

COMPUTING edge

- Human-Computer Interaction (HCI)
- Edge Computing
- Internet of Things (IoT)
- Economics



APRIL 2026

www.computer.org



IEEE Computer Society

Grants for

EMERGING TECHNOLOGY ACTIVITIES

MAKE AN IMPACT | CREATE SOLUTIONS

Are you connecting the computing community with emerging technologies? Help advance emerging tech to create solutions for the betterment of humanity.

Every year, we give up to **US\$50,000** in funding per project for these efforts.

Learn more at

computer.org/communities/emerging-technology-fund



billy

STAFF

Editor

Lucy Holden

Periodicals Portfolio Senior Managers

Carrie Clark and Kimberly Sperka

Director, Publications and Special Projects

Robin Baldwin

Production & Design Artist

Carmen Flores-Garvey

Publications Operations Project Manager

Christine Shaughnessy

Senior Advertising Coordinator

Debbie Sims

Periodicals Portfolio Specialist

Priscilla An

Circulation: *ComputingEdge* (ISSN 2469-7087) is published monthly by the IEEE Computer Society. IEEE Headquarters, Three Park Avenue, 17th Floor, New York, NY 10016-5997; IEEE Computer Society Publications Office, 10662 Los Vaqueros Circle, Los Alamitos, CA 90720; voice +1 714 821 8380; fax +1 714 821 4010; IEEE Computer Society Headquarters, 2001 L Street NW, Suite 700, Washington, DC 20036.

Postmaster: Send address changes to *ComputingEdge*-IEEE Membership Processing Dept., 445 Hoes Lane, Piscataway, NJ 08855. Periodicals Postage Paid at New York, New York, and at additional mailing offices. Printed in USA.

Editorial: Unless otherwise stated, bylined articles, as well as product and service descriptions, reflect the author's or firm's opinion. Inclusion in *ComputingEdge* does not necessarily constitute endorsement by the IEEE or the Computer Society. All submissions are subject to editing for style, clarity, and space.

Reuse Rights and Reprint Permissions: Educational or personal use of this material is permitted without fee, provided such use: 1) is not made for profit; 2) includes this notice and a full citation to the original work on the first page of the copy; and 3) does not imply IEEE endorsement of any third-party products or services. Authors and their companies are permitted to post the accepted version of IEEE-copyrighted material on their own Web servers without permission, provided that the IEEE copyright notice and a full citation to the original work appear on the first screen of the posted copy. An accepted manuscript is a version which has been revised by the author to incorporate review suggestions, but not the published version with copy-editing, proofreading, and formatting added by IEEE. For more information, please go to: http://www.ieee.org/publications_standards/publications/rights/paperversionpolicy.html. Permission to reprint/republish this material for commercial, advertising, or promotional purposes or for creating new collective works for resale or redistribution must be obtained from IEEE by writing to the IEEE Intellectual Property Rights Office, 445 Hoes Lane, Piscataway, NJ 08854-4141 or pubs-permissions@ieee.org. Copyright © 2026 IEEE. All rights reserved.

Abstracting and Library Use: Abstracting is permitted with credit to the source. Libraries are permitted to photocopy for private use of patrons, provided the per-copy fee indicated in the code at the bottom of the first page is paid through the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923.

Unsubscribe: If you no longer wish to receive this *ComputingEdge* mailing, please email IEEE Computer Society Customer Service at help@computer.org and type "unsubscribe *ComputingEdge*" in your subject line.

IEEE prohibits discrimination, harassment, and bullying. For more information, visit www.ieee.org/web/aboutus/whatis/policies/p9-26.html.

2026 IEEE Computer Society Magazine Editors in Chief

Computer

Bret Michael, *Naval Postgraduate School (Interim EIC)*

Computing in Science & Engineering

Jeffrey Carver, *University of Alabama*

IEEE Annals of the History of Computing

Troy Astarte, *Swansea University*

IEEE Computer Graphics and Applications

Pak Chung Wong, *Rocketgraph*

IEEE Intelligent Systems

Jeffrey Voas, *NIST*

IEEE Internet Computing

Weisong Shi, *University of Delaware*

IEEE Micro

Hsien-Hsin Sean Lee, *Intel Corporation*

IEEE MultiMedia

Balakrishnan Prabhakaran, *University of Texas at Dallas*

IEEE Pervasive Computing

Fahim Kawsar, *Omnibuds and University of Glasgow*

IEEE Security & Privacy

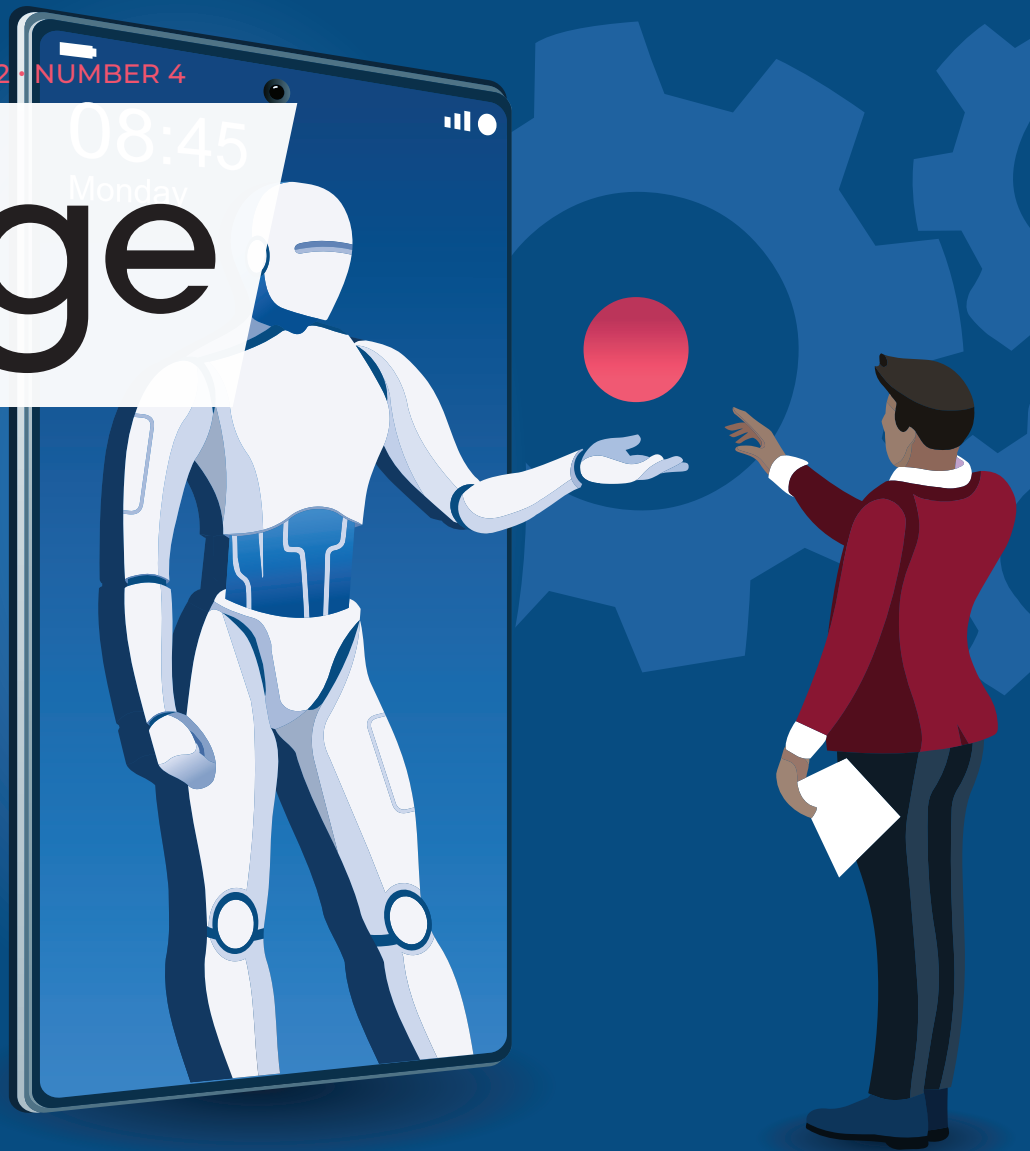
Jeffrey Voas, *NIST*

IEEE Software

Sigrid Eldh, *Ericsson, Mälardalen University, Sweden; Carleton University, Canada*

IT Professional

Charalampos Z. Patrikakis, *University of West Attica*



8

Design Principles
and Challenges
for Gaze + Pinch
Interaction in XR

16

Touch and Feel
Virtual Objects

20

Toward Next-
Generation
Human-Computer
Interface Based on
Earables

Human-Computer Interaction (HCI)

8 Design Principles and Challenges for Gaze + Pinch Interaction in XR

KEN PFEUFFER, HANS GELLERSEN, AND MAR GONZALEZ-FRANCO

16 Touch and Feel Virtual Objects

ADAMOS CHRISTOU AND RAVINDER DAHIYA

20 Toward Next-Generation Human-Computer Interface Based on Earables

YONGJIE YANG, TAO CHEN, AND LONGFEI SHANGGUAN

Edge Computing

24 Computer Education in the Machine Learning Era: Intelligent Systems

BRIAN M. SADLER

30 The Urban Space Information Platform: Opportunities and Challenges

JIABAO LI, RAJIV RANJAN, YUEWEI WANG, XIAOHUI HUANG, PHILIP JAMES, AND SCHAHRAM DUSTDAR

Internet of Things (IoT)

40 Jingxian Wang: "Pushing the Limits of Battery-Free Internet-of-Things"

LAKMAL MEEGAHAPOLA

44 Convenience at a Cost: The Urgent Need for Data Privacy Standards

SYED RIZVI, ANTHONY DEMERI, AND MOHAMMAD R. RIZVI

Economics

50 The Rise of Agentic AI in Finance: Opportunities, Risks, and Human-Centric Integration

NIR KSHETRI

56 Private Returns on Technology Adoption

SHANE GREENSTEIN

Departments

4 Magazine Roundup

7 Editor's Note: Sensory Human-Computer Interactions

62 Conference Calendar

Subscribe to *ComputingEdge* for free at
www.computer.org/computingedge

Magazine Roundup

The IEEE Computer Society's lineup of 12 peer-reviewed technical magazines covers cutting-edge topics ranging from software design and computer graphics to Internet computing and security, from scientific applications and machine intelligence to visualization and microchip design. Here are highlights from recent issues.

Computer

In-Network Collective

Operations: Game Changer or Challenge for AI Workloads?

This article, featured in the January 2026 issue of *Computer*, summarizes the opportunities of in-network collective operations for accelerated collective operations in artificial intelligence (AI) workloads. The authors provide sufficient detail to make this important field accessible to nonexperts in AI or networking, fostering a connection between these communities.

Computing

Designing Future Energy

Systems With Generative AI

Energy systems are experiencing various changes that impact the distribution, use, and reliability of energy. However, planning for and enacting advancements requires significant effort. Emerging generative artificial intelligence (AI) techniques can alleviate pain points and help support the development of the next generation of energy systems. The authors of this October–December 2025 *Computing in Science & Engineering* article highlight

ongoing generative AI work in the areas of atmospheric modeling, building energy management, and distribution network design, and they propose a vision for the role of generative AI that considers opportunities and identifies challenges inherent to this technology.

IEEE Annals

of the History of Computing

Machinery of Ethnic

Cleansing: Punched Card

Machines and the 1920 Greek Population Census

The 1920 Greek population census was conducted using punched card machinery in a period of successive wars and rapid territorial expansion. The authors of this October–December 2025 *IEEE Annals of the History of Computing* article examine this 1920 census, focusing on Macedonia, a newly annexed territory that would soon become the site of the first instance of large-scale ethnic cleansing in modern Europe. The authors approach the census and its machinery as a neglected part of the violent history of the Balkans and bring forward new perspectives on the actual use and potential of punched card machines.

IEEE Computer Graphics and Applications

Toward Agency in

Human–AI Collaboration

Artificial intelligence (AI) is increasingly evolving from a tool for automating repetitive tasks to an intelligent agent actively engaging in dynamic interactions with humans. As AI becomes more integrated into collaborative contexts, it is essential to examine the factors that shape human–AI interaction. Central to this collaboration is AI agency—the capacity for action and effect—a concept that has remained largely peripheral in existing research. In this article featured in the January/February 2026 issue of *IEEE Computer Graphics and Applications*, the authors address this gap by proposing a comprehensive design space for reasoning about agency in human–AI collaboration.

IEEE Intelligent Systems

Considering Sentiment Causes in In-Context Learning for Aspect-Based Sentiment Analysis

Aspect-based sentiment analysis (ABSA) aims to identify aspect terms in texts and determine their sentiment polarities. The in-context



learning paradigm, powered by large language models, has proven effective in low-resource scenarios, where the retrieval of effective demonstration examples is crucial. Existing retrieval methods prioritize semantic and syntactic similarities, overlooking the fact that sentiment is often driven by its underlying causes. Recognizing that similar causes tend to yield similar sentiments, the authors of this article, featured in the November/December 2025 issue of *IEEE Intelligent Systems*, propose the semantic-causal contextual demonstration retrieval (SCCDR), a demonstration retriever that integrates semantic and syntactic information while explicitly modeling sentiment causes.

Internet Computing

Smaller, Smarter, Closer: The Edge of Collaborative Generative Artificial Intelligence

The rapid adoption of generative artificial intelligence (GenAI), particularly large language models, has exposed critical limitations of cloud-centric deployments, including latency, cost, and privacy concerns. Meanwhile, small language models are emerging as viable alternatives for resource-constrained edge environments, although they often lack the capabilities of their larger counterparts. This article from the July/August 2025 issue of

IEEE Internet Computing explores the potential of collaborative inference systems that leverage both edge and cloud resources to address these challenges.

micro

Efficient Disaggregated Cloud Storage for Cold Videos With Neural Enhancement

The rapid growth of video-sharing platforms has driven immense storage demands, with disaggregated cloud storage emerging as a scalable and reliable solution. However, the proportional cost of cloud storage relative to capacity and duration limits the cost-efficiency for managing large-scale video data. This is particularly critical for cold videos, which constitute the majority of video data but are accessed infrequently. To address this challenge, this article from the November/December 2025 issue of *IEEE Micro* proposes neural cloud storage (NCS), leveraging content-aware super-resolution powered by deep neural networks.

MultiMedia

Scalable Neural Light Field With Layer Add-ons of Multilayer Perceptron

Light field (LF) is one of the 3-D image processing techniques that

provides a simple way to generate immersive content. Recently, neural LF (NLF), which incorporates the concept of implicit neural representation into LF, has been introduced, and the difficulties in LF reconstruction have been significantly reduced. Nevertheless, NLF still suffers from the limitations of reusability. To address this issue, the authors of this July–September 2025 *IEEE MultiMedia* article propose scalable NLF (S-NLF), which reconstructs LFs of various qualities via a single multilayer perceptron (MLP). The authors also propose base layer (BL) sharing to further improve the sample-level efficiency.

pervasive COMPUTING MOBILE SYSTEMS | UBIQUITOUS COMPUTING | INTERNET OF THINGS

Ethical Considerations of Extended Reality in the Workplace

As extended reality (XR) technologies are increasingly adopted in workplaces, ethical concerns, such as privacy, equity, and surveillance, must be addressed to ensure their responsible deployment. The authors of this article, featured in the October–December 2025 *IEEE Pervasive Computing* issue, surveyed 39 papers and contextualized them using the widely used National Institute of Standards and Technology artificial

intelligence ethical framework along the dimensions of privacy and surveillance, accountability and governance, transparency and explainability, fairness and inclusivity, and safety and psychological well-being.

IEEE SECURITY & PRIVACY

Differential Privacy in Practice: Lessons Learned From 10 Years of Real-World Applications

Differential privacy (DP) is a widespread data protection mechanism. However, its application in real-world scenarios has been challenging. To shed some light on this, the authors of this article,

featured in the January/February 2026 issue of *IEEE Security & Privacy*, offer a critical analysis of 21 DP deployments by top-tier companies and institutions over the past decade.

IEEE Software

Agentic LMs: Hunting Down Test Smells

In this article, featured in the January/February 2026 issue of *IEEE Software*, the authors explore how agentic language models (Llama-3.2-3B, Gemma-2-9B, DeepSeek-R1-14B, and Phi-4-14B) detect and refactor test smells. Multiagent workflows outperformed single-agent setups in 3 of 5 cases, with

Phi-4-14B achieving the best accuracy and merged open source contributions.

IT Professional

Generative AI in Government Policy Support: A Case Study of Taiwan STPI

This article, featured in the November/December 2025 issue of *IT Professional*, examines the integration of generative AI (GenAI) in government policy making, focusing on Taiwan’s Science & Technology Policy Research and Information Center (STPI). By adopting a boundary perspective, it explores challenges such as the literacy, knowledge, and ownership boundaries encountered during GenAI implementation. The research highlights how STPI leverages retrieval-augmented generation (RAG) technology to overcome these obstacles, enhancing the accuracy and efficiency of policy formulation. 🤖

IEEE DataPort™

STORE, SEARCH & MANAGE RESEARCH DATA

Individual subscriptions to IEEE DataPort are free for all IEEE society members and Young Professionals. Just log in and activate your subscription for unlimited access to datasets, data management tools, dataset storage for your own research, and more.

 Open Access Options	 2 TB of Cloud Storage	 Link to Manuscripts
 Generate Citations	 Reproducible Research	 ORCID Integration
 Host Data Competitions	 DOI Provided	

Join the IEEE Computer Society

computer.org/join



Editor's Note

Sensory Human–Computer Interactions

The field of human–computer interaction (HCI) has existed for decades, since the emergence of personal computing. With the rise of artificial intelligence (AI), HCI has advanced far beyond the level of tapping keyboards and swiping touchscreens. This issue of *ComputingEdge* explores developments in sensory interactive technologies and how they engage all five senses in virtual experiences. The articles also consider changes to edge computing posed by machine learning (ML) and data, as well as changes that could be made to Internet of Things (IoT) technologies to go battery-free and address security risks. The issue concludes with an assessment of how new technology may impact finance.

Researchers are pushing the boundaries of HCI by enhancing its sensory utilization. The authors of “Design Principles and Challenges for Gaze + Pinch Interaction in XR,” from *IEEE Computer Graphics and Applications*, present design principles and issues for the Gaze + Pinch interaction technique for Extended

Reality (XR) headsets, informed by eye-hand research in the HCI field. In “Touch and Feel Virtual Objects,” from *Computer*, the authors outline the current limitations of sensory interactive technologies, as well as future advancements. *IEEE Pervasive Computing* article “Toward Next-Generation Human–Computer Interface Based on Earables” introduces MAF (mobile acoustic field), a human–computer interface that enhances audio services by capturing a wide range of acoustic activities inside and outside of the device.

Engineers need to consider the impacts of ML and increased data when building computing systems. The author of *Computer* article “Computer Education in the Machine Learning Era: Intelligent Systems” expounds upon necessary changes to computer education due to the progression of ML, particularly in its application to computing. In “The Urban Space Information Platform: Opportunities and Challenges,” from *IEEE Internet Computing*, the authors review the USIP (which manages multisource

urban data), analyzing its current and future status as data volume and complexity increases.

Developers are pausing to consider necessary improvements to the IoT. In *IEEE Pervasive Computing* article, “Jingxian Wang: ‘Pushing the Limits of Battery-Free Internet-of-Things,’” Dr. Wang explains his research into a battery-free IoT. *Computer* article “Convenience at a Cost: The Urgent Need for Data Privacy Standards” exposes the huge risks posed by IoT networks’ lack of a universal privacy standard.

New technology has the potential to greatly impact financial services. The article “The Rise of Agentic AI in Finance: Opportunities, Risks, and Human-Centric Integration,” from *IT Professional*, examines the growing role of agentic AI in financial services, focusing on its capabilities and the challenges it poses. *IEEE Micro* article “Private Returns on Technology Adoption” assesses the financial return of consumer computer technologies (CCTs) by analyzing how firms utilize CCTs in their business. 🌍

DEPARTMENT: SPATIAL INTERFACES

Design Principles and Challenges for Gaze + Pinch Interaction in XR

Ken Pfeuffer , Aarhus University, 8000, Aarhus, Denmark

Hans Gellersen , Lancaster University, LA1 4YW, Lancaster, U.K.

Mar Gonzalez-Franco , Google, Seattle, WA, 98103, USA

For Extended Reality (XR) headsets, a key aim is the natural interaction in 3-D space beyond what traditional methods of keyboard, mouse, and touchscreen can offer. With the release of the Apple Vision Pro, a novel interaction paradigm is now widely available where users seamlessly navigate content through the combined use of their eyes and hands. However, blending these modalities poses unique design challenges due to their dynamic nature and the absence of established principles and standards. In this article, we present five design principles and issues for the Gaze + Pinch interaction technique, informed by eye-hand research in the human-computer interaction field. The design principles encompass mechanisms like division of labor and minimalistic timing, which are crucial for usability, alongside enhancements for the manipulation of objects, indirect interactions, and drag & drop. Whether in design, technology, or research domains, this exploration offers valuable perspectives for navigating the evolving landscape of 3-D interaction.

Interaction, innovative control methods, and natural user interfaces (UIs) have long been recognized as significant challenges in achieving a truly immersive and intuitive Extended Reality (XR) experience.^{2,4} However, the input landscape so far remained unsatisfactory, plagued by usability issues such as physical fatigue, ergonomic discomfort, and complex interface designs. The challenge of getting the interface right is amplified by the fragmented input landscape—controllers, hand tracking, eye movements, and voice commands. XR operating systems can courageously strive toward universal support, but it is tricky to unify all possible inputs and combinations in one system. This leaves interaction systems often primitive, often at the exclusion of one modality in favor of another without fully considering their unity.

It is challenging to achieve a harmonious integration of multiple modalities and optimize the effectiveness across various tasks. Yet, the multimodal input trend solidifies with the arrival of the “Gaze + Pinch” XR paradigm: *glance at a UI element with your eyes, then simply pinch with your fingers to activate it.* In the scientific community, researchers have been studying general eye-hand interaction in general since the 1980s⁵ and the specific “Gaze + Pinch” model since 2017.¹⁵ Yet, we see this as a clear innovation with a rapidly growing adoption across XR headsets such as the Microsoft HoloLens 2 or Magic Leap, and especially the OS-wide integration with the Apple Vision Pro.

Comprehensive guidance on multimodal interaction design is rare, with most focusing on one modality or generalizing across input devices.^{1,9} Integrating both raises questions: when to use eyes, hands, or both? How to merge the signals in time and space for optimal ease-of-use and expressiveness? In our article, we aim to highlight our findings, establish principles, and suggest frameworks based on scientific experiments to

guide designers, developers, and researchers in navigating this new interaction paradigm.

We present five design principles and five design issues, drawing insights from our eye-hand research and scientific articles in the area of human-computer interaction. Despite that in our daily lives as scientists we explore all possible future directions, our process of converging and abstraction eventually led to cover this set of principles. This abstraction effort covers the most pressing problem: the specific mechanisms of division of labor and multimodal timing that are key for usability; as well as issues of manipulation of objects with features such as indirect gesture and drag & drop.

We strive for this article to captivate not only the scientific community, but also to offer diverse perspectives that resonate with practitioners across various fields:

- › For designers: This article explores innovative approaches to Gaze + Pinch interaction, offering valuable inspiration.
- › For technologists: This article highlights the technical challenges and opportunities of Gaze + Pinch interaction.
- › For researchers: This article examines the latest research on Gaze + Pinch interaction, offering valuable insights at the intersection between XR, eye-tracking, and multimodal UI.

SUMMARY

The bulk of our work on this article has been to prepare the essence of design principles, which include:

- 1) *Division of labor*: Use a clear separation of tasks: the eyes perform selection tasks, the hands do the actual manipulation work.
- 2) *Minimalistic timing*: One moment matters for the eyes—the moment thumb and index finger contact, the hands take over, relieving the eyes from the explicit motor control tasks.
- 3) *Flexible gesture*: Gaze affords lightweight and flexible gesturing, allowing transitions from one versus two-hands to single versus multitarget control.
- 4) *Infallible eyes*: Eyes operate instantly, with constant accuracy. With good tracking, we cannot miss or overshoot a target when we look at something.
- 5) *Multimodal by design*: Gaze + Pinch complements direct gestures—understand, which tasks can be accomplished with one or another and provide transitions to get the best of both worlds.

We also discuss five behavioral design issues:

- 1) *(Un)Learning*: The Gaze + Pinch interaction challenges conventional action by enabling pointing without physical hand motion, emphasizing the need to unlearn habitual actions for efficiency.
- 2) *Early and late triggers*: The eye-hand timing is key, but ideally UIs are supportive when manual commands precede or lag behind gaze fixation, as errors can result.
- 3) *Input Mappings*: To be efficient across near and far spaces, consider control-display ratios and speed amplification like mouse acceleration to improve the control flexibility.
- 4) *Drag & drop sequences*: There are challenges in re-engaging with dropped objects—UIs can reduce dragging use and provide drop prevention.
- 5) *Continuous eye-input*: Exceptions like pinch-to-zoom showcase natural continuous eye inputs but require careful integration.

GAZE + PINCH

Gaze + Pinch users can directly manipulate objects they look at using familiar gestures like pinch-to-select or two-handed scaling, even when they are far from what would be considered a direct interaction space.

The term “Gaze + Pinch” originated from our 2017 paper,¹⁵ where we studied the foundations for this particular combo of eyes and hands. Gaze + Pinch stems from Pfeuffer’s prior doctoral thesis on the unity of gaze and multitouch gestures,¹⁶ superseding the earlier “Gaze-Touch” technique.¹²

As per Poupyrev et al.’s taxonomy,⁷ Gaze + Pinch is categorized as an egocentric input method, offering a first-person view and employing the virtual pointer metaphor with the eyes as the pointer. Nevertheless, its manipulation closely mirrors direct manipulation (Figure 1), resulting in the adoption of numerous traits from the familiar virtual hand technique. Key distinctions with other main input techniques are as follows:

- › *Gaze + Pinch versus Hand Gesture*: Gaze + Pinch allows users to interact with objects from a distance using the same gesture set, expanding the effective interaction area and maximizing the virtual environment’s vast space.
- › *Gaze + Pinch versus Controller*: Gaze + Pinch liberates users from holding physical devices, enabling them to perform hand gestures on distant objects as if directly manipulating them. This makes the UI highly intuitive, tapping into the inherent spatial manipulation skills of humans.

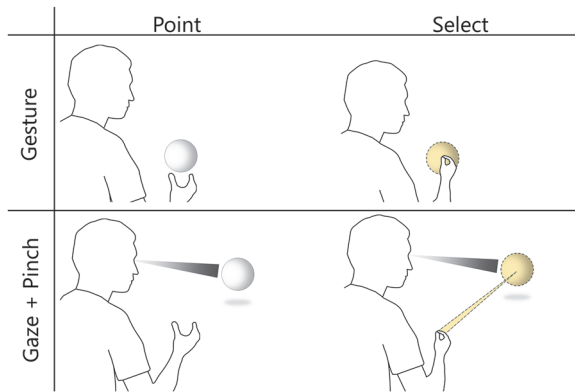


FIGURE 1. Basic operation of Gaze + Pinch in contrast to direct manipulation, demonstrating how similar (and, thus, easy to learn) the gestures are, as well as how the eyes extend reach to any object the user sees. Adapted from Pfeuffer et al.¹⁵

DESIGN PRINCIPLES

We go deeper into each of our key considerations that should guide the design and implementation of UIs based on Gaze + Pinch. In setting up these principles, we strive for a balance between simple yet powerful 3-D manipulation capabilities.

Division of Labor: The Eyes Select, the Hands Manipulate

Our eyes’ natural role involves indicating points of interest, and we can easily look at any point at will. In contrast, the hands are adept at physical manipulation through the interplay of finger movement and hand posture. Use a clear separation of concerns: the eyes perform selection tasks, the hands do the actual confirmation or manipulation work. This avoids the pitfalls of i) overloading the eyes with explicit motor control tasks²⁰—you only actively “use” the eyes to select, ii) physical fatigue²¹—gaze pointing minimizes the hands’ physical motion needs, and iii) supporting the (most) naturalistic roles for each modality.

The hands, then, make indirect gestures. This is similar to a controller in having the ability to interact at a distance, but now with intuitive pinch gestures. Indeed, there are hands-only techniques for selection and manipulation, such as the “handray” raypointing coupled with a pinch-confirmation (e.g., used by the Meta Quest 3 and HoloLens 2). Yet, assigning both selection and manipulation to the hand can be susceptible to hand jitter issues (the Heisenberg problem¹⁸). Our studies showed that Gaze + Pinch (and other eye-hand techniques) leads to improved

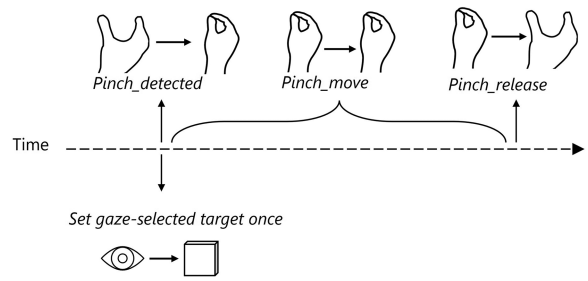


FIGURE 2. Eye-tracking as input is only active the moment that a pinch gesture is registered to avoid erratic behaviors when the eyes are wandering around.

performance and comfort for 3-D selection over gestures alone for interaction over distance.^{8,22}

Minimalistic Timing

There are many ways to mix and match the eye and hand tracking signals. A poorly designed multimodal input fusion can amplify complexity,¹⁰ especially with the eyes that can be wandering around and accidentally select things.²⁰ At the same time, it is important to take advantage of the eyes’ prime strength to offer instant selection.

The primary model for Gaze + Pinch only use a single moment in time for the eyes: the moment that the index finger and thumb have first contact, one has to fixate on the desired target (Figure 2). This instantaneous approach to the interaction is key—but only for the selection. For follow-up object manipulations, such as a drag, pan, or zoom gesture, the hands take over. This affords the freedom to inspect the surroundings independently and avoids accidental actions by eye or hand inputs. For instance in drag & drop, after selection one can freely look around to locate the destination for the dragged object and follow with the hand via indirect control.

In contrast, a hands-only UI typically means you can point with your hand to the target, without continuously monitoring the target. Gaze + Pinch inverts the relationship: the eyes must be on the target but the hands can be anywhere. This is a fundamentally new behavioral pattern that users got to master. Employing gaze minimally, and using the standard gesture set, facilitates a quick adaptation of this new relationship. Beginners may find themselves more attentive to ensure their gaze is on the UI element until receiving the right feedback. More experienced users may swiftly execute a Gaze + Pinch command even without having fully perceived the target and its selection feedback.

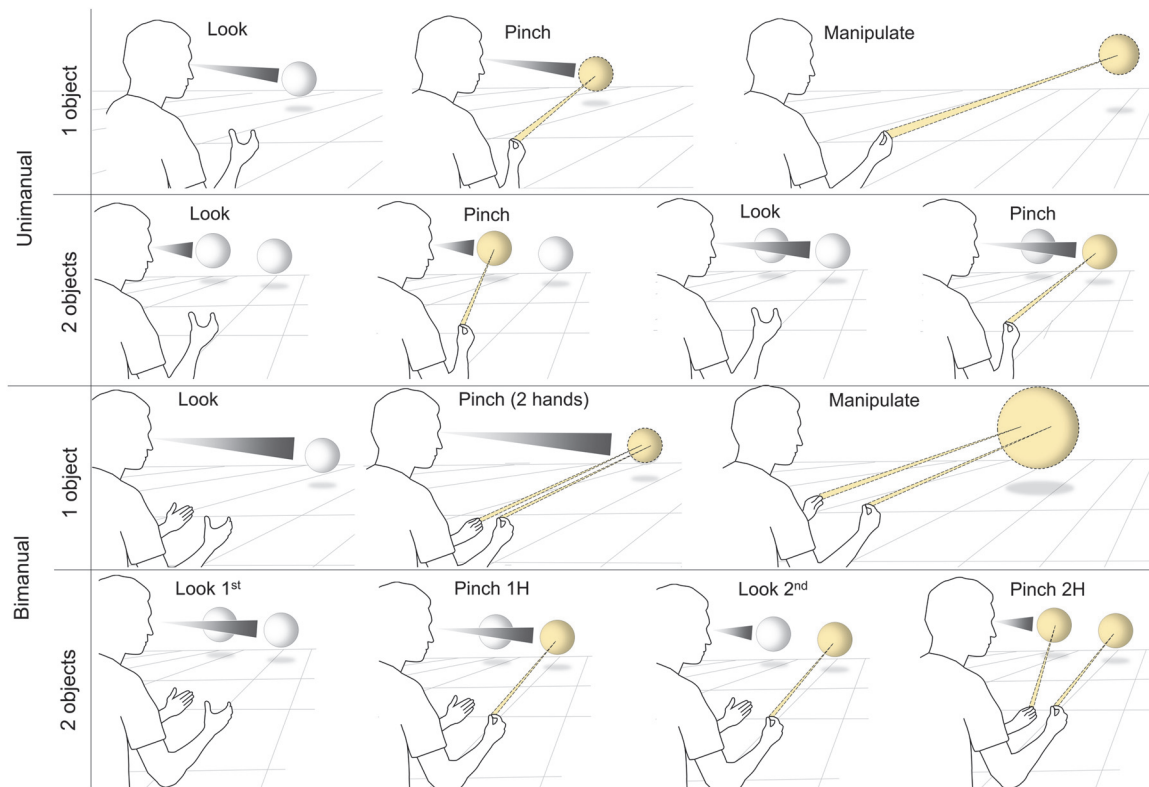


FIGURE 3. Fundamental interaction classes with Gaze + Pinch, natively supporting the atomic classes of one- and two-handed inputs, and single- and multiobject interaction. Taken from Pfeuffer et al.¹⁵

Flexible Gesture

Hand gestures are in control of virtual objects acquired by the eyes. Particularly, the commonly used selection, manipulation, and navigation commands can be covered when employing only pinch-tap and pinch-drag gestures. This flexibly extends to all atomic classes of the hand-based manipulation—one versus two-handed interactions, and single versus multiple target manipulation (Figure 3). Users can seamlessly shift between one or two objects and hands by simply glancing and re-engaging pinch gestures as desired. This is inherited from the default hand-tracking gestures—but the integration with gaze, and the elimination of the manual pointing subtask renders all those basic hand actions extremely lightweight and flexible to use across space.

Infallible Eyes

Our hands adhere to a speed-accuracy tradeoff: faster normally means less accurate when it comes to hands.¹⁹ In contrast, the eyes can fixate on a target in almost instant time, even if the target is in motion, with constant accuracy given by the eye-tracking

sensor. From a user's perspective, the eyes are infallible: hand pointing can miss or overshoot a target, there are even Olympic competitions on target shooting, but the eyes cannot miss as we are either on-target or we look elsewhere. It is crucial for an interaction system that engages eyes to support the simple way of just looking to select, without manual effort. That is of course if sensors and tracking were perfect. However, inaccurate eye-tracking can prompt users to undergo correction measures, such as squinting their eyes or adjust their head position in vain attempts to correct precision limitations, which in turn leads to increased mental exertion and longer selection times. It is perhaps that limitation that has hindered the popularization of Gaze + Pinch until now.

Some of the issues derive from intrinsic of human vision, such as eye dominance, or suboptimal eye strain when looking at targets above the horizon. We are better at looking down than up. In a way eyes are not as symmetric as people might intuitively think. And that means that even the most basic menus, such as home menus, might need to be reconsidered, perhaps they are better off at the bottom of the screen or even on a circular distribution.

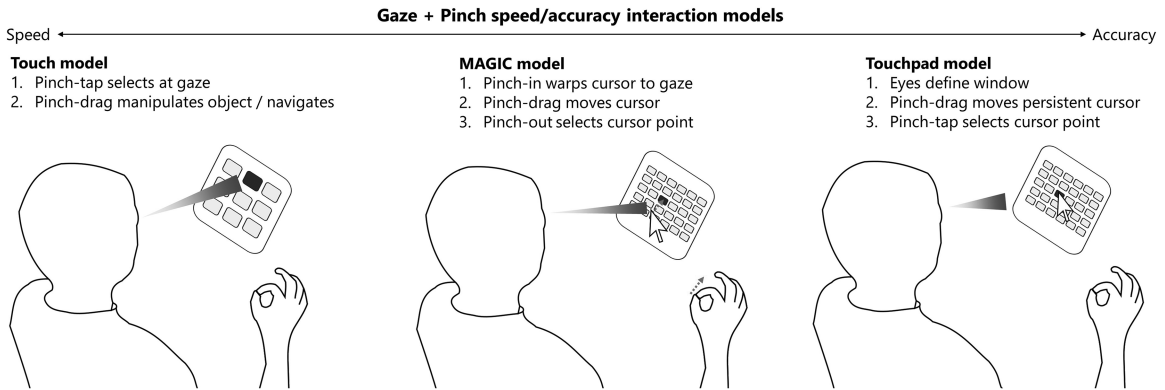


FIGURE 4. Interaction models for Gaze + Pinch, with different speed/accuracy trade-offs. Adopted from Pfeuffer et al.¹⁵

Some examples of design considerations for Gaze + Pinch UIs include somewhat counterintuitive aspects, such as large buttons to achieve a low error rate. In reality, large buttons invite wandering around with the eyes, potentially leading to outliers. Drawing the most salient parts to the center of the button will be welcomed by the selection mechanism, and a generous buffer space around targets makes outliers less impactful. That is, it is better to leave space between buttons than to have large buttons.

For compatibility reasons, not all UIs can be redesigned. If smaller targets are required, hand-refinement and cursor extensions can be integrated with Gaze + Pinch, although departing from the original simplicity. Figure 4 presents interaction models to accommodate both UI requirements and sensor limitations. These methods leverage the Gaze + Pinch signal in different ways to provide flexibility in interaction. The Touch model serves as the default, where a Gaze + Pinch command selects a target and a drag gesture executes manipulation. The MAGIC model (Manual and Gaze Input Cascaded²⁰) enhances precision by introducing a one-time mouse cursor: upon pinch-in, a cursor appears at the gaze position, which is then moved by gestural motion; upon pinch-out, the object under the cursor is selected. The Touchpad model ensures full precision with a persistent mouse in

the window. Here, eye movement is utilized solely for selecting the window, while pinch tap (click) and pinch-drag gestures control the mouse. These models can be selectively implemented by UI systems based on the precision requirements of the application context.

Multimodal by Design

The motto is to get the best of both worlds. Gaze + Pinch can also work together with the direct gestures in nearspace. This is possible via mode-switching methods that use time and space multiplexing of the inputs.¹³ This can have the neat side effect that one can rapidly use direct and indirect inputs at a glance (see Figure 5). For time-multiplexing, imagine a user opens an app with a menu through Gaze + Pinch, which activates the app UI in front of the user. The user switches to direct touch gestures to scroll the app's content. In space multiplexing, picture holding a menu with one hand while using Gaze + Pinch commands with the other, enabling direct and indirect inputs at the same time for on-hand virtual menus. Hand menus usually position menus off the hand to prevent hand-tracking interference. Indirect gestures are spatially separated, avoiding hand overlap (Figure 6). It is one of the intriguing outcomes when UI

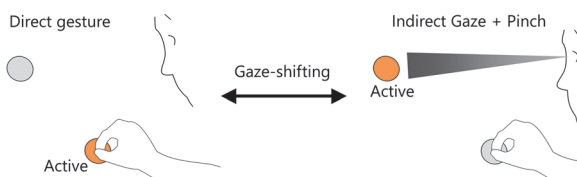


FIGURE 5. Best of both worlds: UIs can aim to support both gesture and Gaze + Pinch control through UI transitions. Adopted from Pfeuffer et al.¹³

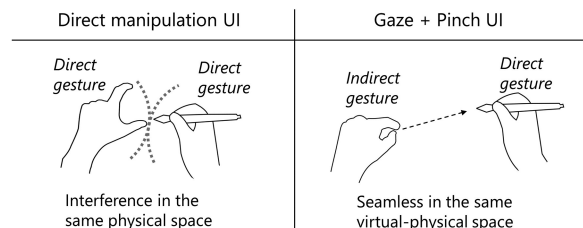


FIGURE 6. Mixing direct and indirect inputs allows for novel bimanual dynamics without physical interference.

systems support transitions between complementary modes of interaction.

BEHAVIORAL ISSUES

Learning a new vocabulary of interactions can feel as challenging as learning a new language from a cognitive perspective. This effect might be even stronger for folks who are experts in current paradigms. But the magic of Gaze + Pinch is that there is little learning to be done. Users employ familiar gestures—such as pinch for selection, translation, rotation, and scaling—reminiscent of direct manipulation, yet distinct in enabling interaction across a broader spatial range. This fusion of familiar and novel features, facilitated by eye gaze, defines the hybrid Gaze + Pinch technique. Perhaps, it is more about the (un)learning and the natural interactions it unlocks beyond controllers or hand gestures.

(Un)Learning

When we want to acquire a target, grab something, we intuitively move toward it. This is partially intuitive because of the physics of the real world, it might even be coded on a very deep layer on our brain. Gaze + Pinch commands do not need this movement anymore, which can be considered almost counterintuitive. Hence, it is important to balance the advantage of using a new way that initially requires rethinking of the action process, over just using the hands without the eyes. In a sense, Gaze + Pinch does not mean learning a new way of interaction—it is about unlearning the common way: do not move your hands, just confirm right where you are with your hand what your eyes have selected. Thus, the learning effort is rather negligible. But changing the nature of action might have consequences we have not fully looked into yet. Transitioning from XR experiences with Gaze + Pinch to those without introduces a perceptible discord. It is open how this big change might affect the behavior of a new generation who may grow up using this UI. What long-term impacts could this have? These questions remain open and necessitate careful consideration.

Early and Late Triggers

A problem of any multimodal interaction, and as such a problem that might arise in Gaze + Pinch, is that the eyes may leave the target before a manual command is registered, or the command is issued just before the gaze lands on the target.⁶ It is possible to have a predictive and generous timing (e.g., using the last fixation) (200–300 ms of a stable gaze), rather than the current gaze coordinate (see fixation detection

methods³). An error in this space of multimodal integration is defined in neuroscience as a body semantic violation.¹¹ If the early or late trigger frequency is known (e.g., through knowledge about user context and application), the timing can be adapted.

Input Mappings

Direct manipulation means a 1:1 control-display mapping from hand to object. With Gaze + Pinch, after selection the user's hand indirectly controls an object. Using a 1:1 mapping between physical hand motion and virtual object makes the interaction feel slow. What one can do is amplify the speed in the transfer function with the object distance. And this is possible because it is naturally a relative interaction paradigm, more akin to the mouse than any other form of natural input. This means, we can even use visual angle to determine dragging speed as a distance-independent metric: if your hand moves by 5° in your FoV, the remote object corresponds with 5° motion. This works well for objects at a distance but may be confusing when targets are near—here the UI can revert to a 1:1 transfer function.

Drag & Drop Sequences

In real life, dropping an object means, hopefully your hand is right there for you to pick it up again. Same happens with direct interactions in XR. But with Gaze + Pinch, you can look away after dropping and finding it again means you have to look back. This can be a hassle, especially if turning your body is involved. So, when designing UIs, it is crucial to think about what type of tasks are supported—ideally, most tasks require only a single action to finish drag & drop, and for sequences of manipulations consider potential enhancements. Hand tracking systems can make sure not to disengage the target from the control of a pinch gesture if just briefly undetected, to keep dragging robust and avoid object loss.

Continuous Eye-Selection

The minimal use principle for eye-based input is a principle rule, but there are cases where it can extend naturally to continuous eye inputs. For example, for zooming into a map. 1) The conservative default would be to set the zooming pivot to the gaze position at the initial pinch-in event, and then allow the hand position to adjust the pivot. This makes zooming like direct manipulation, where the physical input position remains at the virtual position on the underlying map. 2) An alternative model is to use the eyes continuously as input in parallel to the hand gesture. When the eyes

focus on a different area during zooming, the zooming target adjusts accordingly. The online repositioning can lead to more accurate zooming for goal-oriented navigation tasks, and in turn, reduce the need for panning gestures that correct zooming operations afterward.¹⁴ Doing this will put more responsibility on the eyes, which for goal-oriented zooming can feel natural. In addition, the choice between the two models can be informed by the task requirements, i.e., how much the eyes are expected to remain at the area of interest when performing a zoom gesture.

CONCLUSION

The eye-hand interaction design for 3-D experiences is a novel space that is gradually gaining momentum. Thanks to the groundwork laid by scientists and researchers, we are well equipped to explore this field further and anticipate exciting UX developments. Our focus on distilling the essentials of Gaze + Pinch interaction provides a deeper understanding of how to achieve the right balance to achieve a simple-to-use but expressive UI, drawing on basic design principles as well as practical considerations from experience.

Eye-tracking technology can transform how we use our hands, opening up new possibilities for XR interaction. Since the inception of the Gaze + Pinch concept, researchers have been relentlessly advancing the space of multimodal UIs, including our work on advanced selection,^{8,22} and one-handed inputs.¹⁷ We are excited to see how these ideas play out, and what else lies ahead to advance our interactive experience through a symbiosis of our eyes and hands. 😊

REFERENCES

1. Apple, Eyes. Developer Documentation. Accessed: Nov. 3, 24 2023. [Online]. Available: <https://developer.apple.com/design/human-interface-guidelines/eyes>
2. R. T. Azuma, "The most important challenge facing augmented reality," *Presence*, vol. 25, no. 3, pp. 234–238, 2016.
3. R. Andersson, L. Larsson, K. Holmqvist, M. Stridh, and M. Nyström, "One algorithm to rule them all? An evaluation and discussion of ten eye movement event-detection algorithms," *Behav. Res. Methods*, vol. 49, pp. 616–637, 2017.
4. M. Billinghurst, A. Clark, and G. Lee, "A survey of augmented reality," *Found. Trends Hum.-Comput. Interact.*, vol. 8, no. 2–3, pp. 73–272, 2015.
5. R. A. Bolt, "Gaze-orchestrated dynamic windows," in *Proc. 8th Annu. Conf. Comput. Graphics Interact. Techn.*, Assoc. Comput. Mach., New York, NY, USA, 1981, pp. 109–119.
6. M. Kumar, J. Klingner, R. Puranik, T. Winograd, and A. Paepcke, "Improving the accuracy of gaze input for interaction," in *Proc. Symp. Eye Tracking Res. Appl.*, 2008, pp. 65–68.
7. J. J. LaViola Jr, E. Kruijff, R. P. McMahan, D. Bowman, and I. P. Poupyrev, *3D User Interfaces: Theory and Practice*. Boston, MA, USA: Addison-Wesley Professional, 2017.
8. M. N. Lystbæk, P. Rosenberg, K. Pfeuffer, J. E. Grønbæk, and H. Gellersen, "Gaze-hand alignment: Combining eye gaze and mid-air pointing for interacting with menus in augmented reality," in *Proc. ACM Hum.-Comput. Interact.*, 2022, pp. 1–18.
9. Microsoft, "Eye-gaze-based interaction on HoloLens 2." Accessed: Nov. 3, 2023. [Online]. Available: <https://learn.microsoft.com/en-us/windows/mixed-reality/design/eye-gaze-interaction>
10. S. Oviatt, "Ten myths of multimodal interaction," *Commun. ACM*, vol. 42, pp. 74–81, Nov. 1999.
11. G. Padrao, M. Gonzalez-Franco, M. V. Sanchez-Vives, M. Slater, and A. Rodriguez-Fornells, "Violating body movement semantics: Neural signatures of self-generated and external-generated errors," *Neuroimage*, vol. 124, pp. 147–156, 2016.
12. K. Pfeuffer, J. Alexander, M. K. Chong, and H. Gellersen, "Gaze-touch: Combining gaze with multi-touch for interaction on the same surface," in *Proc. 27th Annu. ACM Symp. User Interface Softw. Technol.*, 2014, pp. 509–518.
13. K. Pfeuffer, J. Alexander, M. K. Chong, Y. Zhang, and H. Gellersen, "Gaze-shifting: Direct-indirect input with pen and touch modulated by gaze," in *Proc. 28th Annu. ACM Symp. User Interface Softw. Technol.*, 2015, pp. 373–383.
14. K. Pfeuffer, J. Alexander, and H. Gellersen, "Partially-indirect bimanual input with gaze, pen, and touch for pan, zoom, and ink interaction," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2016, pp. 2845–2856.
15. K. Pfeuffer, B. Mayer, D. Mardanbegi, and H. Gellersen, "Gaze pinch interaction in virtual reality," in *Proc. 5th Symp. Spatial User Interact.*, 2017, pp. 99–108.
16. K. Pfeuffer, "Extending touch with eye gaze input," Doctoral Thesis, Lancaster, LA, U.K.: Lancaster Univ., 2017.
17. K. Pfeuffer et al., "PalmGazer: Unimanual eye-hand menus in augmented reality," in *Proc. ACM Symp. Spatial User Interact.*, Assoc. Comput. Mach., 2023, vol. 10, pp. 1–12.
18. D. Wolf, J. Gugenheimer, M. Combosch, and E. Rukzio, "Understanding the Heisenberg effect of spatial interaction: A selection induced error for spatially tracked input devices," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2020, pp. 1–10.

19. S. Zhai, J. Kong, and X. Ren, "Speed-accuracy tradeoff in Fitts' law tasks—On the equivalency of actual and nominal pointing precision," *Int. J. Hum.-Comput. Stud.*, vol. 61, no. 6, pp. 823–856, 2004.
20. S. Zhai, C. Morimoto, and S. Ihde, "Manual and gaze input cascaded (MAGIC) pointing," in *Proc. SIGCHI Conf. Hum. Factors Comput. Syst.*, 1999, pp. 246–253.
21. J. D. Hincapié-Ramos, X. Guo, P. Moghadasian, and P. Irani, "Consumed endurance: A metric to quantify arm fatigue of mid-air interactions," in *Proc. SIGCHI Conf. Hum. Factors Comput. Syst.*, 2014, pp. 1063–1072.
22. U. Wagner, M. N. Lystbæk, P. Manakhov, J. E. Grønbæk, K. Pfeuffer, and H. Gellersen, "A Fitts' law study of gaze-hand alignment for selection in 3D user interfaces," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2023, pp. 1–15.

KEN PFEUFFER is an assistant professor with the Ubiquitous Computing and Interaction Group, Computer Science Department, Aarhus University, 8000, Aarhus, Denmark. He leads the XI Research Group on topics of human-computer

interaction, in particular AR/VR/XR, eye and hand interaction, UI design, and adaptive UI. He is a member of IEEE. He is the corresponding author of this article. Contact him at ken.pfeuffer@gmail.com.

HANS GELLERSEN is a professor of interactive systems with Lancaster University, LA1 4YW, Lancaster, U.K., and he is also with Aarhus University, Aarhus, Denmark. His recent research interests include eye-tracking, gaze for interaction and multimodal interaction techniques that leverage eye movement in concert with other modalities. Gellersen received his Ph.D. degree in computer science from the University of Karlsruhe, Germany. Contact him at hwg@comp.lancs.ac.uk.

MAR GONZALEZ-FRANCO is a research manager with Google AR & VR, Seattle, WA, 98103, USA, where she works on creating a new generation of Immersive Tech. Contact her at margonzalezfranco@gmail.com.

Contact department editor Kyle Johnsen at kjohnsen@uga.edu or department editor Mark Billingham at mark.billinghurst@auckland.ac.nz or department editor Michele Fiorentino at michele.fiorentino@poliba.it.

ADVERTISER INFORMATION

Advertising Coordinator

Debbie Sims
 Email: dsims@computer.org
 Phone: +1 714-816-2138 | Fax: +1 714-821-4010

Advertising Sales Contacts

Mid-Atlantic US, Northeast, Europe, the Middle East and Africa:
 Dawn Scoda
 Email: dscoda@computer.org
 Phone: +1 732-772-0160
 Cell: +1 732-685-6068 | Fax: +1 732-772-0164

Southwest US, California:
 Mike Hughes
 Email: mikehughes@computer.org
 Cell: +1 805-208-5882

Central US, Northwest US, Southeast US, Asia/Pacific:
 Eric Kincaid
 Email: e.kincaid@computer.org
 Phone: +1 214-553-8513 | Fax: +1 888-886-8599
 Cell: +1 214-673-3742

Midwest US:
 Dave Jones
 Email: djones@computer.org
 Phone: +1 708-442-5633 | Fax: +1 888-886-8599
 Cell: +1 708-624-9901

Jobs Board (West Coast and Asia), Classified Line Ads

Heather Buonadies
 Email: hbuonadies@computer.org
 Phone: +1 623-233-6575

Jobs Board (East Coast and Europe), SE Radio Podcast

Marie Thompson
 Email: marie.thompson@computer.org
 Phone: +1 714-813-5094

DEPARTMENT: PREDICTIONS

Touch and Feel Virtual Objects

Adamos Christou  and Ravinder Dahiya , *Northeastern University*

Artificial and virtual technologies are no longer confined to what we can see or hear. Today's immersive systems are beginning to engage all five senses—introducing touch, scent, and even taste into virtual experiences.

The realm of artificial reality (AR) and virtual reality (VR) has seen remarkable advancements over the past decade. These technologies have transformed from niche applications to mainstream tools used in various sectors. Sensory interactive technologies, which enhance immersive experience by incorporating touch, smell, and other senses, are at the forefront of this transformation. This article explores the current limitations of these technologies, predicts future advancements, and discusses the potential impact, challenges, and risks associated with these predictions.

LIMITATIONS OF EXISTING TECHNOLOGIES

The AR and VR technologies have seen significant progress, despite several limitations that hinder their widespread adoption. Traditional VR headsets are cumbersome, bulky and uncomfortable to wear for long periods of time. They are also not the most fashionable gadgets. Current setups often require a controlled environment and are tethered to powerful computers or rely on standalone devices with limited processing power, restricting mobility and accessibility. Most systems rely on wearable devices like gloves or controllers for haptic feedback, which often fail to deliver a natural experience. The sensation of touch is limited to the areas covered by the devices, and the feedback is usually simplistic, lacking the nuances of real-world interactions. Visual and auditory realism also falls short due to issues like screen door effect,

limited field of view, and latency. While spatial audio has improved, it still does not fully replicate natural soundscapes. Integrating senses beyond sight and sound, such as touch, smell, and taste, remains a significant challenge toward achieving a more immersive experience. Technologists are thus racing to produce viable and affordable alternatives, such as smart glasses. Meta's collaboration with Ray-Ban is perhaps one well-known example.

FUTURE TECHNOLOGIES AND THEIR IMPACT

The future of sensory interactive technologies in AR and VR promises to address current limitations and unlock new possibilities. These are likely to revolve around six main areas:

Advanced haptic feedback systems

Future haptic feedback systems will move beyond wearable devices to incorporate more sophisticated and natural interactions. Several technologies have already provided the glimpse of such implementations, and present good cases for reinvigorating the industry. "Aerohaptics," for example, is a technology which uses controlled jets of air to simulate the sensation of touch (pressure, temperature), allowing users to interact with virtual objects in a more natural and immersive way.¹ In fact, air through controlled jets could also be mixed with artificial fragrance to provide the sense of smell. By delivering mid-air tactile feedback, "aerohaptics" enables users to "touch" virtual objects without the need for gloves or other wearable haptic devices (Figure 1). The technology pairs with 3D or volumetric display systems to create interactive "pseudo-holograms."² Developed using commercially



available components, aerohaptics technology is affordable, and easier to implement. With these attractive features, this futuristic technology marks a significant advancement in haptic sensing for virtual environments.

Another example is ultrahaptics, an approach that uses ultrasound waves to create tactile sensations in mid-air.³ Ultrahaptics enables users to feel virtual objects through focused ultrasound waves that create high-pressure points on the skin. These points can simulate various textures and shapes, providing a realistic sense of touch. However, the pressure from ultrasound waves is much lower than with aerohaptics. Further, the sensation of temperature and smell is not possible with ultrasound waves. Nonetheless, such natural interaction methods still significantly enhance the user experience by making it more intuitive and user-friendly.

Full-spectrum sensory integration

Advancements in sensory technologies will enable the integration of touch, smell, and even taste into AR and VR systems, creating more engaging and convincing experiences. Researchers have already demonstrated a novel wireless olfactory feedback system, designed to incorporate the sense of smell in such scenarios.⁴ This system uses arrays of flexible and miniaturized odor generators that can release various scents on demand. The device, which can be mounted on the upper lip or integrated into a flexible face mask, ensures that scents are delivered close to the user's nose for an ultra-fast response. The odor generators use a subtle heating platform and a mechanical thermal actuator to heat and melt odorous paraffin wax, allowing for programmable odor release with precise control over concentration and combination. The system is wireless, lightweight, and made of soft materials, making it comfortable and unobtrusive for users. This technology significantly enhances the realism and immersion

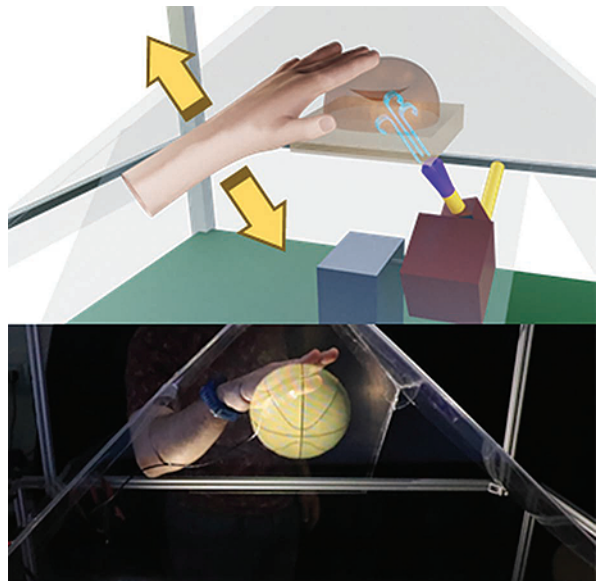


FIGURE 1. The “aerohaptics” haptic feedback system uses directed jets of air to make users feel they are touching a virtual object.

of virtual environments and can find applications beyond entertainment, with potential uses in medical treatment and human emotion control.

THE DEVELOPMENT OF LIGHTWEIGHT, WIRELESS AR AND VR DEVICES WITH POWERFUL PROCESSING CAPABILITIES WILL ENHANCE MOBILITY AND ACCESSIBILITY.

Improved mobility and accessibility

The development of lightweight, wireless AR and VR devices with powerful processing capabilities will enhance mobility and accessibility. These devices will allow users to interact with virtual environments seamlessly in dynamic real-world settings, expanding the applications in remote work, education, and health care.

Enhanced visual and auditory realism

Future AR and VR systems will feature higher resolution displays, wider fields of view, and reduced latency. Advances in spatial audio technology will further enhance the auditory experience, making it indistinguishable from real-world sounds.

Applications in various sectors

The impact of these advancements will be felt across multiple sectors. In health care, for example, clinicians will be able to collaborate on treatments using realistic virtual models, improving patient outcomes. In entertainment, engaging and interactive experiences will transform the way we consume media.

Advanced haptic technologies and full-spectrum sensory integration in AR and VR have the potential to revolutionize education by creating highly immersive and interactive learning environments. These technologies transform abstract concepts into tangible experiences, making learning more engaging and effective. Haptic feedback allows students to physically interact with virtual objects, enhancing their understanding of

ADVANCED HAPTIC TECHNOLOGIES AND FULL-SPECTRUM SENSORY INTEGRATION IN AR AND VR HAVE THE POTENTIAL TO REVOLUTIONIZE EDUCATION.

complex subjects. For example, in Science, Technology, Engineering, and Mathematics education, students can manipulate 3D models of molecules, explore human anatomy, or conduct virtual experiments, providing a hands-on learning experience often not possible in traditional classrooms.⁵ Full-spectrum sensory integration goes beyond visual and auditory stimuli by incorporating touch, smell, and taste into the learning experience, significantly improving memory retention and comprehension.⁶ These technologies also promote inclusivity by catering to diverse learning styles and needs, benefiting students with disabilities through customized sensory inputs. Additionally, haptic and sensory technologies facilitate remote learning, allowing students from different locations to participate in virtual classrooms and collaborate on projects in real time.

Artificial intelligence integration

The convergence of multisensory AR and VR systems with artificial intelligence (AI) opens new frontiers for intelligent, context-aware virtual experiences. AI agents can leverage sensory-rich environments to interpret user behavior, adapt interactions in real time, and provide personalized assistance across different domains. For instance, in an educational setting, an AI-based tutor in a virtual classroom could respond to tactile feedback and emotional cues, dynamically adjusting the lesson.

TECHNOLOGY CHALLENGES

While the future of sensory interactive technologies in AR and VR is promising, several challenges must be addressed to realize the potential they hold. Developing systems that can accurately simulate touch, smell, and other senses is technically complex. Ensuring that these systems are reliable and consistent across different environments adds another layer of complexity. Leveraging AI to enhance the interaction experience requires robust machine learning models for multimodal understanding, and low-latency processing capabilities. High costs associated with advanced AR and VR technologies can limit their accessibility. Developing affordable solutions that can be widely adopted is crucial for the widespread use of these technologies.

Creating a seamless and intuitive user experience is essential for the success of sensory interactive technologies. Users should be able to interact with virtual environments naturally, without the need for extensive training or complex setups. Ensuring that the technology is user-friendly and accessible to people with varying levels of technical expertise is a significant challenge. Lastly, integrating sensory interactive technologies with existing AR and VR systems and applications can be challenging. It requires compatibility with different hardware and software platforms as well as the ability to work seamlessly with other technologies, such as Internet of Things.

RISKS TO PREDICTION

Predicting the future of sensory interactive technologies in AR and VR involves several risks and uncertainties. The pace of technological advancements can be unpredictable. Breakthroughs in one area may lead to rapid progress, while unforeseen challenges in

another area may slow down development. This uncertainty makes it difficult to accurately predict the timeline and impact of future interactive technologies. Even if advanced sensory technologies are developed, their success depends on market adoption. Factors, such as consumer preferences, industry standards, and regulatory frameworks, can influence the adoption and widespread use of these technologies.

The integration of sensory interactive technologies into AR and VR raises ethical and societal concerns too. For example, a virtual avatar of an individual could be created and misused. Issues, such as privacy, data security, and the potential for misuse, must be addressed. Ensuring that these technologies are developed and used responsibly is crucial for their long-term success. The use of advanced AR and VR technologies can have economic and environmental implications. Ensuring that these technologies are sustainable and do not contribute to environmental degradation is essential. Predicting and mitigating these impacts is a significant challenge.

Technologies allowing natural interaction with virtual objects hold immense potential for the future. Advancements in haptic feedback, sensory integration, mobility, and visual and auditory realism will transform the way we interact with virtual environments. This potential could be realized by addressing technical, economic, and societal challenges. By understanding and navigating these challenges, we can unlock the full potential of AR and VR, creating more immersive and impactful experiences in the future. 🌐

REFERENCES

1. A. Christou, R. Chirila, and R. Dahiya, "Pseudo-hologram with aerohaptic feedback for interactive volumetric displays," *Adv. Intell. Syst.*, vol. 4, no. 2, 2022, Art. no. 2100090.
2. Christou, A. Gao, Y. Navaraj, W. T. Nassar, and H. Dahiya, "3D touch surface for interactive pseudo-holographic displays," *Adv. Intell. Syst.*, vol. 4, no. 2, 2022, Art. no. 2000126.
3. T. Carter, S. A. Seah, B. Long, B. Drinkwater, and S. Subramanian, "UltraHaptics: Multi-point mid-air haptic feedback for touch surfaces," in *Proc. 26th Annu. ACM Symp. User Interface Softw. Technol.*, 2013, pp. 505–514.
4. Y. Liu et al., "Soft, miniaturized, wireless olfactory interface for virtual reality," *Nature Commun.*, vol. 14, no. 1, 2023, Art. no. 2297, doi: 10.1038/s41467-023-37678-4.
5. F. Sanfilippo et al., "A perspective review on integrating VR/AR with haptics into STEM education for multi-sensory learning," *Robotics*, vol. 11, no. 2, 2022, Art. no. 41, doi: 10.3390/robotics11020041.
6. W. Ying, "Application of "virtual reality + haptic feedback" in education: Opportunities and challenges," *Frontiers Educational Res.*, vol. 7, no. 8, pp. 33–37, 2024.

ADAMOS CHRISTOU is a postdoctoral researcher at Northeastern University, Boston, MA 02115 USA. Contact him at a.christou@northeastern.edu.

RAVINDER DAHIYA is a professor in the Electrical and Computer Engineering Department at Northeastern University, Boston, MA 02115 USA. Contact him at r.dahiya@northeastern.edu.

The banner features the IEEE Computer Society logo on the left and the 80th Anniversary Celebration logo on the right. Below the logos, the text "FOLLOW US" is centered. Underneath, there are two QR codes: one for LinkedIn (IEEE Computer Society) and one for YouTube (@IEEEComputerSociety). The background is a dark blue space with glowing circuit lines and particles.

DEPARTMENT: RESEARCH BRIEF

Toward Next-Generation Human–Computer Interface Based on Earables

Yongjie Yang , Tao Chen , and Longfei Shangguan , *University of Pittsburgh, Pittsburgh, PA, 15213, USA*

We introduce a new type of human–computer interface—mobile acoustic field (MAF). MAF explores a variety of onboard sensors, such as speaker transducers, inertial measurement unit (IMU), and feedforward and feedback microphones, on earphones/earbuds to capture a wide range of acoustic activities both inside-out and outside-in. We anticipate that the exponential advancement of generative AI, powered by the sheer size of cloud computing resources, will enable personal agents to sense, understand, and further interact with this MAF, offering many exciting mobile services to users. In this article, we summarize our preliminary results, discuss the downstream applications, technical challenges, and our vision on this exciting research field.

MOBILE ACOUSTIC FIELD (MAF)-I: THE POWER AS POWER

Earphones are increasingly being embraced globally for their versatility and convenience. Beyond listening to music and communicating, we found that acoustic waves from bone conduction earphone speakers will transfer the energy into human tissues, producing *surface acoustic waves (SAW)* that propagate along the surface of the human head. Moreover, these SAW will also dissipate into the air as they propagate through the user’s face, producing another signal called *leaky surface acoustic waves (LSAW)*. The combination of these two signals effectively creates an acoustic field surrounding the user’s head—this leads to our initial thought of MAF, *a.k.a., MAF-I*.

Outside-In: User gestures performed on or in the vicinity of the face can perturb the channel of SAW or LSAW signal. By observing how the received SAW and LSAW signals change over time, it is possible to detect and further distinguish these gestures. With this MAF, the human head essentially becomes an independent space for interaction. For instance, it can enable mobile users to define personalized gestures for controlling volume, playback, and muting without the

The “Research Brief” section in *IEEE Pervasive Computing* aims to rethink science communication through ultra-compact papers. In this article, Prof. Shyamnath Gollakota discusses cutting-edge semantic hearing devices drawing on examples from his research group’s work.

From the Editor

need for dedicated sensors. Likewise, by accurately detecting and interpreting the gestures of the user’s hands, the virtual reality (VR) applications allow users to manipulate and control virtual objects and perform various actions, without the need for physical controllers. Users can communicate nonverbally through their over-the-face hand gestures, fostering a more engaging VR gaming experience. Preliminary results can be found in MAF published at CHI’24.³

Inside-Out: Motivated by previous earable-based vital sign (e.g., heart rate) monitoring systems,² we realized that the MAF is not limited to acoustic events that happen outside the human body, the natural contact of the earphone speaker transducer and the wearer’s ear also creates a unique opportunity to capture the fine-grained acoustic event inside the human body from the ear canal. We term this inside-out MAF. In *MobiCom’24*,¹ we demonstrated such feasibility by designing a hardware–software system that turns the wearer’s earphones

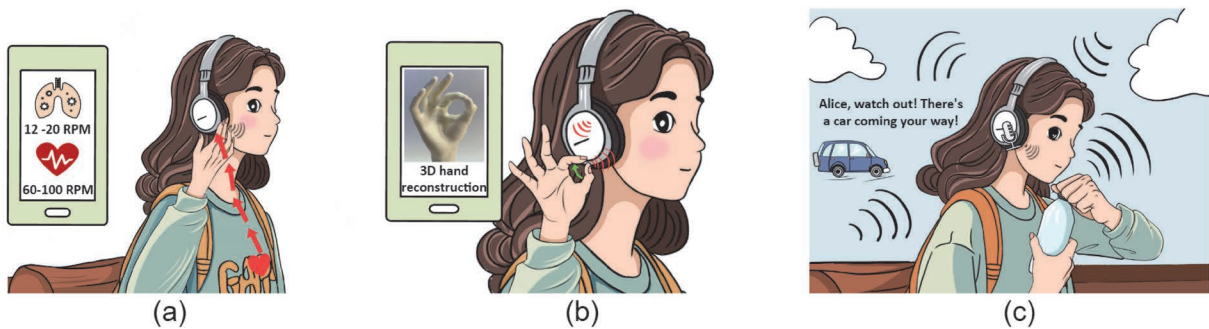


FIGURE 1. (a) MAF-I enables both physiological activity sensing and human gesture recognition. (b) MAF-II leverages leaky waves to reconstruct user 3D gesture for fine-grained gesture control. (c) MAF-III empowers mobile users to interact with sensor-in-the-loop AI agent for just-in-time service.

into a “stethoscope,” allowing the specialist to hear the earphone wearer’s heart sound (i.e., phonocardiogram (PCG) signals) in video clinic visits.

MAF-II: 3D ACOUSTIC FIELD

Many individuals have likely encountered a situation where, at a relatively close distance, they can hear the substantial sound leakage from a nearby person’s earphones, particularly low-end earphones. This phenomenon often arises because these earphones typically lack sound insulation materials, resulting in sound leakage. Even with high-end earphones, our preliminary studies with Sony WF-1000XMS earphones and Google Pixel buds show that acoustic signal leakage still exists.

These leaky signals push us to rethink the MAF. On the one hand, leaky signals reflect off objectives (e.g., human hands or environment) can be received by the feedforward microphones positioned at the back of the earphones. The motions of reflectors will change the length of the reflection paths, resulting in different power delay profiles. This is essentially the SONAR principle, which gives us a fine-grained “field of view” in the MAF. On the other hand, some reflectors, such as gestures (e.g., rubbing fingers, squeezing, and pinching), usually produce unique sounds that can also be captured by the feedforward microphones on earbuds. The combination of active probing and passive listening leaves us a question: *can we build a more generic MAF—one that can capture the earphone wearer’s working environment and behaviors in a unified 3D representation?* We term this *MAF-II*.

The key algorithmic questions are: how to embed signal and their associated spatial information into a unified 3D acoustic field? Are techniques, such as neural radiance field and Gaussian splatting, adequate with minor modifications? Or are clean-slate algorithms necessary? Given that

acoustic signal attenuates severely over distance, how could we expand the “field of view”? Can each earphone wearer exchange their partial or full view with each other to see unseens? We are currently working on these problems by embedding the passive and active acoustic sensing data into a unified feature space and build a transformer-based neural network to reconstruct the fine-grained 3D scenes.

MAF-III: HUMAN-IN-THE-LOOP

When the 3D acoustic field is available, we envision a mobile agent powered by large language models that can monitor this 3D acoustic field, identifying changes in the environment, irregular behaviors, and abnormal physiological activities for safety awareness. Many research questions remain unsolved. For instance, given the success of large language models, how can we adapt these foundation models to our specific scenario? How can we obtain sufficient, high-quality training data for model fine-tuning? Can we create useful synthetic data in the VR space to fine-tune the model?

After constructing the 3D acoustic field and enabling the cyber-agent to comprehend it, our next step is to develop an effective user interface that enables the earphone wearer in the physical world to interact with the cyber-agent. We envision some interactions can be in the form of, for example, “Hi, dude, keep an eye on approaching cars behind me.” The cyber-agent will then pay attention to the rear part of the acoustic field, detecting the presence of approaching cars timely. We term this human-in-the-loop MAF as *MAF-III*. The key algorithmic questions are: once the targeting event is detected, what is the best way to notify the earphone wearer? Would voice reminders plus vibrations be enough? If not, could spatial sound filters, e.g., playing only the car approaching sound event with location embedding, be a good way? Similar to the bounding box

being used in computer vision for visual guidance, can we design an *acoustic bounding box* for acoustic guidance? If yes, what could be the best way of designing an acoustic bounding box? All these problems are crucial to the success of the mobile acoustic field. Figure 1 illustrates the development of three generations of MAF systems. 🌐

ACKNOWLEDGMENTS

This work was supported in part by the National Science Foundation under Grant 2337537 and Grant 2302724 and in part by a Google Research Scholar Award.

REFERENCES

1. T. Chen, Y. Yang, X. Fan, X. Guo, J. Xiong, and L. Shangguan, "Exploring the feasibility of remote cardiac auscultation using earphones," in *Proc. 30th Annu. Int. Conf. Mobile Comput. Netw.*, 2024, pp. 357–372.
2. X. Fan et al., "HeadFi: Bringing intelligence to all headphones," in *Proc. 27th Annu. Int. Conf. Mobile Comput. Netw.*, 2021, pp. 147–159.
3. Y. Yang, T. Chen, Y. Huang, X. Guo, and L. Shangguan, "MAF: Exploring mobile acoustic field for hand-to-face gesture interactions," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2024, pp. 1–20.

YONGJIE YANG is currently working toward the Ph.D. degree with the Department of Computer Science, University of Pittsburgh, Pittsburgh, PA, 15213, USA. His research interests include IoT, acoustic sensing, and mobile computing. Yang received his master's degree in electronics and electrical engineering from Boston University, Boston, MA, USA. Contact him at yoy28@pitt.edu.

TAO CHEN is a postdoctoral associate with the Department of Computer Science, University of Pittsburgh, Pittsburgh, PA, 15213, USA. His research interests include sensory AI, acoustics, and mobile computing. Chen received his Ph.D. degree from the City University of Hong Kong, Hong Kong. Contact him at tac194@pitt.edu.

LONGFEI SHANGUAN is an assistant professor with the Department of Computer Science, University of Pittsburgh, Pittsburgh, PA, 15213, USA. His research interests include IoT, wireless systems, and mobile computing. Shangguan received his Ph.D. degree from the Hong Kong University of Science and Technology, Hong Kong. He is the corresponding author of this article. Contact him at longfei@pitt.edu.



IEEE COMPUTER SOCIETY CALL FOR PAPERS

EXTEND YOUR INFLUENCE

Engage a global network of 350k+ computing professionals through peer-reviewed journals and magazines with the highest impact factors in the industry.


GET PUBLISHED
computer.org/cfp


IEEE
COMPUTER
SOCIETY


IEEE



www.computer.org



PURPOSE: Engaging professionals from all areas of computing, the IEEE Computer Society sets the standard for education and engagement that fuels global technological advancement. Through conferences, publications, and programs, IEEE CS empowers, guides, and shapes the future of its members, and the greater industry, enabling new opportunities to better serve our world.

OMBUDSMAN: Contact ombudsman@computer.org.

CHAPTERS: Regular and student chapters worldwide provide the opportunity to interact with colleagues, hear technical experts, and serve the local professional community.

PUBLICATIONS AND ACTIVITIES

Computer: The flagship publication of the IEEE Computer Society, *Computer*, publishes peer-reviewed technical content that covers all aspects of computer science, computer engineering, technology, and applications.

Periodicals: The IEEE CS publishes 12 magazines, 18 journals

Conference Proceedings & Books: Conference Publishing Services publishes more than 275 titles every year.

Standards Working Groups: More than 150 groups produce IEEE standards used throughout the world.

Technical Communities: TCs provide professional interaction in more than 30 technical areas and directly influence computer engineering conferences and publications.

Conferences/Education: The IEEE CS holds more than 215 conferences each year and sponsors many educational activities, including computing science accreditation.

Certifications: The IEEE CS offers three software developer credentials.

AVAILABLE INFORMATION

To check membership status, report an address change, or obtain information, contact help@computer.org.

IEEE COMPUTER SOCIETY OFFICES

WASHINGTON, D.C.:

2001 L St., Ste. 700,
Washington, D.C. 20036-4928

Phone: +1 202 371 0101

Fax: +1 202 728 9614

Email: help@computer.org

LOS ALAMITOS:

10662 Los Vaqueros Cir.,
Los Alamitos, CA 90720

Phone: +1 714 821 8380

Email: help@computer.org

IEEE CS EXECUTIVE STAFF

Interim Executive Director: Anne Marie Kelly

Director, Governance & Associate Executive Director:
Anne Marie Kelly

Director, Conference Operations: Silvia Ceballos

Director, Information Technology & Services: Sumit Kacker

Director, Marketing & Sales: Michelle Tubb

Director, Membership Development: Eric Berkowitz

Director, Publications & Special Projects: Robin Baldwin

IEEE CS EXECUTIVE COMMITTEE

President: Grace A. Lewis

President-Elect: Joaquim Jorge

Past President: Hironori Washizaki

Vice President: Terry Benzel

Secretary: Yuhong Liu

Treasurer: Fernando Bouche

VP, Member & Geographic Activities: Andrew Seely

VP, Professional & Educational Activities: Cyril Onwubiko

VP, Publications: Charles (Chuck) Hansen

VP, Standards Activities: Darren Galpin

VP, Technical & Conference Activities: Andreas Reinhardt

2025–2026 IEEE Division VIII Director: Cecilia Metra

2026 IEEE Division VIII Director-Elect: Nita Patel

2026–2027 IEEE Division V Director: Leila De Floriani

IEEE CS BOARD OF GOVERNORS

Term Expiring 2026:

Megha Ben, Mrinal Karvir, Sean Peisert, Andreas Reinhardt, Deborah Silver, Yoshiko Yasuda

Term Expiring 2027:

Sven Dickinson, Alfredo Goldman, Daniel S. Katz, Yuhong Liu, Ladan Tahvildari, Damla Turgut

Term Expiring 2028:

Min Chen, Sap Ghosh, Lizy John, Rick Kazman, Carolyn McGregor, Walter Scheirer

IEEE EXECUTIVE STAFF

Executive Director and COO: Sophia Muirhead

General Counsel and Chief Compliance Officer:
Ahsaki Benion

Chief Human Resources Officer: Cheri N. Collins Wideman

Managing Director, Technical Activities: Mojdeh Bahar

Managing Director, IEEE-USA: Russell Harrison

Chief Marketing Officer: Jayne O'Brien

Chief Publication Officer and Managing Director:
Steven Heffner

Chief Governance Officer: Donna Hourican

Managing Director, Member and Geographic Activities:
Cecelia Jankowski

Chief of Staff to the Executive Director: Allison Taylor

Managing Director, Educational Activities: Jamie Moesch

IEEE Standards Association Managing Director: Alpesh Shah

Chief Financial Officer: Kelly Armstrong

Chief Information Digital Officer: Jeff Strohschein

Managing Director, Conferences, Events, and Experiences:
Marie Hunter

IEEE OFFICERS

President & CEO: Mary Ellen Randall

President-Elect: Jill I. Gostin

Past President: Kathleen A. Kramer

Director & Secretary: Jenifer P. Castillo Rodriguez

Director & Treasurer: Gerardo Barbosa

Director & VP, Publication Services & Products: W. Clem Karl

Director & VP, Educational Activities: Timothy P. Kurzweg

Director & VP, Membership and Geographic Activities:
Andrew D. Lowery

Director & President, Standards Association:
Gary R. Hoffman

Director & VP, Technical Activities: F D Tan

Director & President, IEEE-USA: Barry C. Tilton



DEPARTMENT: EDUCATION

Computer Education in the Machine Learning Era: Intelligent Systems

Brian M. Sadler , University of Texas at Austin

In the machine learning (ML) Era computing will become synonymous with intelligence, blending classical and ML-based methods. This points to fundamental changes in computer education.

As profound as was the progression from analog to digital computing, so is the current progression into the era of machine learning (ML). The early ML Era is marked by two key trends, a technology progression to systems, and the use of data-driven compute architectures.

Together, these are enabling all kinds of intelligent systems R&D and evolving the foundations of computer education.

Stemming from the digital era, technology consolidation has progressed over the past few decades into smaller highly capable devices and (mobile) platforms, linked through worldwide communications with dramatically increasing bandwidth and corresponding lower latency. Much of computer science and engineering research was previously focused on component technologies, feeding into system development typically carried out in industry. Today, beyond “system silos,” the combination and convergence of microscale component technologies is commonplace. Sensing, perception, action, and networking are readily combined in secondary school projects, and highly multidisciplinary complex systems are common in graduate research.

A networked consolidation trend also holds at macroscale. Vast on-demand computing is readily available and becoming much more so as compute center infrastructure continues to grow. Cellular communications trends include emerging standards for extremely low

latency messaging, distributed and edge computing, and increasing communications capacity with rapidly maturing autonomous collaborative unmanned aerial vehicles and intelligent reflective surfaces. So, the interconnection of micro- and macroscale systems, into systems of systems, will accelerate.

It is not much of an exaggeration to say that in addition to our own formal disciplines we are all becoming systems scientists and engineers.^a

The second marked trend, data-driven computing, is the defining aspect of this new era. Large and small ML compute architectures have provided entirely new and highly flexible ways to approximate and predict nonlinear functions and dynamic processes. Escaping linear and other explicit modeling techniques, the ML Era is especially delineated by the ability to represent, manipulate, and generate naturally occurring signals and information (text, speech, vision). However, this remarkable advance in learned computing results in systems for which there are no predictive theories, and progress is generally measured by empirical tests and comparisons.^b

^a We may be engaging in the study of systems (system science) or aspects of developing, integrating, and managing systems (systems engineering).

^b Consider a comparison with, say, designing a radar. The physics of electromagnetic propagation and object interaction, and radar signal design, provide an accurate predictive model that enables optimal design criteria and tradeoffs. However, using radar sensors in intelligent systems with data-driven learning (for example, driverless cars, cognitive radars) leans on empirical testing and evaluation.



TEAMS AND TESTBEDS

These two trends have opened doors to intelligent systems R&D. With data-driven elements, progress has a critical dependence on experimental cycles and spiral development. Virtual and physical testbeds provide the proving grounds, combining “real” data collection and simulations at various levels of abstraction.

Developing and evolving new ML system architectures requires flexible software modularity and careful repository management, with the ability to make surgical alterations and testing from the module to the larger system context, to alter and rearrange the system modules, and to explore various forms of feedback and intrasystem interactions.

Exhaustive tests do not scale, so ablation studies and other targeted statistical test measures are needed.^c

It is increasingly easy to pull together open source software to create a new system, and it is increasingly difficult to fully characterize system safety and reliability. The lack of predictive systems theory forces reliance on testing and red teaming and creation of system monitors. Even swapping out a single model-based module (for example, physics-based) for a data-driven one leads to test and verification challenges. Implicit or explicit system checking relies on carefully constructed critics, rules, boundaries, and constraints. Test, verification, and safety will be supported by a growing toolbox that is adapted to the case at hand, so education may focus on gaining familiarity with the tools and preparing students for their use in new contexts.^d

Intelligent systems development requires both compute and domain expertise. Multidisciplinary

teams now readily combine two or more of computer science, any form of engineering, math/stat, linguistics, biology, neuroscience, psychology/sociology, agronomy, and more. This brings a cross-specialty burden that often lands on the computing specialist, requiring sufficient knowledge in a new domain.

Each area of science has evolved its own semantics and core concepts, and building a new mental model is not easy. Although ML tools^e are now taught routinely within the various disciplines, there remains a “concept gap” between disciplines.^f Computer education can

STUDENTS NEED TO BE PREPARED FOR NOT JUST CREATING SAFE AND RELIABLE INTELLIGENT SYSTEMS, BUT ALSO HOW TO MONITOR IN SITU AND GAUGE THEIR ONGOING BENEFITS AND INTERACTIONS.

prepare students by offering cross-discipline coursework, collaborative opportunities, internship placement, double majoring, and university centers focused on major application areas. Of course, these are the bread and butter of flexible higher education.

Looking forward, it is interesting to think about deeper links with specific disciplines and how these may be incorporated into core computer education curricula, such as learning and cognitive science, neuroscience, and autonomy. These are natural combinations, and they naturally lead to systems.

^c Well-crafted ablation studies are the norm in ML research, as are a variety of test statistics and measures of goodness.

^d Continuing the radar example, the physics enable constraint and resource specifications that can be folded into learning with great benefit. However, this is specific to the application and doesn't provide a predictive theory for, say, learning-based cognitive radar.

^e Software tools, libraries, and “platforms” have blossomed in the ML Era. The successful student will gain skills in the understanding and application of tools, up to modifying and creating new ones.

^f Perhaps we might appeal to AI for a useful cross-disciplinary translator just as we might hope, for example, that an AI could smooth patient–physician interactions.

INTELLIGENT FOUNDATION COMPUTING

With expansive ML Era compute centers come massive foundation models^g and new educational challenges. Training large scale models is a demanding task, with layers of pretraining, fine-tuning, and interacting with experts (for example, preference optimization). While remarkably general purpose, large models provide some form of best estimate or statistical prediction.^h Adapting or improving foundation model response may rely on additional external processing, such as feedback-based query refinement, output bias detection and elimination, and use of external memory such as retrieval-augmented generation.

When broadly trained, a foundation model may respond in ways that are unacceptable (for example, for safety or ethical reasons), so input filters (for example, query classification) can be used to avoid undesired prompts.

Foundation models may form a backbone that interacts with supporting models or headsⁱ to create an intelligent system, for example, a driverless car autonomy with heads that provide navigation plans, perception, object motion prediction, and so on.^j This enables separate creation, evolution, and improvement of the specific purpose models within the system context.

System management and user interaction are now relying on foundation models, and these will blend with traditional operating systems. Foundation models are viewed as a potential new hybrid form of operating system, for example, using large language models, providing the user with semantic interaction

^g It seems remarkable that the term “foundation model,” attributed to a 2021 report from Stanford, is so new and yet already so fundamental.

^h In some cases, such as large language models, generation of content not based on fact is undesirable and referred to as a “hallucination.” In other cases, a large model may be used for artistic generation (for example, music, imagery), and novelty is sought after.

ⁱ Now used more generally, head originally refers to the top (that is, output or last) layer in a neural network, such as a classification layer in a convolutional neural network.

^j We can think of this as a single-agent foundation (backbone) model that is capable of interaction and reasoning over various sources of information, evidence, planning, and prediction provided by additional modules.

and linking with tools, generators, and solvers. These will influence and coevolve with specialized integrated circuits designed for mass market computing (for example, smart phones, laptops).^k

The educational challenge is to understand various forms of foundation model training and refinement, querying, measures of risk and reward, human interaction, and how a foundation model can connect with and draw information from heads to serve a specific application. This, combined with a solid grasp of digital era operating systems, will prepare the student for next-gen intelligent foundation computing.

NETWORKED INTELLIGENT COMPUTING

At ultra-scale, the energy and communications resource consumption of foundation models has become very painfully evident, further motivating smaller model approximations that achieve a desired performance–complexity tradeoff. Connecting at a distance is good news (on-demand availability of a variety of models and applications) and bad news (cost, delay, bandwidth limitations).

We can expect a growing proliferation of models in all shapes and sizes, deployed at both micro- and macroscale. Over time, what is commonly called a “big” model will evolve, and models with billions of parameters (or far more, depending on the future time scale we care to predict) will become even more commonplace and embedded into microscale devices.^l With communications, these combine into networked computing and heterogeneous multiagent systems.^m

^k State-of-the-art dense integrated circuit development and manufacturing is very expensive and so relatively few new devices are generated, and these are typically aimed at large scale applications. Analog computing alternatives have classically struggled with accuracy, programmability, and the need for (digital) calibration, yet may offer much higher speed and lower power. Forms of probabilistic and neuromorphic computing potentially handle these issues.

^l The continued evolution of memory and computing circuits will accommodate larger models at micro-scale.

^m We might think of a computer network as a (relatively static) form of distributed computing, whereas a multiagent system consists of independent and collaborating agents with mixed goals. However, there is a spectrum of possibilities, and the definitions are not rigid.

From micro to macro, edge to cloud, computing networks rely on resource allocation that includes communications bandwidth and quality of service, power and energy, and available compute resources. This also includes data representation, processing, and storage, for example, to locally prestore data in a consumer on-demand multimedia application, or compactly represent information optimized for learning an edge task.

When coupled with multiple concurrent tasks or users sharing a network, resource allocation becomes a challenging dynamic multiobjective optimization problem. Widely deployed compute networks, such as controlling the power grid or the Internet of Things, are complex systems of systems, for example, building on cellular systems and the Internet, each of which is highly complex and dynamic.

With stationary demand and dedicated centralized networking infrastructure, a specific multiobjective resource optimization along a Pareto front may be possible. This is more often an ideal and unachievable goal, with evolving real-time conditions requiring abstractions for robust dynamically updated allocations to steer tradeoffs and achieve overall performance metrics. Beyond classical optimization, introducing students to dynamic multiobjective methods is desirable.

ML SYSTEMS PRINCIPLES

ML Era intelligent systems are a blend of classic systems and intelligent computing; see Table 1. Learning-based representation and abstraction, coupled with architecture and connectivity, provide a rich systems design space.

In particular, graph-based architectures are remarkably versatile and can be understood intuitively. They fit communications networks and underlie multiagent interactions. They model the physical world through geometric-semantic maps, enable planning and navigation, and generalize beyond two dimensions. They naturally express dependence and directed graphs model asymmetry (for example, causalityⁿ). Graphs also provide a powerful inference

ⁿ Causality is easy to intuitively understand, but difficult to formally characterize. A graph provides a structure over which causality and other forms of dependence can be learned.

TABLE 1. Intelligent systems combine key elements of classical system theory and ML-based intelligent computing.

State space representation	<ul style="list-style-type: none"> • Discrete and continuous state-space • State-action • Markov decision processes • Model-based and learning-based
Architecture and connectivity	<ul style="list-style-type: none"> • Centralized-decentralized • Synchronous-asynchronous • Hierarchical, graphs, trees • Distributed, worker-server • Feedforward and feedback • Memory and computing resources
Networking	<ul style="list-style-type: none"> • Bandwidth • Routing • Delay • Reliability
Learning	<ul style="list-style-type: none"> • Information representation and abstraction • Tokens and attention • Semantic and neuro-symbolic • Prediction and generation • Encoding-decoding • Explore-exploit • Games • Resource allocation • Optimization (multiobjective, distributed, dynamic, stochastic) • Gradient descent
Metrics	<ul style="list-style-type: none"> • Value, cost, reward, error, and utility functions • Measures of uncertainty • (Formal) model checking
Safety and Security	<ul style="list-style-type: none"> • Privacy and data protection • Constraints • In situ monitoring • Resilience to deception and attack

architecture, for example, graph neural networks. An educational challenge is not whether graphs are fundamentally important, but rather how to avoid confusion among all the many ways they may be employed.

MONITORING INTELLIGENCE

It is apparent that intelligent systems offer much and threaten many. As the ML Era progresses, in various ways, computing will become synonymous with intelligence.

Open source, powerful edge devices, and cloud computing availability all add up to consumer opportunity. Digital playgrounds allow any of us to teach robots to walk and create new interactive artificial intelligence (AI) agents.

We've argued that education should include creating and applying tools for monitoring, safety, test and evaluation. However, a larger challenge is to educate

students on the issues relating to use and application, risks and rewards, and societal gains and losses. Students need to be prepared for not just creating safe and reliable intelligent systems, but also how to monitor in situ and gauge their ongoing benefits and interactions.

General computer education will include secure computing, privacy, and data protection. These are greatly magnified with intelligent computing, and the importance cannot be overstated. These should not be viewed as add-ons, to be taught separately. Rather, they need to be an integral part of any intelligent system, including resilience to adversarial attack. It is our responsibility to embrace this into all aspects of computer education.^o

As the ML Era progresses, the expected baseline knowledge and core curricula will distill and mature.

^o Just as in communications and network engineering education, where security is inherent and inseparable, similar ideas should be part of the culture of computer education.

Many of these concepts become apparent through in-depth study of systems and examples. It is also important to instill the underlying ideas and build understanding. With a foundation in rudiments, the student can mature into an intelligent computing maestro.^p 🎵

BRIAN M. SADLER is a senior research fellow at the University of Texas at Austin in the Oden Institute for Computational Engineering and Sciences, Austin, TX 78712 USA. Contact him at brian.sadler@ieee.org.

^p Beginning percussion students (such as the author, many years ago) often learn drum rudiments, including various drum rolls, diddles, drags, ruffs, and flams. For example, there are 26 U.S. Standard Rudiments, defined in 1933 and commonly taught today. These are short sticking patterns that can be combined to build formal drum parts (for example, orchestral, marching). The rudiments span a rich set of subtle coordination and mastering them provides a solid technical baseline.

The image is a promotional banner for the IEEE Computer Society. The left side features a dark blue background with a circuit-like pattern. At the top left is the IEEE Computer Society logo, and at the top right is the 80th Anniversary Celebration logo. The main headline reads "UNLEASH YOUR POTENTIAL". Below this are several bullet points: "ATTEND WORLD-CLASS CONFERENCES" (Over 195 globally recognized conferences), "EXPLORE THE DIGITAL LIBRARY" (Over 1 million articles covering world-class peer-reviewed content), "ANSWER CALLS FOR PAPERS" (Write and present your ground-breaking accomplishments), "LEARN NEW SKILLS" (Strengthen your resume with the IEEE Computer Society Course Catalog), "LEVEL UP YOUR CAREER" (Search for new positions in the IEEE Computer Society Jobs Board), and "CREATE YOUR NETWORK" (Make connections in local Region, Section, and Chapter activities). The right side of the banner shows a photograph of four people standing in front of a blue wall with the text "Your Profession". Two people in the center are holding framed certificates or awards. Below the photo is a QR code and the text "Sign Up for a Membership at the IEEE Computer Society" with the URL "computer.org/membership". At the bottom right is the IEEE logo.

Get Published in the *IEEE Transactions on Privacy*

This fully open access journal is soliciting papers for review.

IEEE Transactions on Privacy serves as a rapid publication forum for groundbreaking articles in the realm of privacy and data protection. Submit a paper and benefit from publishing with the IEEE Computer Society! With over 5 million unique monthly visitors to the IEEE Xplore® and Computer Society digital libraries, your research can benefit from broad distribution to readers in your field.

Submit a Paper Today!

Visit computer.org/tp to learn more.



The Urban Space Information Platform: Opportunities and Challenges

Jiabao Li , China University of Geosciences, Wuhan 430074, China

Rajiv Ranjan , Newcastle University, NE1 7RU Newcastle upon Tyne, U.K.

Yuewei Wang  and Xiaohui Huang , China University of Geosciences, Wuhan 430074, China

Philip James , Newcastle University, NE1 7RU Newcastle upon Tyne, U.K.

Schahram Dustdar , TU Wien, 1040 Vienna, Austria, and UPF ICREA, 08018, Barcelona, Spain

Urban construction and development rely on the efficient use of geospatial information. The urban spatial information platform (USIP) serves as a centralized hub for managing and applying multisource urban data. However, the increasing volume and complexity of urban data pose significant challenges in management and processing. This article systematically reviews the theory and technology of USIP, analyzing its current status and future directions. We explore the diversity of data sources in urban scenarios and their characteristics. For data management, we review organization and storage technologies, and for data processing, we analyze architectures and mainstream methods. Furthermore, several practical applications of USIP and its impact on urban development are discussed. This study examines USIP's advancements from a full-chain perspective, encompassing data acquisition, technological implementation, and key applications, offering valuable insights for future research and practice.

As global urbanization accelerates, modern cities have become hubs for population, resources, economic activities, and centers for data production and application. However, they also face complex challenges, including traffic congestion, environmental pollution, and resource disparities. The essence of urban operation lies in the intricate interactions of multiple factors, necessitating the use of information technology to achieve comprehensive perception and refined management.¹ In this context, the urban space information platform (USIP) has emerged as an indispensable tool for enhancing urban governance

and improving service quality. USIP integrates multi-source urban data to enable comprehensive perception, dynamic analysis, and intelligent decision making, providing a robust technical foundation for smart city development.² With the rapid expansion of data and the growing complexity of urban systems, addressing the challenges of efficient data management and utilization has become increasingly imperative.

In recent years, advancements in multiple frontier technologies like the Internet of Things (IoT), artificial intelligence (AI), big data, and cloud computing have significantly expanded USIP's capabilities and applications.³ IoT devices and sensor networks enable real-time data collection, while AI enhances predictions and optimizations for complex urban systems.^{4,5} Additionally, 5G and edge computing improve real-time response, and digital twin technology offers innovative approaches for 3-D modeling and dynamic

1089-7801 © 2025 IEEE. All rights reserved, including rights for text and data mining, and training of artificial intelligence and similar technologies.

Digital Object Identifier 10.1109/MIC.2025.3527821

Date of current version 27 August 2025.

urban simulations.^{6,7,8} However, these technologies also pose challenges, such as managing increasingly complex data and meeting high-performance processing and service demands for cross-departmental collaboration.⁹ Addressing these problems will help shape future research directions and guide practical applications.

Research on USIP has advanced significantly in data management, processing, intelligent analysis, and cross-domain applications. For data management, advancements in spatiotemporal databases, cloud storage, and distributed architectures have enhanced data organization and sharing, particularly for complex geospatial and dynamic sensor data.¹⁰ In data processing, methods that leverage deep learning, reinforcement learning, and graph computing offer robust support for spatiotemporal pattern recognition, dynamic prediction, and complex system optimization. Moreover, USIP applications have expanded into critical domains such as traffic management, environmental monitoring, and disaster emergency response, significantly improving urban operational efficiency and driving industrial development.

Despite these advancements, the development and evolution of USIP face a series of challenges that demand further research and practical solutions. One major issue is the lack of unified formats and standards for multisource data, which hampers efficient data integration and sharing. Additionally, meeting the high-performance demands of large-scale, real-time data processing remains a significant challenge, particularly in dynamically evolving urban scenarios. Furthermore, the diverse demands and complex conditions among cities challenge the versatility and scalability of USIP. Addressing these challenges requires continued innovation in processing architectures and

collaboration among urban planners, technologists, and policy makers.

Therefore, in this work, we provide 1) an analysis of urban data sources and the characteristics of different types of data, 2) a discussion of USIP's technical advances and challenges in data management and processing, and 3) an examination of typical applications of USIP in smart cities and other fields, showing its broad prospects for future research directions.

DATA SOURCE IN URBAN SPACE

Data sources form the foundation of USIP, directly influencing its analytical capabilities and application scenarios. Given the diversity, complexity, and heterogeneity of urban data, it is essential to classify and understand the characteristics of these sources.¹¹ Data sources can be broadly categorized, and their specific types and characteristics are summarized in Table 1 for clarity and reference. For instance, sensor-based data include environmental sensors for meteorological, noise, and air quality as well as traffic sensors like cameras and microwave radars, which are widely distributed across urban areas. Additionally, city-related data can be classified based on temporal dimensions and spatial coverage, depending on the specific application scenarios. The inherent temporal and spatial heterogeneity of urban data, especially from dynamic sources such as sensor and traffic data, requires real-time processing to ensure timely and accurate analysis.

The continuous enrichment of urban spatial information data is closely tied to advancements in collection methods, which have become increasingly diverse. With the rapid development of information and communication technologies, common data collection methods now include sensor networks, remote sensing satellites, drones, and big data crawlers. Figure 1 shows the multisource heterogeneous data

TABLE 1. Main categories of urban spatial information data sources.

Data category	Sample data source	Characteristic
Sensor-based data	IoT devices such as weather sensors and smart street lights	Real time and dynamic, covering environmental monitoring, traffic, building monitoring, and so on
Remote sensing image data	Satellite and drone photography	High resolution and multispectral, but with limited timeliness
Socioeconomic data	Government statistics and business activity records	Population and economic information, the data cycle is long but the impact is far-reaching
Traffic data	Bus card swiping records and shared bicycle tracks	Space and time intensive
Social media and crowdsourced data	Weibo geo-tag data and user-generated content	Unstructured, dynamic, and with a strong social orientation

currently available around urban scenes. Each method is tailored to specific application scenarios and data requirements. For example, lidar, as an active sensor, can rapidly collect high-density point cloud data with a high frame rate, making it ideal for high-precision measurements and 3-D modeling.¹² However, urban data collected through different methods often exhibit quality issues, such as inaccuracies, incompleteness, and inconsistencies, posing significant challenges for subsequent data management and processing. Moreover, challenges remain in the widespread adoption of data collection technologies, including the high cost of acquiring high-precision data, the need for stronger data privacy protection and ethical safeguards, and the fragmentation caused by scattered sources and a lack of unified standards.

To meet the increasingly extensive and complex data application needs, the relevant technologies that surround data acquisition must also adapt to the ever-changing real-world foundation. In the future, the evolution and utilization of data sources will progress in several key directions. First, it is essential to break down data silos and promote unified data management standards to achieve more comprehensive data integration. Second, new data collection methods, such as social media, mobile device data, drones, and connected vehicle data, can be introduced. Additionally, intelligent data collection, combined with AI for automated data filtering and classification, presents significant growth potential. Most importantly, establishing secure and reliable privacy protection mechanisms will be crucial to balance data openness and protection.

DATA MANAGEMENT OF USIP

Data management serves as the core of USIP, offering essential technical support for the storage, processing, and utilization of multisource heterogeneous urban data. This section focuses on data organization and storage, emphasizing their key role in enabling efficient storage and retrieval. It further explores the necessity of data quality control and evaluation to ensure reliability and underscores the role of data security in maintaining data integrity and accessibility.

Data Organization and Storage

Data organization and storage form the foundation of USIP data management, directly influencing query efficiency, storage costs, and system scalability. With the increasing scale and diversity of urban data, traditional single-node storage architectures have become insufficient, driving the adoption of distributed, cloud, and edge storage as key solutions. Meanwhile, the

spatiotemporal characteristics of urban data impose new demands on data models and indexing technologies to ensure efficient organization and retrieval. Next, we discuss three key points in data management technology: data model, storage architecture, and spatiotemporal index.

Data Model

The data model describes data, their relationships, and operational rules, providing an abstract framework for representing information and facilitating operations within database systems. In USIP, it specifies the structure, storage, and access methods of data, offering a standardized approach for integrating and managing multisource heterogeneous urban data. By leveraging data models, complex urban data can be processed into formats that are interpretable and operable by computer systems, enabling efficient storage, querying, and analytical operations. There are now various types of data models, including vector models, raster models, spatiotemporal data models, hierarchical data models, and graph data models, each tailored to specific data structures and application requirements.¹³ Selecting the appropriate data model based on data characteristics and application needs is essential as it bridges the gap between physical data storage and logical user requirements, enabling users to interact with complex data structures through simplified interfaces.

The spatiotemporal data model represents data with both spatial and temporal attributes. By integrating spatial and temporal data models, it effectively handles spatial data that change over time. Urban spatial information data are defined not only by their location and distribution in the spatial dimension but

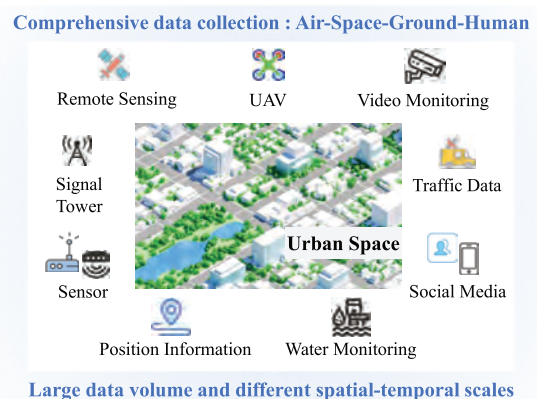


FIGURE 1. Multiple data sources in urban scenarios. UAV: unpiloted aerial vehicle.

also by their changes and dynamics in the temporal dimension. As a key model for organizing and storing spatiotemporal data, they enable the efficient processing, storage, and querying of dynamic information, supporting the analysis of temporal variations and spatial patterns. In the spatiotemporal data model, the spatial dimension is typically represented by geometric data, such as points, lines, and polygons, which define the location and form of urban objects. The temporal dimension captures the time attributes of the data by using time stamps, time intervals, or other time-based representations. The strength of spatiotemporal data models lies in their ability to manage large-scale, dynamic datasets, offering robust support for spatial analysis and time-series analysis. In USIP research, this model applies to various scenarios, including urban traffic flow prediction, environmental change monitoring, and real-time positioning services. Overall, the spatiotemporal data model plays a crucial role in enhancing the intelligence and decision-making capabilities of USIP.

Distributed Storage

The data storage architecture is critical to data management and access efficiency and directly impacts the platform's overall performance, including scalability, fault tolerance, and reliability. Traditional single-node storage, limited by capacity and coverage, is insufficient to meet the demands of managing large-scale urban data. To address these challenges, distributed storage systems have emerged, enabling the storage of data across multiple physical nodes while ensuring high availability, scalability, and fault tolerance. Mainstream distributed storage technologies, such as Hadoop HDFS, Ceph, OpenStack Swift, and Google File System, support the distributed storage and processing of large-scale data, offering efficient read-write performance and robust data backup mechanisms.¹⁴ Furthermore, cloud storage, built on distributed architectures, offers on-demand storage services, effectively addressing dynamic and flexible storage requirements in modern urban applications. In the context of cloud-edge-end collaborative development, edge storage also plays a crucial role. It reduces data transmission delays and supports application scenarios with stringent requirements for real-time data processing and localized storage.

Managing diverse types of urban data ultimately relies on database systems for data archiving, storage, and access interfaces to support operations such as adding, deleting, modifying, and querying data. However, different data types necessitate appropriate data

models for accurate representation and recording, along with suitable database systems to ensure efficient storage and management. Based on data organization and application requirements, database systems can be categorized into relational, nonrelational, and graph. Relational databases, such as MySQL and PostgreSQL, are ideal for storing structured data, including city demographic information and traffic flow records. Nonrelational databases, on the other hand, are better suited for semistructured or unstructured data and offer flexible data models and superior scalability to handle diverse and dynamic urban datasets. To manage the massive and continuously growing volume of urban spatial information data, including satellite imagery, video surveillance, and sensor data, distributed storage systems and cloud storage solutions have been widely adopted. These systems integrate seamlessly with big data processing platforms, such as Hadoop and Spark, enabling efficient large-scale data processing and analysis. In essence, database systems not only provide persistent data storage but also enable the effective utilization of large-scale data through optimized management and operations.

Spatiotemporal Index

Indexing technology is a key method for enhancing data storage and query efficiency. It organizes and manages data through the establishment of efficient data structures. For spatial data, spatial indexes facilitate complex query operations, such as range, proximity, and overlap queries, enabling faster and more accurate data retrieval. However, spatiotemporal data, which integrate both spatial and temporal attributes, present greater challenges. Traditional spatial indexes, such as R-tree and B-tree, and time-stamp-based indexes struggle to efficiently process the combined complexities of spatiotemporal data. Spatiotemporal indexing technology provides an effective solution for organizing and retrieving spatiotemporal data, particularly in large-scale urban spatial information applications. Current mainstream spatiotemporal indexing methods include space-time R-tree, grid-based indexing, and space-time hash indexing. Looking ahead, spatiotemporal indexing technology is expected to advance toward greater query optimization, multidimensional data fusion, and big data processing capabilities.

Data Quality Control and Evaluation

Quality control and evaluation are vital components of data lifecycle management and play a critical role in enhancing the reliability and decision-making capabilities of information application platforms. However,

during the collection, transmission, and storage of multisource urban data, questions surrounding noise, missing values, redundancy, and inconsistency often arise. These issues can undermine data accuracy and, in turn, mislead urban planning and emergency decision-making processes. Effective data quality control entails cleaning erroneous data, imputing missing values, removing redundancy, and ensuring consistency and integrity. These processes provide reliable data support for critical applications, including traffic management, environmental monitoring, and urban planning. For data quality assessment, integrating AI and big data technologies can significantly enhance both the efficiency and accuracy of the evaluation process.

However, the increasing volume and diversity of sensitive data have heightened concerns about privacy breaches and data misuse. Data leaks or tampering can result in severe social and economic consequences. To safeguard data during storage, transmission, and use, common security technologies include data encryption, access control, and multifactor authentication. Furthermore, with the rise of distributed storage systems, blockchain technology has emerged as a decentralized security solution, enhancing data storage security through tamper-proof distributed ledgers. Properly addressing data security challenges in the design and implementation of USIP is critical to ensuring a robust foundation for smart city development.

Data Security

As a platform for urban data integration and utilization, USIP handles vast numbers of geospatial and real-time dynamic data, which are closely related to urban traffic management, emergency response, and public safety.

DATA PROCESSING OF USIP

The effective processing of massive urban data depends on the adoption of suitable processing architectures and methods, directly influencing the efficiency of data analysis and application. The key technologies

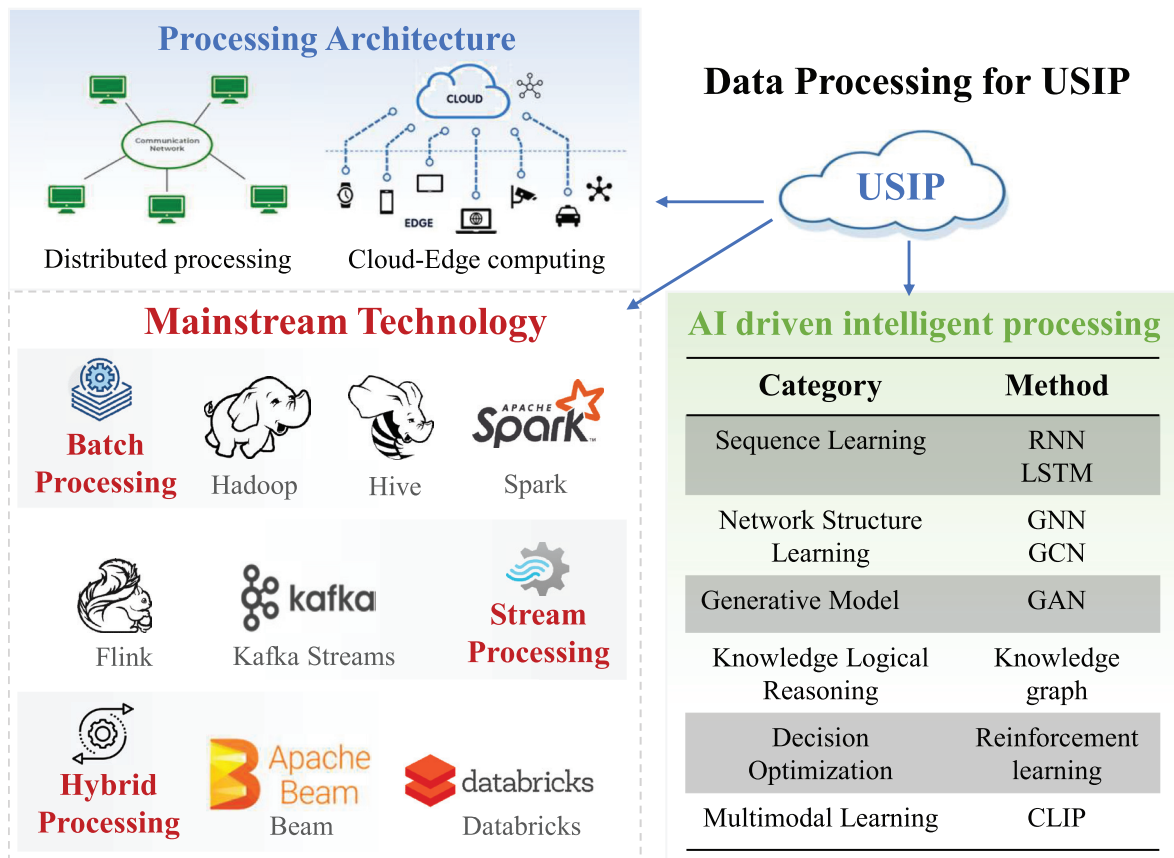


FIGURE 2. Key technologies for data processing of urban spatial information. RNN: recurrent neural network; LSTM: long short-term memory; GNN: graph neural network; GAN: generative adversarial network; CLIP: contrastive language–image pretraining; GCN: graph convolutional network.

for data processing of urban spatial information are shown in Figure 2. This review examines three key aspects: data processing architectures, mainstream processing technologies, and AI-driven computational analysis.

In USIP, the data processing architecture defines the flow and computation of data, directly shaping computation patterns and influencing the platform's performance and efficiency. These architectures continuously evolve to address the challenges posed by increasing data volume, real-time demands, and growing complexity. Traditional stand-alone architecture has become inadequate, leading to the adoption of distributed architecture and cloud-edge collaborative solutions as mainstream approaches. Distributed architecture achieves high scalability and fault tolerance by distributing data and tasks across multiple nodes for parallel computation. By integrating in-memory computing, task scheduling optimization, and storage-computation decoupling technologies, this architecture overcomes the performance limitations of the stand-alone system, enabling efficient processing of large-scale data.¹⁵ Additionally, cloud-edge collaborative architecture serves as a critical processing model to address the conflict between real-time responsiveness and

scalability. By offloading real-time processing tasks to edge devices and leveraging cloud systems for data storage and in-depth analysis, it plays a pivotal role in balancing real-time response with a centralized analysis. The cloud-edge-end architecture, through the integration of hierarchical task scheduling, edge intelligence, and low-latency network technologies, has been successfully applied to scenarios such as real-time traffic control and emergency response.

For specific data processing technologies, depending on the real-time requirements and application scenarios, mainstream approaches include offline batch processing, real-time stream processing, and hybrid processing. Offline batch processing is designed primarily for handling static, large-scale datasets, making it well suited for historical data analysis and computationally intensive tasks. Representative technologies like MapReduce, Flink, and Spark have been extensively applied to scenarios such as urban development trend modeling and environmental change assessment.

Specific data processing technologies vary based on real-time requirements and application scenarios, with mainstream approaches that include offline batch processing, real-time stream processing, and

TABLE 2. Comparison of mainstream big data processing frameworks.

Framework	Description	Characteristic	Applicable scenarios
MapReduce	The distributed computing model proposed by Google; simplifies parallel programming to efficiently process large-scale datasets.	Easy to use, highly scalable, reliable; follows data locality.	Designed for handling static, large-scale batch processing tasks, such as log analysis and extract, transform, load tasks.
Spark	A memory-based distributed computing framework that supports batch and stream processing and uses memory to cache intermediate results to improve processing speed.	The Resilient Distributed dataset is used as the core component; high performance and strong ecosystem.	Applicable to a variety of data processing scenarios, such as large-scale data statistics and social media analysis.
Flink	A distributed stream processing framework that is focused on high-performance, high-throughput real-time stream data processing.	Provides true stream processing; event-time-based processing engine; low latency, supports stateful computing.	Suitable for scenarios that require high concurrency, complex event processing, and batch-stream integration.
Storm	A distributed real-time stream processing framework that is designed for low-latency, high-throughput real-time data processing.	Processes data with millisecond latency; provides horizontal scalability; lightweight deployment and operation.	Suitable for scenarios with extremely high real-time requirements and relatively simple data stream processing logic.
Beam	A unified programming model for the definition and execution of batch and streaming tasks.	Provides an abstraction layer to write code once and run it on different distributed processing engines.	Suitable for scenarios that require multiengine compatibility and batch and stream unification.

hybrid processing. Table 2 lists the current mainstream spatiotemporal big data processing frameworks, including their characteristics and applicable scenarios. Offline batch processing is tailored for handling static, large-scale datasets, making it ideal for historical data analysis and computationally intensive tasks. Computational frameworks such as MapReduce and Spark have been widely utilized in applications like urban development trend modeling and environmental change assessment.¹⁶ Real-time stream processing, characterized by low latency and high concurrency, is well suited for dynamic monitoring and rapid response scenarios. Applications include real-time traffic flow optimization, air quality monitoring, and early warning systems, which are often implemented by using frameworks like Flink. Hybrid processing combines the strengths of both batch and stream processing, enabling efficient analysis of historical data while providing rapid responses to real-time data. For example, the Lambda architecture achieves efficient processing of historical and real-time data by separating batch and stream processing modules. In USIP, hybrid processing technologies can be applied to scenarios that require multidimensional analysis. For instance, in environmental monitoring, it integrates historical meteorological data with real-time sensor data for multilevel pollution source analysis. By leveraging these technologies, USIP provides comprehensive support, seamlessly bridging historical data analysis and real-time responsiveness. This integration enhances the intelligence and precision of urban management, facilitating more informed decision making across various application domains.

The rapid development and widespread application of AI in data processing and analysis are fundamentally transforming the technological landscape of USIP. With its powerful capabilities in data mining, pattern recognition, and prediction, AI offers innovative solutions for efficiently processing multisource heterogeneous data. Currently, machine learning and deep learning methods are continuously evolving, with tailored approaches being applied to specific scenarios, providing robust support for diverse urban applications.¹⁷ Specifically, long short-term memory-based spatiotemporal prediction models can be used for traffic flow forecasting and environmental trend analysis. Graph neural networks enable efficient path optimization and resource allocation by modeling complex urban spatial networks. Moreover, multimodal data integration facilitates geospatial information fusion and multisource data analysis, while generative adversarial networks are employed to reconstruct missing

data and enhance low-quality data. In addition, knowledge-driven approaches such as knowledge graph reasoning and analysis, as well as reinforcement learning for strategy optimization, are other AI methods gaining traction in USIP. Despite challenges such as limited interpretability and high resource demands, AI continues to play a pivotal role in enhancing USIP's intelligence and decision-support capabilities.

Although USIP has adopted advanced architectures, AI, and other technical methods in data processing, it still faces significant challenges, including balancing data scale with real-time performance, integrating multi-source heterogeneous data, and optimizing resources and task allocation. Looking ahead, USIP's data processing capabilities are expected to advance toward greater efficiency, intelligence, and distributed collaboration. Moreover, the integration of large-scale pre-trained models with green computing holds promise for breakthroughs in urban scenario forecasting and real-time responsiveness, providing robust technical support for intelligent and resilient urban governance.

APPLICATIONS OF USIP

As a foundational infrastructure for advancing urban digitalization and refined management, USIP plays a vital role in multisource data integration, dynamic analysis, and intelligent decision support, as shown in Figure 3. It has been successfully applied to management and decision making across various domains. This section uses USIP applications in three key areas, smart cities, environmental monitoring, and emergency management, to demonstrate its transformative impact and extensive applications across various urban domains.

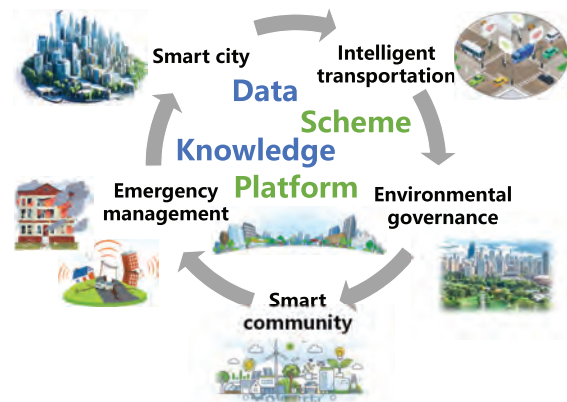


FIGURE 3. Examples of various applications of USIP.

Smart City

Smart cities strive to enhance urban operational efficiency, improve residents' quality of life, and promote sustainable development through the integration of advanced technologies. Leveraging the data integration, efficient processing, analysis, and intelligent decision-making capabilities provided by USIP, the platform has been effectively applied to various domains, including traffic management, smart communities, and social governance.

In traffic management scenarios, USIP enables the integration of multisource data from cameras and sensors to support traffic flow prediction, dynamic signal control, and route optimization. By mining traffic big data and optimizing real-time traffic strategies, it improves road efficiency while reducing energy consumption. With the support of USIP, community services can benefit from optimized resource allocation, enhanced safety monitoring, and elder care services. In addition, USIP plays a significant role in social governance and public safety by integrating data from social media, cameras, and sensors to dynamically monitor and manage public events. For example, by analyzing real-time security data across urban areas, it can identify crime hot spots and deploy targeted measures to address them effectively. Overall, USIP has become a cornerstone of smart city development, fostering more efficient, intelligent, and inclusive urban environments.

Environmental Monitoring

Currently, cities face significant challenges in environmental management, including monitoring inefficiencies, inaccurate predictions, and unequal resource allocation. USIP provides effective solutions to address these issues. As a data-driven platform, USIP facilitates real-time monitoring and dynamic analysis of critical environmental indicators such as air quality, water resource distribution, and noise levels.¹⁸ With the integration of AI, USIP enables the prediction of pollution dispersion trends and the issuance of early warnings, effectively guiding rapid pollution source management and policy adjustments. At the same time, USIP plays a pivotal role in sustainable development planning by supporting regional carbon emission tracking, green energy planning, and low-carbon policy evaluation, providing scientific evidence for achieving urban sustainability goals. Looking ahead, with the deeper application of AI, large-scale pretrained models, and green computing technologies, USIP is poised to unlock greater potential in environmental

protection and sustainable development, offering innovative solutions for the green transformation of global urban governance.

Emergency Management

Emergency management is a vital pillar of urban safety and resilience. In scenarios such as natural disasters, public health crises, and sudden security incidents, the speed and quality of response directly impact the scale of losses and recovery efficiency. Information platforms utilize IoT sensors and remote sensing technologies for dynamic monitoring, integrating historical and real-time data to accurately predict event trends.¹⁹ For instance, in flood management, USIP leverages water level and rainfall data collected by sensors to dynamically simulate flood spread and optimize evacuation routes, significantly enhancing rescue efficiency. Simultaneously, the platform employs digital twin technology to visually represent resource distribution in affected areas, aiding government agencies in efficiently allocating rescue efforts. During public health emergencies, USIP combines population mobility trajectories with social media data to trace epidemic spread, identify high-risk areas, and provide scientific support for effective containment measures.

Looking ahead, with the deep integration of AI and federated learning technologies, USIP is poised to facilitate intelligent and collaborative solutions in emergency management, significantly enhancing urban emergency response capabilities and public safety governance. By leveraging rich data together with efficient and intelligent analytical processes, USIP serves as a cornerstone for disaster prevention, response, and mitigation.

CONCLUSION

This study systematically explored the key components and practical applications of USIP by leveraging advances in computer science theory and cutting-edge technologies, focusing on four aspects: data sources, management, processing, and applications. First, we analyzed the diversity of urban spatial information sources, categorized the data, and discussed their respective characteristics. Next, we focused on data organization, storage, quality control, and security, examining the core technologies and challenges in USIP's data management. In addition, we delved into the evolution and current state of data processing technologies, including processing frameworks, mainstream techniques, and the accelerating impact of AI on data

processing. Finally, through an in-depth discussion of typical scenarios such as smart city construction, environmental governance, and emergency management, we demonstrated the immense potential of USIP to enhance urban resilience, quality of life, and sustainable development.

Although USIP has already been applied in various urban scenarios and domains, delivering economic value and significant advances, many challenges remain for its development. Future efforts must address bottlenecks such as integrating data from multiple sources, processing large amounts of data in real time, and balancing privacy with data sharing. By bridging the gap between complex urban challenges and data-driven solutions, USIP offers a solid technological foundation for advancing smart, resilient, and human-centered cities. 🌐

ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China under Grant U21A2013.

REFERENCES

1. A. Camero and A. Enrique, "Smart city and information technology: A review," *Cities*, vol. 93, pp. 84–94, Apr. 2019, doi: 10.1016/j.cities.2019.04.014.
2. D. Jiang, "The construction of smart city information system based on the Internet of Things and cloud computing," *Comput. Commun.*, vol. 150, pp. 158–166, Jan. 2020, doi: 10.1016/j.comcom.2019.10.035.
3. D. Singh and C. K. Reddy, "A survey on platforms for big data analytics," *J. Big Data*, vol. 2, no. 1, pp. 1–20, Dec. 2015, doi: 10.1186/s40537-014-0008-6.
4. W. Han et al., "Geological remote sensing interpretation using deep learning feature and an adaptive multisource data fusion network," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022, doi: 10.1109/TGRS.2022.3183080.
5. Z. Allam and Z. A. Dhunny, "On big data, artificial intelligence and smart cities," *Cities*, vol. 89, pp. 80–91, Jun. 2019, doi: 10.1016/j.cities.2019.01.032.
6. A. Gohar and G. Nencioni, "The role of 5G technologies in a smart city: The case for intelligent transportation system," *Sustainability*, vol. 13, no. 9, May 2021, Art. no. 5188, doi: 10.3390/su13095188.
7. K. Cao, Y. Liu, G. Meng, and Q. Sun, "An overview on edge computing research," *IEEE Access*, vol. 8, pp. 85,714–85,728, 2020, doi: 10.1109/ACCESS.2020.2991734.
8. E. Shahat, C. T. Hyun, and C. Yeom, "City digital twin potentials: A review and research agenda," *Sustainability*, vol. 13, no. 6, Mar. 2021, Art. no. 3386, doi: 10.3390/su13063386.
9. M. Talebkhah, A. Sali, M. Margani, M. Gordan, S. J. Hashim, and F. Z. Rokhani, "IoT and big data applications in smart cities: Recent advances, challenges, and critical issues," *IEEE Access*, vol. 9, pp. 55,465–55,484, 2021, doi: 10.1109/ACCESS.2021.3070905.
10. P. Repette, J. S. Marques, T. Yigitcanlar, D. Sell, and E. Costa, "The evolution of city-as-a-platform: Smart urban development governance with collective knowledge-based platform urbanism," *Land*, vol. 10, no. 1, Jan. 2021, Art. no. 33, doi: 10.3390/land10010033.
11. L. Wang, B. Zuo, Y. Le, Y. Chen, and J. Li, "Penetrating remote sensing: Next-generation remote sensing for transparent earth," *Innovation*, vol. 4, no. 6, Sep. 2023, Art. no. 100519, doi: 10.1016/j.xinn.2023.100519.
12. B. Huang and J. Wang, "Big spatial data for urban and environmental sustainability," *Geo-Spatial Inf. Sci.*, vol. 23, no. 2, pp. 125–140, Jun. 2020, doi: 10.1080/10095020.2020.1754138.
13. A. F. Tollefsen, H. Strand, and H. Buhaug, "PRIO-GRID: A unified spatial data structure," *J. Peace Res.*, vol. 49, no. 2, pp. 363–374, Apr. 2012, doi: 10.1177/0022343311431287.
14. M. Tang, Y. Yu, Q. M. Malluhi, M. Ouzzani, and W. G. Aref, "LocationSpark: A distributed in-memory data management system for big spatial data," *Proc. VLDB Endowment*, vol. 9, no. 13, pp. 1565–1568, Sep. 2016, doi: 10.14778/3007263.3007310.
15. L. Wang, Y. Ma, A. Zomaya, R. Ranjan, and D. Chen, "A parallel file system with application-aware data layout policies for massive remote sensing image processing in digital earth," *IEEE Trans. Parallel Distrib. Syst.*, vol. 26, no. 6, pp. 1497–1508, Jun. 2014, doi: 10.1109/TPDS.2014.2322362.
16. L. Cheng, L. Wang, R. Feng, and J. Yan, "Remote sensing and social sensing data fusion for fine-resolution population mapping with a multimodel neural network," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 5973–5987, 2021, doi: 10.1109/JSTARS.2021.3086139.
17. G. Nguyen et al., "Machine learning and deep learning frameworks and libraries for large-scale data mining: A survey," *Artif. Intell. Rev.*, vol. 52, no. 1, pp. 77–124, Jan. 2019, doi: 10.1007/s10462-018-09679-z.
18. S. Manfreda et al., "On the use of unmanned aerial systems for environmental monitoring," *Remote Sens.*, vol. 10, no. 4, Apr. 2018, Art. no. 641, doi: 10.3390/rs10040641.
19. C. P. Gonzalez, M. Colebrook, J. R. Garcia, and C. R. Remedios, "Developing a data analytics platform to support decision making in emergency and security management," *Expert Syst. Appl.*, vol. 120, pp. 167–184, Apr. 2019, doi: 10.1016/j.eswa.2018.11.023.

JIABAO LI is pursuing his Ph.D. degree in geoscience information engineering with the School of Computer Science, China University of Geosciences, Wuhan 430074, China. Contact him at wutonglbj@cug.edu.cn.

RAJIV RANJAN is the University Chair Professor for the Internet of Things research with the School of Computing, Newcastle University, NE1 7RU Newcastle upon Tyne, U.K. Contact him at rajanranjan@ncl.ac.uk.

YUEWEI WANG is an associate professor with the School of Computer Science, China University of Geosciences, Wuhan 430074, China. Contact him at yuewei.w@cug.edu.cn.

XIAOHUI HUANG is a lecturer with the School of Computer Science, China University of Geosciences, Wuhan 430074, China. Contact him at xhhuang@cug.edu.cn.

PHILIP JAMES is a senior lecturer in the Urban Observatory, School of Civil Engineering, Newcastle University, NE1 7RU Newcastle upon Tyne, U.K. Contact him at philip.james@ncl.ac.uk.

SCHAHRAM DUSTDAR is a professor of computer science and head of the Distributed Systems Group at TU Wien, 1040 Vienna, Austria, and UPF ICREA, 08018, Barcelona, Spain. Contact him at dustdar@dsg.tuwien.ac.at.



computer.org/pervasive

IEEE
pervasive
COMPUTING

Call for Articles

IEEE Pervasive Computing publishes accessible, useful peer-reviewed papers on the latest developments in pervasive, mobile, and ubiquitous computing. Topics include hardware technology, software infrastructure, real-world sensing and interaction, human-computer interaction, and systems considerations, including deployment, scalability, security, and privacy.



Author Guidelines
bit.ly/3UoGDT7

 **IEEE**
 **IEEE COMPUTER SOCIETY**

Jingxian Wang: “Pushing the Limits of Battery-Free Internet-of-Things”

Lakmal Meegahapola , ETH Zurich, 8092, Zurich, Switzerland

FROM THE EDITOR

The “Emerging Rockstar” segment in *IEEE Pervasive Computing* highlights rising stars in the field of pervasive computing through captivating interviews. In the series’ inaugural article, Dr. Lakmal Meegahapola, a

member of *IEEE Pervasive Computing*’s editorial board, conducts an interview with Dr. Jingxian Wang, an Assistant Professor at the National University of Singapore.



Jingxian Wang is an Assistant Professor with the National University of Singapore’s Computer Science Department. Previously, he was a Researcher with Microsoft Research, Redmond, WA and has a Ph.D. from Carnegie Mellon University. His research is in the area of wireless systems. In 2023, he was honored with the ACM SIGMOBILE Dissertation Award. His contributions to wireless systems and ubiquitous computing have earned best paper awards at top-tier venues, including ACM UbiComp/ISWC 2020 and ACM/IEEE IPSN 2021. Furthermore, his innovations have twice been featured as Research Highlights in *Communications of the ACM*. He is also the recipient of the Microsoft Research Ph.D. Fellowship.

IEEE Pervasive: Could you provide an overview of your research in battery-free Internet-of-Things, some of the key challenges in the area, and how you have addressed these challenges in your work?

Jingxian Wang: The dream of the Internet of Things (IoTs), or pervasive computing, is to have everyday objects such as clothing and keys smartly connected to the Internet. However, there is a challenge: reliance on batteries is a major drawback. Batteries require frequent replacements, add weight, and increase costs. My research is focused on creating a battery-free Internet of Things. We have observed how tiny and cost-effective RFID tags can

be, setting a precedent for future IoT platforms. However, these technologies currently face limitations in communication range and sensing capabilities. Overcoming these limitations is exactly what my work aims to achieve.^{1,2,3,4,5,6,7}

In battery-free IoT systems, the main issue is their short communication range. This is because their energy sources blindly disperse wireless power in all directions. What we need is a more focused energy transfer, like beamforming in wireless communication. However, the challenge with battery-free devices is that they do not provide feedback before receiving power, unlike battery-powered devices such as mobile phones. That is where our innovation comes in—“blind wireless beamforming.” By precisely delivering the wireless energy without prior feedback from the device, we have significantly extended the communication range of these systems.⁸ Interestingly, this concept has inspired our team to apply wireless beaming

in a microwave oven, such as for evenly heating a pizza slice.⁹

Pervasive: Regarding your work with electronics-free soft robots, could you detail the specific aspects you focus on? In addition, what are the ultimate objectives or aspirations you have for this area of your research?

Wang: Electronics-free soft robots, composed entirely of stretchable, shape-memory polymers, represent a fascinating area of research.¹⁰ These robots are also battery-free and are actuated by heat. Traditionally, tethered connections to a dc power supply provide the necessary heat, but our aim is to use beamforming techniques to wirelessly heat and move the robot. This involves two primary challenges: precisely directing wireless energy to the robot as it moves and developing a new material architecture for the robot to efficiently absorb and utilize the wireless energy. Our experiments with a 2-watt WiFi system have been promising,^{11,12} demonstrating the robot's movement through wireless actuation. The long-term goal is to develop these robots as autonomous, battery-free sensors or wireless computing platforms that can navigate complex environments.

IT IS ACTUALLY OUR NETWORKS FOR WIRELESS COMMUNICATION THAT CONSUME THE MOST ENERGY.

Pervasive: Considering the nature of your work, which spans fields such as artificial intelligence, robotics, and materials science, what challenges do you encounter in conducting such interdisciplinary research?

Wang: Jumping into fields such as robotics, materials science, and AI from wireless systems is like moving to a new country—exciting, but also a bit intimidating. At first, it is tough to step out of your comfort zone. But once you do, you start seeing how your own expertise, such as understanding signal propagation, communication, and wireless networks, can actually shake things up in these other areas. It is all about connecting the dots between different fields to create something totally new and groundbreaking.

Pervasive: What advice would you offer to upcoming researchers interested in pursuing a career in a similar research field as you, based on your experiences and insights?

Wang: I encourage researchers to broaden their horizons beyond their specific areas of study. Reading papers and engaging in activities outside your primary field can be incredibly enlightening. This approach has been instrumental in helping me generate unique ideas and solve complex problems by integrating perspectives from various disciplines. It not only enriches your research, but also brings in unexpected yet valuable insights. In addition, my research is highly problem driven. I focus on the practical challenges and everyday problems, which has been a key driver in developing solutions that are both theoretically sound and practically impactful. This problem-driven approach, combined with a broad, interdisciplinary perspective, is truly essential.

Pervasive: Finally, what are your future goals?

Wang: Looking forward, my focus is on developing wireless systems that are both sustainable and efficient. At first glance, one might assume that the most energy-consuming sector of the digital world is the “Cloud”—data centers storing vast amounts of data and running computationally intense AI algorithms. Surprisingly, that is not the case. It is actually our networks for wireless communication that consume the most energy. This is especially true with the rise of emerging technologies like 5G, IoT, and Industry 4.0, which have been overlooked in this regard. In fact, they account for over 55% of worldwide energy consumption in computing systems. The critical question then becomes: How can we build wireless connectivity that is both energy-efficient and sustainable? This is likely to be the next big challenge in achieving next-generation wireless ubiquity. By enhancing the scalability, longevity, and energy efficiency of wireless systems, my research stands to make a significant contribution toward a more connected and sustainable future. 🌍

REFERENCES

1. J. Wang, J. Zhang, K. Li, C. Pan, C. Majidi, and S. Kumar, “Locating everyday objects using NFC textiles,” in *Proc. 20th Int. Conf. Inf. Process. Sensor Netw.*, 2021, pp. 15–30.
2. J. Wang et al., “Rfid tattoo: A wireless platform for speech recognition,” *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 3, no. 4, pp. 1–24, 2019.
3. C. Pan et al., “Silver-coated poly (dimethylsiloxane) beads for soft, stretchable, and thermally stable conductive elastomer composites,” *ACS Appl. Mater. Interfaces*, vol. 11, no. 45, pp. 42561–42570, 2019.

4. H. Jin, J. Wang, Z. Yang, S. Kumar, and J. Hong, "RF-wear: Towards wearable everyday skeleton tracking using passive RFIDs," in *Proc. ACM Int. Joint Conf. Int. Symp. Pervasive Ubiquitous Comput. Wearable Comput.*, 2018, pp. 369–372.
5. H. Jin, J. Wang, Z. Yang, S. Kumar, and J. Hong, "Wish: Towards a wireless shape-aware world using passive RFIDs," in *Proc. 16th Annu. Int. Conf. Mobile Syst., Appl., Serv.*, 2018, pp. 428–441.
6. W. Sun et al., "Microfluid: A multi-chip RFID tag for interaction sensing based on microfluidic switches," *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 6, no. 3, pp. 1–23, 2022.
7. J. Wang, V. Ranganathan, J. Lester, and S. Kumar, "Ultra low-latency backscatter for fast-moving location tracking," *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 6, no. 1, pp. 1–22, 2022.
8. J. Wang, J. Zhang, R. Saha, H. Jin, and S. Kumar, "Pushing the range limits of commercial passive {RFIDs}," in *Proc. 16th USENIX Symp. Netw. Syst. Des. Implementation*, 2019, pp. 301–316.
9. H. Jin, J. Wang, S. Kumar, and J. Hong, "Software-defined cooking using a microwave oven," in *Proc. 25th Annu. Int. Conf. Mobile Comput. Netw.*, 2019, pp. 1–16.
10. X. Huang, M. Ford, Zach J. Patterson, M. Zarepoor, C. Pan, and C. Majidi, "Shape memory materials for electrically-powered soft machines," *J. Mater. Chem. B*, vol. 8, no. 21, pp. 4539–4551, 2020.
11. J. Wang et al., "Wireless actuation for soft electronics-free robots," in *Proc. 29th Annu. Int. Conf. Mobile Comput. Netw.*, 2023, pp. 1–16.
12. Y. Song et al., "Navigating soft robots through wireless heating," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2023, pp. 2598–2605.

LAKMAL MEEGAHAPOLA is a postdoctoral researcher at ETH Zurich, 8092, Zurich, Switzerland. He earned his Ph.D. in electrical engineering from EPFL, Lausanne, Switzerland. His research interests lie at the intersection of mobile and wearable sensing, digital health, machine learning, deep learning, and human-computer interaction. Contact him at Imeegahapola@ethz.ch.



IEEE TRANSACTIONS ON SUSTAINABLE COMPUTING

ADVANCING ENERGY-EFFICIENT AND
RESOURCE-CONSCIOUS TECHNOLOGY



For more information visit computer.org/tsusc





CALL FOR SPECIAL ISSUE PROPOSALS

Computer solicits special issue proposals from leaders and experts within a broad range of computing communities. Proposed themes/issues should address important and timely topics that will be of broad interest to *Computer's* readership. Special issues are an essential feature of *Computer*, as they deliver compelling research insights and perspectives on new and established technologies and computing strategies.

Please send us your high-quality proposals for the 2026 - 2027 editorial calendar. Of particular interest are proposals centered on:

- 3D printing
- Robotics
- LLMs
- AI safety
- Dis/Misinformation
- Legacy software
- Microelectronics



Proposal Guidelines Are Available at:

computer.org/csdl/magazine/co/write-for-us/15911

DEPARTMENT: INTERNET OF THINGS

Convenience at a Cost: The Urgent Need for Data Privacy Standards

Syed Rizvi , The Pennsylvania State University – Altoona

Anthony Demeri , The Pennsylvania State University

Mohammad R. Rizvi, Accenture

The widespread use of Internet of Things (IoT) devices offers life-changing advantages, but this convenience comes at a steep cost. Since IoT networks lack a universal privacy standard, users are susceptible to unprecedented risks.

Since the dawn of time, humanity has marveled at the thunderous roars and static shocks of electricity. From Benjamin Franklin's oft-debated kite experiment to the engineering of household lighting, this life-changing energy has served as the spark of innovation for centuries. Today, with the capacity to open doors with our voices, scry across the globe in the palm of our hands, and merchandise without lifting a finger, we are privy to services our ancestors would be unable to distinguish from magic.

THE INTERNET OF THINGS

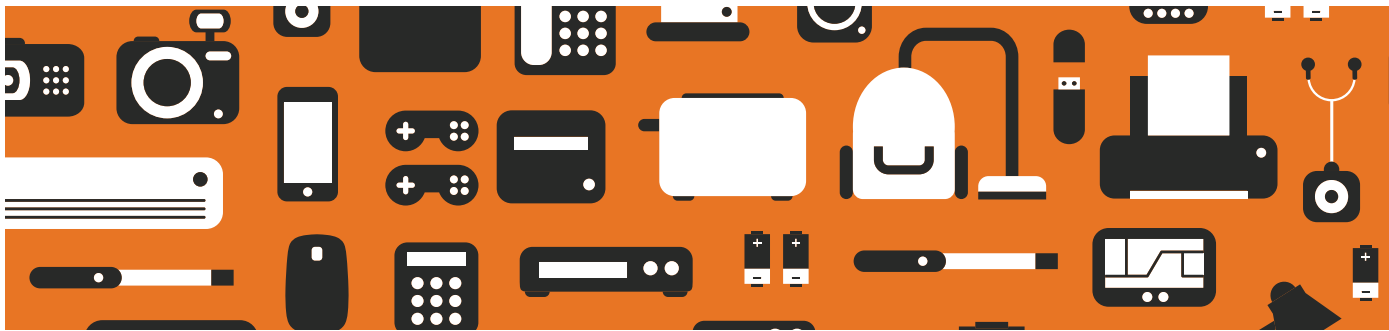
These conveniences, brought to us through the vast network of the Internet of Things (IoT), have revolutionized the way we interact with the world around us. Unfortunately, in offering their services, IoT devices often collect a significant amount of sensitive data, including medical records, financial transactions, and the personal data of their end users (for example, location, preferences, habits, etc.). Fundamentally, a user's convenience comes at the cost of their privacy, for which there are no refunds. By exploring real-world examples of IoT applications and associated data privacy challenges, we will highlight the urgent need for universal data privacy standards across IoT networks. Ultimately, the development of data privacy standards specifically designed for IoT will encourage the

widespread adoption of these devices, benefiting businesses, providing conveniences, and enhancing overall public trust—this time, without the cost of privacy.

AN IoT WORLD

In general, IoT devices are used across various industries and domains, offering numerous benefits to businesses, the economy, and end users. To offer these services with a high degree of accuracy and responsiveness, IoT devices must often collect a vast range of high-resolution data in real time. From health care and finance to smart homes and industrial automation, these devices enable seamless data collection, real-time monitoring, and automation, ultimately improving efficiency and enhancing user experiences.

In health care, IoT devices continuously monitor patient data, transmitting and/or evaluating real-time health data to medical professionals. For example, wearable devices like smartwatches can detect irregular heart rhythms and send immediate alerts to users or doctors, potentially preventing life-threatening conditions such as cardiac arrest. A well-documented case is the Apple Watch's electrocardiogram feature, which has detected atrial fibrillation in users and prompted them to seek medical attention before severe complications occurred.¹ A particularly crucial application of IoT devices occurred during the COVID-19 pandemic, where smartwatch oxygen sensors played a vital role in detecting early signs of respiratory distress, significantly improving treatment outcomes.² Similarly, a closed-loop insulin delivery system connects with



a continuous glucose monitor to adjust insulin delivery in real time, preventing potentially fatal hypoglycemic events for diabetic patients.³ Together, these innovations demonstrate how IoT devices not only collect data but actively contribute to saving human lives through high-frequency real-time monitoring.

In financial services, IoT devices also operate in real time, monitoring transactions and detecting fraudulent activities as they occur. Many banks and payment services use IoT-enabled biometric authentication—such as fingerprint or facial recognition—captured instantly during each transaction to enhance security and user convenience.⁴ At the consumer level, smart home devices like security cameras, smart thermostats, and voice assistants collect and transmit data continuously or at set intervals, depending on their function.⁵ Security cameras record video footage in real time, motion sensors detect activity the moment it occurs, and smart assistants listen for voice commands at all times.

Clearly, IoT devices offer significant utility, but the highly sensitive private data collected by these devices (medical records, biometric data, financial transactions, real-time location tracking, and more) raise concerns regarding how personal information is stored, used, and protected. Before developing an adequate data privacy model, however, we must first understand why such private data are necessarily collected, from a technical perspective.

OBJECTIVES BEHIND DATA COLLECTION

Naturally, given the personalized solutions offered by IoT devices, it is logical that personal data are collected. Less obvious, however, is the extensive quantity, quality, and retention of said data, often collected

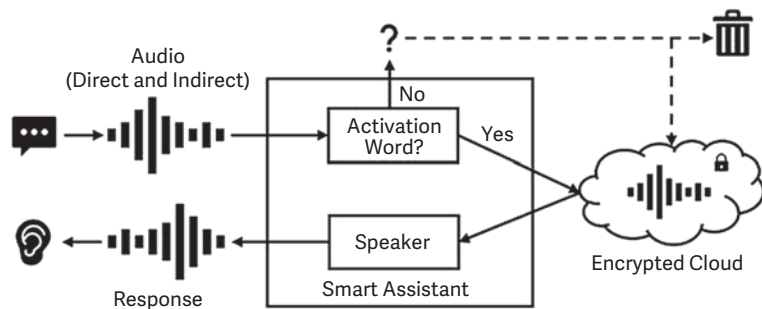


FIGURE 1. A high-level overview of a traditional smart assistant.

without the user’s full awareness. For example, if we consider the case of smart assistants (for example, Apple’s Siri or Amazon’s Alexa), we might imagine that for these assistants to actually be convenient, they need to be able to promptly respond to user requests, which can happen at any time. Of course, without knowing when a user will need their services, these

ULTIMATELY, THE DEVELOPMENT OF DATA PRIVACY STANDARDS SPECIFICALLY DESIGNED FOR IoT WILL ENCOURAGE THE WIDESPREAD ADOPTION OF THESE DEVICES.

devices are left with only one option: to *always* listen.⁶ While individual products may have differences in optimizations and voice processing techniques, fundamentally, they must all process input (that is, your voice) and act accordingly (that is, answer your prompt), as shown in Figure 1.

Notably, although IoT devices are often marketed as upholding data security practices, these practices have been shown to be insufficient, even allowing the remote compromise of IoT home security cameras.⁷ Regardless, even if such security flaws could be patched, there would still be considerable risk to

user data—critically, security does *not* imply privacy. For example, by observing only the encrypted network traffic of IoT home security cameras (considered “secure” by the data security model), researchers were able to “infer the working state of cameras and even the daily activity routine” around said camera,⁸ demonstrating how data security alone is insufficient.

INTEGRATING PRIVACY AND SECURITY BY DESIGN

Unlike data *security*—which ensures that data remain protected from breaches and unauthorized modifications—data *privacy* focuses on how data are collected, used, and shared—a particularly pressing concern within IoT networks. For IoT systems, the Confidentiality, Integrity, and Availability (CIA) model is considered to be a widely accepted data security standard. Simply put: if data are confidential, they cannot be *read* by unauthorized parties; if data have integrity, they cannot be *changed* by unauthorized parties; if data have availability, they are accessible by authorized parties. Data privacy, however, implies a degree of anonymity, such that data cannot even be *associated* with a specific individual or entity—a critical absence from the CIA model. Unfortunately, the nuanced differences between security and privacy are more than theoretical, requiring practical urgency across all domains.

Practical privacy

Thankfully, present advances in IoT technology have paved the way for incredible improvements to health care through the integration of robotic-assisted surgery, task automation, patient monitoring, and effective remote evaluation and treatment.⁹ Unfortunately, these advances, often provided through the collection, storage, and sharing of sensitive patient data, are seldom accompanied by sufficient privacy practices, resulting in real-world consequences.

Notably, in November 2023, the Epic Software and MyChart Patient Portals were compromised, causing ambulance diversion and suspending access to patient services.¹⁰ Shortly thereafter, in February 2024, the cyberthreat actor known as BlackCat compromised the sensitive patient data and services of Change Healthcare in perhaps the most severe cyberattack in health care to date, resulting in more than US\$1 billion in damages and critical impacts to patient care.¹⁰ A

month later, in March 2024, the telehealth organization Cerebral was compromised by malign actors, exposing “dates of birth, contact information, self-assessment response, treatment details, and other clinical information”¹⁰ for more than 3 million people.

Clearly, without stringent privacy protections, there is an enormous risk for individuals’ sensitive information, personal behaviors, activities, and even private conversations to be misused or exposed. Ultimately, a growing reliance on IoT technology, combined with the sensitive nature of IoT-collected data, makes privacy one of the most pressing concerns of the digital age.

But we can’t put all of the blame on IoT device manufacturers; after all, the market is highly competitive, and accounting for all of the complex considerations of privacy can be costly. To truly solve this problem, manufacturers need to be told explicitly what is (or is not) required to consider a given IoT device as privacy adequate. In industry, this is accomplished with a *standard* (like the CIA model), which may be either regulatory (lawfully required) or guiding (recommended). Once a universal data privacy model exists for IoT devices, consumers can make educated decisions on which devices meet their individual needs. While such a globally recognized model does not yet exist for IoT devices, there is hope.

Privacy by design

Today, researchers and organizations are proactively seeking to bridge the widening gap between technological advancement, privacy, and security. The National Institute of Standards and Technology (NIST) is one such organization actively devoted to improving privacy and security, providing guidelines for manufacturers, including “NIST Special Publication (SP) 800-53, which recommends changing default passwords as best practice for implementations of networking equipment, including routers and switches and the servers and services running on them,”¹¹ encouraging manufacturers to play a more active role in promoting such industry best practices. While changing your password is one particularly straightforward means of *reducing* risk, poor privacy and security design principles result in *natively* vulnerable system components, particularly among software publishers.¹¹ In the future, manufacturers and organizations must take care to explicitly incorporate data

TABLE 1. Privacy considerations and their respective design implications.^{13,14,15,16}

Right to access	This refers to the right to view any and all data related to one's person, obtaining complete and structured copies upon demand. In practice, this necessitates organized technical and managerial design principles. As an example, a user of a home assistant should be able to download, access, and replay all collected data in a convenient format.
Right to correct	This refers to the right to correct any inaccuracies within one's collected data, requiring systems capable of rapid retroactive corrections.
Right to delete	This refers to the right to have one's personal data deleted on demand, necessitating explicit organizational and technological procedures for data erasure. Today, for example, social media companies provide users with the option to delete their respective profiles and all associated data; this must be extended to IoT-integrated devices, as well.
Right to object (opt-out)	This refers to the right to opt out of data collection. In tandem with "default privacy," personal data must be, by default, not collected, instead requiring individuals to opt in to such services (with easy opt-out capability). For example, today, when individuals in the EU visit certain websites, they are immediately provided with the option to opt out of data collection; this must be extended to IoT devices on a granular basis, allowing users to configure the type, quantity, and quality of any data collection.
Minimal collection	Systems and organizations must collect only data explicitly required to provide a user-requested service. In other words, IoT devices must be designed with transparent data collection processes. For example, if a user opts to wear a personal heart rate monitor to track their real-time heart rate, they should reasonably expect said device to not collect location data without explicit authorization.
Breach notification	Organizations must be required to notify individuals impacted by data breaches of any scale within an explicit timeline (for example, within three days of detection).
Limited retention	By default, user data must be explicitly deleted in its entirety after a fixed period of time (for example, one month) unless a user provides explicit retention consent. This principle requires systems to be designed with built-in data deletion mechanisms, requiring explicit override consent. Today, for example, some messaging services provide the option to automatically delete message history after a certain point in time.
Antidiscrimination	A user cannot be discriminated against based on their privacy needs. This principle requires regulatory intervention, giving individuals the legal right to action over antiprivacy discrimination.
Default privacy	By default, organizations and systems must presume that data collection and storage consent is not given and, when given, is ephemeral, necessarily expiring after a fixed time period (for example, one month) unless reacquired.

privacy considerations into the design and engineering of IoT-based systems and devices.

LOOKING AHEAD

Fortunately, recognizing an urgent lapse in privacy standards, several governing bodies have started implementing specialized laws, such as the European Union's (EU's) General Data Privacy Regulation (GDPR), to address how consumers' personal data are collected, used, and shared by businesses. For example, California's Consumer Privacy Act (CCPA) applies to all personal data collected by businesses that meet a minimum threshold. Virginia's Consumer Data Protection Act (VCDPA) and the Colorado Privacy Act (CPA) similarly share many common elements with the CCPA. Other states, such as Utah, Connecticut, and Maryland, have also enacted laws to offer privacy

protection to consumers. Additionally, programs like the Wearables and Medical IoT Interoperability and Intelligence program (WAMIII) from IEEE are making an incredible effort to secure privacy within the Clinical IoT, successfully publishing the *IEEE/UL Standard for Clinical Internet of Things (IoT) Data and Device Interoperability with TIPSSTrust, Identity, Privacy, Protection, Safety, and Security* in September 2024.¹² In Tables 1 and 2, we define, compare, and contrast the differences among leading regulations.

On a high level, existing laws and standards share a number of common themes; however, as IoT adoption continues to expand, the need for a general and globally recognized IoT data privacy standard becomes increasingly more urgent. This standard must establish internationally recognized principles such as transparency, user consent, data minimization, ethical data

TABLE 2. A comparison of privacy-focused acts and programs.

	GDPR ¹³	CCPA ¹⁴	VCDPA ¹⁵	CPA ¹⁶
Scope	All EU citizens and organizations	California big business	Virginia big business	Colorado big business
Right to access	✓	✓	✓	✓
Right to correct	✓	Added in 2023	45-day response time	✓
Right to delete	✓	✓	✓	✓
Right to object (opt-out)	Lacking opt-in requirement	Lacking opt-in requirement	✓	Lacking opt-in requirement
Minimal collection	✓	✓	✓	✓
Breach notification	Within 72 h	No explicit timeline	No explicit timeline	Within 30 days
Limited retention	Moderate (lacking autodelete)	Moderate (lacking autodelete)	Moderate (lacking autodelete)	Moderate (lacking autodelete)
Antidiscrimination	×	✓	✓	✓
Default privacy	Yes—but not ephemeral consent	Yes—but not ephemeral consent	Requires opt-in consent	Yes—but not ephemeral consent

handling, and other considerations listed in Table 1, holding groups of all sizes and scopes responsible for failing to adhere to respective regulations.

Ultimately, without an adequate data privacy model, co-engineered with modern security principles, users will remain exposed to sensitive data violations, and IoT ecosystems will continue to operate without clear accountability for protecting personal information. Moving forward, we call on industry, academia, and regulatory bodies to rapidly promote the explicit intertwining of both security and privacy standards in future systems, particularly IoT devices. With such models, the future can be convenient—without the cost. 🌐

REFERENCES

- S. Shahid et al., "Diagnostic accuracy of apple watch electrocardiogram for atrial fibrillation: A systematic review and meta-analysis," *JACC Adv.*, vol. 4, no. 2, Feb. 2025, Art. no. 101538, doi: 10.1016/j.jacadv.2024.101538.
- X. Ding et al., "Wearable sensing and telehealth technology with potential applications in the coronavirus pandemic," *IEEE Rev. Biomed. Eng.*, vol. 14, pp. 48–70, 2021, doi: 10.1109/RBME.2020.2992838.
- A. M. Joshi, P. Jain, and S. P. Mohanty, "iGLU 3.0: A secure noninvasive glucometer and automatic insulin delivery system in IoMT," *IEEE Trans. Consum. Electron.*, vol. 68, no. 1, pp. 14–22, Feb. 2022, doi: 10.1109/TCE.2022.3145055.
- H. U. Khan, M. Z. Malik, S. Nazir, and F. Khan, "Utilizing bio metric system for enhancing cyber security in banking sector: A systematic analysis," *IEEE Access*, vol. 11, pp. 80,181–80,198, 2023, doi: 10.1109/ACCESS.2023.3298824.
- D. Buil-Gil et al., "The digital harms of smart home devices: A systematic literature review," *Comput. Human Behav.*, vol. 145, Aug. 2023, Art. no. 107770, doi: 10.1016/j.chb.2023.107770.
- A. Waseem, "Is Google listening to you? Yes, and here's how to stop it," *Private Internet Access*, Feb. 01, 2025. Accessed: Apr. 20, 2025. [Online]. Available: <https://www.privateinternetaccess.com/blog/google-chrome-listening-in-to-your-room-shows-the-importance-of-privacy-defense-in-depth/>
- O. Almazrouei, P. Magalingam, M. Kamrul Hasan, M. Almehrzi, and A. Alshamsi, "Penetration testing for IoT security: The case study of a wireless IP security CAM," in *Proc IEEE 2nd Int. Conf. AI Cybersecurity (ICAIC)*, Piscataway, NJ, USA: IEEE Press, 2023, pp. 1–5, doi: 10.1109/ICAIC57335.2023.10044176.
- J. Li, Z. Li, G. Tyson, and G. Xie, "Characterising usage patterns and privacy risks of a home security camera service," *IEEE Trans. Mobile Comput.*, vol. 21, no. 7, pp. 2344–2357, Jul. 2022, doi: 10.1109/TMC.2020.3039787.
- J. F. DeFranco and M. J. Metro, "Internet of telemedicine," *Computer*, vol. 55, no. 4, pp. 56–59, Apr. 2022, doi: 10.1109/MC.2022.3143625.

10. "HealthSec USA summit 2024 annual report," HealthSec Cyber Security for Healthcare, Boston, MA, USA, May 2024. Accessed: May 20, 2025. [Online]. Available: <https://healthsec.cs4ca.com/wp-content/uploads/HealthSec-2024-Annual-Report.pdf>
11. J. F. DeFranco and B. Maley, "Closing the security agility gap," *Computer*, vol. 55, no. 8, pp. 100–102, Aug. 2022, doi: 10.1109/MC.2022.3169400.
12. *IEEE/UL Standard for Clinical Internet of Things (IoT) Data and Device Interoperability with TIPPSS–Trust, Identity, Privacy, Protection, Safety, and Security*, IEEE/UL 2933–2024, Sep. 30, 2024. [Online]. Available: <https://standards.ieee.org/ieee/2933/7592/>
13. "Regulation (EU) 2016/679 of the European Parliament and of the council," European Union, Brussels, Belgium, Apr. 05, 2016. Accessed: May 31, 2025. [Online]. Available: <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679>
14. "California Consumer Privacy Act (CCPA)." California Office of the Attorney General, 2025. Accessed: May 31, 2025. [Online]. Available: <https://www.oag.ca.gov/privacy/ccpa>
15. "Chapter 53. Consumer Data Protection Act." Virginia State Law Portal, 2025. Accessed: May 31, 2025. [Online]. Available: <https://law.lis.virginia.gov/vacodefull/title59.1/chapter53/>
16. "Colorado Privacy Act (CPA)." Colorado Office of the Attorney General, 2025. Accessed: May 31, 2025. [Online]. Available: <https://coag.gov/resources/colorado-privacy-act/>

SYED RIZVI is a professor of Information Sciences and Technology at The Pennsylvania State University, Altoona, PA 16601 USA. Contact him at srizvi@psu.edu.

ANTHONY DEMERI is an active software engineer and doctoral student at The Pennsylvania State University, University Park, PA 16803 USA. Contact him at akd6327@psu.edu.

MOHAMMAD R. RIZVI is a senior manager at Accenture, Dallas, TX 75039 USA. Contact him at m.rizvi@accenture.com.

IT Professional

CALL FOR ARTICLES

IT Professional seeks original submissions on technology solutions for the enterprise. Topics include

- Emerging Technologies
- Cloud Computing
- Web 2.0 And Services
- Cybersecurity
- Mobile Computing
- Green IT
- RFID
- Social Software
- Data Management And Mining
- Systems Integration
- Communication Networks
- Datacenter Operations
- IT Asset Management
- Health Information Technology

We welcome articles accompanied by web-based demos.

For more information, see our author guidelines at bit.ly/4faGdch

computer.org/itpro



The Rise of Agentic AI in Finance: Opportunities, Risks, and Human-Centric Integration

Nir Kshetri , University of North Carolina at Greensboro, Greensboro, NC, 27412, USA

This paper looks at the growing role of agentic artificial intelligence (AI) in financial services, focusing on its capabilities to enhance operational efficiency, automate complex tasks, and deliver personalized customer experiences. It also analyzes the challenges related to data quality, regulatory compliance, security risks, and the importance of human oversight to ensure trustworthy and effective AI integration.

Although agentic artificial intelligence (AI) is still in its “trial phase,” its core capabilities—such as language processing, predictive analytics, and reasoning—are already well developed.¹ Despite the challenges of responsible integration in finance, these hurdles are likely to be overcome, driving rapid adoption. Wolters Kluwer N.V., a Dutch information services company, surveyed 392 finance leaders in May 2025, finding that 6% had adopted agentic AI and 38% plan to do so within the next 12 months—bringing expected utilization to 44% by 2026.²

Top firms in accounting are also deploying agentic AI to streamline operations—from standard compliance tasks and routine tax work to advanced financial forecasting and decision-making. In March 2025, Deloitte and Ernst & Young (EY) announced the use of autonomous AI systems to assist clients with tasks like document uploads and financial statement analysis. Deloitte’s Zora AI offers deployable agents that “perceive, reason, and act,” while EY plans to integrate AI agents across over 30 million tax processes. EY adopted autonomous AI tools in its tax practice to leverage its vast in-house tax data and address diverse global client needs, while Deloitte’s early agents, including Zora AI—already in use by Hewlett Packard Enterprise—support finance teams with tasks such as expense management, financial analysis, and scenario

modeling. PricewaterhouseCoopers (PwC) has also developed similar agentic capabilities (see “Exhibit 1: Agentic AI in Action: How PwC Is Transforming Tax Operations with Intelligent Agents”).³

Building on this momentum, specialized startups are also advancing agentic AI tailored to industry-specific needs. In February 2025, Zurich-based Unique secured €28.7 million in Series A funding to advance its agentic AI platform for financial institutions. The platform supports 25 ready-to-use applications or customizable agents, ensuring regulatory compliance and data security while integrating seamlessly into back and middle-office operations.⁴

These trends underscore a rapid shift toward AI-driven processes in the finance sector, and some analysts contend that agentic AI will surpass the Internet era in economic impact⁵; Citigroup even predicts it could usher in a “do it for me” Economy that outpaces previous digital revolutions.⁶ For instance, Accenture was reported to be helping global banks, including Westpac, develop AI agents for customer service, compliance, and sales. Applications include handling calls, writing credit memos, and assessing risks. According to Accenture, banks are seeing 20%–30% productivity gains.⁷

Agentic AI is transforming traditional financial processes by enabling more adaptive and intelligent decision-making. For example, unlike conventional credit scoring models that rely on static data, agentic AI allows banks and fintechs to perform continuous credit assessments using real-time transaction data, behavioral trends, and economic indicators—resulting in faster approvals and more precise risk management.⁸

1520-9202 © 2025 IEEE. All rights reserved, including rights for text and data mining, and training of artificial intelligence and similar technologies.

Digital Object Identifier 10.1109/MITP.2025.3585227

Date of current version 15 August 2025.

This article examines the transformative impact of agentic AI on the financial services industry, highlighting its applications in automation, decision-making, and customer engagement. It also addresses the key challenges and considerations—such as regulatory compliance, data quality, and security—that organizations must navigate for successful and responsible adoption.

TRANSFORMING FINANCE WITH AGENTIC AI: FROM BACK-OFFICE EFFICIENCY TO PROACTIVE CUSTOMER CARE

Agentic AI in banking has two main focuses: internal operational efficiency, such as automating routine tasks and building predictive models for trading and risk management, and customer-facing services, like automated help desks and personalized investment advice. Both aim to improve efficiency and reduce costs.⁶

Regarding internal operational efficiency, agentic AI is often likened to an unlimited pool of intelligent interns, capable of automating tedious back-end tasks and streamlining workflows.⁹ Agentic AI is already boosting productivity and cutting human error by automating repetitive tasks such as data entry, compliance checks, and transaction processing—freeing employees to focus on more strategic work. Key operational gains include real-time compliance monitoring for regulatory adherence, intelligent fraud detection, dynamic risk assessment, and streamlined “know your customer” (KYC)/anti-money laundering (AML) workflows, all powered by autonomous decision-making capabilities.¹⁰ As fraudsters increasingly exploit AI, agentic AI offers a powerful defense by detecting anomalies, tracing threats across systems, and autonomously initiating preventive actions such as freezing transactions or triggering audits. In compliance, these systems can interpret evolving regulations and autonomously adjust internal protocols, transforming compliance from a reactive burden into a proactive, streamlined service.¹¹ Additionally, agentic AI offers significant cost savings by automating tasks traditionally outsourced to contractors or third parties.⁵

As to customer-facing services, agentic AI moves beyond traditional chatbot functionality by recognizing intent, remembering prior interactions, and delivering proactive, customized financial guidance.¹¹ In this evolving landscape, users will have personal AI agents assisting in product selection and transaction execution.⁵ Agentic AI’s most immediate impact lies in real-time needs assessment—such as identifying wellness concerns during a transaction—and autonomously negotiating tailored solutions like payment deferrals or debt restructuring. Rather than offering

static product options, it dynamically crafts customized financial packages by integrating diverse banking products and services.¹⁰ Traditional banking is reactive, with customers initiating interactions. Agentic AI reverses this by anticipating needs and offering proactive solutions. For example, a farmer could receive pre-approved credit before a drought, insurance recommendations before the growing season, or investment advice based on crop yields. This enhances financial inclusion, literacy, and personalized support.¹²

THE EXPANDING ROLE OF AGENTIC AI SYSTEMS IN FINANCE

Table 1 illustrates how agentic AI is revolutionizing the financial services industry through diverse and impactful applications. From optimizing automated trading and enhancing fraud detection to delivering personalized financial advice and streamlining regulatory compliance, these use cases demonstrate the broad potential of agentic AI to improve efficiency, accuracy, and customer experience across the sector. Real-world examples highlight how leading firms are leveraging these technologies to drive measurable gains and innovation.

In automated trading, these systems use advanced algorithms to execute trades with remarkable speed and precision, optimizing returns by analyzing market trends and timing.¹³ By processing massive volumes of real-time data—like price movements, news sentiment, and geopolitical developments—agentic AI bots empower hedge funds to act faster and smarter. Singapore-based investment firm QuantEdge, for instance, leverages these tools to make hyper-timely trades, producing dynamic portfolios that consistently outperform manual approaches by as much as 15% per year.¹⁴

For fraud detection, they identify suspicious activities by scrutinizing transaction patterns and spotting anomalies in real time, thereby enhancing security.¹³ Mastercard leverages AI agents to strengthen fraud defense and support consumer dispute resolution by providing end-to-end security. These agents enable robust authentication using on-device biometrics and assist in explaining unfamiliar or questionable transactions to users.¹⁵ By analyzing transaction data across multiple cards and merchants, Mastercard’s system doubles the detection of compromised cards, decreases false positives by as much as 200%, and speeds up spotting at-risk merchants threefold.¹⁴

In terms of personalized financial advice, agentic AI can analyze vast customer data, providing tailored investment strategies that align with individual risk profiles.¹³ Era, an AI-powered personal finance platform, raised \$6.2 million in February 2025 to advance its agentic AI capabilities for real-time, personalized

EXHIBIT 1: AGENTIC AI IN ACTION: HOW PwC IS TRANSFORMING TAX OPERATIONS WITH INTELLIGENT AGENTS

Across PricewaterhouseCoopers's (PwC's) tax operations, teams of AI agents are actively deployed, enhancing efficiency by automating interpretation and mapping tasks and marking a progression from experimental stages to routine implementation.^{S1}

PwC has joined forces with Microsoft and Salesforce to create AI agents for clients and recently extended this effort through a new alliance with Oracle³. PwC trained its agents to accurately identify and handle various types of expenses and disclosures using historical data, then continuously enhanced their performance through ongoing feedback. Much like managing human teams, PwC reinforced correct decisions, addressed edge cases, and regularly retrained the agents with real-world scenarios to steadily improve accuracy and confidence.^{S1}

Processing Schedule K-1s—a federal tax form that reports each individual's share of income, deductions, and other items from pass-through entities like partnerships, S corporations, and trusts^{S2}—has historically been labor-intensive. Tax teams had to interpret varying formats, map line items (e.g., deductible expenses), and often collaborate to resolve unclear footnotes. Agentic AI has transformed this workflow by interpreting footnotes, determining appropriate treatments, and assigning confidence levels to its decisions. Ambiguities are flagged for human review, and when users make corrections, the agent learns from them. This human-in-the-loop (HITL) approach reduces manual effort while improving consistency and efficiency in tax reporting.^{S1} The hybrid AI-human tax teams have reduced K-1 production from nearly two weeks to just one day.^{S3}

Like human team members, each AI agent leverages a specialized set of tools, coordinated by PwC's Agent OS to efficiently handle tasks such as data mapping and document extraction.^{S1} Introduced in March 2025, Agent OS functions as an enterprise AI command center that streamlines and scales multiagent workflows up to ten times faster than conventional approaches.

By offering a unified, scalable platform for building and integrating AI agents across diverse systems, it accelerates the move from pilot projects to full enterprise adoption.^{S4} PwC manages over 800 custom GPTs and 250+ AI agents, including in tax, and developed Agent OS to better organize, orchestrate, and monitor them efficiently.^{S1}

AI agents are expected to deliver similar efficiency gains across tax functions—from data processing to deal due diligence—by understanding context, taking proactive actions, learning from errors, and integrating seamlessly with enterprise platforms.^{S3} However, AI use remains limited in areas like policy interpretation and high-touch consulting due to the complex reasoning and client interaction that still require a human touch.^{S1}

REFERENCES

- S1. A. Zaki, "How PwC's tax team is using agentic AI," *CFO*, Jun. 6, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.cfo.com/news/how-pwcs-dom-megna-tax-team-is-using-agentic-ai-agent-os-cpa-cfo-/750062/>
- S2. "Schedule K-1 federal tax form: What is it and who is it for?" *Investopedia*, Sep. 9, 2024. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.investopedia.com/terms/s/schedule-k-1.asp>
- S3. D. Megna, "AI agents: Transforming the tax experience," *PwC Blog*, Mar. 27, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.pwc.com/us/en/services/tax/library/tax-ai-agents.html>
- S4. M. Wood, "PwC launches AI agent operating system to revolutionize AI workflows for enterprises," *PwC Newsroom*, Mar. 27, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.pwc.com/us/en/about-us/newsroom/press-releases/pwc-launches-ai-agent-operating-system-enterprises.html>

wealth management. Unlike traditional apps, Era autonomously manages finances—optimizing balances, transferring funds, and offering tailored insights—using automation enhanced through a partnership with Cerebras for institutional-grade intelligence.¹⁶

Furthermore, these systems power portfolio management tools that continuously monitor and adjust investment strategies to ensure maximum returns in dynamic market conditions.¹³ Acting as intelligent financial companions, these systems retain long-term

TABLE 1. How agentic AI is revolutionizing financial services: Use cases and examples.

Capability	Description	Example and Impact
Automated trading	Uses advanced algorithms to execute trades quickly and precisely by analyzing market trends and timing.	QuantEdge uses agentic AI for hyper-timely trades, producing portfolios that outperform manual strategies by ~15% annually.
Fraud detection and security	Identifies suspicious transaction patterns and anomalies in real time to enhance security.	Mastercard employs AI agents for fraud defense, enabling biometric authentication and dispute resolution. Detection of compromised cards doubled, false positives reduced by 200%, and risk identification sped up 3×.
Personalized financial advice and wealth management	Analyzes extensive customer data to provide tailored investment strategies and automated financial management.	Era is advancing agentic AI for real-time personalized wealth care using institutional-grade automation.
Autonomous portfolio management and risk modeling	Continuously monitors and adjusts investment strategies based on evolving goals and market conditions, synthesizing macroeconomic data for risk mitigation.	Agentic AI autonomously adapts asset allocations, proactively rebalances portfolios, and issues hedge recommendations for optimal returns.
Regulatory compliance automation	Streamlines compliance processes like AML and KYC by automating investigations and decision support.	Oracle Financial Services integrated agentic AI into its Investigation Hub Cloud Service for efficient financial crime investigations and compliance.

goals, anticipate evolving needs, and deliver personalized advice seamlessly across digital channels and at scale. Agentic AI enables autonomous portfolio management by continuously adjusting asset allocations in real time based on evolving investor goals and risk tolerance, using adaptive, long-horizon decision-making rather than static rules. It also transforms risk modeling by synthesizing macroeconomic and geopolitical data into live scenarios, proactively mitigating risk through portfolio rebalancing or hedge recommendations.¹¹

Finally, regulatory compliance is streamlined as agentic AI automates processes like AML and KYC checks, ensuring financial institutions adhere to regulatory standards efficiently.¹³ In March 2025, Oracle Financial Services announced the integration of agentic AI and workflow automation into its Investigation Hub Cloud Service to streamline financial crime investigations. These AI agents analyze alert data and sanctions matches to generate summaries that assist compliance teams in decision-making.¹⁷

BALANCING PROMISE AND PITFALLS: THE REAL-WORLD CHALLENGES OF AGENTIC AI DEPLOYMENT

While agentic AI holds great promise for enhancing decision-making and risk management across industries, its practical deployment faces significant challenges. These include high operational costs, complex infrastructure demands, opaque return on investments (ROI), strict

regulatory requirements, data security risks, and the critical need for high-quality data management.¹⁸

Building and maintaining agentic AI requires a costly, specialized talent ecosystem—beyond data scientists—including machine learning operations, security, and ethics experts. Infrastructure demands like GPU clusters and compliance auditing add further expense, often rivaling the AI's deployment cost. For many mid-sized firms, the total cost of ownership quickly outweighs the marginal gains, turning quick wins into long-term burdens.¹⁹

It is worth noting that the ROI of agentic AI can be deceptive: Cloud providers often offer free credits to jumpstart adoption, masking the true costs of running large-scale systems. Once these credits expire, enterprises may face steep expenses and vendor lock-in, undermining long-term return on investment.¹⁹

The challenges of implementing agentic systems are further highlighted when we consider regulatory compliance, which poses risks for autonomous AI systems, especially when performing tasks such as approving loans without sufficient transparency. Failure to meet oversight standards can result in fines, penalties.²⁰ Under the EU AI Act, systems used for credit scoring—such as those that may deny individuals access to loans—are classified as high-risk due to their potential impact on fundamental rights. Before deployment, these systems must comply with strict requirements, including robust risk mitigation frameworks, the use of high-quality, nondiscriminatory data, and appropriate human oversight mechanisms.²¹

Agentic AI systems in finance also pose heightened data security risks due to their broad access to sensitive information and system controls. A breach could expose private financial data, leading to identity theft, fraud, and targeted cyberattacks—making them attractive targets for malicious actors.²²

Finally progress in agentic AI hinges on a critical but often overlooked factor: data quality. Without clean, well-governed data, even advanced AI initiatives face serious limitations, making prior investment in data infrastructure essential for meaningful innovation.²³ Agentic AI systems rely on high-quality data, but many organizations suffer from “data debt” due to legacy systems, silos, and outdated records. When flawed data feeds autonomous agents, mistakes scale rapidly—leading to erroneous payments, poor customer experiences, or operational failures at speed and scale.¹⁹

Before reaping the full benefits of agentic AI, financial institutions must lay a solid foundation that ensures responsible deployment, organizational alignment, and risk-aware implementation. To successfully adopt agentic AI in finance, organizations should first assess readiness by identifying high-impact areas, such as manual bottlenecks or rising compliance risks. Starting with targeted use cases—like fraud detection or invoice processing—allows for measured impact and scalable implementation. Investing in AI training ensures finance teams can effectively collaborate with these tools. Finally, aligning AI solutions with regulatory and data privacy requirements is essential to maintain compliance and trust.²⁴

Embracing a human-in-the-loop (HITL) model allows AI agents to propose actions while humans retain final decision-making, balancing automation with oversight. This approach builds trust, improves models through real-world feedback, and ensures human expertise enhances rather than competes with AI¹⁹—exemplified by PwC’s tax operations, where agents handling Schedule K-1 footnotes learn from human corrections to boost future accuracy.

To address privacy challenges posed by agentic AI, organizations should cultivate a privacy-first culture by educating employees on risks and safe practices. They must also establish clear policies for auditing AI systems and hold teams accountable. Finally, deploying strong security controls—such as role-based access, encryption, and real-time monitoring—can help prevent breaches and unauthorized access.²²

CONCLUSION

The advent of agentic AI is set to revolutionize the financial services sector, enhancing user experiences and operational efficiency. Agentic AI is rapidly moving

beyond experimental phases to become a transformative force, promising significant improvements in efficiency, decision-making, and customer engagement. Leading accounting firms such as Deloitte, EY, and PwC are already deploying AI agents to automate complex tax processes, financial analysis, and client services, demonstrating measurable gains in productivity and accuracy. With capabilities ranging from autonomous trading to personalized wealth management and fraud detection, agentic AI is poised to reshape how financial institutions operate, while also enhancing financial inclusion through proactive, tailored solutions.

However, realizing the full potential of agentic AI requires navigating significant challenges, including the high costs of talent and infrastructure, data quality issues, and stringent regulatory compliance demands—especially in high-risk applications like loan approvals. The integration of HITL models and robust privacy and security frameworks, as exemplified by PwC’s approach, will be essential to balance automation with oversight, build trust, and ensure ethical deployment. As financial institutions strategically assess readiness and focus on scalable, compliant use cases, agentic AI stands to usher in a new era of innovation, efficiency, and personalized service across the industry. 🌟

REFERENCES

1. A. Kenney, “What agentic AI will mean for finance,” *Financial Manage.*, Jun. 21, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.fm-magazine.com/issues/2025/jun/what-agentic-ai-will-mean-for-finance/>
2. “Finance leaders plan sixfold increase in agentic AI: Wolters Kluwer,” *Int. Accounting Bull.*, May 29, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.internationalaccountingbulletin.com/news/wolters-kluwer-agentic-ai-adoption-in-finance/?cf-view>
3. K. Sibayan, “Big four now using agentic AI to boost staff productivity,” *The Trusted Prof.*, Mar. 24, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.nysscpa.org/news/publications/the-trusted-professional/article/big-four-now-using-agentic-ai-to-boost-staff-productivity-032425>
4. P. Pawar, “Agentic AI market to reach USD 196.6 billion by 2034,” *London Daily News*, Jun. 25, 2025. [Online]. Available: <https://www.londondaily.news/agentic-ai-market-to-reach-usd-196-6-billion-by-2034/>
5. “Agentic AI & the ‘do it for me’ economy,” *Citigroup*, Jan. 17, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.citigroup.com/global/insights/agentic-ai>

6. L. Browning, "How agentic AI will transform banking (and banks)," *The Financial Brand*, Feb. 5, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://thefinancialbrand.com/news/artificial-intelligence-banking/agentic-ai-the-next-big-innovation-for-banks-186428>
7. J. Eyers, "Westpac works with Accenture to deploy AI agents," *Financial Rev.*, Feb. 24, 2025. Accessed: Jun. 23, 2025. <https://www.afr.com/companies/financial-services/westpac-works-with-accenture-to-deploy-ai-agents-20250127-p5l7hh>
8. B. Zhang and K. Garvey, "Agentic AI will be the real banking disruptor," *The Banker*, Feb. 25, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.thebanker.com/content/886b880f-fc01-458d-81a5-4ad4c27815da>
9. T. Pathe, "Agentic AI and the future of fintech and banking automation," *FinTech Futures*, Feb. 10, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.fintechfutures.com/2025/02/agentic-ai-and-the-future-of-fintech-and-banking-automation/>
10. J. Marous, "How agentic AI will disrupt retail banking's future," *The Financial Brand*, Jun. 9, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://thefinancialbrand.com/news/artificial-intelligence-banking/how-agentic-ai-will-disrupt-retail-bankings-future-189846>
11. D. Nagasubramanian, "Why agentic AI could be more disruptive for finance than the internet," *Forbes*, Jun. 12, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.forbes.com/councils/forbestechcouncil/2025/06/12/why-agentic-ai-could-be-more-disruptive-for-finance-than-the-internet/>
12. M. Nkosi, "Exclusive: Absa's CTO discusses the future of AI & the rise of Agentic AI," *IT News Africa*, Feb. 3, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.itnewsafrica.com/2025/02/exclusive-absas-cto-discusses-the-future-of-ai-the-rise-of-agentic-ai/>
13. V. Dugar, "The dawn of agentic AI systems: Revolutionizing financial services," *Forbes Bus. Develop. Council*, Feb. 6, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.forbes.com/councils/forbesbusinessdevelopmentcouncil/2025/02/06/the-dawn-of-agentic-ai-systems-revolutionizing-financial-services/>
14. S. Aliakbar, "How agentic AI is transforming financial services: Four game-changing use cases," *SP Edge Blog*, Apr. 3, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://blog.sp-edge.com/post/how-agentic-ai-is-transforming-financial-services-four-game-changing-use-cases>
15. "Mastercard unveils Agent Pay, pioneering agentic payments technology to power commerce in the age of AI," *Mastercard Newsroom*, Apr. 29, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.mastercard.com/news/press/2025/april/mastercard-unveils-agent-pay-pioneering-agentic-payments-technology-to-power-commerce-in-the-age-of-ai/>
16. A. Tardif, "FundingEra raises \$6.2M to advance agentic AI for personalized wealth-care," *Unite.AI*, Feb. 12, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.unite.ai/era-raises-6-2m-to-pioneer-ai-powered-wealth-care-with-agentic-ai/>
17. "Oracle brings AI agents to the fight against financial crime," *Oracle Newsroom*, Mar. 13, 2025. [Online]. Available: <https://www.oracle.com/news/announcement/oracle-brings-ai-agents-to-the-fight-against-financial-crime-2025-03-13/>
18. "Agentic AI: Financial services disruptor?" *Financial Times*, Apr. 22, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://agenticai.live.ft.com/>
19. R. Goswami, "Agentic AI in the enterprise: Cost of autonomy vs. ROI reality check," *CTO Mag.*, Jun. 7, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://ctomagazine.com/agentic-ai-in-enterprise/>
20. A. Chopra, "Agentic AI in banking: The future and the challenges," *Forbes Technol. Council*, May 5, 2025. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.forbes.com/councils/forbestechcouncil/2025/05/05/agentic-ai-in-banking-the-future-and-the-challenges/>
21. "AI act," European Commission. Accessed: Jun. 23, 2025. [Online]. Available: <https://digital-strategy.ec.europa.eu/en/policies/regulatory-framework-ai>
22. E. Kron, "Five privacy concerns around agentic AI," *SC Media*, Jun. 18, 2024. Accessed: Jun. 23, 2025. [Online]. Available: <https://www.scworld.com/perspective/five-privacy-concerns-around-agentic-ai>
23. N. Alexander, "How agentic AI and regulated stablecoins are reshaping banking," *Bobsguide*, Jun. 18, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.bobsguide.com/agentic-ai-and-regulated-stablecoins-are-reshaping-banking/>
24. R. Gupta, "Automation to intelligence: Agentic AI and the finance industry," *Forbes*, Jun. 4, 2025. Accessed: Jun. 22, 2025. [Online]. Available: <https://www.forbes.com/councils/forbesfinancecouncil/2025/06/04/automation-to-intelligence-agentic-ai-and-the-finance-industry/>

NIR KSHETRI is *IT Professional's* associate editor in chief and IT Economics editor, as well as a professor at the Bryan School of Business and Economics, University of North Carolina at Greensboro, Greensboro, NC, 27412, USA. Contact him at nbkshetr@uncg.edu

DEPARTMENT: MICRO ECONOMICS

Private Returns on Technology Adoption

Shane Greenstein , Harvard Business School, Boston, MA, 02163, USA

From a firm’s perspective, the emergence of a new technology wave is a new opportunity to generate a financial return. The question is precisely how. That topic remains as salient today in the era of artificial intelligence as it was in the era when firms first encountered smartphones, the commercial Internet, and personal computers.

Before we fully embrace this new era, let me suggest that we review lessons from the most recent era of technology adoption. In particular, let’s focus on consumer computer technologies (CCTs)—the mix of the mobile ecosystem and the widely used Internet, enhanced by Web 2.0.

Some research partners and I recently analyzed returns from co-invention in CCTs. By co-invention, we mean the invention of new applications by firms that utilize CCTs in their business. That investment typically involves developing business processes and practices to complement the adoption of CCTs and support the introduction of new services and business models.

This column provides a somewhat brief overview of high-level points from the study. Some of the implications arising from these points might seem obvious, but it was a surprise to us that all the consequences would come from the same framework. For more information, please see the reference at the end of the column.

FRAMEWORK

Consider a straightforward model of the two types of projects undertaken by firms, where one utilizes the new technology as an intermediate input and involves *incremental* co-invention. The other also uses the technology as an intermediate input, but requires something more ambitious and innovative for the firm. Referred to by many names in everyday speech, for brevity, we will label it as *novel* co-invention.

Incremental co-invention in CCTs was exceedingly common. That is because, when faced with a low-cost

expansion of existing services, most firms choose to invest in it. For example, should a firm build an app? This was a comparatively straightforward economic decision, especially when the benefits were directly measurable in terms of web traffic that enabled a marketing funnel into sales leads or ad revenue. Accordingly, most firms did build apps.

Novel co-invention differed. Success was difficult, costly, and rare, but you knew a successful novel co-invention effort when you saw it. Many users found the service compelling. It established a new category of uses or altered leadership in an existing end-market. Think of Netflix in its early years, or the elaborate efforts required by your pharmacy to send you text reminders to refill a prescription.

For a firm involved in novel co-invention, adopting CCTs was usually less straightforward, often because the payoffs were several years out from the initiation of investment. For example, when Netflix first began transitioning from mail-order DVD rentals to streaming,

SOME OF THE IMPLICATIONS ARISING FROM THESE POINTS MIGHT SEEM OBVIOUS, BUT IT WAS A SURPRISE TO US THAT ALL THE CONSEQUENCES WOULD COME FROM THE SAME FRAMEWORK.

they had to confront considerable uncertainty about how to do so at scale. It took them more than seven years to find a predictable and reliable process.

In an entrepreneurial setting, the economic constraints bind in different ways. For example, should an entrepreneur build an app for a use case that no other firm covers? Consider Snap, who introduced ephemeral messaging with Snapchat. This novel social experience drove rapid adoption among young users, creating network effects that fueled explosive early growth before competitors caught up with similar features. Indeed, Snap reached a place that few entrepreneurs ever attain. Alas for them, their time at the top was brief.

0272-1732 © 2025 IEEE. All rights reserved, including rights for text and data mining, and training of artificial intelligence and similar technologies.

Digital Object Identifier 10.1109/MM.2025.3615287

Date of current version 24 December 2025.

Here is my point: Novel co-invention was risky in the CCT era. Private returns were uncertain. Yet, as we all know, many firms did build these. Some succeeded, and most failed, sometimes spectacularly.

IMPLICATIONS

What does all that imply? For one, comparisons between investors in incremental co-invention should reveal a specific pattern of returns, one proportional to the preexisting customer revenue.

Here is what that means. If one observes a single industry as it evolves into an online market over time, and if every firm made incremental changes, then everybody would have received a 2%, 3%, or 4% return, and nobody should miraculously gain a 50% return on their website.

We tested that prediction. We rolled up our sleeves and collected information from the Internet Archive about the online experiences of terrestrial radio and newspapers, in each case in 2013. We compared that with their listenership in 1993, 20 years earlier. In radio, we were able to track the experiences of over 2000 stations. In newspapers, we examined approximately 100.

The framework predicts that radio stations with more listeners in 1993 would have a larger online audience in 2013 than stations that started with a smaller number of listeners in 1993. That is what we see. In other words, radio appears to be an industry without any novel successful co-invention as of 2013. We make a similar prediction for the audience size for newspapers, but see a slightly different outcome. Newspapers with larger circulations gained a more than proportionate online readership compared to those with smaller circulations.

COMPARE INCREMENTAL AND NOVEL

Here is another implication: Incremental co-invention should be as widespread as the businesses that use it. In contrast, the returns on novel co-invention should be skewed with only a few big winners and many losers. Novel co-inventions should also be tied to the labor market for technical talent or industry specialties.

Testing this contrast is no small task. It requires an examination of all industries. We rolled up our sleeves again and found a way to collect data on private returns for 2400 leading online properties in 2013, roughly split between smartphone apps and websites. Through extensive (and sometimes tedious) work, we classified all of these into incremental and novel co-invention efforts that led to their production. Some of

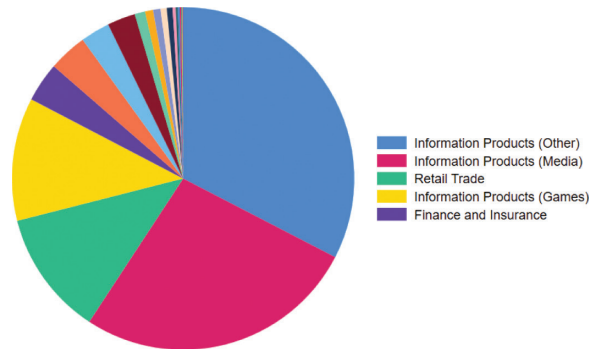


FIGURE 1. Incremental co-invention distribution across industries.

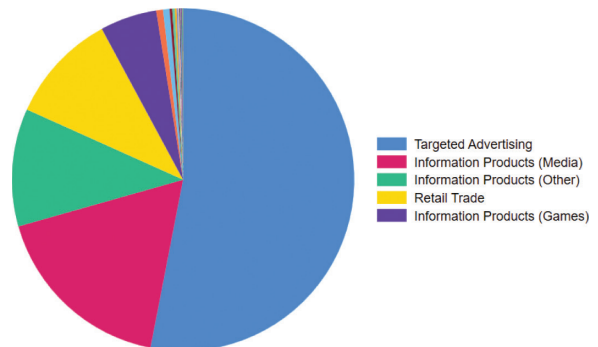


FIGURE 2. Novel co-invention distribution across industries.

these industries are so new that we also had to build a new product classification to capture the main categories in the data.

Figures 1 and 2 illustrate the distribution of value creation across firms. These should differ between incremental and novel innovations. Indeed, we find that it does. Incremental co-invention grows from firm assets that support existing business, while novel co-invention grows from an entirely different origin. Local labor markets with appropriate technical and commercial domain knowledge supported investment in novel uses of CCTs.

We then examined a related, open question. There is no theoretical reason why the total value created by novel co-inventors should be bigger or smaller than that made by incremental co-inventors. A small number of big hits could exceed the total of a large volume of small contributions, or vice versa. Only data can settle the question.

As it turned out, the aggregate private value created by incremental co-invention is smaller than that made by novel co-invention, and by a considerable amount. Empirically, we estimate that incremental co-invention accounts for approximately 6% of the total observed value among the top 2400 firms. If every

bit of value we do not observe in the long tail of small firms is attributed to some incremental investment, it still accounts for less than 18% of the total value.

In other words, novel co-invention drove the vast majority of new private value creation in CCTs, and in a small set of industries where novel co-invention thrived; like it or not, Google, Facebook, Netflix, Amazon, and other top hundred CCT-based companies owned a significant portion of the profits from adapting CCTs.

To be clear, sometimes these companies came up with the novel co-inventions themselves (e.g., the Google Search Engine and Ad network), and sometimes they bought parts of it (e.g., Google bought YouTube after they tried and failed to develop something similar), absorbing it in the big company. Sometimes these received enormous investments, and sometimes these purchases went nowhere.

That result begs questions about the origins of why novel co-inventions sometimes originate in entrepreneurial firms and sometimes within leading firms. That is a bigger topic for another day. Today's column focuses on the lessons from characterizing these fact patterns.

LIKE IT OR NOT, GOOGLE, FACEBOOK, NETFLIX, AMAZON, AND OTHER TOP HUNDRED CCT-BASED COMPANIES OWNED A SIGNIFICANT PORTION OF THE PROFITS FROM ADAPTING CCTS.

CONCENTRATION

There is one more subtle prediction. If all co-inventions were incremental, there would be little impact on competition and market structure. Similarly, the industrial and geographical distribution of economic activity would change little, as they are based on the existing choices of firms and their existing assets. However, novel co-inventions should not be widespread. It should concentrate on a few industries and in a few places.

As expected, four related media and entertainment industries, as well as retail trade, collectively account for 97% of the value from novel co-invention, with targeted advertising being the most significant. Most of the manufacturing and several additional significant sectors, such as health care, are absent.

To be sure, incremental co-invention is more widespread industrially. However, it is also less impactful in creating value.

Novel co-invention should also concentrate geographically. To check this, we undertook another forensic investigation into the location of the co-inventing firm's headquarters for each product.

Again, as expected, we find that the location of incremental and novel co-invention differs. Because incremental co-invention builds on existing business, it is geographically distributed widely, following the existing location patterns of the co-inventing firms. In contrast, the geographic outcome of novel co-invention differs sharply, with four prominent regions emerging for novel innovation: the greater San Francisco region, New York, Seattle, and Los Angeles.

San Francisco and Seattle are no surprise, as they have long been centers for technical talent, and many firms there zig-zagged their way into media and entertainment markets. New York and Los Angeles firms had the opportunity to join them during this technology wave because their cities possessed the entertainment industry's human capital, and many firms there had to zig-zag their way into technology.

CONCLUSION

We are not the first commentators to notice the bifurcation of results in the commercial Internet. However, most commentators emphasize the winner-take-all competitive dynamics that result in a small number of massive providers.

We took another approach. You might call it reductionist for eschewing elaborate explanations about why things happened, but simple has advantages too. It boils results down to either/or, incremental/novel, which makes it easy to observe. If that trend continues into the new AI era, expect to be able to track project outcomes this way. 🤖

REFERENCE

1. T. Bresnahan, S. Greenstein, and P.-L. Yin, "New economic forces behind the value distribution of innovation," Nat. Bur. of Econ. Res., Cambridge, MA, USA, Working Paper 34090, 2025. [Online]. Available: <https://www.nber.org/papers/w34090>

SHANE GREENSTEIN is a professor with Harvard Business School, Boston, MA, 02163, USA. Contact him at sgreenstein@hbs.edu.

Get Published in the *IEEE Open Journal of the Computer Society*

Get more citations by publishing with the *IEEE Open Journal of the Computer Society*

Your research on computing and informational technology will benefit from 5 million unique monthly users of the *IEEE Xplore*[®] Digital Library. Plus, this journal is fully open and compliant with funder mandates, including Plan S.



Submit your paper today!
Visit www.computer.org/oj to learn more.



IEEE SECURITY & PRIVACY

IEEE Security & Privacy is a bimonthly magazine communicating advances in security, privacy, and dependability from the top thinkers in the field.

computer.org/security

Find the latest research and practical articles alongside case studies, surveys, tutorials, columns, and in-depth interviews. Topics include:

- Internet, software, hardware, and systems security
- Legal and ethical issues and privacy concerns
- Privacy-enhancing technologies
- Data analytics for security and privacy
- Usable security
- Integrated security design methods
- Security of critical infrastructures
- Pedagogical and curricular issues in security education
- Security issues in wireless and mobile networks
- Real-world cryptography
- Emerging technologies, operational resilience, and edge computing
- Cybercrime and forensics, and much more



Join the IEEE Computer Society for
subscription discounts today!

computer.org/product/csdl-full-access





stay connected.

Join our online community! Follow us to stay connected wherever you are:



| @ComputerSociety



| facebook.com/IEEEComputerSociety



| IEEE Computer Society



| youtube.com/IEEEComputerSociety



| instagram.com/ieee_computer_society

Conference Calendar

IEEE Computer Society conferences are valuable forums for learning on broad and dynamically shifting topics from within the computing profession. With over 200 conferences featuring leading experts and thought leaders, we have an event that is right for you. Questions? Contact conferences@computer.org.

MAY

4 May

- HOST (IEEE Int'l Symposium on Hardware Oriented Security and Trust), Washington, DC, USA
- MOST (IEEE Int'l Conf. on Mobility, Operations, Services and Technologies), Detroit, USA

8 May

- BigDataSecurity (IEEE Conf. on Big Data Security on Cloud), New York City, USA
- CAI (IEEE Int'l Conf. on Artificial Intelligence), Granada, Spain
- HPSC (IEEE Int'l Conf. on High Performance and Smart Computing), New York City, USA
- IDS (IEEE Int'l Conf. on Intelligent Data and Security), New York City, USA
- SmartCloud (IEEE Int'l Conf. on Smart Cloud), New York City, USA

11 May

- SenSys (ACM/IEEE Int'l Conf. on Embedded Artificial Intelligence and Sensing Systems), Saint Malo, France

12 May

- RTAS (IEEE Real-Time and Embedded Technology and Applications Symposium), Saint Malo, France

13 May

- FCCM (IEEE Annual Int'l Symposium on Field-Programmable Custom Computing Machines), Atlanta, USA

15 May

- ICES (Int'l Conf. on Energy Storage), Shenyang, China

18 May

- CCGrid (IEEE Int'l Symposium on Cluster, Cloud and Internet Computing), Sydney, Australia
- ICFEC (IEEE Int'l Conf. on Fog and Edge Computing), Sydney, Australia
- ICST (IEEE Int'l Conf. on Software Testing, Verification and Validation), Daejeon, Korea
- S&P (IEEE Symposium on Security and Privacy), San Francisco, USA

19 May

- ICDE (IEEE Int'l Conf. on Data Eng.), Hong Kong, China
- ISMVL (IEEE Int'l Symposium on Multiple-Valued Logic), Sendai, Japan

25 May

- FG (IEEE Int'l Conf. on Automatic Face and Gesture Recognition), Kyoto, Japan
- IPDPS (IEEE Int'l Parallel and Distributed Processing Symposium), New Orleans, USA

JUNE

1 June

- ICHI (IEEE Int'l Conf. on Healthcare Informatics), Minneapolis, USA

3 June

- CBMS (IEEE Int'l Symposium on Computer-Based Medical Systems), Limassol, Cyprus
- CVPR (IEEE/CVF Conf. on Computer Vision and Pattern Recognition), Denver, USA

10 June

- SVCC (Silicon Valley Cybersecurity Conf.), San Jose, USA

16 June

- WoWMoM (IEEE Int'l Symposium on a World of Wireless, Mobile and Multimedia Networks), Bologna, Italy

22 June

- DCOSS-IoT (Int'l Conf. on Distributed Computing in Smart Systems and the Internet of Things), Reykjavik, Iceland
- DSN (Annual IEEE/IFIP Int'l Conf. on Dependable Systems and Networks), Charlotte, USA
- ICDCS (IEEE Int'l Conf. on Distributed Computing Systems), Seoul, Korea
- ICSA (IEEE Int'l Conf. on Software Architecture), Amsterdam, Netherlands



- SmartComp (IEEE Int'l Conf. on Smart Computing), Messina, Italy

23 June

- ISCC (IEEE Symposium on Computers and Communications), Vila Moura, Portugal

25 June

- IEEE Cloud Summit, Washington, DC, USA

27 June

- ISCA (ACM/IEEE Annual Int'l Symposium on Computer Architecture), Raleigh, USA

28 June

- ARITH (IEEE Symposium on Computer Arithmetic), Fulda, Germany

29 June

- MDM (IEEE Int'l Conf. on Mobile Data Management), Athens, Greece

JULY

1 July

- IOLTS (IEEE Int'l Symposium on On-Line Testing and Robust System Design), Polignano a Mare, Italy

6 July

- EuroS&P (IEEE European Symposium on Security and Privacy), Lisbon, Portugal
- ICALT (IEEE Int'l Conf. on Advanced Learning Technologies), Hung Yen, Vietnam
- ICME (IEEE Int'l Conf. on Multimedia and Expo), Bangkok, Thailand

7 July

- COMPSAC (IEEE Annual Computers, Software, and Applications Conf.), Madrid, Spain
- ISVLSI (IEEE Computer Society Annual Symposium on VLSI), Kolkata, India

8 July

- ICEDEG (Int'l Conf. on eDemocracy & eGovernment), Lisbon, Portugal

13 July

- ICCP (IEEE Int'l Conf. on Computational Photography), Princeton, USA
- SERVICES (IEEE World Congress on Services), Sydney, Australia

20 July

- FiCloud (Int'l Conf. on Future Internet of Things and Cloud), Granada, Spain

26 July

- CSF (IEEE Computer Security Foundations Symposium), Lisboa, Portugal

27 July

- CISOSE (IEEE Int'l Congress on Intelligent and Service-Oriented Systems Eng.), Fukuoka, Japan

31 July

- IRI (IEEE Int'l Conf. on Information Reuse and Integration for Data Science), Seattle, USA

AUGUST

3 August

- SCC (IEEE Int'l Space Computing Conf.), Pasadena, USA

- SMC-IT (IEEE Int'l Conf. on Space Mission Challenges for Information Technology), Pasadena, USA

9 August

- MIPR (IEEE Int'l Conf. on Multimedia Information Processing and Retrieval), Bangkok, Thailand

11 August

- RTCSA (IEEE Int'l Conf. on Embedded and Real-Time Computing Systems and Applications), Qingdao, China

15 August

- PCDS (IEEE Int'l Conf. on Privacy Computing and Data Security), Jeju, Korea

17 August

- RE (IEEE Int'l Requirements Eng. Conf.), Montreal, Canada

19 August

- HOTI (IEEE Symposium on High-Performance Interconnects), Virtual

21 August

- SmartIoT (IEEE Int'l Conf. on Smart Internet of Things), Shenyang, China

Learn more about
IEEE Computer
Society conferences
computer.org/conferences

Career Accelerating Opportunities

Explore new options—upload your resume today

careers.computer.org



Changes in the marketplace shift demands for vital skills and talent. The **IEEE Computer Society Career Center** is a valuable resource tool to keep job seekers up to date on the dynamic career opportunities offered by employers.

Take advantage of these special resources for job seekers:



JOB ALERTS



TEMPLATES



WEBINARS



CAREER
ADVICE



RESUMES VIEWED
BY TOP EMPLOYERS

No matter what your career level, the IEEE Computer Society Career Center keeps you connected to workplace trends and exciting career prospects.



IEEE
COMPUTER
SOCIETY

