mathbf{f}^*) + \exp(-\log^C m)$  for any constant  $C > 0$ , where  $\mathbf{f}^*$  minimizes  $\text{cost}(\mathbf{f}^*)$  over flows with demands  $\mathbf{d}$ .*

This generalization gives the first almost-linear-time algorithms for solving entropy-regularized optimal transport (equivalently, matrix scaling),  $p$ -norm flow problems, and  $p$ -norm isotonic regression for  $p \in [1, \infty]$ .

## 2. OVERVIEW

### 2.1. Computing minimum-cost flows via undirected min-ratio cycles

Recall the linear program for minimum-cost flow given in Equation (1). We assume that this LP has a unique optimal solution at  $\mathbf{f}^* \in \mathbb{R}^E$  and let  $F^* = \mathbf{c}^\top \mathbf{f}^*$  (this can be achieved by a negligible perturbation using the famous Isolation Lemma).

For our algorithm, we use a potential-reduction interior point method,<sup>20</sup> where in each iteration we measure progress by reducing the value of the potential function

$$\begin{aligned}\Phi(\mathbf{f}) &\stackrel{\text{def}}{=} 20m \log(\mathbf{c}^\top \mathbf{f} - F^*) \\ &+ \sum_{e \in E} ((\mathbf{u}_e^+ - \mathbf{f}_e)^{-\alpha} + (\mathbf{f}_e - \mathbf{u}_e^-)^{-\alpha})\end{aligned}$$

for  $\alpha = 1/(1000 \log mU)$ . The reader can think of the barrier  $x^{-\alpha}$  as the more standard  $-\log x$  for simplicity instead. We use it for technical reasons beyond the scope of this paper.

Using standard techniques, one can add  $O(n)$  additional, artificial edges of large capacity and cost to the graph  $G$  such that the optimal solution to the minimum-cost flow problem remains unchanged (and in particular is not supported on the artificial edges) and such that one can easily find a feasible flow  $\mathbf{f}$  on the artificial edges such that  $\mathbf{B}^\top \mathbf{f} = \mathbf{d}$  with bounded potential, that is,  $\Phi(\mathbf{f}) = O(m \log m)$ .

Given the current feasible solution  $\mathbf{f}$ , the potential reduction interior point method asks to find a circulation  $\Delta$ , that is, a flow that satisfies  $\mathbf{B}^\top \Delta = \mathbf{0}$ , such that  $\Phi(\mathbf{f} + \Delta) \leq \Phi(\mathbf{f}) - m^{-o(1)}$ . Given  $\Delta$ , it then sets  $\mathbf{f} \leftarrow \mathbf{f} + \Delta$  and repeats. When  $\Phi(\mathbf{f}) \leq -200m \log mU$ , we can terminate because then  $\mathbf{c}^\top \mathbf{f} - F^* \leq (mU)^{-10}$ , at which point standard techniques let us round to an exact optimal flow.<sup>13</sup> Thus if we can reduce the potential by  $m^{-o(1)}$  per iteration, the method terminates in  $m^{1+o(1)}$  iterations.

Let us next describe how to find a circulation  $\Delta$  that reduces the potential sufficiently. Given the current flow  $\mathbf{f}$ , defining the gradient and lengths  $\mathbf{g}(\mathbf{f}) \stackrel{\text{def}}{=} \nabla \Phi(\mathbf{f})$  and  $\ell(\mathbf{f})_e \stackrel{\text{def}}{=} (\mathbf{u}_e^+ - \mathbf{f}_e)^{-1-\alpha} + (\mathbf{f}_e - \mathbf{u}_e^-)^{-1-\alpha}$ , and we let  $\mathbf{L} \stackrel{\text{def}}{=} \text{diag}(\ell)$  be the matrix with these lengths on the diagonal and zeros elsewhere. We can then define the *min-ratio cycle* problem

$$\begin{aligned}\text{minimize } & \mathbf{g}^\top \Delta / \|\mathbf{L}\Delta\|_1 \\ \text{subject to } & \mathbf{B}^\top \Delta = \mathbf{0}\end{aligned}\tag{2}$$

Given any solution  $\Delta$  to this problem with  $\mathbf{g}(\mathbf{f})^\top \Delta / \|\mathbf{L}\Delta\|_1 \leq -\kappa$  for some  $\kappa < 1/100$ , scaled so that  $\|\mathbf{L}\Delta\|_1 = \kappa/50$ . Then a direct Taylor expansion shows that  $\Phi(\mathbf{f} + \Delta) \leq \Phi(\mathbf{f}) - \kappa^2/500$ . Hence, it suffices to show that such a  $\Delta$  exists with  $\kappa = \tilde{\Omega}(1)$ , because then a data structure that returns an  $m^{o(1)}$ -approximate solution still has  $\kappa = m^{-o(1)}$ , which suffices. Fortunately, the *witness circulation*  $\Delta(\mathbf{f})^* = \mathbf{f}^* - \mathbf{f}$  satisfies  $\mathbf{g}(\mathbf{f})^\top \Delta / \|\mathbf{L}\Delta\|_1 \leq -\tilde{\Omega}(1)$ .

We emphasize that it is essential for our data structure that the witness circulation  $\mathbf{f}^* - \mathbf{f}$  yields a sufficiently good solution. This assumption ensures that good solutions to the min ratio cycle instances do not change arbitrarily between iterations. Further, even though the algorithm does not have access to the witness circulation  $\mathbf{f}^* - \mathbf{f}$ , it still knows how it changes between iterations as it can track changes in  $\mathbf{f}$ . We are able to leverage this guarantee to ensure our data structure succeeds for the updates coming from the IPM.

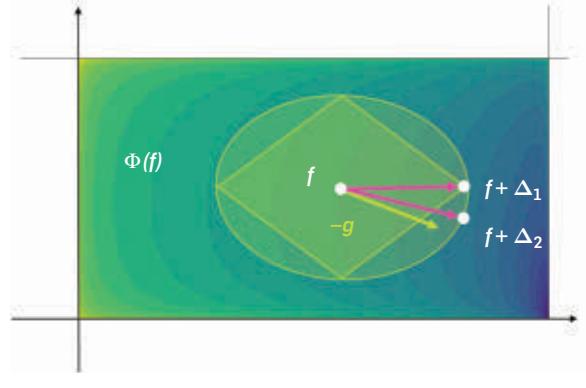
Finally, let us contrast our approach with previous approaches: previous analyses of IPMs solved  $\ell_2$  problems, that is, problems of the form given in Equation (2) but with the  $\ell_1$  norm replaced by a  $\ell_2$  norm (see Figure 1), which can be solved using a linear system. Karmarkar<sup>20</sup> shows that repeatedly solving  $\ell_2$  subproblems, the IPM converges in  $\tilde{O}(m)$  iterations. Later analyses of path-following IPMs<sup>31</sup> showed that a sequence of  $\tilde{O}(\sqrt{m}) \ell_2$  subproblems suffices to give a high-accuracy solution. Surprisingly, we are able to argue that a solving sequence of  $\tilde{O}(m) \ell_1$  minimizing subproblems of the form in Equation (2) suffice to give a high-accuracy solution to Equation (1). In other words, changing the  $\ell_2$  norm to an  $\ell_1$  norm does not increase the number of iterations in a potential reduction IPM. The use of an  $\ell_1$ -norm-based subproblem gives us a crucial advantage: Problems of this form must have optimal solutions in the form of simple cycles—and our new algorithm finds approximately optimal cycles vastly more efficiently than any known approaches for  $\ell_2$  subproblems.

### 2.2. High-level overview of the data structure for dynamic min-ratio cycle

As discussed in the previous section, our algorithm computes a minimum-cost flow by solving a sequence of  $m^{1+o(1)}$  min-ratio cycle problems  $\min_{\mathbf{B}^\top \Delta = 0} \mathbf{g}^\top \Delta / \|\mathbf{L}\Delta\|_1$  to  $m^{o(1)}$  multiplicative accuracy. Because our IPM ensures stability for lengths and gradients, and is even robust to approximations of lengths and gradients, we can show that over the course of the algorithm, we only need to update the entries of the gradients  $\mathbf{g}$  and lengths  $\ell$  at most  $m^{1+o(1)}$  total times.

**Warm-up:** A simple, static ALGORITHM. A simple approach to finding an  $\tilde{O}(1)$ -approximate min-ratio cycle is the following: given our graph  $G$ , we find a probabilistic low stretch spanning tree  $T$ , that is, a tree such that for each edge  $e = (u, v) \in G$ , the stretch of  $e$ , defined as  $\text{st}_e^{T, \ell} \stackrel{\text{def}}{=} \frac{\sum_{f \in T[u, v]} \ell(f)}{\ell(e)}$  where  $T[u, v]$  is the unique path from  $u$  to  $v$  along the tree  $T$ , is  $\tilde{O}(1)$  in expectation. Such a tree can be found in  $\tilde{O}(m)$  time.<sup>2,3</sup> This fact will allow us to argue that with probability

**Figure 1.** The IPM takes steps to minimize the potential  $\Phi(\mathbf{f})$  by updating  $\mathbf{f}$  to  $\mathbf{f} + \Delta$ . Previous approaches suggest obtaining  $\Delta$  by solving an  $\ell_2$  subproblem (here finding  $\Delta_2$  as the optimal step on an ellipsoid), but our approach obtains  $\Delta$  by minimizing an  $\ell_1$  problem (illustrated by finding  $\Delta_1$  as the optimal step in a box). While the new strategy possibly makes less progress at a step, this allows us to find the step  $\Delta$  more efficiently.



at least  $\frac{1}{2}$ , one of the tree cycles is an  $\tilde{O}(1)$ -approximate solution to Equation (2).

Let  $\Delta^*$  be the optimal circulation that minimizes Equation (2), and assume w.l.o.g. that  $\Delta^*$  is a cycle that routes one unit of flow along the cycle. We assume for convenience, that edges on  $\Delta^*$  are oriented along the flow direction of  $\Delta^*$ , that is,  $\Delta^* \in \mathbb{R}_{\geq 0}^E$ . Then, for each edge  $e = (u, v)$ , the *fundamental tree cycle* of  $e$  in  $T$ , denoted  $e \oplus T[v, u]$ , is formed by edge  $e$  concatenated with the path in  $T$  from its endpoint  $v$  to  $u$ . To work again with vector notation, we denote by  $p(e \oplus T[v, u]) \in \mathbb{R}^E$  the vector that sends one unit of flow along the cycle  $e \oplus T[v, u]$  in the direction that aligns with the orientation of  $e$ . A classic fact from graph theory now states that  $\Delta^* = \sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot p(e \oplus T[v, u])$  (note that the tree paths used by adjacent off-tree edges cancel out, see Figure 2). In particular, this implies that  $g^\top \Delta^* = \sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot g^\top p(e \oplus T[v, u])$ .

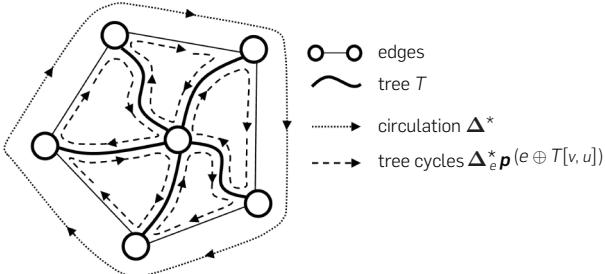
From the guarantees of the low-stretch tree distribution, we know that the circulation  $\Delta^*$  is not stretched by too much in expectation. Thus, by Markov's inequality, with probability at least  $\frac{1}{2}$ , the circulation  $\Delta^*$  is not stretched by too much. Formally, we have that  $\sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot \|L p(e \oplus T[v, u])\|_1 \leq \gamma \|L \Delta^*\|_1$  for  $\gamma = \tilde{O}(1)$ . Combining these insights, we can derive that

$$\begin{aligned} \frac{g^\top \Delta^*}{\|L \Delta^*\|_1} &\geq \frac{1}{\gamma} \cdot \frac{\sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot g^\top p(e \oplus T[v, u])}{\sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot \|L p(e \oplus T[v, u])\|_1} \\ &\geq \frac{1}{\gamma} \min_{e: \Delta_e^* > 0} \frac{g^\top p(e \oplus T[v, u])}{\|L p(e \oplus T[v, u])\|_1} \end{aligned}$$

where the last inequality follows from the fact that  $\min_{i \in [n]} \frac{x_i}{y_i} \leq \frac{\sum_{i \in [n]} x_i}{\sum_{i \in [n]} y_i}$  (recall also that  $g^\top \Delta^*$  is negative). This tells us that for the edge  $e$  minimizing the expression on the right, the tree cycle  $e \oplus T[v, u]$  is a  $\gamma$ -approximate solution to Equation (2), as desired. We can boost the probability of success of the above algorithm by sampling  $\tilde{O}(1)$  trees  $T_1, T_2, \dots, T_s$  independently at random and conclude that w.h.p. one of the fundamental tree cycles approximately solves Equation (2).

Unfortunately, after updating the flow  $f$  to  $f'$  along such a fundamental tree cycle, we cannot reuse the set of trees  $T_1, T_2, \dots, T_s$  because the next solution to Equation (2) has to be found with respect to gradients  $g(f')$  and lengths  $\ell(f')$  depending on  $f'$  (instead of  $g = g(f)$  and  $\ell = \ell(f)$ ). But  $g(f')$  and  $\ell(f')$  depend on the randomness used in trees  $T_1, T_2, \dots, T_s$ . Thus, naively, we have to recompute all trees, spending

**Figure 2. Illustrating the decomposition  $\Delta^* = \sum_{e: \Delta_e^* > 0} \Delta_e^* \cdot p(e \oplus T[v, u])$  of a circulation into fundamental tree cycles.**



again  $\Omega(m)$  time. But this leads to run-time  $\Omega(m^2)$  for our overall algorithm which is far from our goal.

**A dynamic approach.** Thus we consider the data structure problem of maintaining an  $m^{o(1)}$  approximate solution to Equation (2) over a sequence of at most  $m^{1+o(1)}$  changes to entries of  $g, \ell$ . To achieve an almost linear time algorithm overall, we want our data structure to have an amortized  $m^{o(1)}$  update time. Motivated by the simple construction above, our data structure will ultimately maintain a set of  $s = m^{o(1)}$  spanning trees  $T_1, \dots, T_s$  of the graph  $G$ . Each cycle  $\Delta$  that is returned is represented by  $m^{o(1)}$  off-tree edges and paths connecting them on some  $T_i$ .

To obtain an efficient algorithm to maintain these trees  $T_i$ , we turn to a recursive approach. Each level of the recursion will *partially* construct a tree, by choosing a spanning forest of the vertices, and contracting the connected components of the forest. We obtain a tree by repeating forest selection-and-contraction until only a single vertex is left. Then, we compose the forest edges obtained at different levels, yielding a spanning tree of the original graph. At each level of the construction, our forest is probabilistic and only succeeds with constant probability at preserving the hidden witness circulation well enough. To preserve the witness with high probability, we construct  $O(\log n)$  different forests at each level and recurse on them separately.

In each level of our recursion, we first reduce the number of vertices using a forest contraction and then the number of edges by making the contracted graph sparse. To reduce the number of vertices, we produce a *core graph* (the result of contracting forest components) on a subset of the original vertex set, and we then compute a *spanner* of the core graph which reduces the number of edges. The edge-reduction step is important to ensure the overall recursion reduces the graph size in each step, which is essential to obtaining almost linear running time in our framework.

Both the *core graph* and *spanner* at each level need to be maintained dynamically, and we ensure they are very stable under changes in the graphs at shallower levels in the recursion. In both cases, our notion of stability relies on some subtle properties of the interaction between data structure and hidden witness circulation.

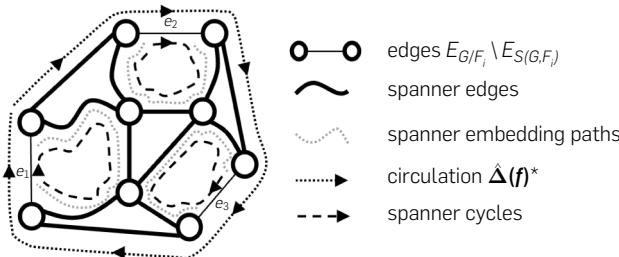
We maintain a recursive hierarchy of graphs. At the top level of our hierarchy, for the input graph  $G$ , we produce  $B = O(\log n)$  core graphs. To obtain each such core graph, for each  $i \in [B]$ , we sample a (random) forest  $F_i$  with  $\tilde{O}(m/k)$  connected components for some size reduction parameter  $k$ . The associated core graph is the graph  $G/F_i$  which denotes  $G$  after contracting the vertices in the same components of  $F_i$ . We can define a map that lifts circulations  $\hat{\Delta}$  in the core graph  $G/F_i$  to circulations  $\Delta$  in the graph  $G$  by routing flow along the contracted paths in  $F_i$ . The lengths in the core graph  $\hat{\ell}$  (again let  $\hat{L} = \text{diag}(\hat{\ell})$ ) are chosen to upper bound the length of circulations when mapped back into  $G$  such that  $\|\hat{L}\hat{\Delta}\|_1 \geq \|L\Delta\|_1$ . Crucially, we must ensure these new lengths  $\hat{\ell}$  do not stretch the witness circulation  $\Delta^*$  when mapped into  $G/F_i$  by too much, so we can recover it from  $G/F_i$ . To achieve this goal, we choose  $F_i$  to be a low-stretch forest, that is, a forest with properties similar to those of a low-stretch tree. In Section 2.3, we summarize the central aspects of our core graph construction.

While each core graph  $G/F_i$  now has only  $\tilde{O}(m/k)$  vertices, it still has  $m$  edges which are too large for our recursion. To overcome this issue we build a spanner  $S(G, F_i)$  on  $G/F_i$  to reduce the number of edges to  $\tilde{O}(m/k)$ , which guarantees that for every edge  $e = (u, v)$  that we remove from  $G/F_i$  to obtain  $S(G, F_i)$ , there is a  $u$ -to- $v$  path in  $S(G, F_i)$  of length  $m^{o(1)}$ . Ideally, we would now recurse on each spanner  $S(G, F_i)$ , again approximating it with a collection of smaller core graphs and spanners. However, we face an obstacle: removing edges could destroy the witness circulation so that possibly no good circulation exists in any  $S(G, F_i)$ . To solve this problem, we compute an explicit embedding  $\Pi_{G/F_i \rightarrow S(G, F_i)}$  that maps each edge  $e = (u, v) \in G/F_i$  to a short  $u$ -to- $v$  path in  $S(G, F_i)$ . We can then show the following dichotomy: Let  $\hat{\Delta}^*$  denote the witness circulation when mapped into the core graph  $G/F_i$ . Then, either one of the edges  $e \in E_{G/F_i} \setminus E_{S(G, F_i)}$  has a spanner cycle consisting of  $e$  combined with  $\Pi_{G/F_i \rightarrow S(G, F_i)}(e)$  which is almost as good as  $\hat{\Delta}^*$ , or re-routing  $\hat{\Delta}^*$  into  $S(G, F_i)$  roughly preserves its quality. Figure 3 illustrates this dichotomy. Thus, either we find a good cycle using the spanner, or we can recursively find a solution  $S(G, F_i)$  on that almost matches  $\hat{\Delta}^*$  in quality. To construct our dynamic spanner with its efficient updates and strong stability guarantees under changes in the input graph, we design a new approach that diverges from other recent works on dynamic spanners.

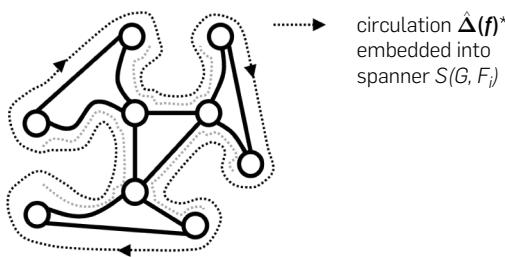
Our recursion uses  $d$  levels, where we choose the size reduction factor  $k$  such that  $k^d \approx m$  and the bottom level graphs have  $m^{o(1)}$  edges. Note that since we build  $B$  trees on  $G$  and recurse on the spanners of  $G/F_1, G/F_2, \dots, G/F_B$ , our recursive hierarchy has a branching factor of  $B = O(\log n)$  at each level of recursion. Thus, choosing  $d \leq \sqrt{\log n}$ , we get  $B^d = m^{o(1)}$  leaf nodes in our recursive hierarchy. Now, consider the

**Figure 3. Illustration of a dichotomy:** either one of the edges  $e \in E_{G/F_i} \setminus E_{S(G, F_i)}$  has a spanner cycle consisting of  $e$  combined with  $\Pi_{G/F_i \rightarrow S(G, F_i)}(e)$  which is almost as good as  $\hat{\Delta}(f)^*$ , or re-routing  $\hat{\Delta}(f)^*$  into  $S(G, F_i)$  roughly preserves its quality.

$G/F_i$  and spanner  $S(G, F_i)$



spanner  $S(G, F_i)$



forests  $F_{i_1}, F_{i_2}, \dots, F_{i_d}$  on the path from the top of our recursive hierarchy to a leaf node. We can patch these forests together to form a tree associated with the leaf node. For each of these trees, we maintain a link-cut tree data structure. Using this data structure, whenever we find a good cycle, we can route flow along it and detect edges where the flow has changed significantly. The cycles are either given by an off-tree edge or a collection of  $m^{o(1)}$  off-tree edges coming from a spanner cycle. We call the entire construction a *branching tree chain*. In Section 2.4, we elaborate on the overall composition of the data structure.

What have we achieved using this hierarchical construction compared to our simple, static algorithm? First, consider the setting of an *oblivious adversary*, where the gradient and length update sequences and the optimal circulation after each update are fixed in advance (i.e., the adversary is oblivious of the algorithm's random choices). In this setting, we can show that our spanner-of-core graph construction can survive through  $m^{1-o(1)}/k^i$  updates at level  $i$ . Meanwhile, we can rebuild these constructions in time  $m^{1+o(1)}/k^{i-1}$ , leading to an amortized cost per update of  $km^{o(1)} \leq m^{o(1)}$  at each level. This gives the first dynamic data structure for our undirected min-ratio problem with  $m^{o(1)}$  query time against an oblivious adversary.

However, our real problem is harder: the witness circulation in each round is  $\Delta(f)^* = f^* - f$  and depends on the updates we make to  $f$ , making our sequence of subproblems adaptive. While we cannot show that our data structure succeeds against an adaptive adversary, we give a data structure that works against a restricted adaptive adversary. We show that the witness circulation  $f^* - f$  lets us express the IPM as such a restricted adversary.

### 2.3. Building core graphs

In this section, we describe our core graph construction, which maps our dynamic undirected min-ratio cycle problem on a graph  $G$  with at most  $m$  edges and vertices into a problem of the same type on a graph with only  $\tilde{O}(m/k)$  vertices and  $m$  edges, and handles  $\tilde{O}(m/k)$  updates to the edges before we need to rebuild it.

**Forest routings and stretches.** To understand how to define the stretch of an edge  $e$  with respect to a forest  $F$ , it is useful to define how to *route* an edge  $e$  in  $F$ . Given a spanning forest  $F$ , every path and cycle in  $G$  can be mapped to  $G/F$  naturally (where we allow  $G/F$  to contain self-loops). On the other hand, if every connected component in  $F$  is rooted, where  $\text{root}_u^F$  denotes the root corresponding to a vertex  $u \in V$ , we can map every path and cycle in  $G/F$  back to  $G$  as follows. Let  $P = e_1, \dots, e_k$  be any (not necessarily simple) path in  $G/F$  where the preimage of every edge  $e_i$  is  $e_i^G = (u_i^G, v_i^G) \in G$ . The preimage of  $P$ , denoted  $P^G$ , is defined as the following concatenation of paths:

$$P^G \stackrel{\text{def}}{=} \bigoplus_{i=1}^k F[\text{root}_{u_i^G}^F, u_i^G] \oplus e_i^G \oplus F[v_i^G, \text{root}_{v_i^G}^F],$$

where we use  $A \oplus B$  to denote the concatenation of paths  $A$  and  $B$ , and  $F[a, b]$  to denote the unique  $ab$ -path in the forest  $F$ . When  $P$  is a circuit (that is, a not necessarily simple cycle),  $P^G$  is a circuit in  $G$  as well. One can extend these maps linearly

to all flow vectors and denote the resulting operators as  $\Pi_F : \mathbb{R}^{E(G)} \rightarrow \mathbb{R}^{E(G/F)}$  and  $\Pi_F^{-1} : \mathbb{R}^{E(G/F)} \rightarrow \mathbb{R}^{E(G)}$ . Since we let  $G/F$  have self-loops, there is a bijection between the edges of  $G$  and  $G/F$ , and thus  $\Pi_F$  acts like the identity function. Related routing schemes date back to Spielman-Teng<sup>34</sup> and are generally known as *portal routing*.

To make our core graph construction dynamic, the key operation we need to support is the dynamic addition of more root nodes, which results in forest edges being deleted to maintain the invariant each connected component has a root node. Whenever an edge is changing in  $G$ , we ensure that  $G/F$  approximates the changed edge well by forcing both its endpoints to become root nodes, which in turn makes the portal routing of the new edge trivial and this guarantees its stretch is 1.

For any edge  $e^G = (u^G, v^G)$  in  $G$  with image  $e$  in  $G/F$ , we set  $\hat{\ell}_e$ , the edge length of  $e$  in  $G/F$ , to be *an upper bound* on the length of the *forest routing* of  $e$ , that is, the path  $F[\text{root}_{u^G}^{F_e}, u^G] \oplus e^G \oplus F[v^G, \text{root}_{v^G}^{F_e}]$ . Meanwhile, we define  $\text{str}_e \stackrel{\text{def}}{=} \hat{\ell}_e / \ell_e$ , as an overestimate on the stretch of  $e$  w.r.t. the forest routing. A priori, it is unclear how to provide a single upper bound on the stretch of every edge, as the root nodes of the endpoints are changing over time. Providing such a bound for every edge is important for us as the lengths in  $G/F$  could otherwise be changing too often when the forest changes. We guarantee these bounds by a scheme that makes auxiliary edge deletions in the forest in response to external updates, with the resulting additional roots chosen carefully to ensure the length of upper bounds.

Now, for any flow  $f$  in  $G/F$ , its length in  $G/F$  is at least the length of its pre-image in  $G$ , that is,  $\|\mathbf{L}\Pi_F^{-1}f\|_1 \leq \|\hat{\mathbf{L}}^F f\|_1$ . Let  $\Delta^*$  be the optimal solution to Equation (2). We will show later how to build  $F$  such that with constant probability  $\|\hat{\mathbf{L}}^F \Delta^*\|_1 \leq \gamma \|\mathbf{L}\Delta^*\|_1$  holds for some  $\gamma = m^{o(1)}$ , solving Equation (2) on  $G/F$  with edge length  $\hat{\ell}$  and properly defined gradient  $\hat{g}$  on  $G/F$  yields an  $\frac{1}{\gamma}$ -approximate solution for  $G$ . The gradient  $\hat{g}$  is defined so that the total gradient of any circulation  $\Delta$  on  $G/F$  and its preimage  $\Pi_F^{-1}\Delta$  in  $G$  is the same, that is,  $\hat{g}^\top \Delta = g^\top \Pi_F^{-1}\Delta$ . The idea of incorporating gradients into portal routing was introduced in Kyng et al.<sup>24</sup>; our version of this construction is somewhat different to allow us to make it dynamic efficiently.

**Collections of low stretch decompositions (LSD).** The first component of the data structure is constructing and maintaining forests of  $F$  that form a *Low Stretch Decomposition (LSD)* of  $G$ . Informally, a  $k$ -LSD is a rooted forest  $F \subseteq G$  that decomposes  $G$  into  $O(m/k)$  vertex disjoint components. Given some positive edge weights  $\mathbf{v} \in \mathbb{R}_{>0}^E$  and reduction factor  $k > 0$ , we compute a  $k$ -LSD  $F$  and length upper bounds  $\hat{\ell}^F$  of  $G/F$  that satisfy two properties:

1.  $\text{str}_e = \hat{\ell}_e / \ell_e = \tilde{O}(k)$  for any edge  $e^G \in G$  with image  $e$  in  $G/F$ , and
2. The weighted average of  $\text{str}_e$  w.r.t.  $\mathbf{v}$  is only  $\tilde{O}(1)$ , that is,  $\sum_{e^G \in G} \mathbf{v}_{e^G} \cdot \text{str}_e \leq \tilde{O}(1) \cdot \|\mathbf{v}\|_1$ .

Item 1 guarantees that the solution to Equation (2) for  $G/F$  yields a  $\tilde{O}(k)$ -approximate one for  $G$ . However, this guarantee is not sufficient for our data structure, as our

$B$ -branching tree chain has  $d \approx \log_k m$  levels of recursion and the quality of the solution from the deepest level would only be  $\tilde{O}(k)^d \approx m^{1+o(1)}$ -approximate.

Instead, we compute  $k$  different edge weight assignments  $\mathbf{v}_1, \dots, \mathbf{v}_k$  via multiplicative weight updates so that the LSDs  $F_1, \dots, F_k$  have  $\tilde{O}(1)$  an average stretch on every edge in  $G : \sum_{j=1}^k \text{str}_{e_j}^{F_j} = \tilde{O}(k)$ , for all  $e^G \in G$  with image  $e$  in  $G/F$ .

By Markov's inequality, for any fixed flow  $f$  in  $G$ ,  $\|\hat{\mathbf{L}}^F f\|_1 \leq \tilde{O}(1) \|\mathbf{L}f\|_1$  holds for at least half the LSDs corresponding to  $F_1, \dots, F_k$ . Taking  $O(\log n)$  samples uniformly from  $F_1, \dots, F_k$ , say  $F_1, \dots, F_B$  for  $B = O(\log n)$  we get that with high probability

$$\min_{j \in [B]} \left\| \widetilde{\text{str}}_{e_j}^{F_j} \circ \mathbf{L}\Delta^* \right\|_1 \leq \tilde{O}(1) \|\mathbf{L}\Delta^*\|_1. \quad (3)$$

That is, it suffices to solve Equation (2) on  $G/F_1, \dots, G/F_B$  to find an  $\tilde{O}(1)$ -approximate solution for  $G$ .

## 2.4. Maintaining a branching tree chain

Our branching chain is constructed as follows:

1. Sample and maintain  $B = O(\log n)$   $k$ -LSDs  $F_1, F_2, \dots, F_B$ , and their associated core graphs  $G/F_i$ . Across  $O(m/k)$  updates at the top level, the forests  $F_i$  are *decremental*, that is, only undergo edge deletions (from root insertions), and will have  $\tilde{O}(m/k)$  connected components.
2. Maintain spanners  $\mathcal{S}(G, F_i)$  of the core graphs  $G/F_i$ , and embeddings  $\Pi_{E(G/F_i) \rightarrow \mathcal{S}(G, F_i)}$ , say with length increase  $\gamma_\ell = m^{o(1)}$ .
3. Recursively process the graphs  $\mathcal{S}(G, F_i)$ , that is, maintain LSDs and core graphs on those, and spanners on the contracted graphs, etc., for  $d$  total levels, with  $k^d = m$ .
4. Whenever a level  $i$  accumulates  $m/k^i$  total updates, hence doubling the number of edges in the graphs at that level, we rebuild levels  $i, i+1, \dots, d$ .

Recall that on average, the LSDs stretch lengths by  $\tilde{O}(1)$ , and the spanners  $\mathcal{S}(G, F_i)$  stretch lengths by  $\gamma_\ell$ . Hence, the overall data structure stretches lengths by  $\tilde{O}(\gamma_\ell)^d = m^{o(1)}$  (for appropriately chosen  $d$ ).

We now discuss how to update the forests  $G/F_i$  and spanners  $\mathcal{S}(G, F_i)$ . Intuitively, every time an edge  $e = (u, v)$  is changed in  $G$ , we will delete  $\tilde{O}(1)$  additional edges from  $F_i$ . This ensures that no edge's total stretch/routing-length increases significantly due to the deletion of  $e$ . As the forest  $F_i$  undergoes edge deletions, the graph  $G/F_i$  undergoes *vertex splits*, where a vertex has a subset of its edges moved to a newly created vertex. Thus, a key component of our data structure is to maintain spanners and embeddings of graphs undergoing vertex splits (and edge insertions/ deletions). Importantly, the amortized recourse (number of changes) to the spanner  $\mathcal{S}(G, F_i)$  is  $m^{o(1)}$  independent of  $k$ , even though the average degree of  $G/F_i$  is  $\Omega(k)$ . Hence, on average  $\Omega(k)$  edges will move in  $G/F_i$  per vertex split.

Overall, let every level have recourse  $\gamma_r = m^{o(1)}$  (independent of  $k$ ) per tree. Then each update at the top level induces  $O(B\gamma_r)^d$  (as we branch using  $B$  forests/core graphs at each level of the recursion) updates in the data structure overall. Intuitively, for the proper choice of  $d = o(1)$ , both the total

recourse  $O(B\gamma)^d$  and approximation factor  $\tilde{O}(\gamma_\ell)^d$  are  $m^{o(1)}$  as desired.

## 2.5. Going beyond oblivious adversaries by using IPM guarantees

The precise data structure in the previous section only works for *oblivious adversaries*, because we used that if we sampled  $B = O(\log n)$  LSDs, then w.h.p. there is a tree whose average stretch is  $\tilde{O}(1)$  with respect to a *fixed flow*  $f$ . However, since we are updating the flow along the circulations returned by our data structure, we influence future updates, so the optimal circulations our data structure needs to preserve are not independent of the randomness used to generate the LSDs. To overcome this issue, we leverage the key fact that the flow  $f^* - f$  is a good witness for the min-ratio cycle problem at each iteration.

To simplify our discussion, we focus on the role of the witness in ensuring the functioning of a single layer of core graph construction, which already captures the main ideas. We can prove that for any flow  $f$ ,  $g(f)^\top \Delta(f) / (100m + \|L(f)\Delta(f)\|_1) \leq -\tilde{\Omega}(1)$  holds where  $\Delta(f) = f^* - f$ . Then, the best solution to Equation (2) among the LSDs  $G/F_1, \dots, G/F_B$  maintains an  $\tilde{O}(1)$ -approximation of the quality of the witness  $\Delta(f) = f^* - f$  as long as

$$\min_{j \in [B]} \|\hat{L}^{F_j} \Delta(f)\|_1 \leq \tilde{O}(1) \|L(f)\Delta(f)\|_1 + \tilde{O}(m). \quad (4)$$

In this case, let  $\hat{\Delta}$  be the best solution obtained from graphs  $G/F_1, \dots, G/F_B$ . We have

$$\frac{g(f)^\top \hat{\Delta}}{\|L(f)\hat{\Delta}\|_1} \leq \frac{g(f)^\top \Delta(f)}{\tilde{O}(1) \|L(f)\Delta(f)\|_1 + \tilde{O}(m)} = -\tilde{\Omega}(1).$$

The additive  $\tilde{O}(m)$  term is there for technical reasons that can be ignored for now. We define the *width*  $w(f)$  of  $\Delta(f)$  as  $w(f) = 100 \cdot 1 + |L(f)\Delta(f)|$ . The name comes from the fact that  $w(f)_e$  is always at least  $|\ell(f)_e(f^*_e - f_e)|$  for any edge  $e$ . We show that the width is also slowly changing across IPM iterations, in that if the width changed by a lot, then the residual capacity of  $e$  must have changed significantly. This gives our data structure a way to predict which edges' contribution to the length of the witness flow  $f^* - f$  could have significantly increased.

Observe that for any forest  $F_j$  in the LSD of  $G$ , we have  $\|\hat{L}^{F_j} \Delta(f)\|_1 \leq \|\tilde{L}^{F_j} \circ w(f)\|_1$ . Thus, we can strengthen Equation (4) and show that the IPM potential can be decreased by  $m^{-o(1)}$  if

$$\min_{j \in [B]} \|\tilde{L}^{F_j} \circ w(f)\|_1 \leq \tilde{O}(1) \|w(f)\|_1. \quad (5)$$

Equation (5) also holds with w.h.p. if the collection of LSDs is built after knowing  $f$ . However, this does not necessarily hold after augmenting with  $\Delta$ , an approximate solution to Equation (2).

Due to the stability of  $w(f)$ , we have  $w(f + \Delta)_e \approx w(f)_e$  for every edge  $e$  whose length does not change a lot. For other edges, we update their edge length and force the stretch to be 1, that is,  $\tilde{str}_e^{F_j} = 1$  via the dynamic LSD maintenance, by shortcircuiting the routing of the edge  $e$  at its endpoints. To distinguish between the earlier stretch values and those after updating edges, let us denote the former by  $str_{old}$  and

the latter by  $\tilde{str}_{new}^{F_j}$ . This gives that for any  $j \in [B]$ , the following holds:

$$\|\tilde{str}_{new}^{F_j} \circ w(f + \Delta)\|_1 \lesssim \|\tilde{str}_{old}^{F_j} \circ w(f)\|_1 + \|w(f + \Delta)\|_1.$$

Using the fact that  $\min_{j \in [B]} \|\tilde{str}_{old}^{F_j} \circ w(f)\|_1 \leq \tilde{O}(\|w(f)\|_1)$ , we have the following:

$$\min_{j \in [B]} \|\tilde{str}_{new}^{F_j} \circ w(f + \Delta)\|_1 \lesssim \tilde{O}(\|w(f)\|_1) + \|w(f + \Delta)\|_1.$$

Thus, solving Equation (2) on the updated  $G/F_1, \dots, G/F_B$  yields a good enough solution for reducing IPM potential as long as the width of  $w(f + \Delta)$  has not decreased significantly, that is,  $\|w(f)\|_1 \leq \tilde{O}(1) \|w(f + \Delta)\|_1$ .

If the solution on the updated graphs  $G/F_1, \dots, G/F_B$  does not have a good enough quality, we know by the above discussion that  $\|w(f + \Delta)\|_1 \leq 0.5 \|w(f)\|_1$  must hold. Then, we re-compute the collection of LSDs of  $G$  and solve Equation (2) on the new collection of  $G/F_1, \dots, G/F_B$  again. Because each recomputation reduces the  $\ell_1$  norm of the width by a constant factor, and all the widths are bounded by  $\exp(\log^{o(1)} m)$  (which can be guaranteed by the IPM), there can be at most  $\tilde{O}(1)$  such recomputations. At the top level, this only increases our runtime by  $\tilde{O}(1)$  factors.

The full construction is along these lines, but more complicated since we recursively maintain the solutions on the spanners of each core graph  $G/F_1, \dots, G/F_B$ . In the full version of the paper, we describe and analyze a multi-level rebuilding scheme that extends the above reasoning to our full data structure.

## 3. CONCLUSION

In this paper, we presented an almost-linear time algorithm for minimum-cost flow, maximum flow, and more generally, all convex single-commodity flows. Our work essentially settles the complexity of several fundamental and intensely studied problems in algorithms design. We hope that the ideas introduced in this work will spur further research in several directions, including simpler and faster algorithms for flow problems; hopefully resulting in a significant impact on algorithms in practice.

## Acknowledgments

Li Chen was supported by NSF Grant CCF-2106444. Rasmus Kyng and Maximilian Probst Gutenberg have received funding from the grant “Algorithms and complexity for high-accuracy flows and convex optimization” (no. 200021 204787) of the Swiss National Science Foundation. Yang P. Liu was supported by NSF CAREER Award CCF-1844855 and NSF Grant CCF-1955039. Richard Peng was partially supported by NSF CAREER Award CCF-1846218 and NSERC Discovery Grant RGPIN-2022-03207. Sushant Sachdeva’s research is supported by a Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant RGPIN-2018-06398 and an Ontario Early Researcher Award (ERA) ER21-16-284.

The authors thank the 2021 Hausdorff Research Institute for Mathematics Program Discrete Optimization. The authors are very grateful to Yin Tat Lee and Aaron Sidford for several useful discussions.  $\square$

## References

1. Abboud, A., Krauthgamer, R., Li, J., Panigrahi, D., Saranurak, T., Trabelsi, O. Breaking the cubic barrier for all-pairs max-flow: Gomory-hu tree in nearly quadratic time. In *63rd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2022* (Denver, CO, USA, October 31–November 3, 2022), IEEE, NY, 884–895.
2. Abraham, I., Neiman, O. Using petal-decompositions to build a low stretch spanning tree. *SIAM J. Comput.* 48, 2 (2019), 227–248.
3. Alon, N., Karp, R.M., Peleg, D., West, D. A graph-theoretic game and its application to the k-server problem. *SIAM J. Comput.* 24, 1 (1995), 78–100.
4. Axiotis, K., Mądry, A., Vladu, A. Circulation control for faster minimum cost flow in unit-capacity graphs. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)* (2020), IEEE, NY, 93–104.
5. Axiotis, K., Mądry, A., Vladu, A. Faster sparse minimum cost flow by electrical flow localization. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)* (2022), IEEE, NY, 528–539.
6. Bernstein, A., Gutenberg, M.P., Saranurak, T. Deterministic decremental ssssp and approximate min-cost flow in almost-linear time. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)* (2022), IEEE, NY, 1000–1008.
7. Bernstein, A., Nanongkai, D., Wulff-Nilsen, C. Negative-weight single-source shortest paths in near-linear time. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)* (2022), IEEE, NY, 600–611.
8. Brand, J.v.d., Lee, Y.T., Liu, Y.P., Saranurak, T., Sidford, A., Song, Z., et al. Minimum cost flows, mdps, and  $\ell_1$ -regression in nearly linear time for dense instances. In *STOC* (2021), ACM, NY, 859–869.
9. Brand, J.v.d., Lee, Y.-T., Nanongkai, D., Peng, R., Saranurak, T., et al. Bipartite matching in nearly-linear time on moderately dense graphs. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)* (2020), IEEE, NY, 919–930.
10. Brand, J.v.d., Lee, Y.T., Sidford, A., Song, Z. Solving tall dense linear programs in nearly linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing (STOC)* (2020), ACM, NY, 775–788.
11. Cen, R., Li, J., Nanongkai, D., Panigrahi, D., Saranurak, T., Quanrud, K. Minimum cuts in directed graphs via partial sparsification. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)* (2021), IEEE, Denver, CO, 1147–1158.
12. Christiano, P., Kelner, J.A., Mądry, A., Spielman, D.A., Teng, S. Electrical flows, Laplacian systems, and faster approximation of maximum flow in undirected graphs. In *Proceedings of the 43rd ACM Symposium on Theory of Computing (STOC)* (2011), ACM, NY, 273–282.
13. Daitch, S.I., Spielman, D.A. Faster approximate lossy generalized flow via interior point algorithms. In *Proceedings of the forty ninth annual ACM symposium on Theory of computing* (2008), ACM, NY, 451–460.
14. Dantzig, G.B. Application of the simplex method to a transportation problem. In *Activity Analysis and Production and Allocation*, T.C. Koopmans (ed.), (1951), John Wiley and Sons, New York, 1862–1867.
15. Even, S., Tarjan, R.E. Network flow and testing graph connectivity. *SIAM J. Comput.* 4, 4 (1975), 507–518.
16. Gao, Y., Liu, Y.P., Peng, R. Fully dynamic electrical flows: Sparse maxflow faster than goldberg-rao. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)* (2022), IEEE, NY, 516–527.
17. Goldberg, A., Tarjan, R. Finding minimum-cost circulation by successive approximation. *Math. Oper. Res.* 15 (1990), 430–466.
18. Goldberg, A.V., Rao, S. Beyond the flow decomposition barrier. *J. ACM* 45, 5 (1998), 783–797. Announced at FOCS’97.
19. Hopcroft, J.E., Karp, R.M. An  $n^{5/2}$  algorithm for maximum matchings in bipartite graphs. *SIAM J. Comput.* 2, 4 (Dec 1973), 225–231.
20. Karmarkar, N. A new polynomial-time algorithm for linear programming. In *STOC* (1984), ACM, NY, 302–311.
21. Karzanov, A.V. On finding maximum flows in networks with special structure and some applications. *Matematicheskie Voprosy Upravleniya Proizvodstvom* 5, (1973), 81–94.
22. Kathuria, T., Liu, Y.P., Sidford, A. Unit capacity maxflow in almost  $O(m^{4/3})$  time. In *61st IEEE Annual Symposium on Foundations of Computer Science (FOCS)* (2020), IEEE, NY, 119–130.
23. Kelner, J.A., Lee, Y.T., Orecchia, L., Sidford, A. An almost-linear-time algorithm for approximate max flow in undirected graphs, and its multicommodity generalizations. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)* (2014), Society for Industrial and Applied Mathematics, Portland, OR, 217–226.
24. Kyng, R., Peng, R., Sachdeva, S., Wang, D. Flows in almost linear time via adaptive preconditioning. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing* (2019), ACM, NY, 902–913.
25. Li, J., Nanongkai, D., Panigrahi, D., Saranurak, T., Yingchareonthawornchai, S. Vertex connectivity in poly-logarithmic max-flows. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing* (2021), ACM, NY, 317–329.
26. Li, J., Panigrahi, D. Approximate gomory-hu tree is faster than  $n-1$  max-flows. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing* (2021), ACM, NY, 1738–1748.
27. Liu, Y.P., Sidford, A. Faster energy maximization for faster maximum flow. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing (STOC 2020)* (2020), ACM, NY, 803–814.
28. Mądry, A. Navigating central path with electrical flows: From flows to matchings, and back. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science* (2013), IEEE, NY, 253–262.
29. Mądry, A. Computing maximum flow with augmenting electrical flows. In *57th IEEE Annual Symposium on Foundations of Computer Science (FOCS)* (2016), IEEE Computer Society, NY, 593–602.
30. Peng, R. Approximate undirected maximum flows in  $O(m\text{polylog}(n))$  time. In *Proceedings of the twenty-seventh annual ACM-SIAM symposium on Discrete algorithms* (2016), SIAM, Philadelphia, PA, 1862–1867.
31. Renegar, J. A polynomial-time algorithm, based on newton's method, for linear programming. *Math. Program.* 40, 1 (1988), 59–93.
32. Sherman, J. Nearly maximum flows in nearly linear time. In *54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)* (2013), IEEE Computer Society, NY, 263–269.
33. Sherman, J. Area-convexity,  $\ell_\infty$  regularization, and undirected multicommodity flow. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing* (2017), ACM, NY, 452–460.
34. Spielman, D.A., Teng, S. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing (STOC)* (2004), ACM, NY, 81–90.
35. van den Brand, J., Gao, Y., Jambulapati, A., Lee, Y.T., Liu, Y.P., Peng, R., et al. Faster maxflow via improved dynamic spectral vertex sparsifiers. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing (STOC 2022)* (2022), ACM, NY, 543–556.

**Li Chen** (lichen@gatech.edu), Georgia Institute of Technology, Atlanta, GA, USA.

**Rasmus Kyng and Maximilian Probst Gutenberg** ([kyng, maxprobst ]@inf.ethz.ch), ETH Zurich, Zurich, Switzerland.

**Yang P. Liu** (yangpliu@stanford.edu), Stanford University, Stanford, CA, USA.

**Richard Peng** (ly5peng@uwaterloo.ca), University of Waterloo, Waterloo, ON, Canada.

**Sushant Sachdeva** (sachdeva@cs.toronto.edu), University of Toronto, Toronto, ON, Canada.

# CAREERS

## Hofstra University

### **Two Assistant/Associate Computer Science Faculty Positions (Computer Science & Computer Science Specializing in Cybersecurity)**

The Department of Computer Science, at Hofstra University, invites applications for two tenure-track positions: one in Computer Science with no area restrictions and one in Computer Science specializing in Cybersecurity. The positions are at the Assistant Professor or Associate Professor level and begin in Fall 2024. Preference for the Cybersecurity position will be given to those able to teach courses in system security, ethical hacking, cryptography or digital forensics. We value faculty's commitment to teaching and sustained research productivity and applicants are expected to have a Ph.D. in Computer Science or a closely related field.

Our academic programs include undergraduate and masters level degrees. We have ten full time faculty housed in the new Science & Innovation Center which boasts state-of-the-art academic facilities.

Please upload the following with your application for consideration for the Computer Science position at <https://hofstra.peopleadmin.com/postings/2519> or the Cybersecurity specialization position at <https://hofstra.peopleadmin.com/postings/2551>:

- ▶ Cover letter
- ▶ Title, name, address, email and phone of three references
- ▶ Curriculum vitae
- ▶ Statement of teaching experience and philosophy
- ▶ Statement of current research interests and anticipated future research directions
- ▶ Two sample publications in a single PDF

All applications received by January 15, 2024 will receive full consideration. Inquiries should be directed to [SEAS4@hofstra.edu](mailto:SEAS4@hofstra.edu).

Hofstra University is an equal opportunity employer, committed to fostering diversity in its faculty, administrative staff and student body, and encourages applications from the entire spectrum of a diverse community.

---

## Rutgers University

### **Assistant Professor and Associate Professor Positions (Tenure-Track/Tenured)**

The Computer Science Department at Rutgers University, New Brunswick NJ, invites applications for multiple tenure-track/tenured positions at the Assistant Professor and Associate Professor levels. We will consider outstanding candidates at the Professor level as well.

We invite applications from candidates making research contributions in any area of CS, and welcome applicants with interdisciplinary approaches. We are especially interested in Algorithms, Machine Learning and Data Science,

High-performance Computing and Scalable Systems, Quantum Computing, and other subfields with broad potential for collaborative impact across the department and the university.

Rutgers is committed to invest significant resources to promote interdisciplinary research and education in Data Science and Artificial Intelligence. CS expects to make one appointment this year as part of a larger interdisciplinary cluster hire at Rutgers. The successful candidate will be expected to contribute to cluster initiatives as well as to Rutgers research and education programs in Data Science, Machine Learning, and AI.

Rutgers CS has a broad portfolio of research, innovation, undergraduate and graduate education, and outreach, with strengths in AI and robotics, foundations of computer science, and systems. Rutgers stands among America's oldest and most diverse public research universities and is the top-ranked public university in the New York/New Jersey metropolitan area.

Fostering diversity and inclusion in computing is part of the department's culture and mission, and the department values faculty efforts to broaden the field through scholarship, teaching, mentoring, and professional service. We welcome applications from all qualified candidates, including those with non-traditional career paths or who have achieved excellence in careers outside academia.

Responsibilities will include research, supervision of Ph.D. students, and teaching undergraduate- and graduate-level courses in Computer Science. Pursuit of external research funding is expected.

Requirements: Successful completion of a Ph.D. in Computer Science or a closely related field is required by the start date.

Timeline: The appointment will start September 1, 2024. Applications received by January 1, 2024, will be given priority.

How to Apply: Applicants should submit their cover letter, CV, a research statement addressing both past and future work, a diversity statement outlining accomplishments and approach for broadening participation in computing, a teaching statement, and contact information for at least three references. For details, including application procedure, please see <https://go.rutgers.edu/CSTTFaculty>.

Contact Info: [hiring-committee@cs.rutgers.edu](mailto:hiring-committee@cs.rutgers.edu).

Rutgers Policies: Offer is contingent upon successful completion of all pre-employment screenings. Rutgers is an equal opportunity employer: see <http://ehr.rutgers.edu/non-discrimination-statement>.

---

## Rutgers University

### **Instructional, Non-Tenure Track Position**

The Department of Computer Science at Rutgers University invites applications for an instructional, nontenure track position. A PhD degree

in Computer Science or related field is required. The teaching faculty position targets Rutgers's instructional needs in the areas of Introductory Topics in Computer Science, Algorithms, Computer Systems, and Data Science. Rutgers has recently launched a new major in Data Science and is committed to teaching excellence at the undergraduate and graduate levels. Depending on experience, suitable candidates will be invited to teach in our Masters programs or to develop and manage applied learning opportunities such as internships.

The appointment will start as early as January 1, 2024, and the length of the contract will be a minimum of one year, and may be longer contingent on the candidate's experience. Reappointments are typically for multiple years. We are searching for an Assistant Teaching Professor, but will consider excellent applications for all ranks, including Associate Teaching Professor, Teaching Professor, and Distinguished Teaching Professor.

Responsibilities include teaching Computer Science undergraduate classes, interviewing and hiring recitation instructors and graders (typically senior undergraduate or Masters students), coordinating and supervising recitation sections, creating exams, homework, and programming assignments, possibly in collaboration with other instructors, coordinating and supervising grading, and curriculum development. Other duties, such as managing applied learning opportunities such as internships, are also possible depending on experience.

Consideration of candidates will begin immediately and continue until positions are filled. For questions regarding this position please contact: [nt-hiring@cs.rutgers.edu](mailto:nt-hiring@cs.rutgers.edu).

For more information about CS at Rutgers go to <https://cs.rutgers.edu> and for Data Science see <https://mps.rutgers.edu/data-science>.

To apply for the position, go to: <https://jobs.rutgers.edu/postings/211853> and submit your CV and contact information for three references.

Rutgers subscribes to the value of academic diversity and encourages applications from individuals with varied experiences, perspectives, and backgrounds. Women, minorities, and persons with disabilities are encouraged to apply. Rutgers is an affirmative action/equal opportunity employer. Offer is contingent upon successful completion of all pre-employment screenings.

---

## University of Illinois Chicago (UIC)

### **Full-Time Teaching Faculty (Lecturer or Clinical Professor)**

The Computer Science Department at the University of Illinois Chicago (UIC) seeks to hire full-time teaching faculty (Lecturer or Clinical Professor). Candidates would work alongside 24 full-time teaching faculty with over 150 years of experience and 14 awards for excellence. Standard teaching load is three course sections per se-



At the University of the Bundeswehr Munich, the Cyber Defence and Smart Data (CODE) research institute brings together experts from different scientific disciplines as well as competencies from the private and public sectors. The common goal is joint research in the fields of cyber security, artificial intelligence, and quantum technologies. As part of its commitment to education, CODE offers a master's degree in cyber security and contributes to the undergraduate and postgraduate computer science curricula offered by the university's department of computer science.

Research and teaching at CODE are supported by world-class equipment, labs, and computing resources. A dedicated building to house the institute is in its planning phase and will be constructed over the coming years. To further strengthen the institute, the Department of Computer Science wishes to fill the following position as soon as possible:

### Full Professor (W3) in Reverse Engineering

We are looking for an exceptional candidate with an outstanding scientific track record and didactic skills. She or he should be internationally renowned for at least one current direction within the broad research area of reverse engineering, usually demonstrable through publications at top-level international conferences and journals.

The advertised full professorship focuses on methods and applications for the analysis of software whose design is not or only partially disclosed, for example, in order to analyse malicious software used by attackers, unknown data formats, or manipulations in the context of software supply chains, and also includes specifics of the executing hardware, for example, security vulnerabilities in the design of microprocessors. The focus may be on the advancement and use of static and dynamic code analysis methods, the application of machine learning approaches for this purpose, or the development of reverse engineering countermeasures.

The appointee is expected to contribute to CODE and joint initiatives and activities. Experience in the acquisition of research funding is required.

The successful candidate is expected to offer courses in her or his area of specialization as part of the cyber security master's program and should be willing to contribute to the computer science undergraduate and postgraduate degrees. Duties will also include participation in the academic self-governance of the university and interaction across departments.

Formal requirements for appointment are a university degree, pedagogical ability, an ability for scientific excellence demonstrated through a doctoral degree, and additional scientific achievements (habilitation or successful evaluation of a junior/assistant professor position, or equivalent scientific achievements). Further, the candidate will be expected to assume a leadership role based on gender equality and diversity.

The Universität der Bundeswehr München offers academic programs directed primarily at officer candidates and officers, who can obtain bachelor's and master's degrees in a trimester system. Study programs are complemented by an integrated program entitled "studium plus", which consists of interdisciplinary seminars, tutorials and training in key professional qualifications.

The recruitment requirements and the employment status of professors are governed by the Federal Civil Service Act (Bundesbeamtengegesetz). To be appointed a civil servant (Beamter), candidates must be below the age of 50 on the date of their appointment.

The university seeks to increase the number of female professors and thus explicitly invites women to submit applications. Severely handicapped candidates with equal qualifications will receive special consideration.

Please submit your application documents **31.12.2023** via email to berufungen.informatik@unibw.de or via mail, marked as confidential personnel matter, to **the Dean of the Department of Computer Science, Prof. Dr. Oliver Rose, Universität der Bundeswehr München, D-85577 Neubiberg.**

By submitting your application, you agree to the processing of your data by the responsible parties. More information about data privacy is available on the website of UniBw München.



mester. Additionally, our teaching faculty participates and often lead efforts in shared governance and committees, curriculum decisions and development, and computer science education at a local and national level.

#### Minimum Qualifications:

The Clinical Professor track is a long-term career track that starts with Clinical Assistant Professor and offers advancement to Clinical Associate and Clinical Full Professor. Minimum qualifications include a PhD in Computer Science or a closely related field. Teaching is the primary focus for clinical faculty, but there are also opportunities for research.

The Lecturer track is a long-term career track that starts with Lecturer and offers opportunities for advancement to Senior Lecturer. Minimum qualifications include an MS in Computer Science and significant teaching experience.

The department seeks candidates interested in all areas of computer science, but in particular systems and software engineering. Submit applications online at <https://jobs.uic.edu>. Include:

- ▶ A curriculum vitae,
- ▶ Contact information for at least three references,
- ▶ One-page statement on your teaching philosophy and how it is inclusive to a diverse student population.

For more information, send an email to [shanon@uic.edu](mailto:shanon@uic.edu). For fullest consideration, apply by 10/18/23. Applications will be accepted and reviewed until the positions are filled.

The University of Illinois at Chicago is a Minority Serving Institution, an HSI, and an AANAPISI. UIC is an affirmative action, equal opportunity employer, dedicated to the goal of building a culturally diverse and pluralistic faculty and staff committed to teaching and working in a multicultural environment. We are committed to equal employment opportunity regardless of race, color, national origin, sex, religion, age, sexual orientation, gender identity, Veteran, or disability status.

The University of Illinois may conduct background checks on all job candidates upon acceptance of a contingent offer. Background checks will be performed in compliance with the Fair Credit Reporting Act.

The University of Illinois System requires candidates selected for hire to disclose any documented finding of sexual misconduct or sexual harassment and to authorize inquiries to current and former employers regarding findings of sexual misconduct or sexual harassment. For more information, see the UI HR website: <https://www.hr.uillinois.edu/cms/One.aspx?portalId=4292&mp;pageId=1411899>

The University of Illinois System is an equal opportunity employer, including but not limited to disability and/or veteran status, and complies with all applicable state and federal employment mandates.

Please visit Required Employment Notices and Posters to view our non-discrimination statement and find additional information about required background checks, sexual harassment/misconduct disclosures, COVID-19 vaccination requirement, and employment eligibility review through E-Verify.

The university provides accommodations to applicants and employees.

[CONTINUED FROM P. 96] And it's not clear that you could build the same kind of dual-license company today—not just because we got lucky, but because software architecture is fundamentally different. We were able to leverage the fact that using Berkeley DB required linking our code with a customer's code, and our license relied on that. This architecture appears less frequently today.

**Let's talk about your work on data provenance, which suggests various practical uses for collecting information about where data comes from and how it is manipulated and stored. However, arriving at useful applications has been something of a rocky journey.**

The history of data provenance is frustrating. The field dates back to a paper by Peter Buneman, who pointed out that because relational databases have rigorous, formal mathematical underpinning, you could state precisely where data comes from and why it was being produced as the result of a query. From there, the workflow community said, "If we use this to track scientific workflows, it will make things reproducible." And then my group were the fringe lunatics who said, "Actually, you can't make things reproducible unless you have provenance at the system level to understand if people installed new libraries or upgraded the compiler or changed their system in some other way."

**The community then issued a series of provenance challenges to better understand the different representations used for provenance and explore their expressive power.**

Yes, to their credit. However, while the scenarios were realistic, there wasn't a user community who cared. So, as a result, in my opinion, the standards that emerged did not help people who want to use provenance, who often end up having to implement their own language on top of the standard.

The other thing that's held provenance back is that the places where it's most useful, at least so far, are the kinds of places that I liken to insurance. Nobody goes out and buys insurance because they want it; they buy insurance because they have to. People don't want to "pay for" provenance

**"People don't want to 'pay for' provenance; they only discover they *should* have when something bad happens—  
their system gets infiltrated, they end up with malware on it."**

(that is, add any overhead to their work); they only discover they *should* have when something bad happens—their system gets infiltrated, they end up with malware on it.

**Without compelling applications, it's difficult to say people should collect this type of data.**

One of my "aha" moments came during a chat with an ecologist who had been a computer scientist. I described what we were building, and she was very receptive. I asked if she would use those types of tools. She said, "Well, what would I have to do to use them?"

I said, "Basically, once the package is installed, all you really have to do is type one command before you run your analysis." And that was a deal breaker.

**You've spoken a lot about how the field has been partitioned over the years, with different micro-communities that don't talk to one another. What sort of reception do you get?**

I think people largely agree, and at the same time, they don't know how to push back, because if you want to get tenure, the easy thing to do is pick your little community and do really well in that community.

**I detected a familiar note of desperation in some of the follow-up questions you've received—like, 'how can you expect everyone to keep up with adjacent fields when we struggle to stay on top of our own?' To which your response was,**

'you don't need to be an expert in everything, just to have a vague sense of what people are working on so you can draw ideas when the opportunity arises.'

I think it's really hard for junior people to walk into a project and realize they can't control everything. By the same token, how do I know my students are doing the right things? It's hard.

**You've made a similar point about the separation between academic and product communities.**

Back when both communities were smaller, there was less of a sharp divide between the people who built products and the people who do research. And I understand why that's shifted, but I do worry that our research is less informed by real problems, as opposed to problems that academics find interesting.

And yet industry has very real problems, and most companies aren't large enough to have their own internal research team to solve them. What can we do to bridge that gap? Longing for the good old days when industry and academia met at conferences may not be the right solution. But I do think there's a problem, and we need to find a 2023 solution instead of the 1980 solution.

**It also seems the research community's definition of practical may be a little different than the product community's. Is the goal to build something that's forward-looking or something that can be commercialized within six months?**

I think that's really key. If your research could be commercialized in six months, that assumes you're building a new thing. Academics like that kind of research. But the research that has real impact on industry, that solves a problem faced by every single company, might not be a brand-new thing that's going to make you rich. And that's the kind of research that is harder to do now. In some sense, this brings us full circle. Every single data scientist has really good problems to solve, and if we can find the right way to trick them into capturing data provenance by giving them tools that make their job easier, then maybe we can do work that has real impact. □

**Leah Hoffmann** is a technology writer based in Piermont, NY, USA.

## Q&A

# Ulterior Motives

*2023–2024 ACM Athena Lecturer Margo Seltzer recalls the motivations behind the development of the Berkeley DB database software library, and other achievements during her career.*

MARGO SELTZER, THE Canada 150 Research Chair in Computer Systems at the University of British Columbia and 2023–2024 ACM Athena Lecturer, is the kind of researcher who stands out not just for her accomplishments, but for her tirelessness. After building a database software library that underpinned many first-generation Internet services, she worked on topics that range from file systems and storage to capturing and accessing data provenance. Here, she speaks with Leah Hoffmann about finding impactful research projects—and keeping up with everything that's going on in the field.

The story of Berkeley DB, the database software library that you built with Keith Bostic and Mike Olson, has been told before at greater length, but let me see if I can summarize. Your work on packages such as hash and B-tree was released with Berkeley Unix as the DB 1.85 library. Then, as a side project, you and Mike created a transactional storage system—which Netscape later wanted to integrate into the LDAP directory server it was building on a Berkeley DB core. That prompted you and Keith to launch Sleepycat Software and create a production-quality transactional library. Before Netscape approached you, did you guys ever think, “maybe we should commercialize, maybe there’s something there”?

I had always wanted to build a database with the tool-based Unix-style philosophy. And at one of my first jobs, I was involved in a development effort that started me thinking along these lines.



Keith had also been thinking about something similar; he was reimplementing the vi editor, and he felt the right way to do it was to build it on top of a record management package. Finally, Mike Olson had built B-trees at all three of the jobs that he had before going back to grad school, so his motivation was, “I’m going to build one more B-tree, release it, and I will never have to do this again.”

So, in some sense, all three of us had ulterior motives to do this.

**The Netscape deal essentially provided seed money for Sleepycat—you did not raise additional funds.**

People seem to think the only way to build a company is to raise \$1 billion

and go big, and I do not understand that. An organically grown startup is a much saner lifestyle—and if you genuinely care about the product, as opposed to getting famous and/or rich, that is absolutely the way to do it. We controlled our destiny, start to stop. And we had a really good time. To be fair, we worked our butts off, but many of our employees still describe it as the best job they ever had.

**Sleepycat also pioneered the first open-source dual license.**

Today, you look at open source and everybody says, “Oh, we know what that is.” But that was not always the case. We had to educate the market to a large extent. [CONTINUED ON P. 95]



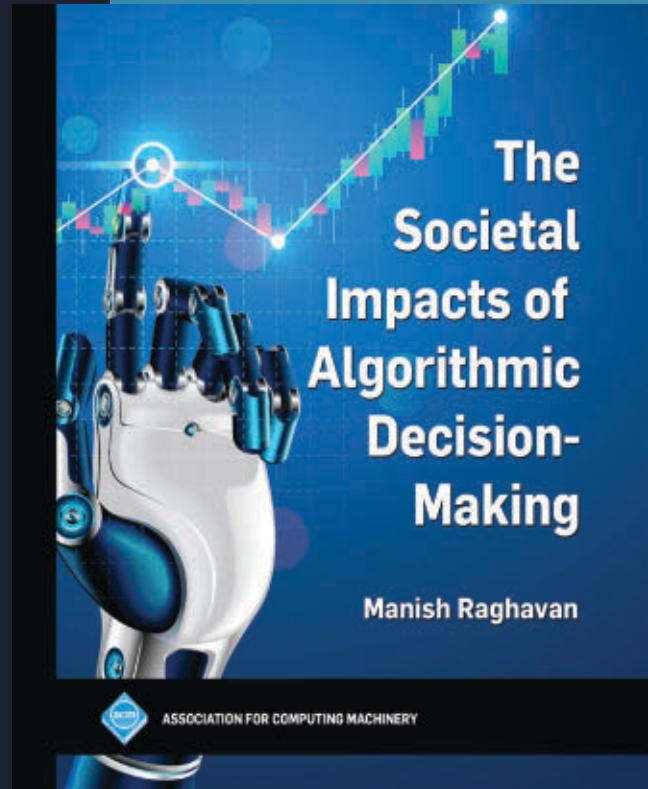
## ACM BOOKS

### Collection III

This book demonstrates the need for and the value of interdisciplinary research in addressing important societal challenges associated with the widespread use of algorithmic decision-making. Algorithms are increasingly being used to make decisions in various domains such as criminal justice, medicine, and employment. While algorithmic tools have the potential to make decision-making more accurate, consistent, and transparent, they pose serious challenges to societal interests. For example, they can perpetuate discrimination, cause representational harm, and deny opportunities.

*The Societal Impacts of Algorithmic Decision-Making* presents several contributions to the growing body of literature that seeks to respond to these challenges, drawing on techniques and insights from computer science, economics, and law. The author develops tools and frameworks to characterize the impacts of decision-making and incorporates models of behavior to reason about decision-making in complex environments. These technical insights are leveraged to deepen the qualitative understanding of the impacts of algorithms on problem domains including employment and lending.

The social harms of algorithmic decision-making are far from being solved. While easy solutions are not presented here, there are actionable insights for those who seek to deploy algorithms responsibly. The research presented within this book will hopefully contribute to broader efforts to safeguard societal values while still taking advantage of the promise of algorithmic decision-making.



## The Societal Impacts of Algorithmic Decision-Making

**Manish Raghavan**

ISBN: 979-8-4007-0859-6  
DOI: 10.1145/3603195

# MATLAB FOR AI

Boost system design and simulation with explainable and scalable AI. With MATLAB and Simulink, you can easily train and deploy AI models.

[mathworks.com/ai](https://mathworks.com/ai)

