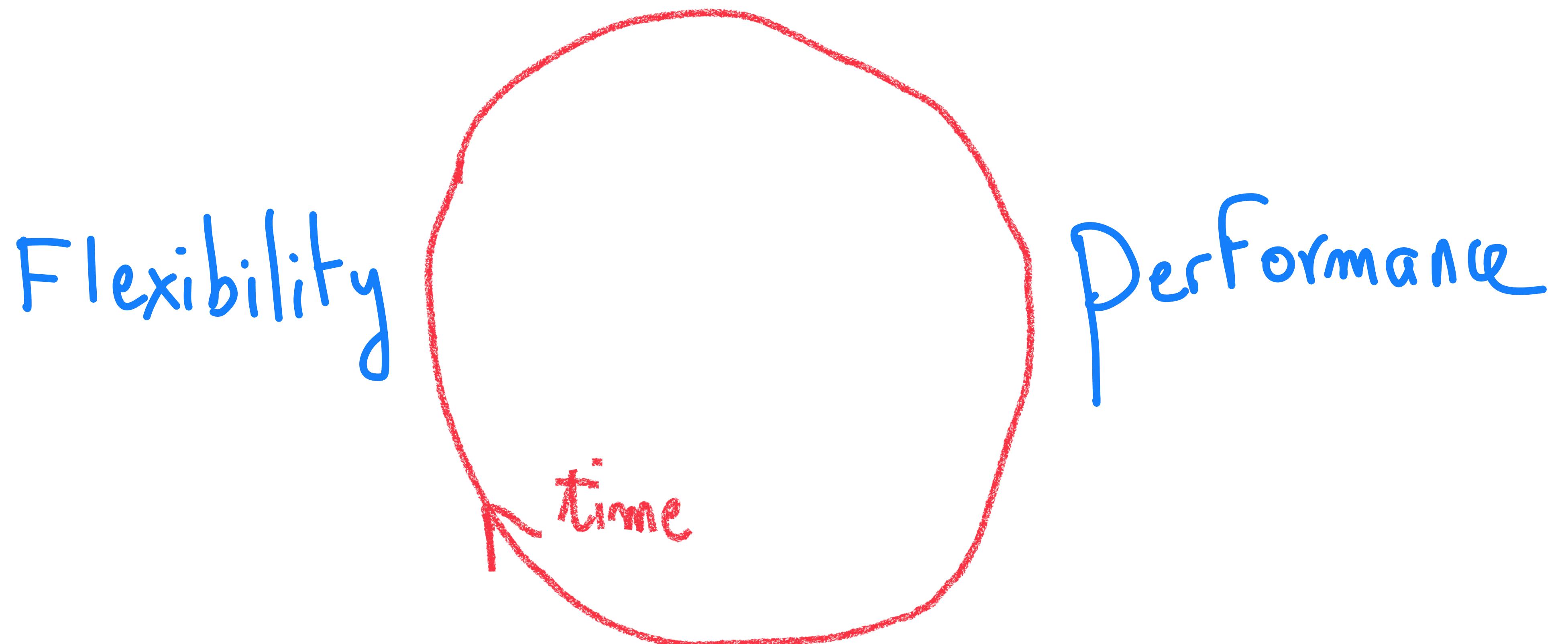


Hardware-Software Codesign

Thomas - POCS - Fall 2023 - EPFL

Outline



Outline

Some historical perspective on the HW/SW dance.

Cost of Abstraction - Value of Specialization

Chasing the Unicorn

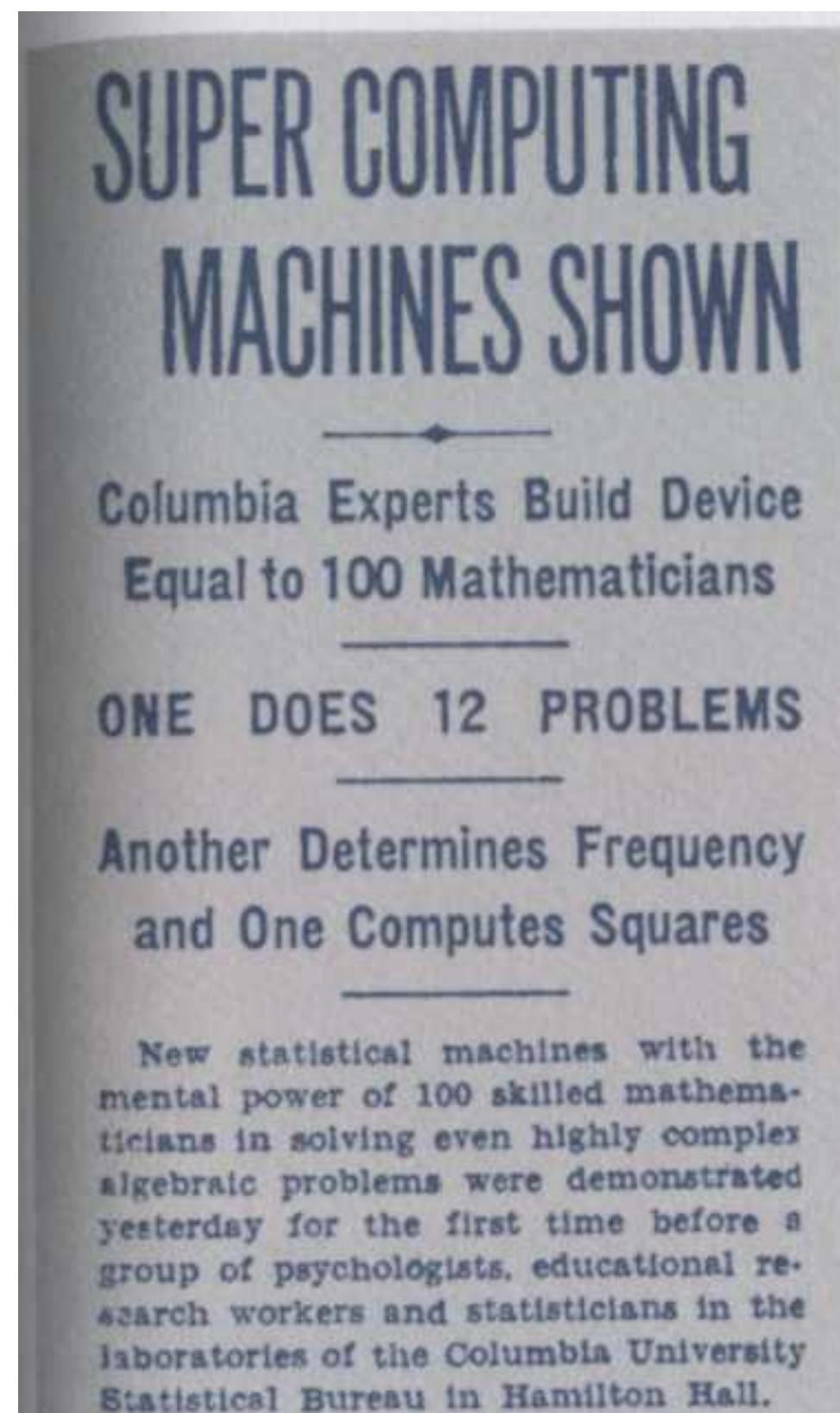
Communication and Interfaces - 3 case studies

Principles of Acceleration

Cost of Specialization - Value of Abstraction

Hardware/Human Codesign

Pre-Turing Era - 1930s



Fancy tabulating machine could High Value Computations: $\sum_i x_i^2$ or even $\sum_i w_i \cdot x_i$

Technology was electromechanical relays

Typically, quite a few “telescoping” (pipelining) tricks (example)

Humans interact with the kernel performed by the machine:

Machine Input: Feed it punchcards with data

Machine Outputs:

- New punchcards (write once medium, for partially accumulated data)
- Human readable summary on fan-fold paper

Cool computations: SELECT/WHERE/GROUP BY !

Performance: Could ingest 150 punch cards per minutes

1945 - Vacuum Tube's Era

Von Neumann's magic - Death of hardware, birth of Software

Reports on Building an “Automatic Computing Device” (C[entral], M[emory], R[ecording] (for IO))

The orders which are received by CC come from M, i.e. from the same place where the numerical material is stored.

Von Neumann's Insight/Suggestion (Mechanical - ms, Vacuum tube - us)

Thus it seems worthwhile to consider the following viewpoint: The device should be as simple as possible, that is, contain as few elements as possible. This can be achieved by never performing two operations simultaneously, if this would cause a significant increase in the number of elements required.

It is also worth emphasizing that up to now all thinking about high speed digital computing devices has tended in the opposite direction: Towards acceleration by telescoping processes at the price of multiplying the number of elements required. It would therefore seem to be more instructive to try to think out as completely as possible the opposite viewpoint

(<https://web.mit.edu/STS.035/www/PDFs/edvac.pdf>)

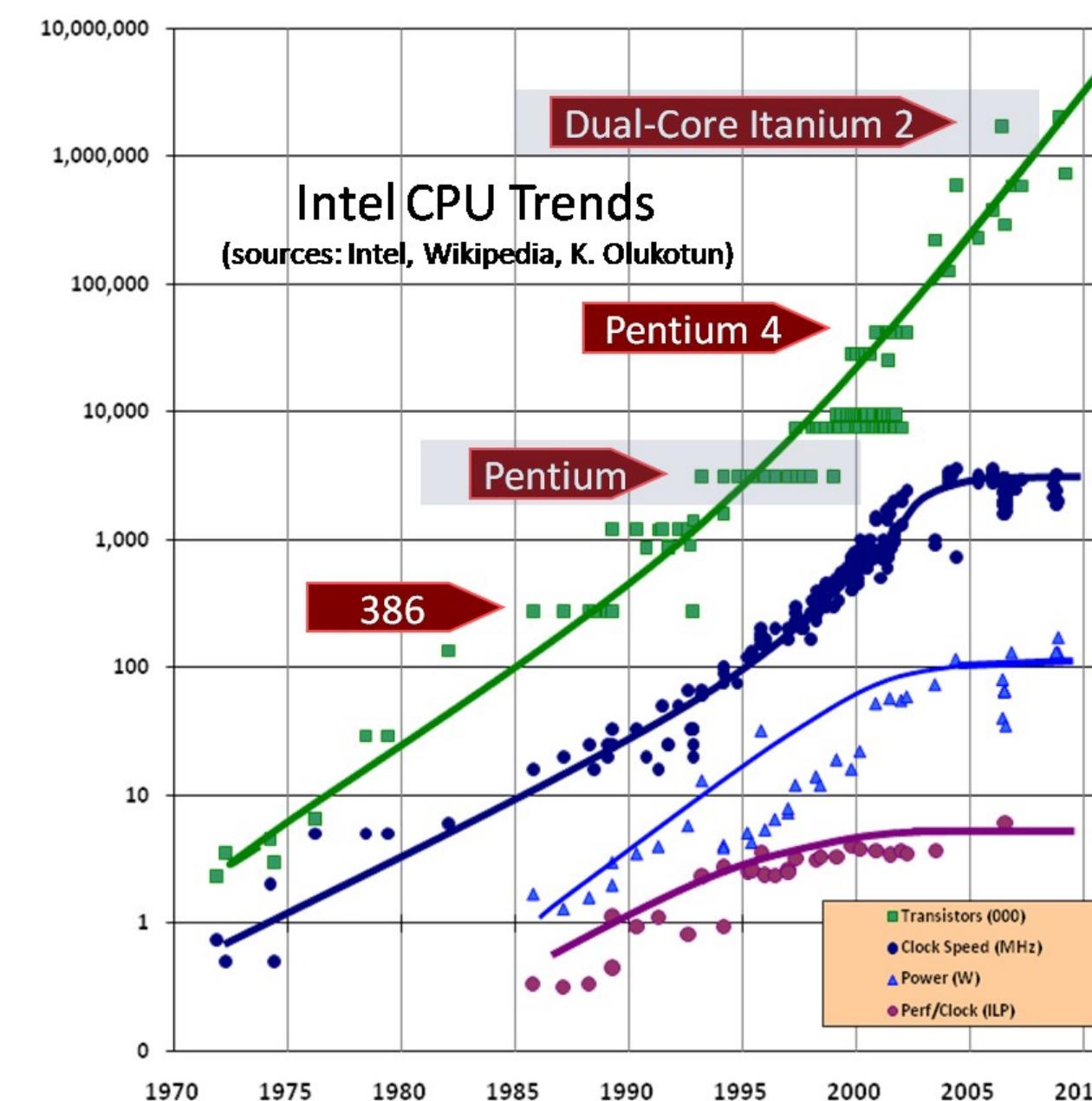
Very rich paper: why binary is best, thoughts about errors in computing, JIT, biomimetics, brain VS computer, synchrony/asynchrony ...

Feynman's observation

Plenty of Room at the Bottom

1959 Feynman mentally explores the future of miniaturisation: in the future (present) we should be able to do insanely small machines!

Research and industry makes room at the bottom: Moore's law



https://web.pa.msu.edu/people/yang/RFeynman_plentySpace.pdf

Feynman was right

Transistor Free Lunch Party!

Side-effects:

Smaller transistor ~ faster clock ~ no effort, same design go faster

Smaller transistor ~ more transistors on a given chip, what do we do with them?

Von Neumann's suggestion of simple machine is dead:

- And now it is not just arithmetic tricks
- New optimization opportunities!

So everything got so much faster?

danluu.com

computer	latency (ms)	year	clock	# T
apple 2e	30	1983	1 MHz	3.5k
ti 99/4a	40	1981	3 MHz	8k
custom haswell-e 165Hz	50	2014	3.5 GHz	2G
commodore pet 4016	60	1977	1 MHz	3.5k
sgi indy	60	1993	.1 GHz	1.2M
custom haswell-e 120Hz	60	2014	3.5 GHz	2G
thinkpad 13 chromeos	70	2017	2.3 GHz	1G
imac g4 os 9	70	2002	.8 GHz	11M
custom haswell-e 60Hz	80	2014	3.5 GHz	2G
mac color classic	90	1993	16 MHz	273k
powerspec g405 linux 60Hz	90	2017	4.2 GHz	2G
macbook pro 2014	100	2014	2.6 GHz	700M
thinkpad 13 linux chroot	100	2017	2.3 GHz	1G
lenovo x1 carbon 4g linux	110	2016	2.6 GHz	1G
imac g4 os x	120	2002	.8 GHz	11M
custom haswell-e 24Hz	140	2014	3.5 GHz	2G
lenovo x1 carbon 4g win	150	2016	2.6 GHz	1G
next cube	150	1988	25 MHz	1.2M
powerspec g405 linux	170	2017	4.2 GHz	2G
packet around the world	190			
powerspec g405 win	200	2017	4.2 GHz	2G
symbolics 3620	300	1986	5 MHz	390k

**"All problems in computer science can be solved by
another level of indirection ...
except for the problem of too many levels of indirection"**

~David J. Wheeler (Maybe ?)

**Mid 2010s, after 70 years of indirections:
time for a spring cleaning**

Cost of Abstraction

There is Plenty of Room at The Top

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	—	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45
TPU v1 (2016)		<0.001	~90,000.000	~15,000,000	~270	100% (???)

Source: <https://www.microsoft.com/en-us/research/uploads/prod/2020/11/Leiserson-et-al-Theres-plenty-of-room-at-the-top.pdf>

TPUv1/Processor: 28/22nm, ~350mm², ~700MHz/~3GHz

Remarks

Flame war between C/Java is a fight between 0.01% efficiency and 0.03% efficiency (This depends from program to program, compiler to compiler!)

Modern processors already embed a whole bunch of accelerators: Vector unit, often vastly underused

Algorithm complexities make simplifying assumptions (RAM model) and hides real cost
size(Workforce) usually decreases with speedup :(

Surprisingly, there are also performance costs to specialisation (end of the lecture)

Pragmatically: Modern AI is enabled by codesign

GPT3 evaluation ~

>> Read all the paper books in the EPFL library

Compute a nontrivial (arithmetic intensive) function of all the words present in all the books and some context

Output 1 token

Rinse and repeat, with the updated context of adding the produced token

Hundreds of GB per token! It is a pharaonic amount of compute

What about other computations?

Answer: the quantitative approach

A Characterization of Processor Performance in the VAX-11/780

Joel S. Emer

Douglas W. Clark

This paper reports the results of a study of VAX-11/780 processor performance using a novel hardware monitoring technique. A micro-PC histogram monitor was built for these measurements. It keeps a count of the number of microcode cycles executed at each microcode location. Measurement experiments

What about other computations?

Answer: the quantitative approach

Average VAX Instruction Timing (Cycles per Instruction)

	Compute	Read	R-Stall	Write	W-Stall	IB-Stall	Total
Decode	1.000					0.613	1.613
Spec1	0.895	0.306	0.364				1.565
Spec2-6	1.052	0.148	0.116	0.161	0.192	0.102	1.771
B-Disp	0.221					0.005	0.226
Simple	0.870	0.029	0.017	0.033	0.027		0.977
Field	0.482	0.049	0.058	0.007	0.002		0.600
Float	0.292	0.000	0.000	0.008	0.001		0.302
Call/Ret	0.937	0.133	0.074	0.130	0.184		1.458
System	0.434	0.015	0.031	0.014	0.028		0.522
Character	0.318	0.039	0.099	0.046	0.004		0.506
Decimal	0.026	0.002	0.000	0.001	0.002		0.031
Int/Except	0.055	0.002	0.005	0.004	0.006		0.071
Mem Mngmt	0.555	0.061	0.200	0.004	0.003		0.824
Abort	0.127						0.127
TOTAL	7.267	0.783	0.964	0.409	0.450	0.720	10.593

Performance evaporation (Knuth's challenge)

Intuitions about the cost of software abstractions

Knuth's challenge (1989): “Make a thorough analysis of everything your computer does during one second of computation.”

Performance evaporation (Knuth's challenge)

Intuitions about the cost of software abstractions

Knuth's challenge (1989): “Make a thorough analysis of everything your computer does during ~~one second~~ of computation.”

50 μ

Takeaway 1

There are no intuition for the performance of any complicated systems. A computer is a complicated system.

Measure, characterize, model and then organize/specify your findings

Chasing the unicorn

(Chasing the Unicorn)

The Ideal of Codesign

Serpent de mer of languages and compiler

Write your algorithm in SuperLang

```
# superc --powerConstraint=1W --areaConstraint=100um2 --clock=3GHz --  
IPblockLibraries armA72 matmulGoogle128 ether10GBBroadcom -O3 myProgram.super
```

-> Produce hardware description (Maybe FPGA configuration or directly ASICs description) + software for the CPUs + software glue code + user app that uses all that

Next Level Unicorn : SuperLang is actually C/Python/Haskell/Scala, so we can just reuse existing code

Escape Hatch:

"Super" does not work for all programs

"Superc" might produce suboptimal code

Progress made while looking for the unicorn

Behavioral VS Structural Verilog, and others High-Level Hardware Languages

Modern High Level Synthesis

System-Level Design Tools

Communication and Interfaces

Partitioning and Mapping, Design Space Exploration

Communication and Interfaces

by examples

Case Study 1: AES instructions - Tightly coupled accelerators

What is AES

Symmetric key crypto Block Cipher

Compiling with -O3 (for RISCV):

~300 straight-line instructions per AES-round (Same for X86-clang)

Encrypting 128 bits requires 10 rounds -> ~3000 instructions

Question:

If I send a file by email, how many times will the content of the file be AES-ed before it is opened by the recipient?

AES is a “High-Value” Application

AES is a part of most commonly used cipher suites for TLS.

TLS protocol massively used on the internet

WIFI also uses AES

FileVault on Mac

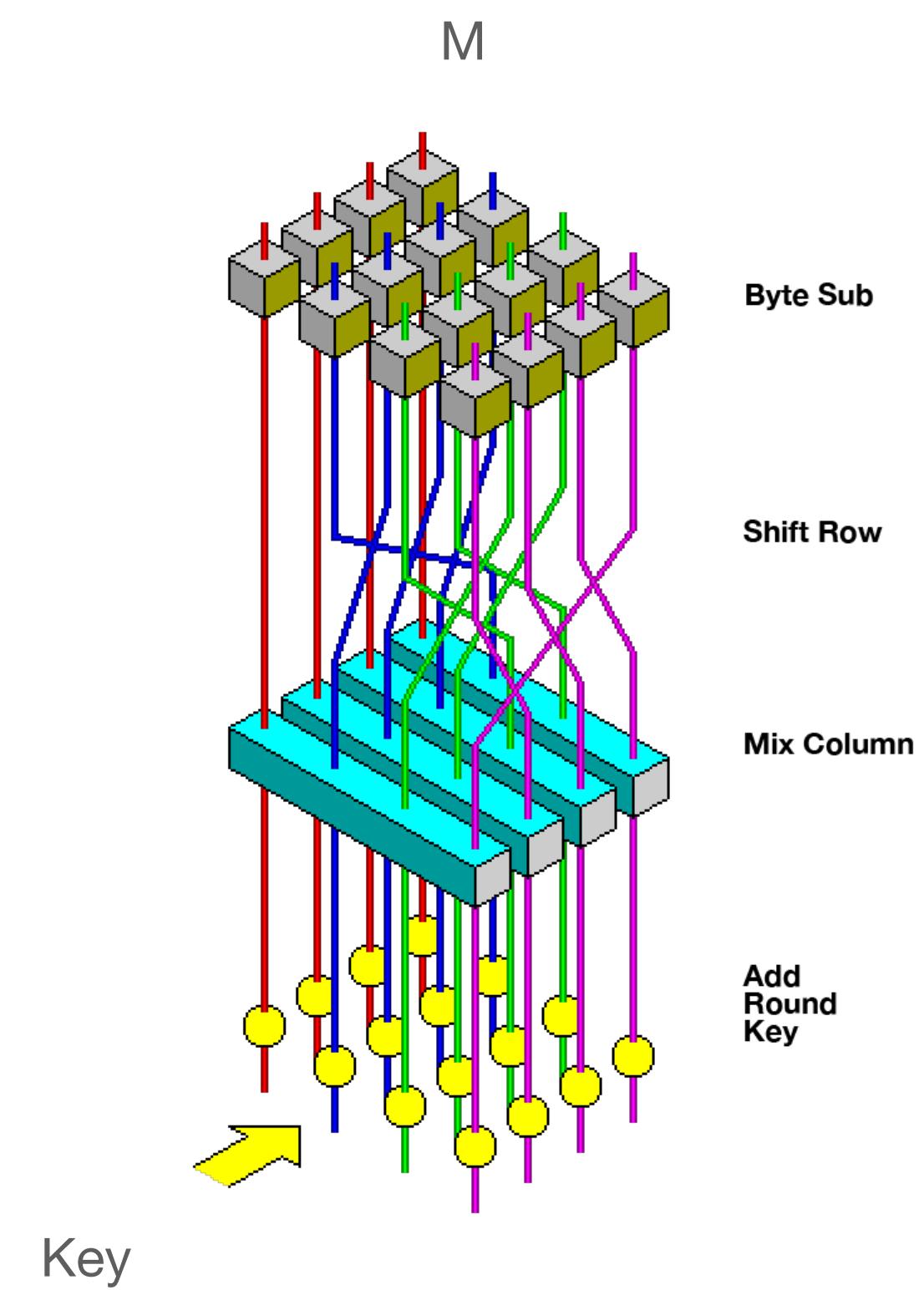
-> quite Important for server workloads

There are other reasons to accelerate AES (Security)

AES encryption

M and Key are 128 bits ($16 * 8$ bits)

$\text{AES_encrypt}(M, \text{Key}) = 128$ new bits



AESENC

From <https://eprint.iacr.org/2016/299.pdf>

AESENC xmm1, xmm2

“Perform one round of an *AES encryption* flow, using one 128-bit data (state) from xmm1 with one 128-bit round key from xmm2”

AESENC can be thought of like ADD, a new arithmetic instruction. Can nicely integrate it in the pipeline of the processor, easy!

Latency: 4 cycles, Throughput: 1 token per cycle

Deployment “challenge”:

People should use the AESENC instruction and not their own C implementation anymore!

Scope and Limitations of AESENC-style of acceleration

There are quite a few applications that can be enabled by low-latency instructions

Intel Vector Extensions (MMX -> SSE -> AVX) originally for multimedia and then more general

Compiler are still struggling hard to properly use those extensions

Limitations:

Many computations act on more data than just data that can be held in registers

- Need access to memory

Many computations need hundred of cycles to complete

Why is it a limitation?

Why is it a problem if instruction needs to access a lot of memory?

Accelerators are just devices!

How do devices work?

Examples of successful “accelerators”

DL accelerators (example: Gemmini, NPU, TPU)

Network Interface Cards (NIC)

TCP/IP stack runs on some NIC

GPU, GPGPU

Sound Card

Basic observations

Wide discrepancy in programmability

Programmable:

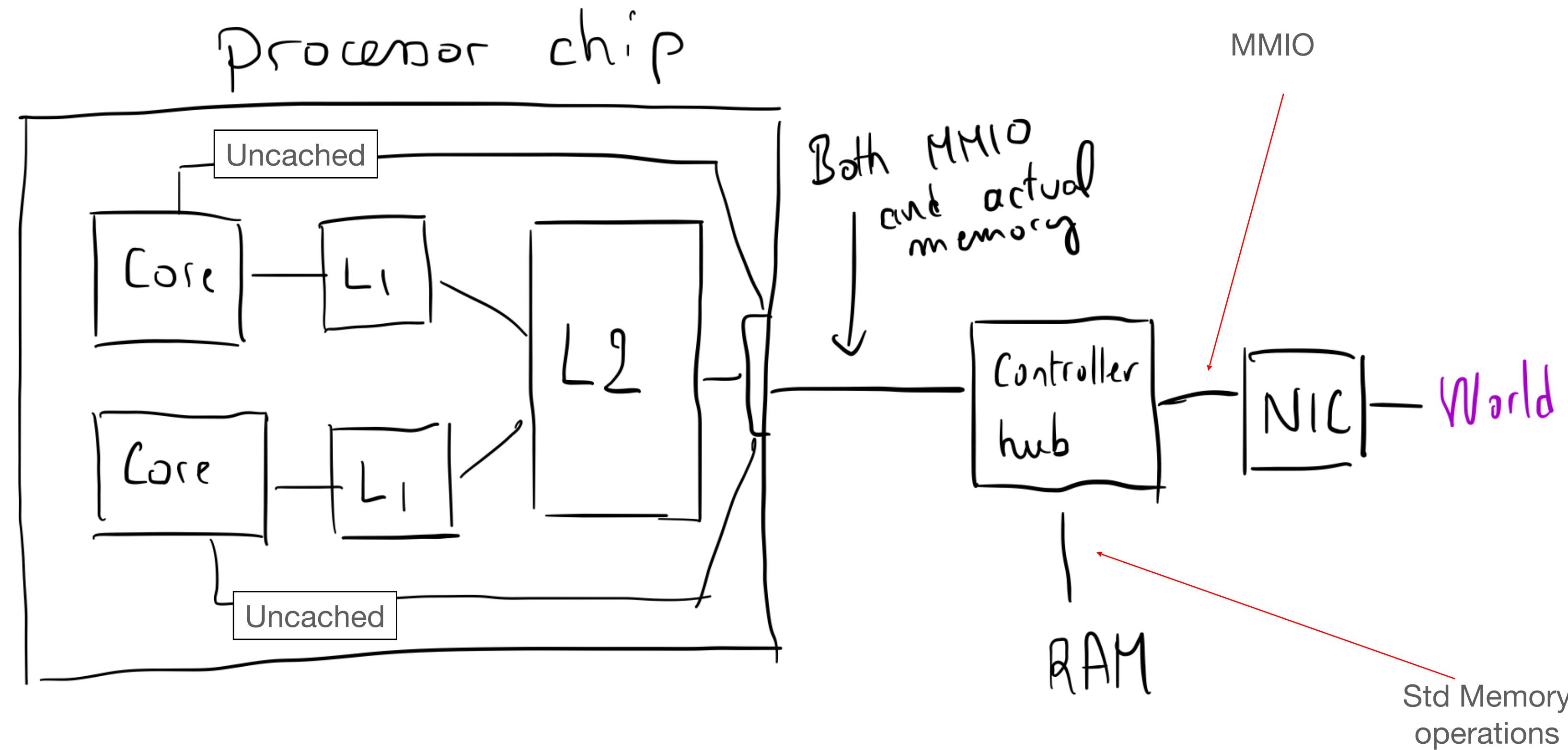
- Modern GPU - runs more or less arbitrary Cpp code
- Modern smartNIC can run programs too

Limited programmability:

- my old ethernet card
- my old GPU
- my old sound card

Case study 2: a simplified-NIC

Approximate High-level picture – NIC



Simplified journey of a packet in a fancy NIC

Processor tells NIC:

Address in RAM for SEND packets

Address where NIC should put the RECEIVED packets

Application generate data to send:

“GET / HTTP/1.1

Host: www.example.com

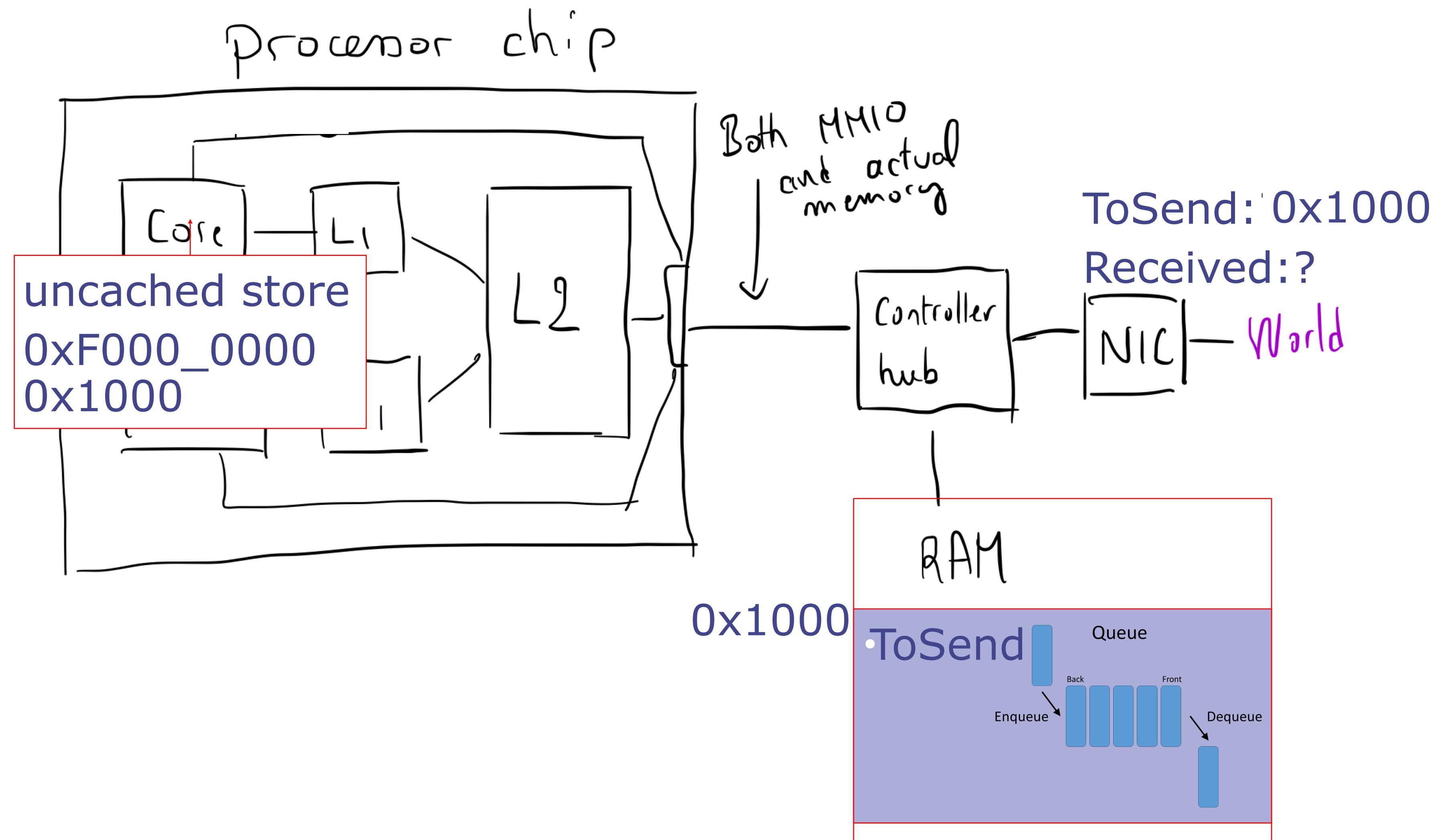
User-Agent: Mozilla/5.0

...”

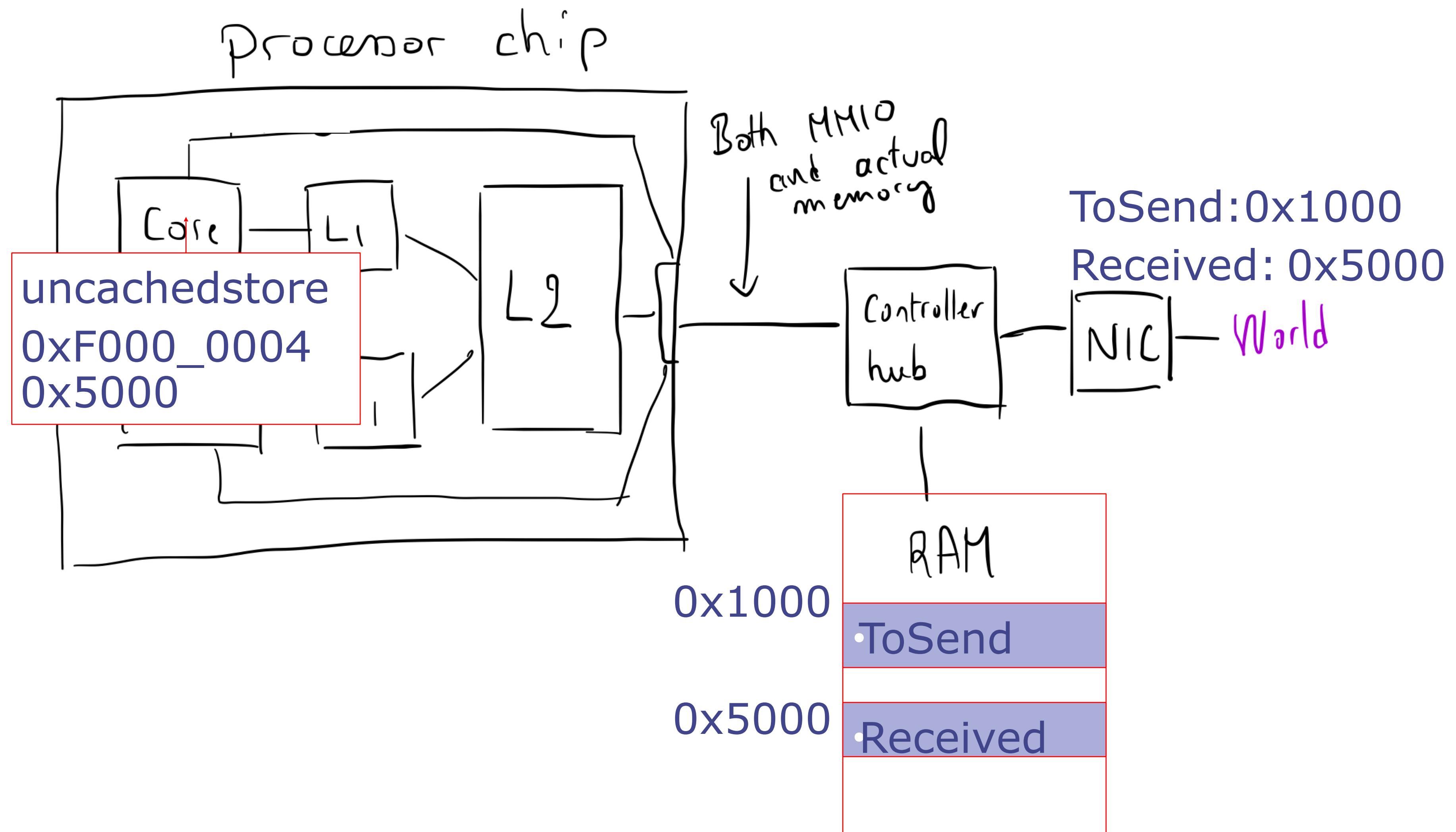
Wrap the string in an envelope with an address

Put the envelope in the SEND location

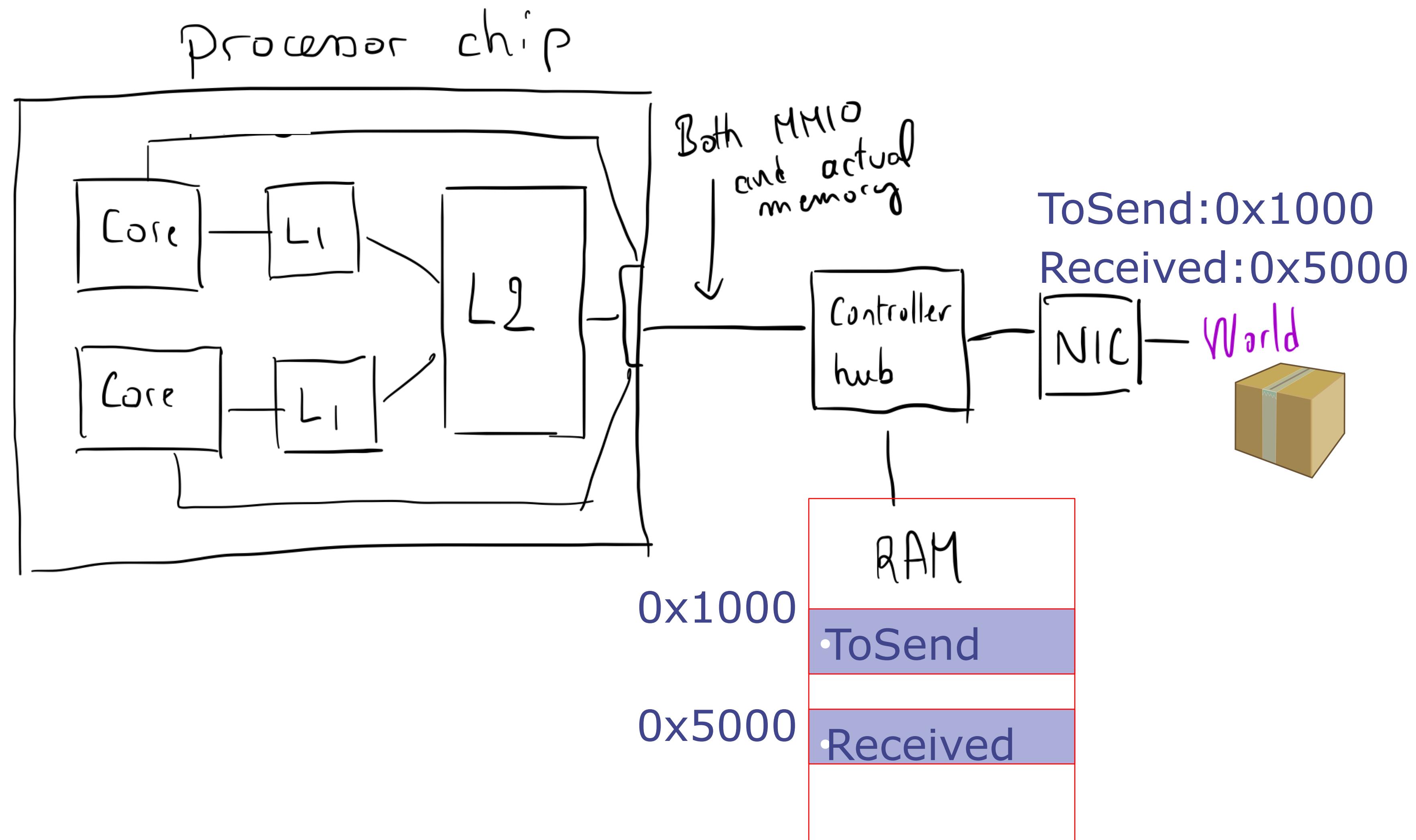
Agreeing on send locations



Agreeing on receive location



Simplified journey of a received packet



The NIC-CPU system

ToSend and Receive are in-mem circular buffers:

CPU push into ToSend and pull from Received

NIC push into Received and pull from ToSend

The NIC, once configured, operate independently of the core

The NIC performs computations – networking cost little CPU compute

Arrivals and departures

The challenge of synchronisation

How does NIC know when something should be sent?

How does CPU know when something arrived?

Two solutions:

Active polling – check every 1ms

Doorbell mechanism:

CPU -> NIC: store to special location

NIC -> CPU: Interrupts

Characteristic time: ~1us

MMIO vs Shared memory

Through MMIO the processor sends commands and pointers

The data is not sent directly through MMIO, but through shared memory

Why MMIO to NIC must be uncached?

Why stores to ToSend must be uncached?

Using Devices from SW

Write a program to manage the shared buffers with the device:

- Allocate buffers

- Recycle buffers

- Send address of buffers to device

- Doorbell code, etc...

Such a program is named a device driver!

Takeaway

A large class of accelerators interact with the rest of the systems like if they were devices.

Devices interact with the host through shared-memory.

Performance considerations - When accelerators can't do miracles

Latency considerations

On-chip latency:

1 cycle (registers memory) -~10 clock cycles (last level cache)

Off-chip (e.g. PCIE) latency:

1us (~4000 cycles)

A whole lot of CPU work!

Throughput considerations

PCIE throughput (gen 2):

500MB/s per lane (up to 16 lanes)

Reality: a few GB/s

CPU/DRAM bandwidth ~5-10x more (60-150GB/s)

Internal GPU bandwidth: ~1TB/s

- GPU companies work hard to avoid the PCIe bottleneck

More modern PCIe: 1GB/s (resp. 2GB/s) per lane for Gen3 (resp. Gen4)

Throughput takeaways

FPGA being disappointing

If a CPU already manages to saturate memory bandwidth, an accelerator won't go faster!

FPGA will typically go slower when because PCIE bandwidth < DRAM bandwidth ?

**Of the importance of the memory system:
“It’s the memory, stupid!”**

Richard L. Sites, Digital Equipment Corporation (Microprocessor Reports, 1996)

Some papers promise a better future for PCIE

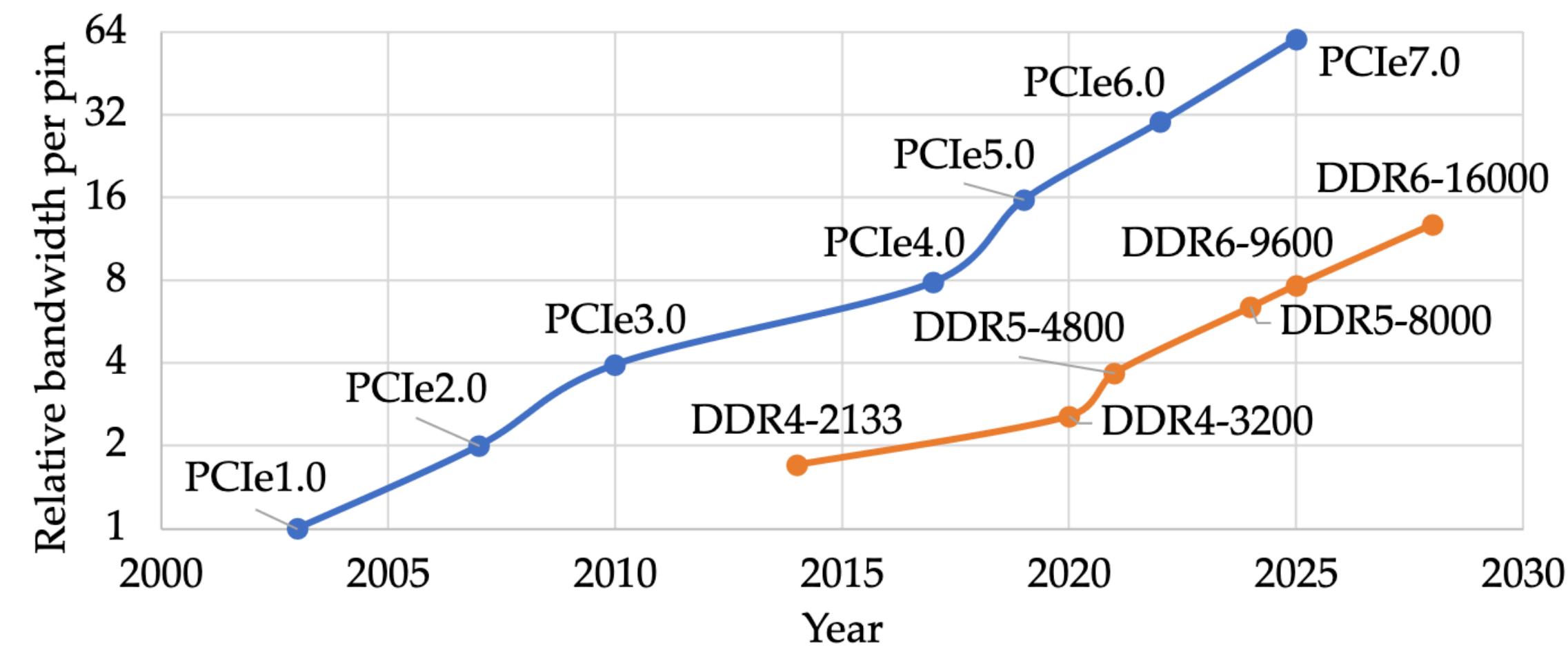


Figure 1: Bandwidth per processor pin for DDR and CXL (PCIe) interface, norm. to PCIe-1.0. Note that y-axis is in log scale.

Not everything is lost for accelerators

Sometimes crossing the pcie bus is already embedded in the nature of the task!

Bandwidth savings:

- offload database operator/analytics operators to the storage (ssd) itself
- compress on the fly to increase bandwidth

Latency/x86 resource savings:

- offload part of the network stack to the NIC itself (programmable NIC)
- offload RPC scheduling - bypass kernel

Why? In the datacenter, Jeff Dean and Luis Barroso said so in datacenter “Tail at Scale”

Principles of Accelerators

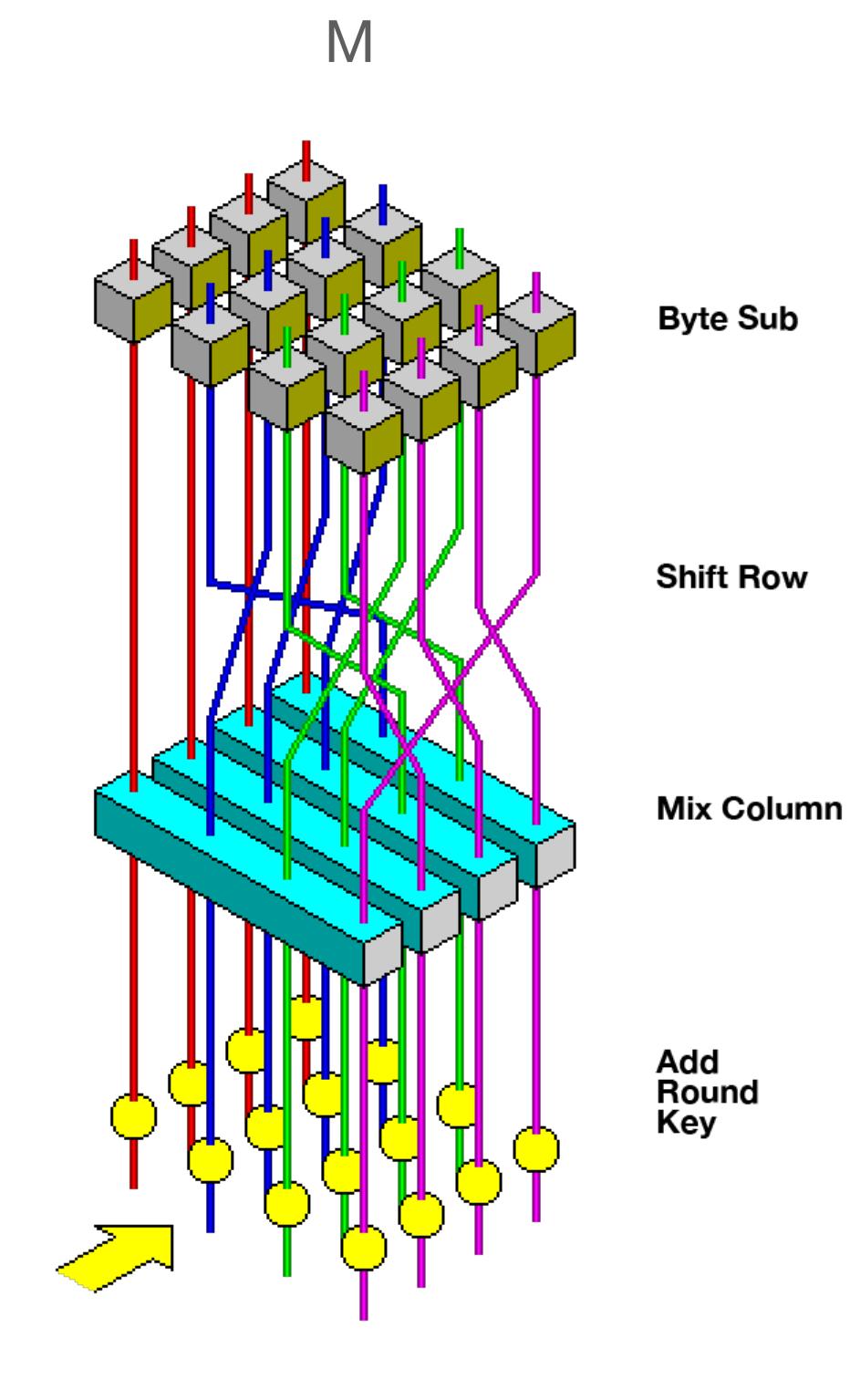
Guidelines for accelerators

(Borrowed from Bill Dally's talk)

1. Parallelism
2. Locality
3. Optimize Memory Orchestration and Control
4. Custom datatypes and operations

Parallelism

1. Look at the dataflow graph of your computation, at a fine granularity
2. What is the critical path of your computation, compared to the amount of compute nodes?
3. Does it make sense to pipeline the computation?



Locality - Arithmetic Intensity

Reuse the data you bring from memory multiple times!

Arithmetic Intensity is the ratio of arithmetic operations to data movement (bytes)

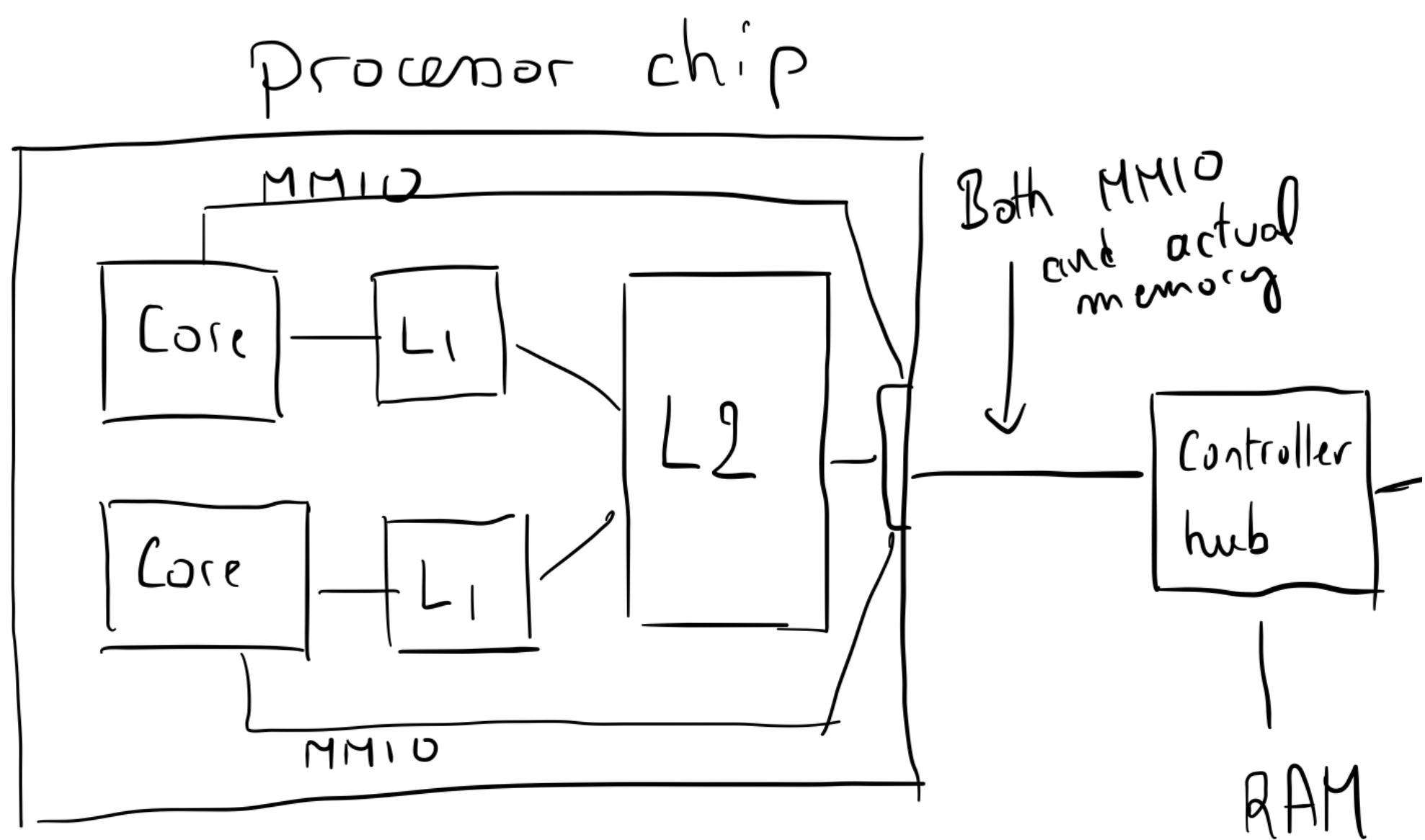
Low Arithmetic Intensity <-> Memory bottlenecked!

Dot product: not so high arithmetic density $\sim \frac{m}{2m} = O(1)$

Matrix multiply: much higher arithmetic density $\sim \frac{m^3}{2m^2} = O(m)$

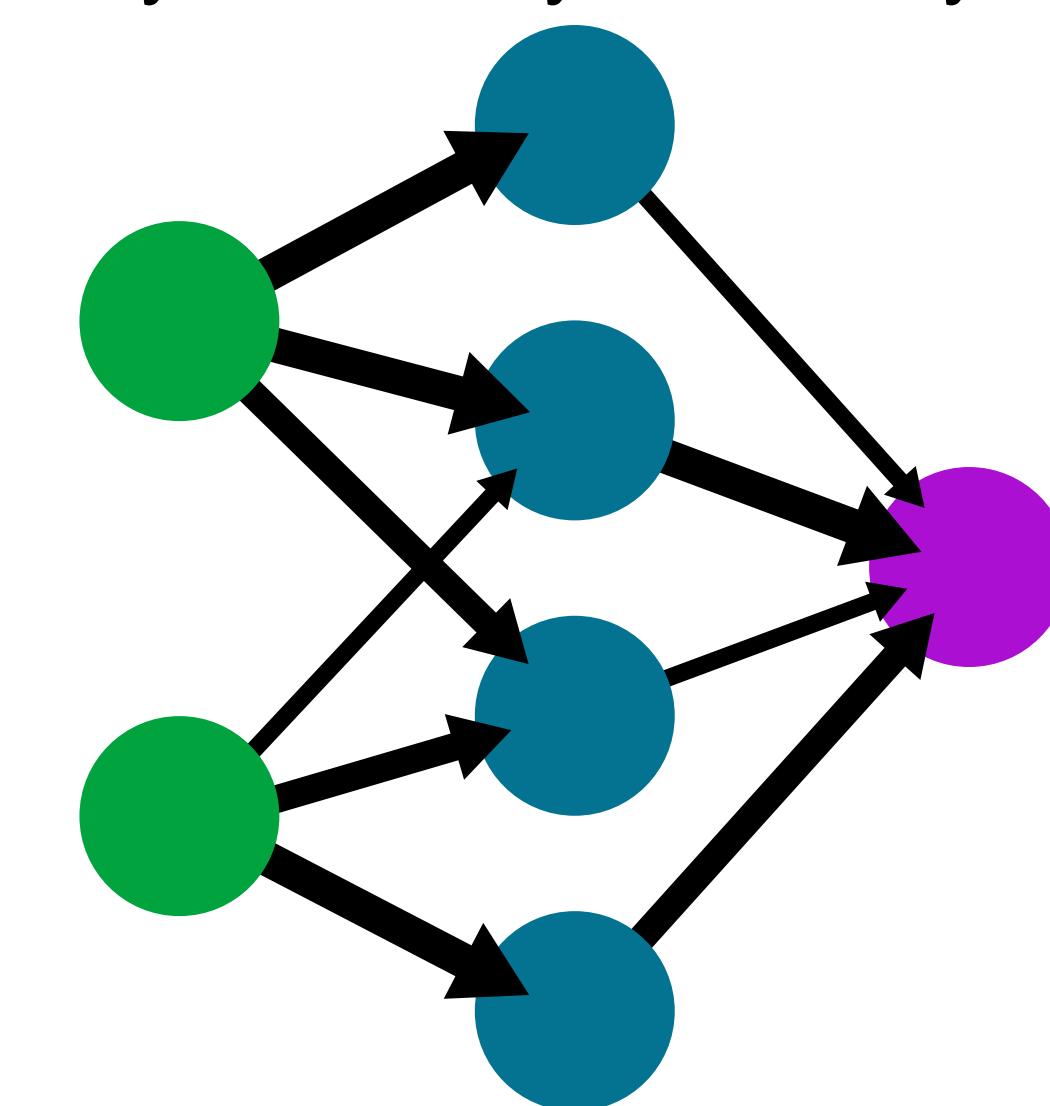
Optimize memory orchestration

Application specific “memory pipelining” / Decoupled execution



A simple neural network

input layer hidden layer output layer

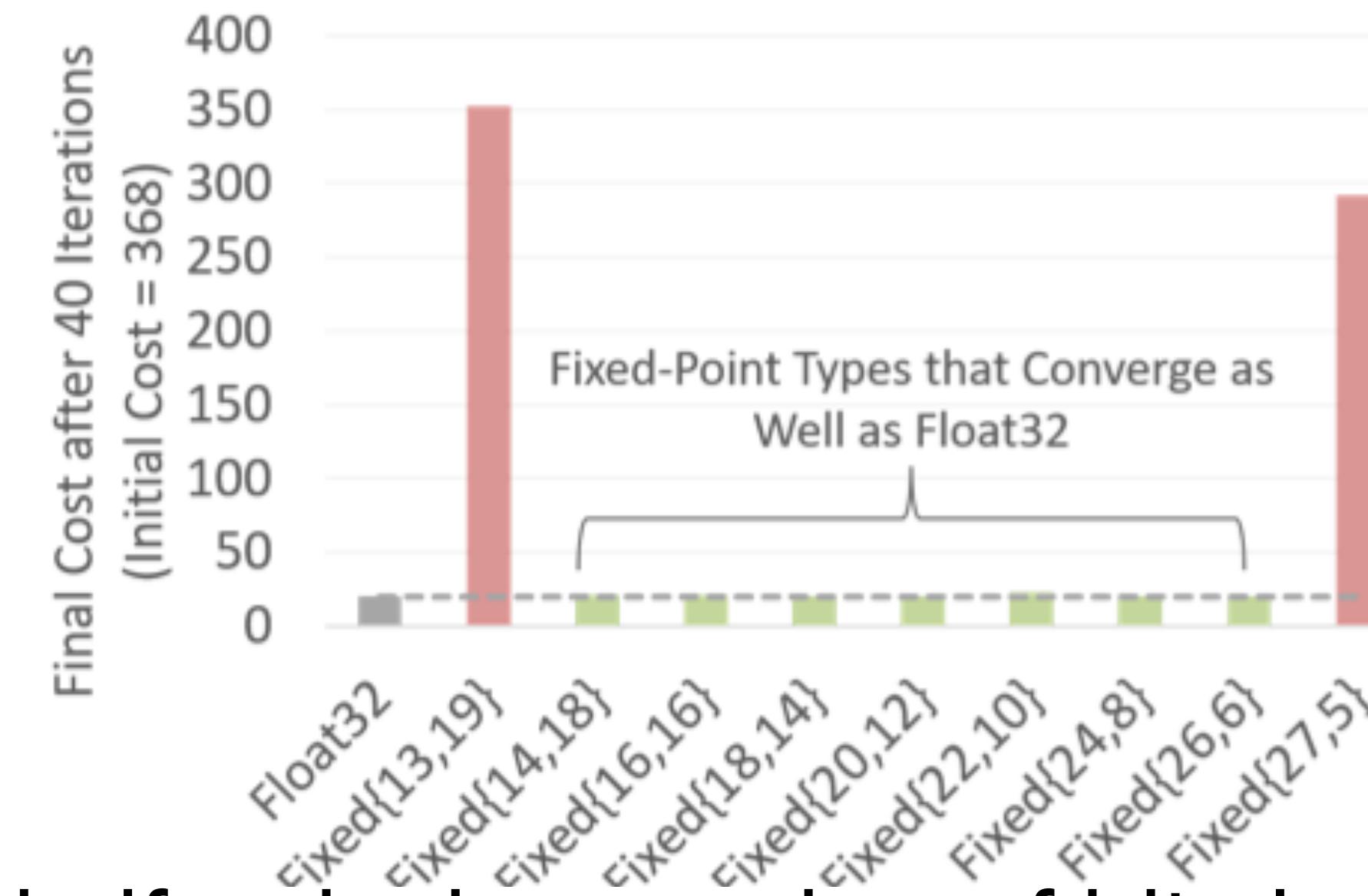


Specialized Datatypes

Software easily goes bloat when working on custom Datatypes with custom operations

Remember AES is expensive in SW because does weird things to 128 bits!

Domain specific analysis of structure of computation:



Modern ML -> lot to gain if reducing number of bits in representation

Joker:
CPU can achieve max speed (random access in large memory/streaming from large memory),
Accelerator = Energy Savings

The Cost of Specialization

Software Integration

The elephants are in the room

Having defined new instructions is not the end of the story:

If I add AESENC, my libssl library won't suddenly start using it

Have to worry about cross-platform

If I have matmul 16x16, or matmul (nxm forall n,m<64),

512x512 matmul?

convolution?

Arbitrary linear/tensor algebra operator

Drawback of Old School HW/SW codesign

Take application X, handwrite code that leverage Accelerator W (V1), performs great

[...] Time elapse

10 years later, Accelerator W (V13)

-> Handwritten code performs poorly (when it works)

Compilation challenges

The elephants are in the room

All things considered standard compilation is easy:

Usually not too many different ways to compile

With Domain Specific, program space is typically very complicated

$A^*(B^*C)^*A$ or $(A^*B) * (C * A)$...

Compiler must find good sequence of instructions modulo rewrites!

Every domain has its own set of rewrite - every time requiring a new compiler

Compiler must consider an accurate cost model of memory!

Compilation – Sync issues

The elephants are in the room

Sync issues:

Need to add explicit data movement instructions if I want to use the data computed on CPU

Hot research ideas:

Decoupling functionality and scheduling/performance [Halide/Exo]

Maybe not a fully push-button compilation, more an “assistant-compiler”

Maybe enable user to augment the compiler - easily add transformations with triggers

Conclusion

“Accelerators” and Codesigns are not new - Yesterday they enabled census, today they are enabling AI

There are costs to abstraction

There are good guidelines/models for what can be accelerated

There are costs to specialization

There is still work to do to better understand those costs in general

Accelerators 100 years later

SIGCOMM 2023

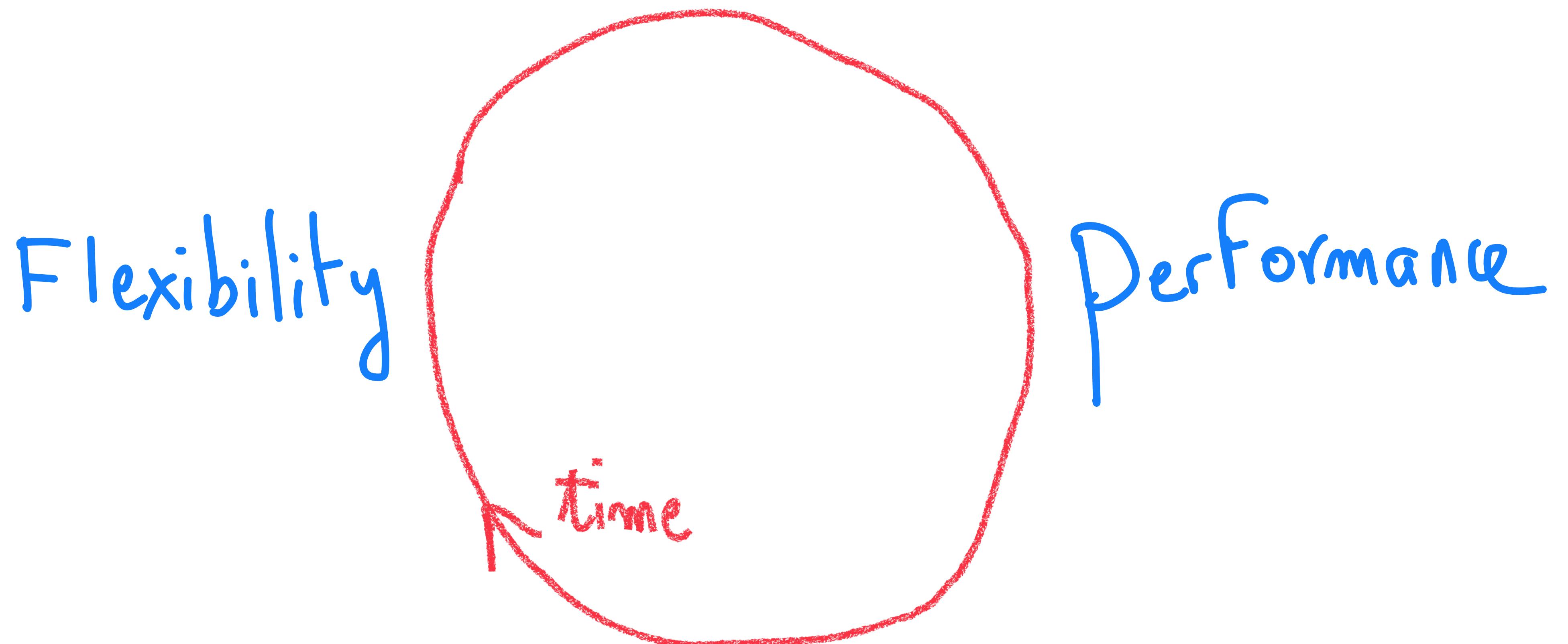
LIGHTNING: A Reconfigurable Photonic-Electronic SmartNIC for Fast and Energy-Efficient Inference

Zhizhen Zhong Mingran Yang Jay Lang Christian Williams Liam Kronman
Alexander Sludds Homa Esfahanizadeh Dirk Englund Manya Ghobadi

2.1 Photonic Vector Dot Product

Same good old $\sum_i w_i \cdot x_i$, with photons

Outro



The Missing Quote

Von Neumann's condition

5.6 Accelerating these arithmetical operations does therefore not seem necessary —at least not until we have become thoroughly and practically familiar with the use of very high speed devices of this kind, and also properly understood and started to exploit the entirely new possibilities for numerical treatment of complicated problems which they open up. Furthermore it seems questionable whether

Sources and Inspiration

Bill Dally :

<https://www.youtube.com/watch?v=fnd05AeeFN4>

David Patterson/John Hennessy Turing Award Speed:

<https://www.youtube.com/watch?v=3LVeEjsn8Ts>

Chris Lattner & al. , LLVM Dev Mtg (10 days ago):

<https://www.youtube.com/watch?v=SEwTjZvy8vw>