

Building the Infinite Brain

COMP 590/790

Raghavendra Pradyumna Pothukuchi



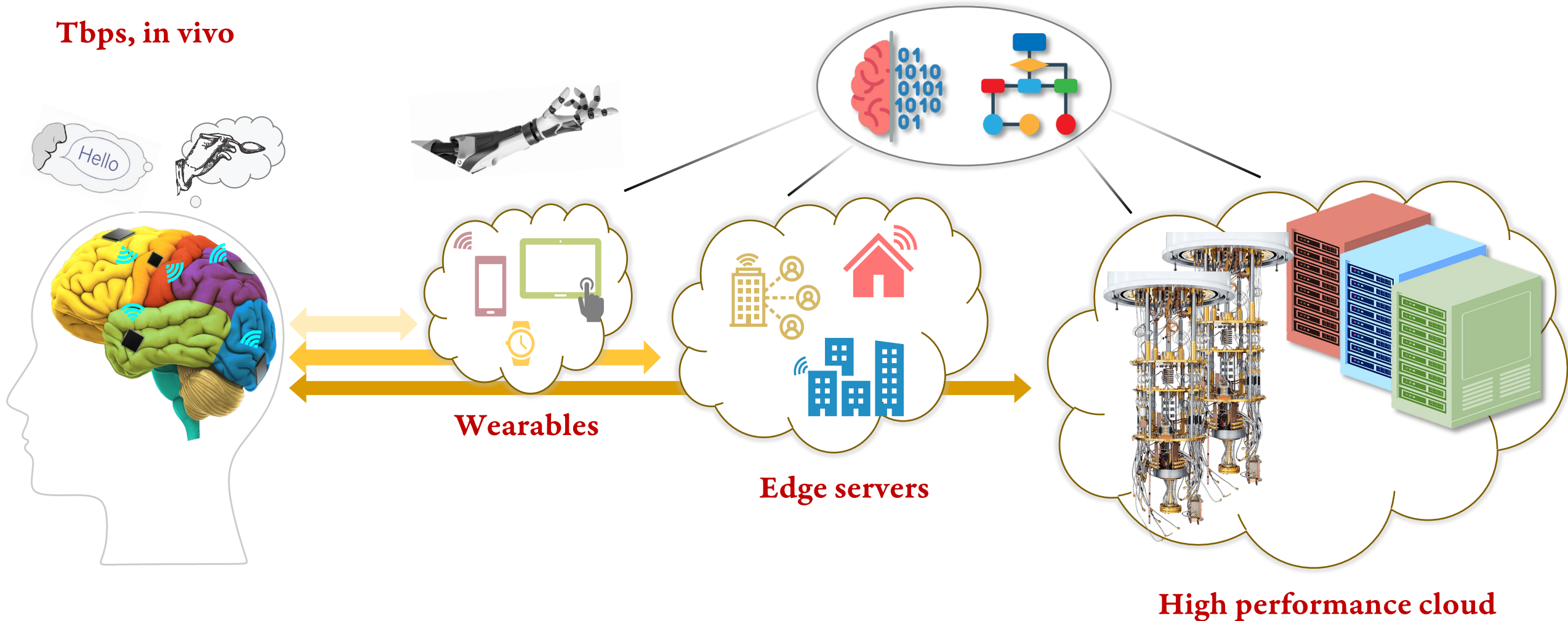
THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

✉ raghav@cs.unc.edu

Building the Infinite Brain

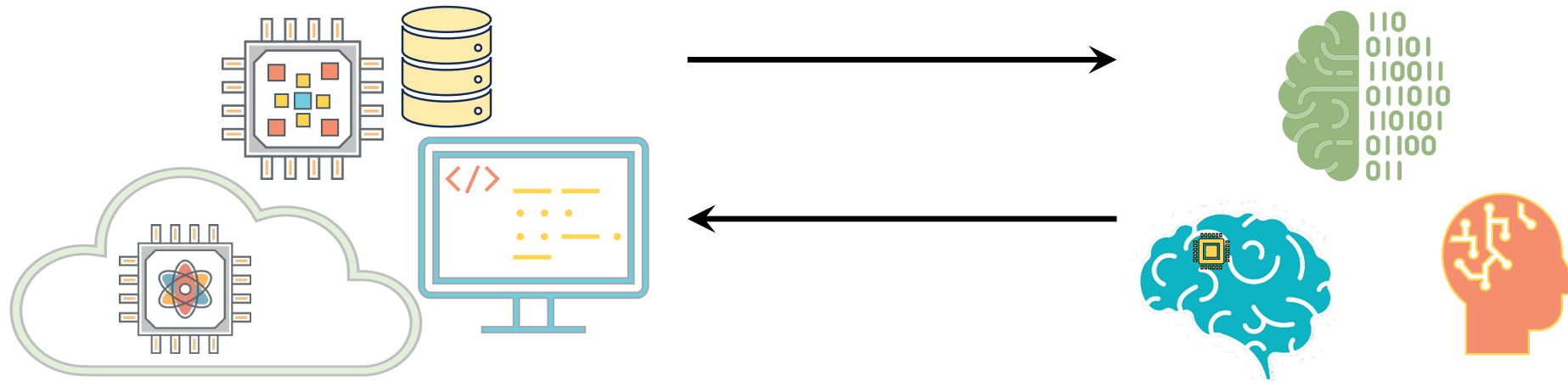
Real-time cognitive frameworks

Tbps, in vivo



Central Theme: A Virtuous Cycle of Innovation

Architectures and systems for the brain sciences



Constraints and inspiration from the brain sciences for architectures

BCI Processing Challenge: Efficiency **vs** Flexibility

Efficiency

High throughput

100-1000 Mbps

Real-time performance

Tens of ms

Safety

$\leq 1\text{ }^{\circ}\text{C}$

$\leq \text{few mW}$



Flexibility

Complex processing

Signal processing, Machine learning

Customization

Personalized treatment,
target multiple conditions, adaptation
to the brain, support evolving methods

Summary: Computer Architecture

What is computer architecture?

Historically, the ISA; but now encompasses organization

What are the goals?

Mnemonic: Simple Timely Efficient Adaptable Dependable Yummy

How to estimate impact of fixing bottlenecks?

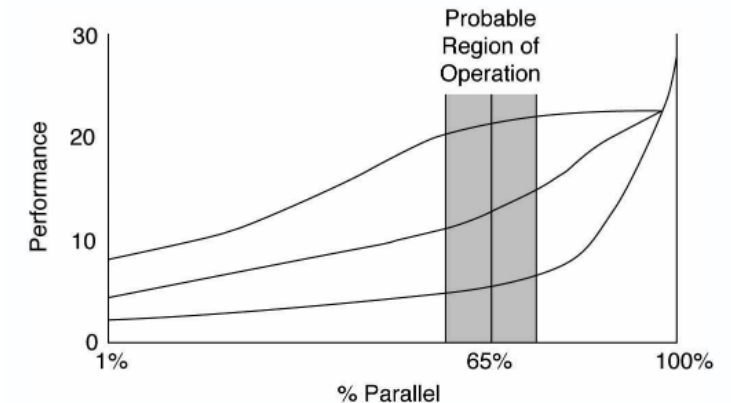
Amdahl's law

How to fix bottlenecks?

Algorithms, adding a fast path

How to improve efficiency?

Technology, approximation, locality, regroup



Pipelining Summary

What is pipelining?

Splitting work into many components executed independently, and passed in a chain

Why pipelining?

Increases efficiency via parallelism

What are the challenges in pipelining?

Hazards: structural, data, control

How to fix hazards?

Concurrency (structural), eager (forwarding data), eager or speculative (branch delays, prediction)

What systems design principles does pipelining touch?

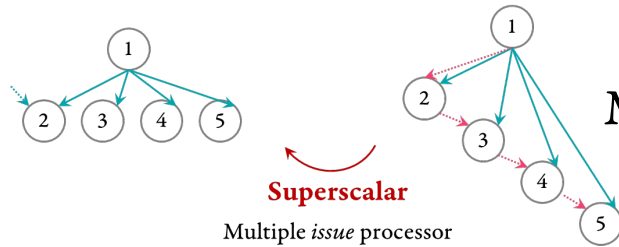
Tradeoff simplicity for efficiency, leveraging parallelism through fragmenting and regrouping



Superscalar and Dynamic Scheduling Summary

What are superscalar processors?

Multiple issue processors (send multiple instructions to execute)



Why dynamic scheduling?

To be resource efficient in exploiting instruction level parallelism

What are hardware methods for dynamic scheduling?

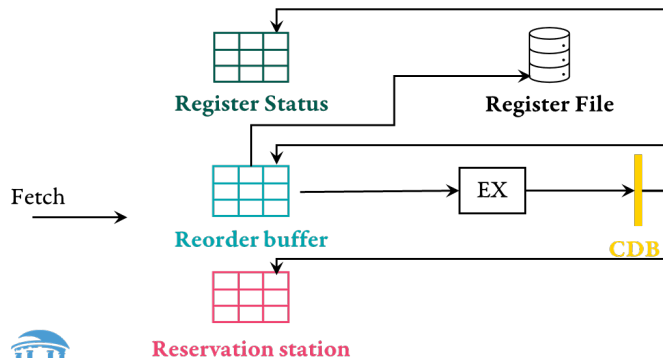
Scoreboarding, Tomasulo's algorithm, speculation

What systems principles does dynamic scheduling involve?

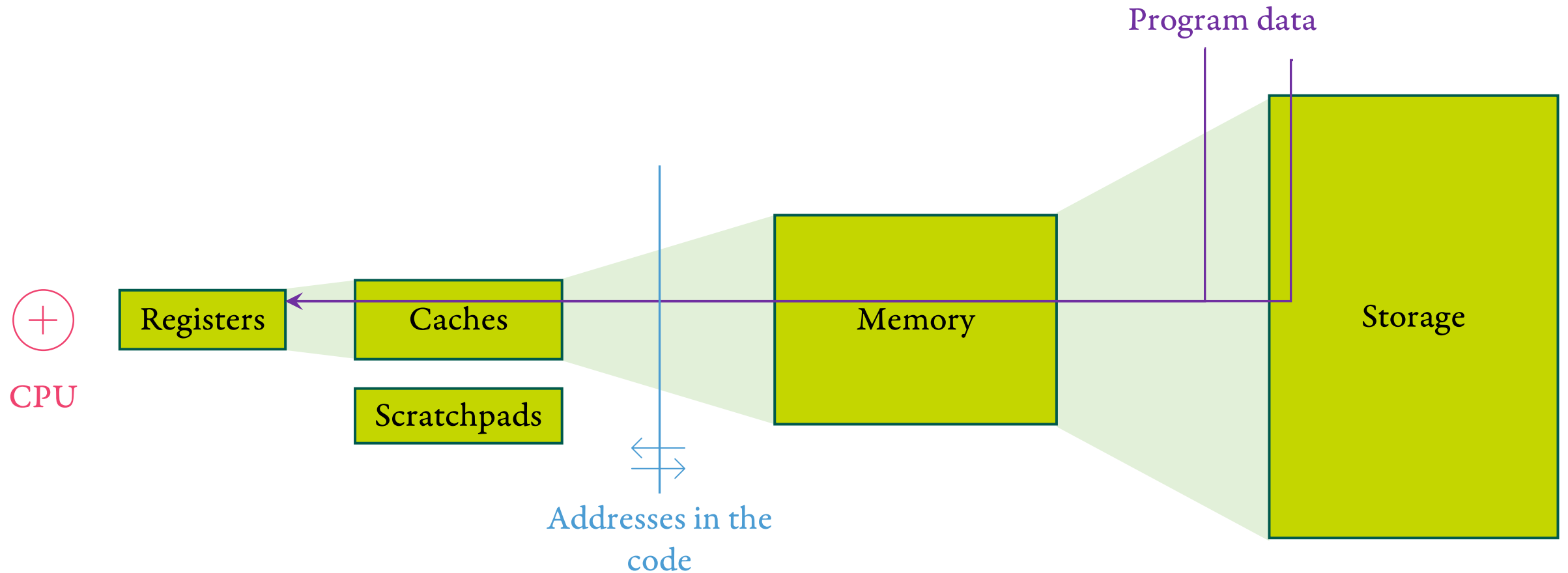
Eager, speculative, and concurrency

What are some other methods to exploit ILP?

Software pipelining, unrolling, which target static ILP



Memory Hierarchy



Cache Summary

What is caching?

A technique to minimize the impact of long memory latencies

What are the basic cache parameters?

Associativity (placement), block size, capacity, write through/back, write-allocate or no

What are the types of cache misses?

Compulsory, capacity, conflict

What are some ways in which cache performance can be improved?

Non-blocking, banking, software or hardware prefetching etc.

How to choose the right memory hierarchy design?

Tailor it to the access patterns and layout: general purpose to domain specific



Memory Summary

What are some goals in designing the memory subsystem?

Capacity, program abstraction, protection, and sharing

What is virtual memory?

Abstraction to deliver the goals

What is a bottleneck in virtual memory systems?

Address translation

How can addressing bottlenecks be tackled?

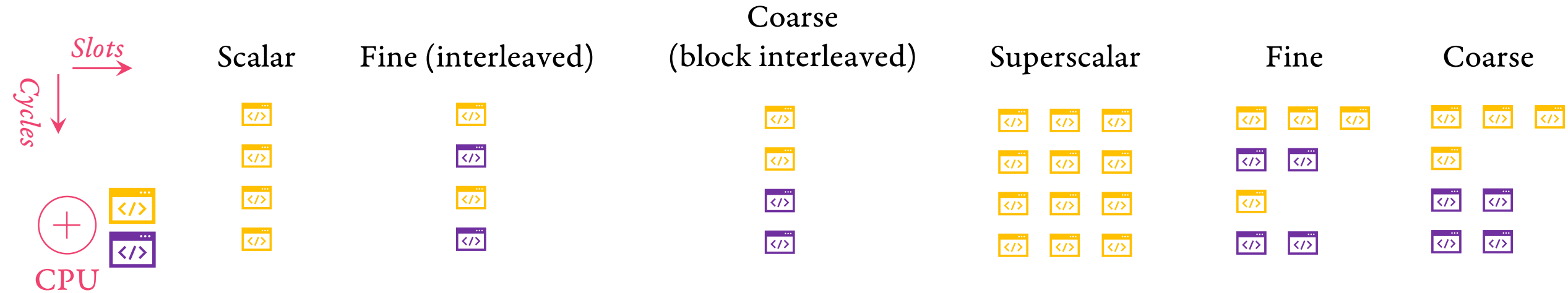
Through fast paths: TLBs

How to choose the right memory hierarchy design?

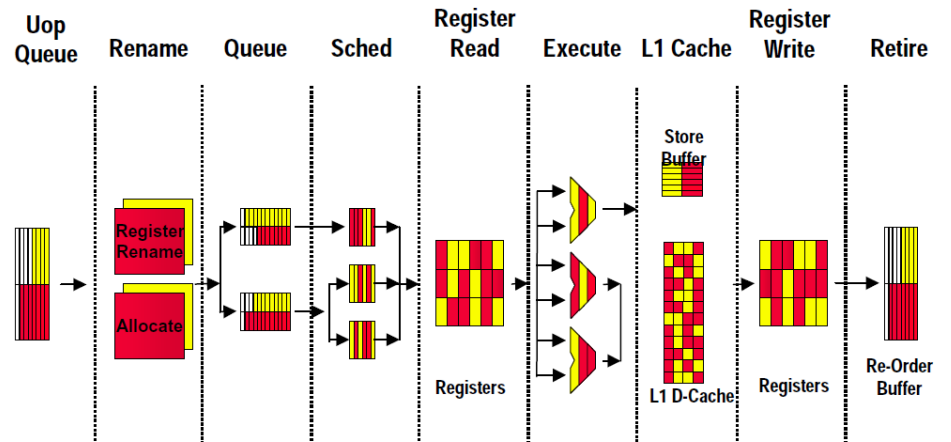
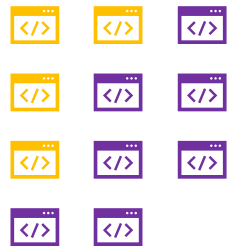
Tailor it to the access patterns and layout: general purpose to domain specific



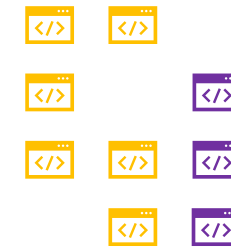
Multithreading



Simultaneous



Very Large Instruction Word (VLIW)
Static



“Hyper-Threading Technology Architecture and Microarchitecture”, Intel Technology Journal Q1, 2002

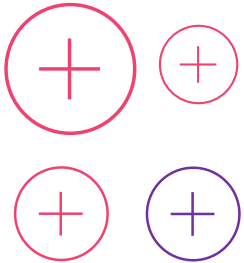


Multicore Organization

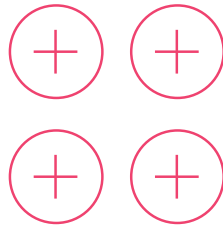
Compute

Memory

Asymmetric



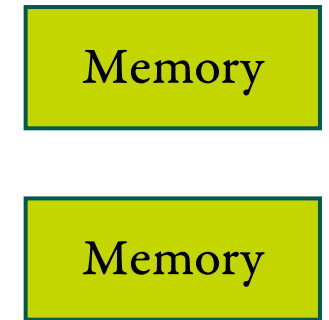
Symmetric



Centralized



Distributed



Uniform or Non-uniform
memory access

Multicore and Multiprocessor Summary

What are the ways to improve throughput beyond single instruction stream?

Multithreading, chip multiprocessing, and combinations

What is the difference between multithreading and chip multiprocessing?

Multithreading: multiple instructions on the same core/processor/CPU

Chip multiprocessing: using multiple cores on the same chip

What are the aspects of multicore organization?

Compute (symmetric or asymmetric) and memory (centralized or distributed)

How are hardware and software parallelism related?

Abstractions are different—software (threads, processes), hardware (threads/contexts)

Hardware: implicit (single instruction stream) or explicit (multiple instruction streams)

Software: sequential (via automatic parallelism) or parallel (shared memory or message passing)



Coordinating Memory Accesses

Define correct behavior!

“Write the spec”

Shared data block

Cache coherence

Also called cache consistency

All accesses

Memory consistency

Applies to sequential programs too!



Multicore Memory Summary

What is the challenge with (multicore/distributed) system design and programming?

Implementing and reasoning *correctly* about it

What are two specifications that help the above?

Cache coherence: about writes to the same location across caches

Memory consistency: about ordering of memory accesses

What is cache coherence?

Each write is eventually visible and writes to the same location are serialized

What is memory consistency?

Specification of what values a read can return and when

What are the principles at play in defining and implementing coherence and consistency?

Abstraction, Efficiency, “Yummy”, Concurrency, Approximate



Parallel Architectures

Low-level parallelism

Pipelining, superscalar, CPU+I/O separation

Parallel architecture

Multiple processors executing concurrently to solve a problem, with an explicit high-level framework

Various classifications and taxonomies

Flynn: Single/Multiple Instruction and Single/Multiple Data

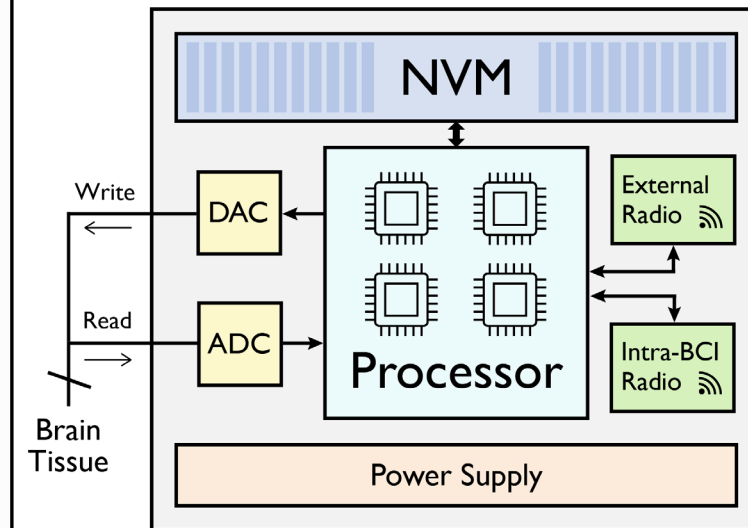
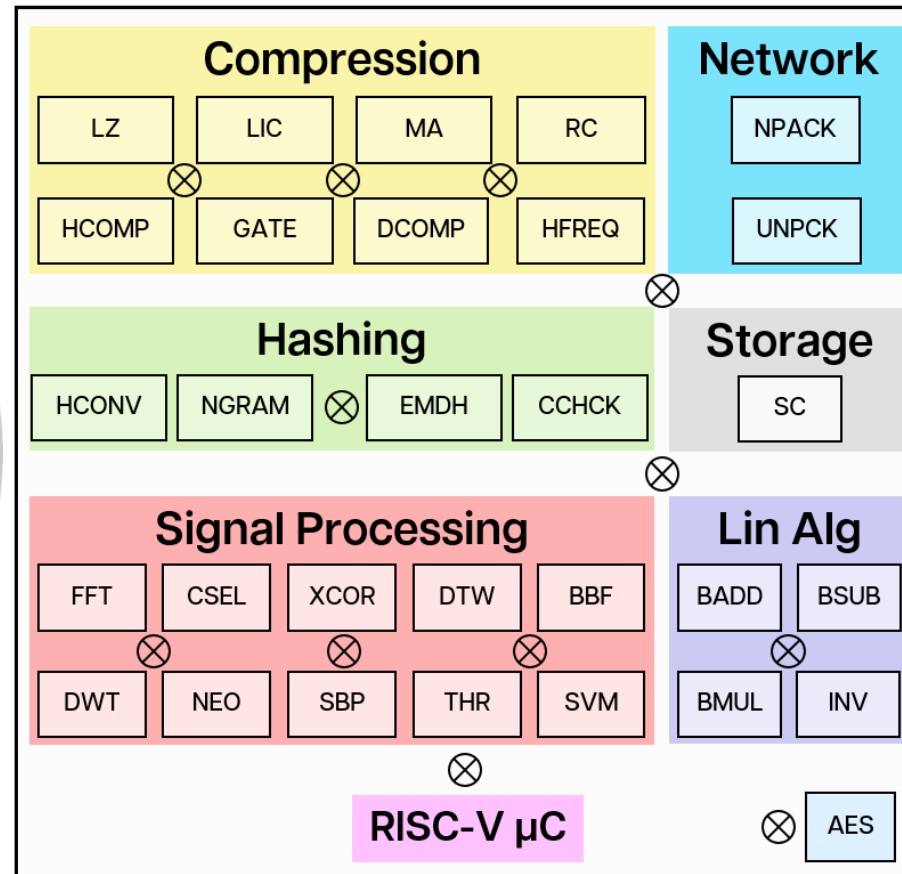
Vector processors, SIMD units, Associative processors, Systolic arrays,

MIMD (Shared and distributed memory)

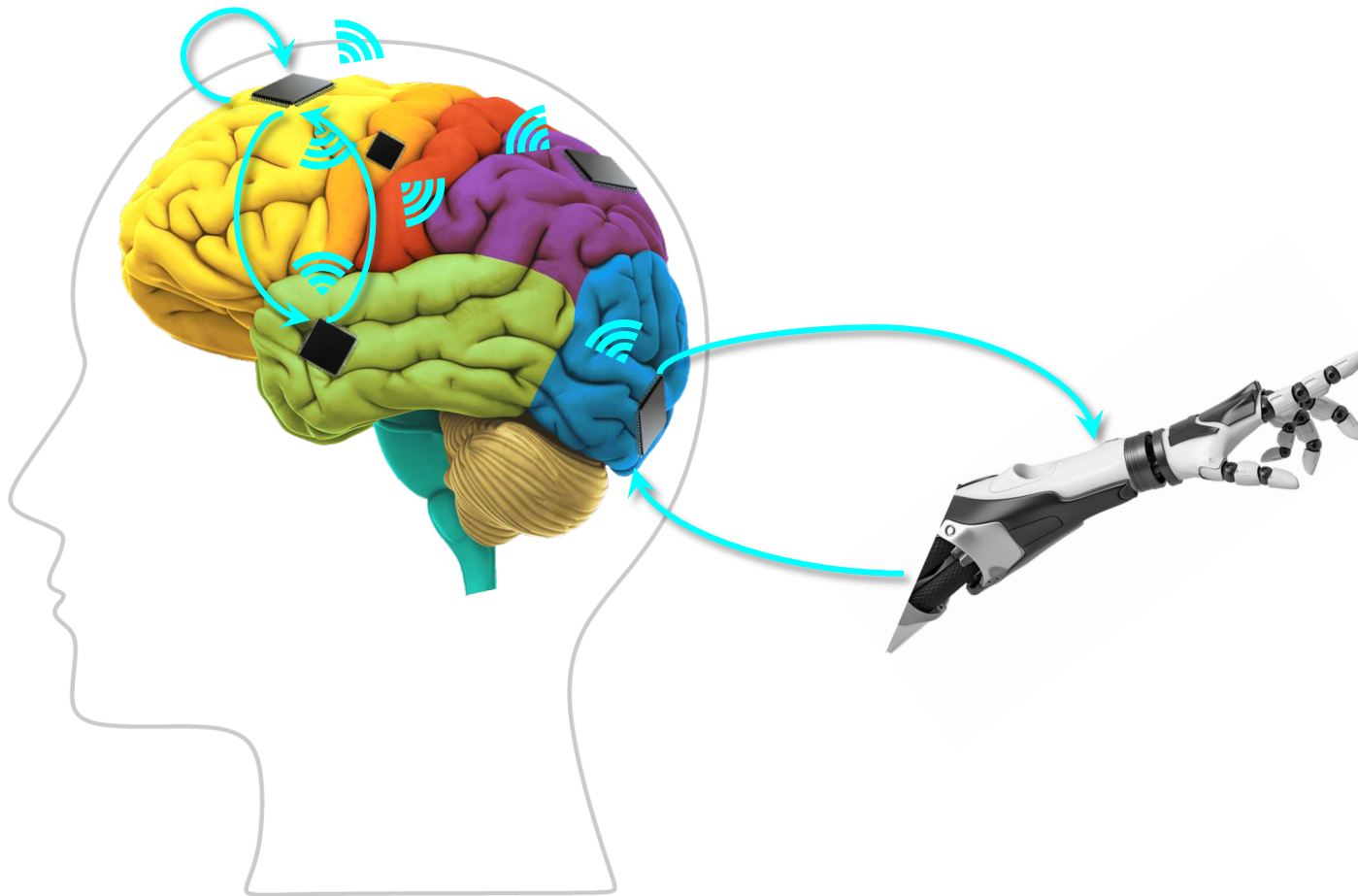
Dataflow processors



Example with a BCI Processor: SCALO



Example with a BCI Processor: SCALO



Dataflow

Static resource sharing and ordering

Software PE scheduling

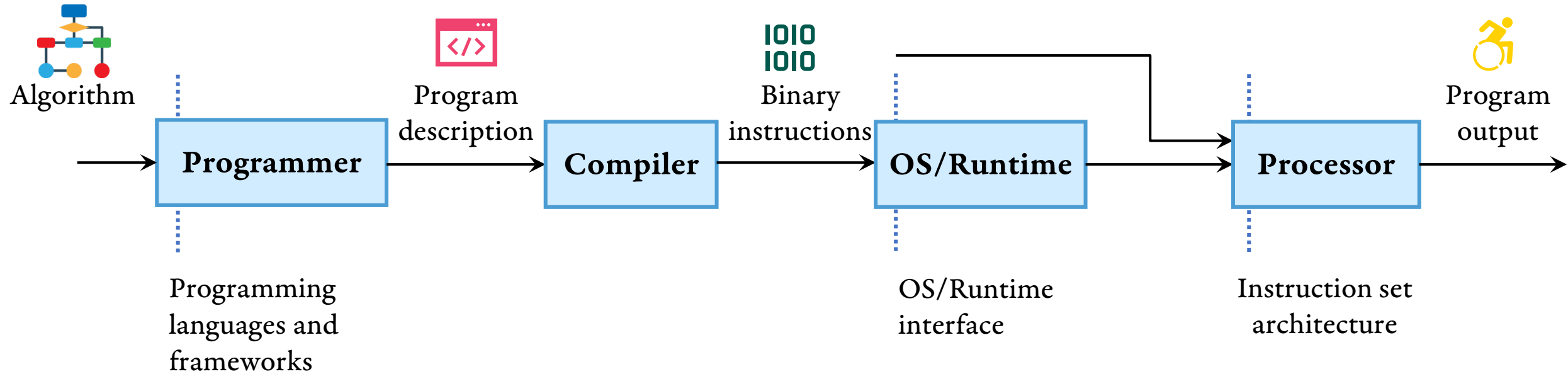
Co-designed fast paths

Explicit communication

No caching

Non-blocking writes with merging

BCI Computing is a Full Stack Problem



Emerging user trends (multimodal, wireless, portability)

Emerging algorithms (DNNs, Online learning, RNNs)

Emerging hardware trends (Beyond ASIC)

Co-design (Infinimind, Marple, Noema, KalmMind)

Abstractions, Programming or Development Platforms (BRAND, xDev, OpenVIBE)

Full stack constraints (Neuralink)

What You Should Have Learned

Computer architecture concepts

Recognizing and understanding them

BCI applications and algorithms

Understand what BCIs can offer, and how they work

BCI computational needs

Understanding them and the sources of these needs and constraints

BCI state of the art

Understand latest developments in the broad BCI space across the stack, limitations, future directions

Apply computer system design to BCIs

Demonstrate understanding of BCIs and computer system design through a working project



Feedback on Paper Reading

Read in detail

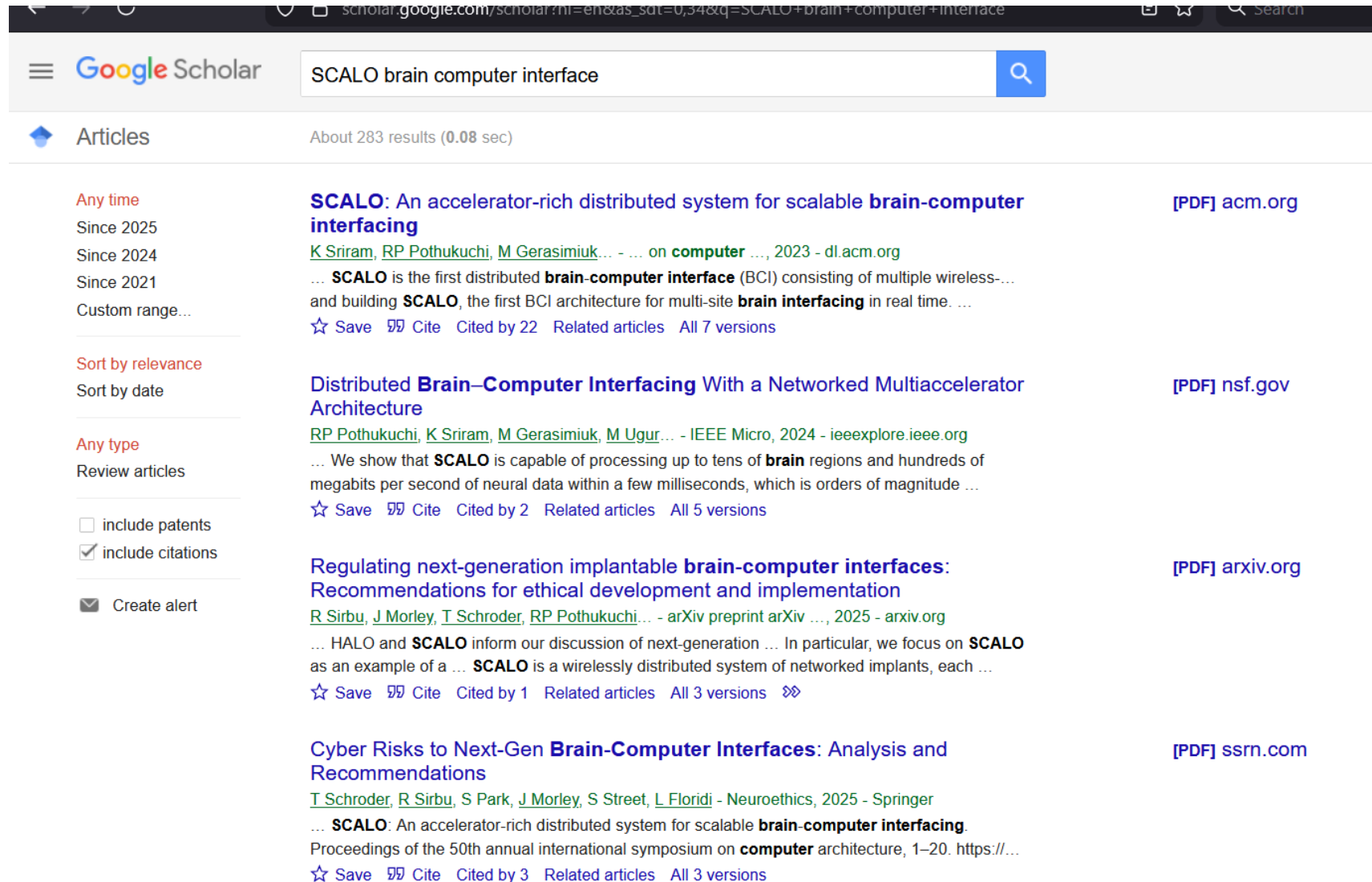
Multiple passes

Use manual search first

You'll discover many things—helps with breadth and depth



Feedback on Paper Reading



The screenshot shows a Google Scholar search results page. The search query is "SCALO brain computer interface". The results are sorted by relevance. The first result is "SCALO: An accelerator-rich distributed system for scalable brain-computer interfacing" by K Sriram, RP Pothukuchi, M Gerasimiuk, et al., published in 2023 in the DL.acm.org journal. The second result is "Distributed Brain-Computer Interfacing With a Networked Multiaccelerator Architecture" by RP Pothukuchi, K Sriram, M Gerasimiuk, M Ugur, et al., published in 2024 in the IEEE Micro journal. The third result is "Regulating next-generation implantable brain-computer interfaces: Recommendations for ethical development and implementation" by R Sirbu, J Morley, T Schroder, RP Pothukuchi, et al., published in 2025 as an arXiv preprint. The fourth result is "Cyber Risks to Next-Gen Brain-Computer Interfaces: Analysis and Recommendations" by T Schroder, R Sirbu, S Park, J Morley, S Street, L Floridi, published in 2025 in the Neuroethics journal. The left sidebar contains filters for time, sort, and type, as well as checkboxes for patents and citations, and a "Create alert" button.

Google Scholar

SCALO brain computer interface

Articles About 283 results (0.08 sec)

Any time
Since 2025
Since 2024
Since 2021
Custom range...

Sort by relevance
Sort by date

Any type
Review articles

☐ include patents
☒ include citations

☒ Create alert

SCALO: An accelerator-rich distributed system for scalable brain-computer interfacing [PDF] acm.org
K Sriram, RP Pothukuchi, M Gerasimiuk... - ... on computer ..., 2023 - dl.acm.org
... SCALO is the first distributed brain-computer interface (BCI) consisting of multiple wireless-... and building SCALO, the first BCI architecture for multi-site brain interfacing in real time. ...
☆ Save Cite Cited by 22 Related articles All 7 versions

Distributed Brain-Computer Interfacing With a Networked Multiaccelerator Architecture [PDF] nsf.gov
RP Pothukuchi, K Sriram, M Gerasimiuk, M Ugur... - IEEE Micro, 2024 - ieeexplore.ieee.org
... We show that SCALO is capable of processing up to tens of brain regions and hundreds of megabits per second of neural data within a few milliseconds, which is orders of magnitude ...
☆ Save Cite Cited by 2 Related articles All 5 versions

Regulating next-generation implantable brain-computer interfaces: Recommendations for ethical development and implementation [PDF] arxiv.org
R Sirbu, J Morley, T Schroder, RP Pothukuchi... - arXiv preprint arXiv ..., 2025 - arxiv.org
... HALO and SCALO inform our discussion of next-generation ... In particular, we focus on SCALO as an example of a ... SCALO is a wirelessly distributed system of networked implants, each ...
☆ Save Cite Cited by 1 Related articles All 3 versions

Cyber Risks to Next-Gen Brain-Computer Interfaces: Analysis and Recommendations [PDF] ssrn.com
T Schroder, R Sirbu, S Park, J Morley, S Street, L Floridi - Neuroethics, 2025 - Springer
... SCALO: An accelerator-rich distributed system for scalable brain-computer interfacing. Proceedings of the 50th annual international symposium on computer architecture, 1-20. https://...
☆ Save Cite Cited by 3 Related articles All 3 versions



Feedback on Paper Reading

Read in detail

Multiple passes

Use manual search first

You'll discover many things—helps with breadth and depth

Mozart: Taming taxes and composing accelerators with shared-memory

[PDF] acm.org

[V Suresh](#), [B Mishra](#), [Y Jing](#), [Z Zhu](#), [N Jin...](#) - Proceedings of the ..., 2024 - dl.acm.org

... [42] propose a **brain-computer interface** (BCI) system that includes processing units for small kernels frequently used in BCI pipelines. Recent works [43, 83] also highlight the ...

☆ Save  Cite Cited by 3 Related articles All 6 versions



Stimulates broad domain learning, abstract reasoning, and multi-task learning



Feedback on Presentations

Different types of presentations

Tease them (motivate the interaction), Show them (storytelling), Throw at them (detailed review)

Presentation is an aid to an immersive experience

Draw attention where you want through color, font, objects, graphics, and *absence*
Reading notes or long text stimulates reading or hearing—not listening, and not engaging



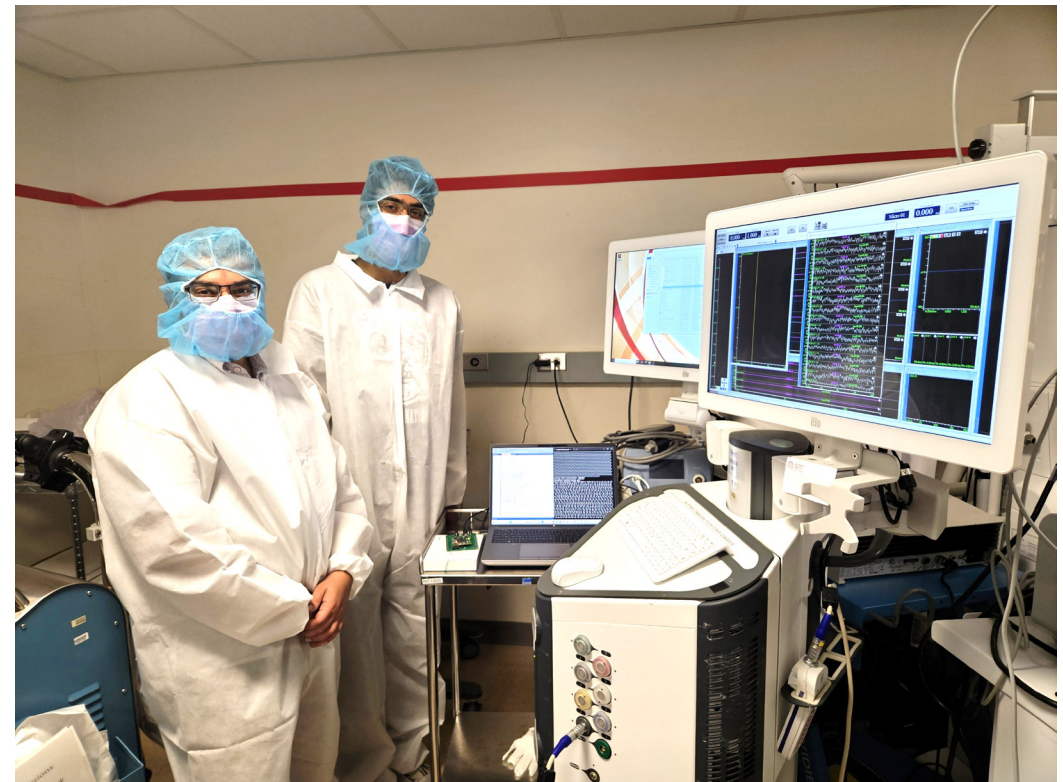
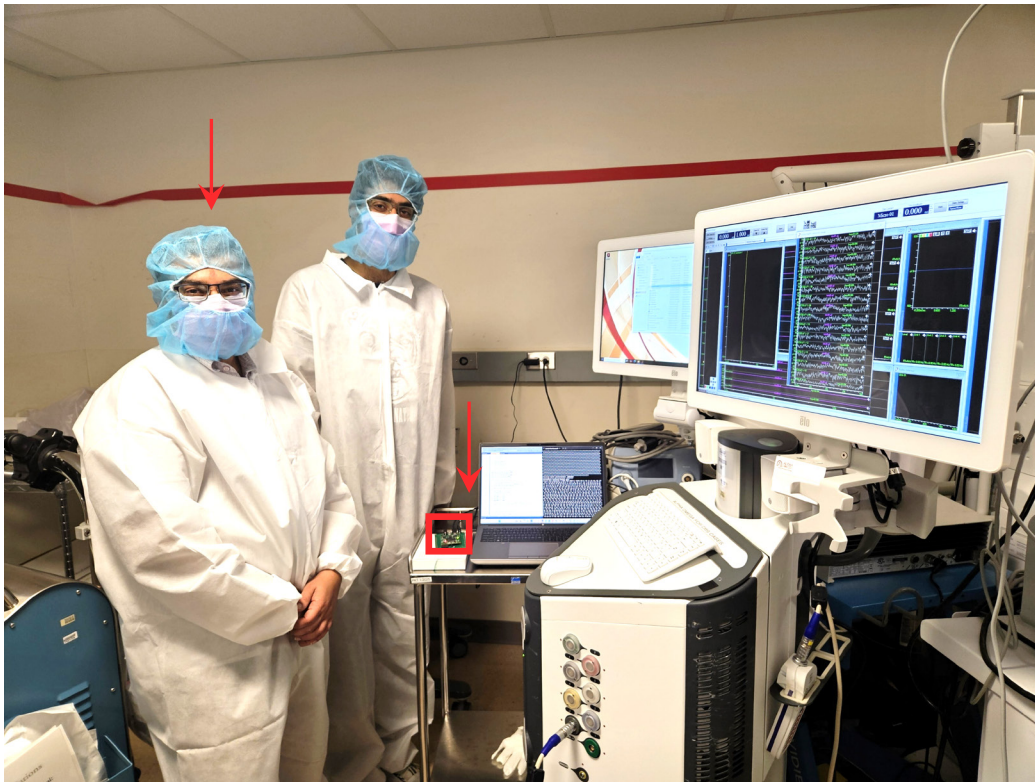
Feedback on Presentations

Different types of presentations

Tease them (motivate the interaction), Show them (storytelling), Throw at them (detailed review)

Presentation is an aid to an immersive experience

Draw attention where you want through color, font, objects, graphics, and *absence*



Feedback on Presentations

Different types of presentations

Tease them (motivate the interaction), Show them (storytelling), Throw at them (detailed review)

Presentation is an aid to an immersive experience

Draw attention where you want through color, font, objects, graphics, and *absence*
Reading notes or long text stimulates reading or hearing—not listening, and not engaging

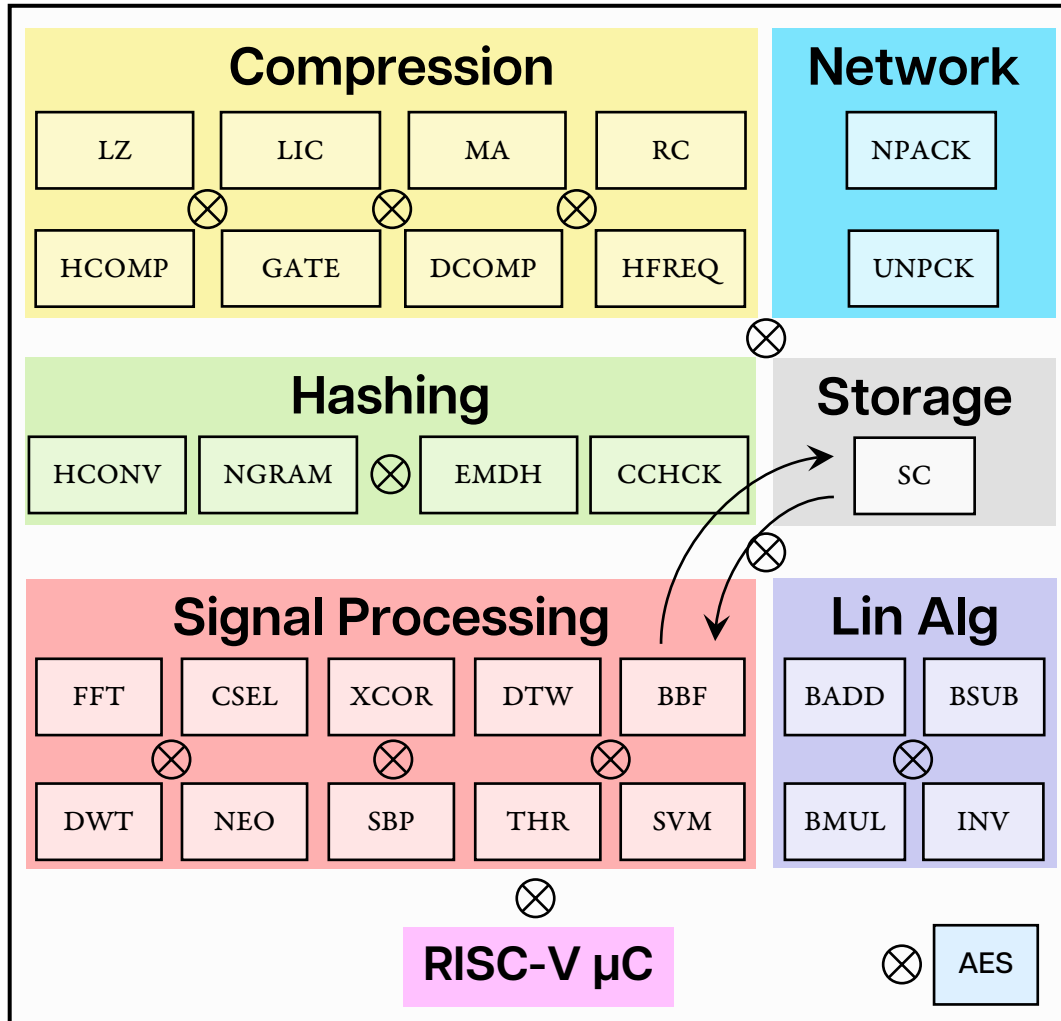
Lead the audience

Move from a high-level introduction to the nitty gritty and go back up with the wisdom

Example slide follows



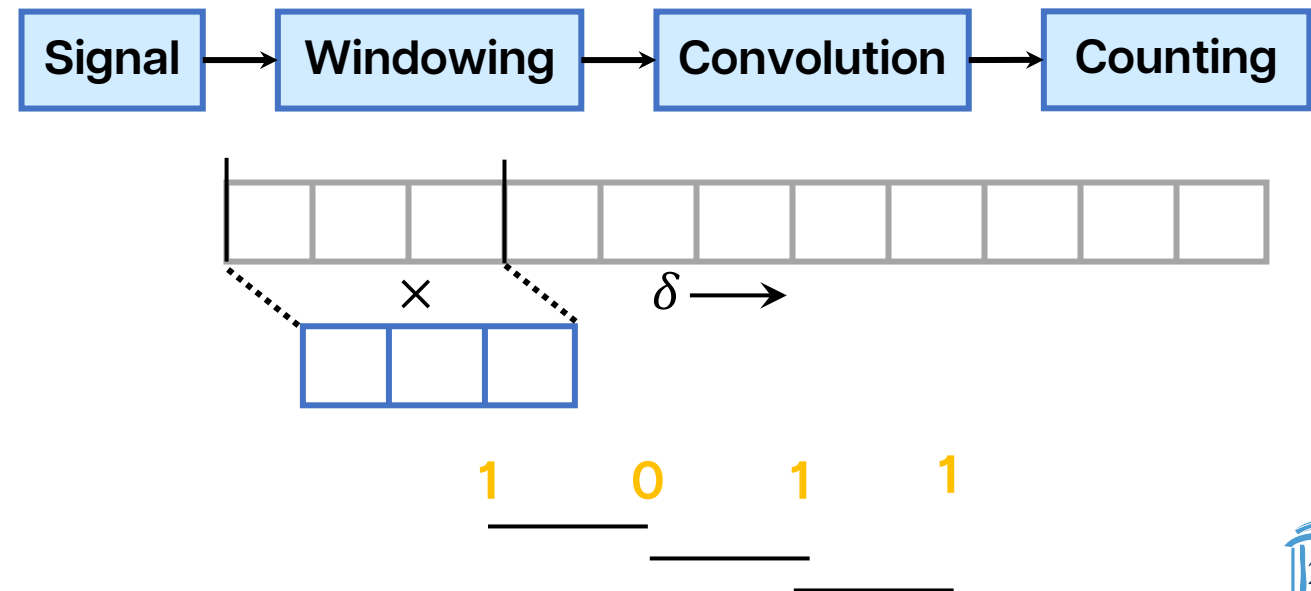
Accelerator-Rich Architecture



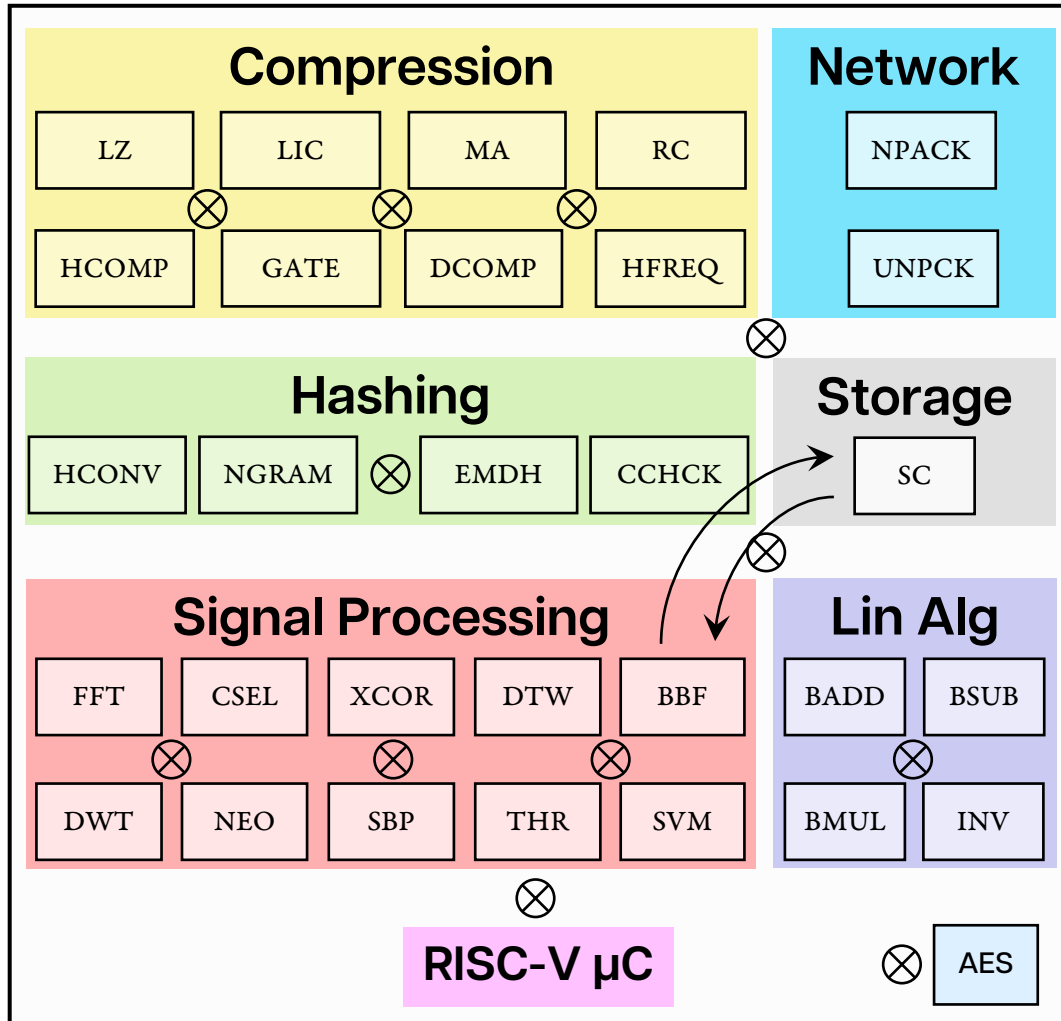
Optimize processing elements (PEs) for reuse

Locality Sensitive Hashing

Dynamic Time Warping, Euclidean Distance, Correlation, Earth-Mover's Distance



Accelerator-Rich Architecture



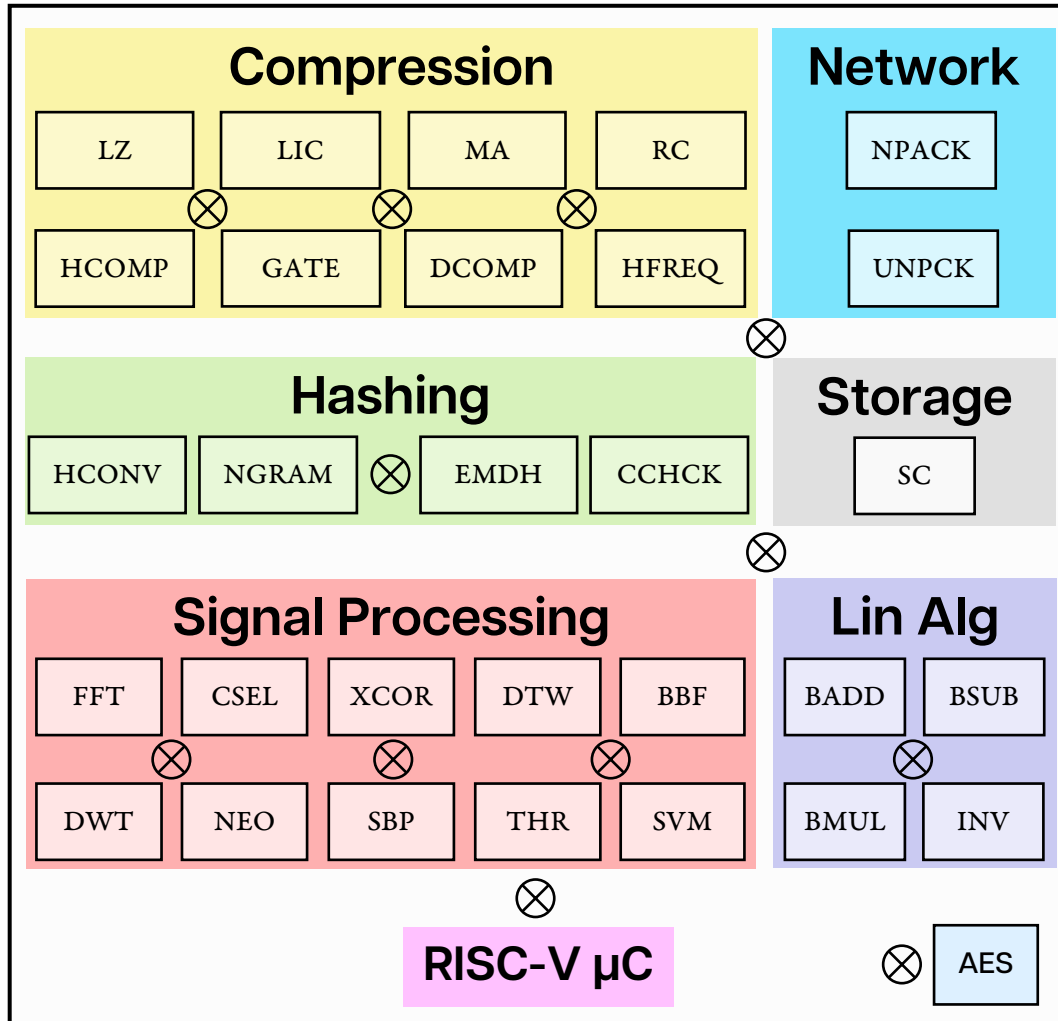
Optimize processing elements (PEs) for reuse

Per-PE clock, and frequency scaling

Configurable interconnect

Predictable power and performance

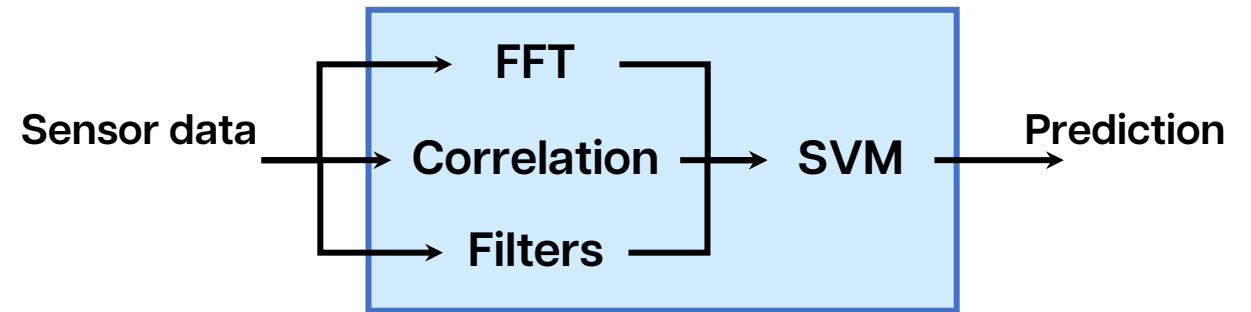
Accelerator-Rich Architecture



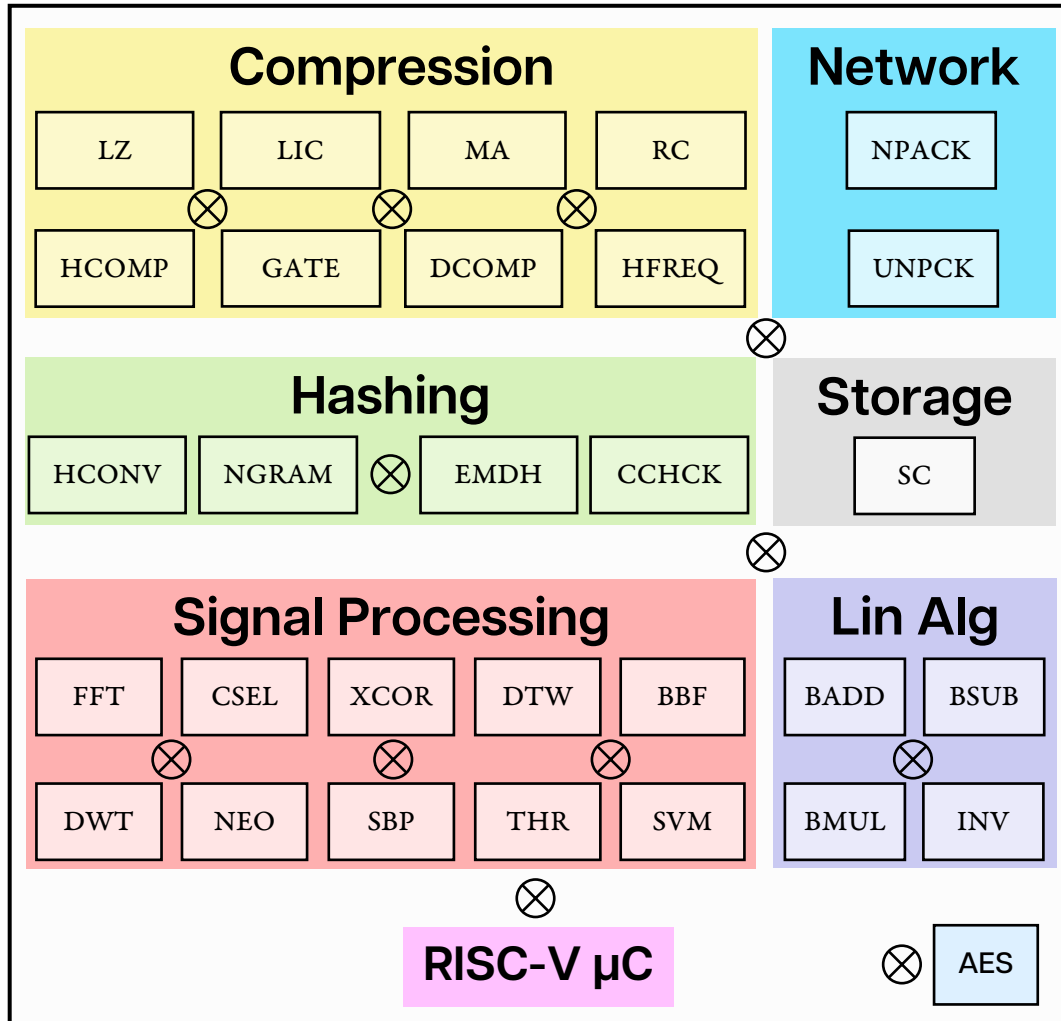
ISA, but with accelerators, for BCI applications

Accelerator as the instruction

Efficiency with flexibility



Accelerator-Rich Architecture

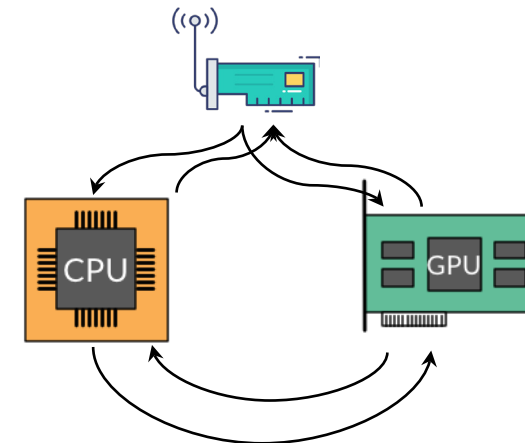


ISA, but with accelerators, for BCI applications

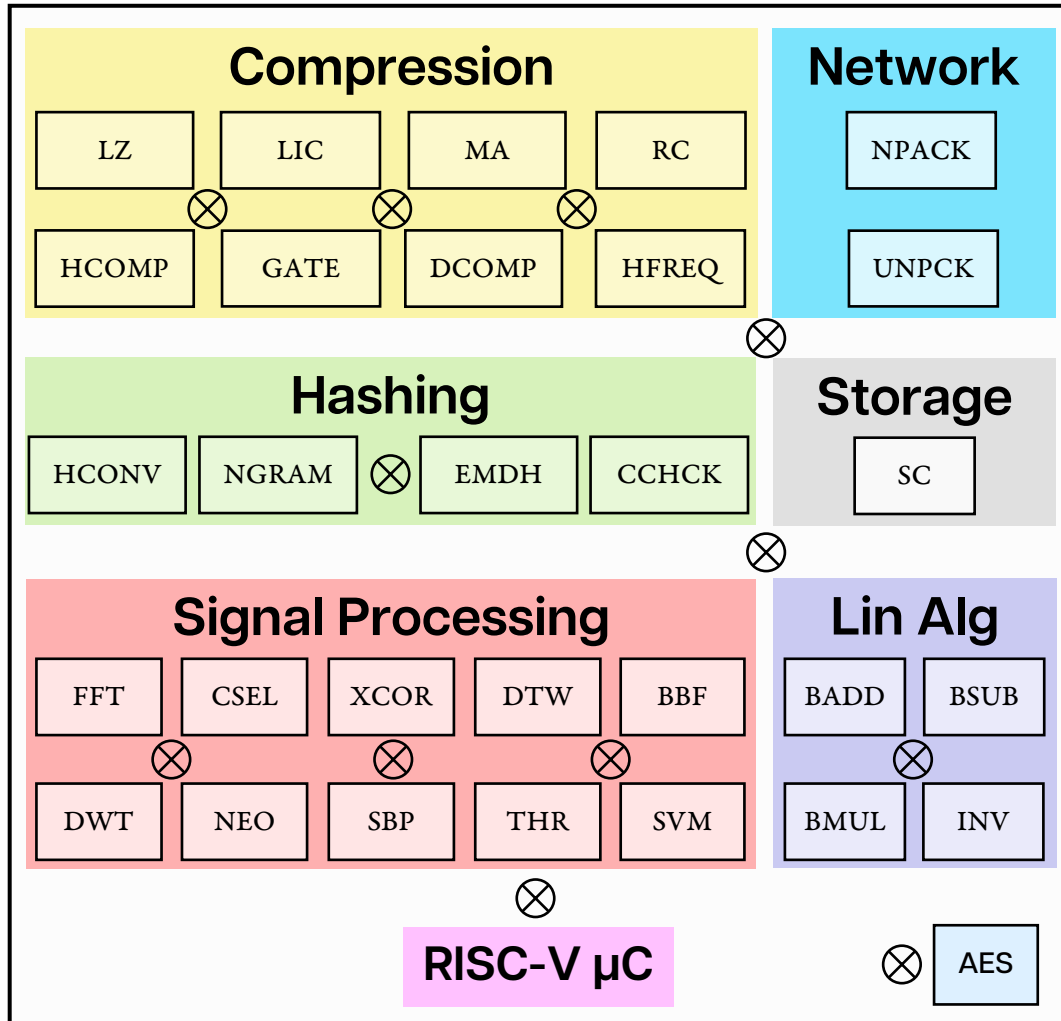
Accelerator as the instruction

Efficiency with flexibility

Minimizing CPU-led overheads



Accelerator-Rich Architecture



ISA, but with accelerators, for BCI applications

Accelerator as the instruction

Efficiency with flexibility

Minimizing CPU-led overheads

Hash-based communication for scalability

Feedback on Presentations

Different types of presentations

Tease'em (motivate the interaction), Show'em (storytelling), Throw at'em (detailed review)

Presentation is an aid to an immersive experience

Draw attention where you want through color, font, objects, graphics, and *absence*
Reading notes or long text stimulates reading or hearing—not listening, and not engaging

Lead the audience

Move from a high-level introduction to the nitty gritty and go back up with the wisdom

Make the audience think, and be excited

Show connections that weren't obvious, discuss significance, tie to real examples that matter

Practice!

Perfect practice makes perfect—to deliver the best experience



Feedback?

Content

Delivery

Organization

Alignment with expectations

Knowledge and skill building



Image Credits (Educational, Fair Use)

- Title image: VLADGRIN, https://www.istockphoto.com/vector/human_-machine-gm147409511-16840728 (Educational fair use)
- Infinite brain: Science wonder stories, May 1930, Illustrator: Frank R Paul, Editor: Hugo Gernsback
- Brain color, ICs, cloud server, black rat: No attribution required (Hiclipart)
- Hand with spoon: public domain freepng
- Signals: <https://www.nature.com/articles/nrn3724>
- Thought clouds: F. Willett et al./*Nature* 2021/Erika Woodrum, <https://med.stanford.edu/neurosurgery/news/2022/bci-award>. <https://www.the-scientist.com/news-opinion/brain-computer-interface-user-types-90-characters-per-minute-with-mind-68762>
- Picture of scientists: <https://www.cs.auckland.ac.nz/~brian/rutherford8.html> (original: Pierre de Latil), Bush (Carnegie Science), Others (Wikipedia, National Academies, IEEE, and university profile images)
- Flowchart: Pause08 – flaticon.com; Digital brain: Smashicons – flaticon.com; Quantum processor icons created by Paul J. - Flaticon
- Server rack: upklyak – freepik.com
- Arm, Lotus: Adobe stock
- Quantum processor: Rigetti computing
- Images of implanted users: Top: Case Western Reserve University (<https://thedaily.case.edu/man-quadruplegia-employs-injury-bridging-technologies-move-just-thinking/>), Bottom: Jan Scheuermann (University of Pittsburgh/UPMC; <https://www.upmc.com/media/news/bci-press-release-chocolate>)
- Images of wearable BCIs: Cognixion, NextMind
- Types of BCIs: “Brain–computer interfaces for communication and rehabilitation,
- Illustrative BCI: Neuralink
- Electrodes: “Electrochemical and electrophysiological considerations for clinical high channel count neural interfaces”, Vatsyayan et al.
- Form factors: Neuropace, Medtronic, Bloomberg, “Fully Implanted Brain–Computer Interface in a Locked-In Patient with ALS” by Vansteensel et al., Blackrock Neurotech
- Jose Delgado’s video: Online, various sources (CNN, Youtube)
- Video of Kennedy and Ramsey: Online, various sources (Youtube, Neural signals)
- Code snippet inspiration: ECE 252 slides at Duke (Dan Sorin et al.)
- Apple processor pipeline: <https://dougallj.github.io/applecpu/firestorm.html>

Logos, trademarks are all properties of respective owners

Not to be shared outside the course

