



Building Abstractions at the  
Hardware-software Boundary

HELLO CPPNorth

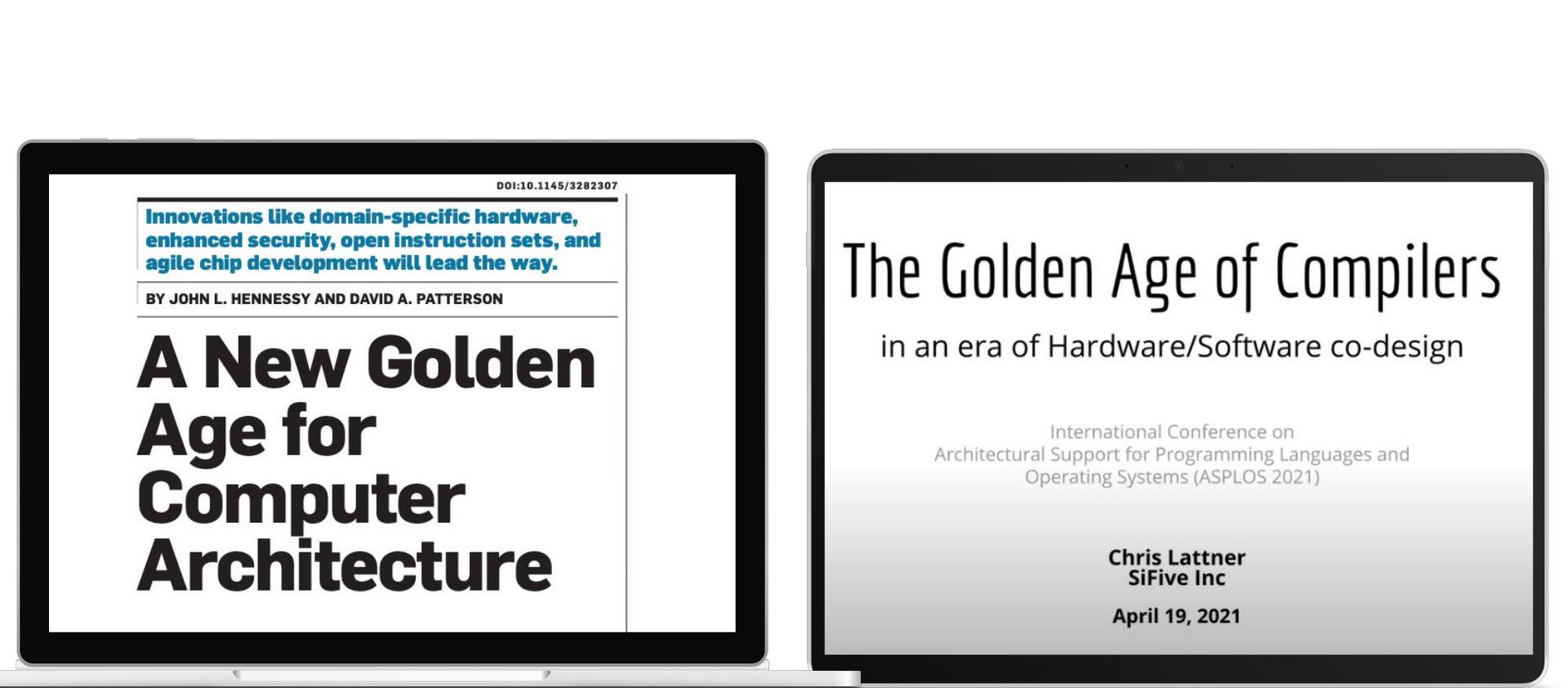
# **Andrew Bitar**

Compiler Technical Lead & Manager

[abitar@groq.com](mailto:abitar@groq.com)

## MY GOAL TODAY

How should we approach  
building **abstractions** in a  
computing era beyond CPUs?



DOI:10.1145/3282307

**Innovations like domain-specific hardware,  
enhanced security, open instruction sets, and  
agile chip development will lead the way.**

BY JOHN L. HENNESSY AND DAVID A. PATTERSON

# A New Golden Age for Computer Architecture

# The Golden Age of Compilers

in an era of Hardware/Software co-design

International Conference on  
Architectural Support for Programming Languages and  
Operating Systems (ASPLOS 2021)

Chris Lattner  
SiFive Inc

April 19, 2021

turing lecture

In Computing,  
Abstractions Are  
**Ubiquitous**

# Abstractions Are Ubiquitous

**Applications** Functions, types, inheritance, template metaprogramming

**Programming Languages** Domain-specific languages (e.g. PyTorch, TensorFlow, P4)

Functional programming languages (e.g. Haskell)

Scripting languages (e.g. Python)

Augmenting lower-level languages with higher-level constructs (C → C++ → C++20)

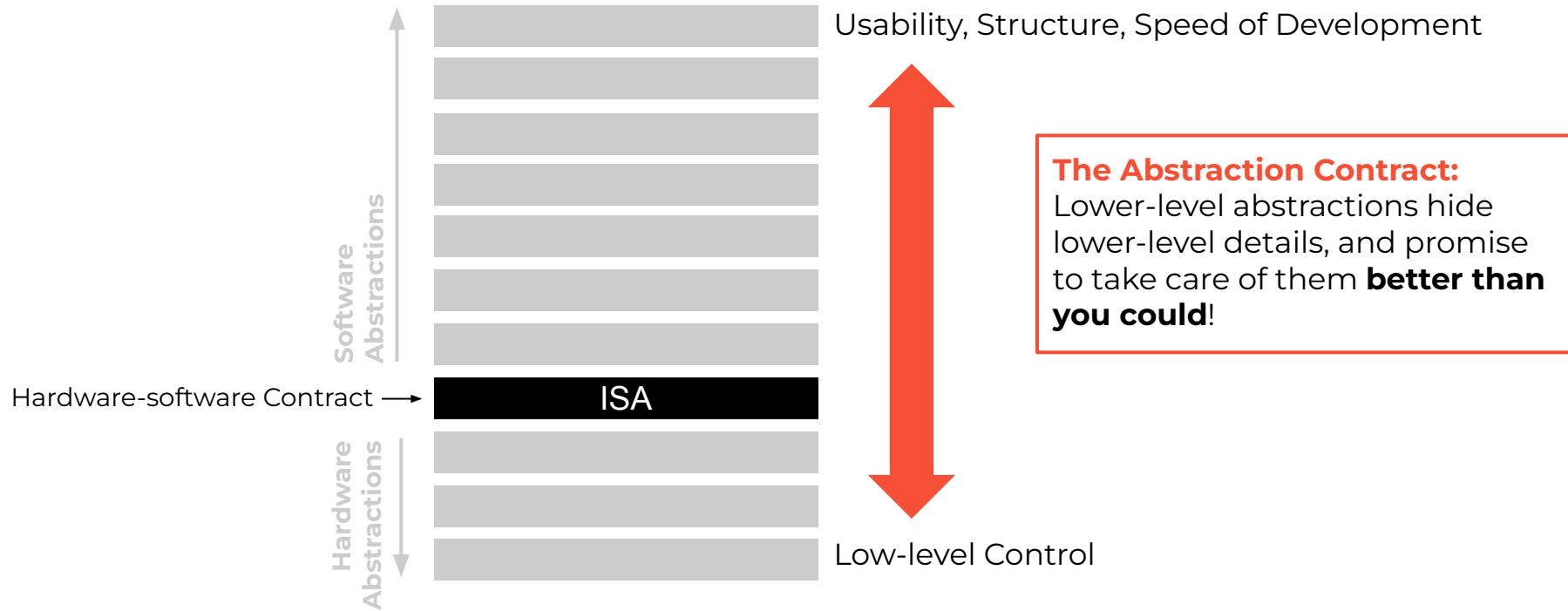
**Compilers** Multi-pass compilers, LLVM, MLIR

Gradual IR “lowering” from higher-level representations down to machine language

**Hardware** CAD flows (HDL → Tech-Mapped Netlist → P&R)

Logic design (block hierarchies, interfaces)

# The Abstraction Trade-off



# The Abstraction Trade-off

Abstractions rarely come free

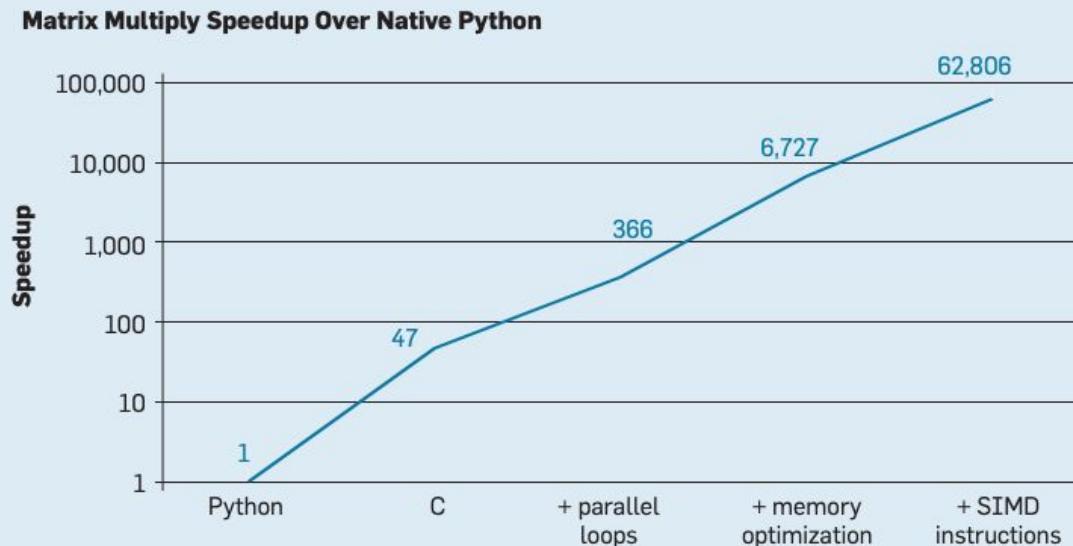
“There are no zero-cost abstractions.”

Chandler Carruth, CppCon’19

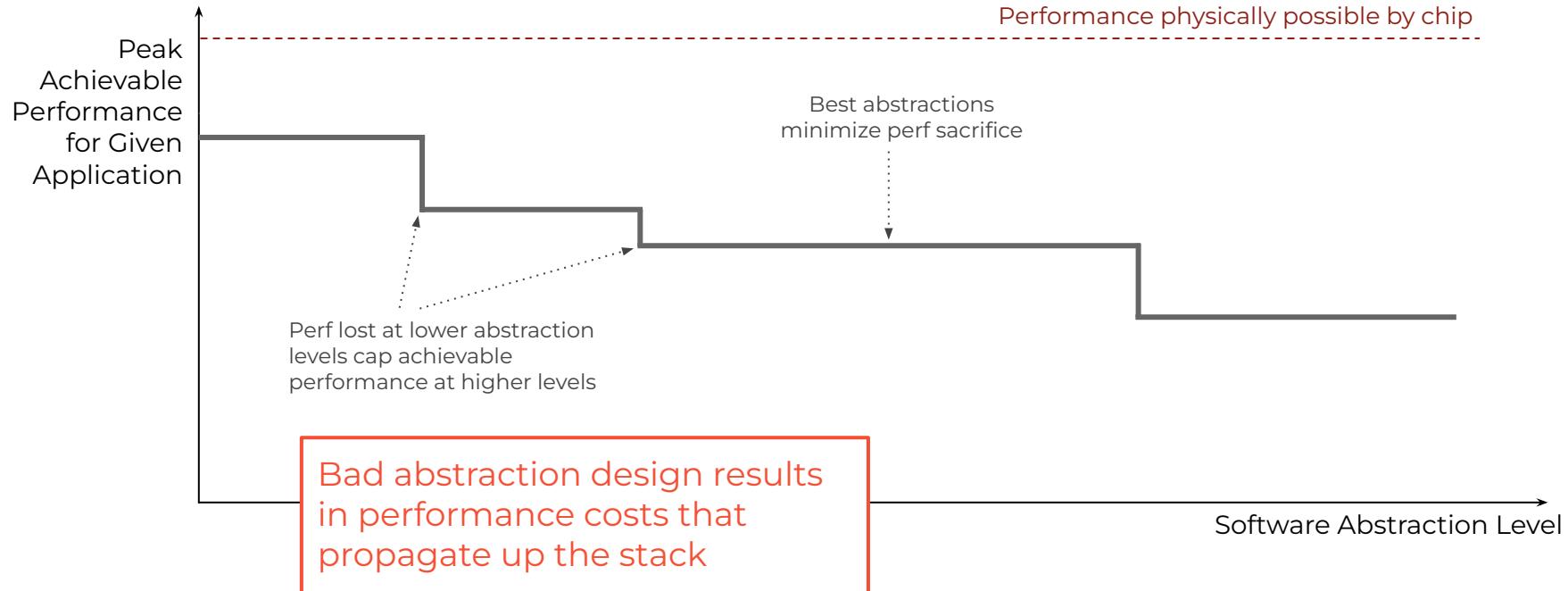
Even leveraging an abstraction as simple as `std::unique_ptr` can come at a runtime cost!

Hennessy and Patterson. A New Golden Age for Computer Architecture. ISCA’18.

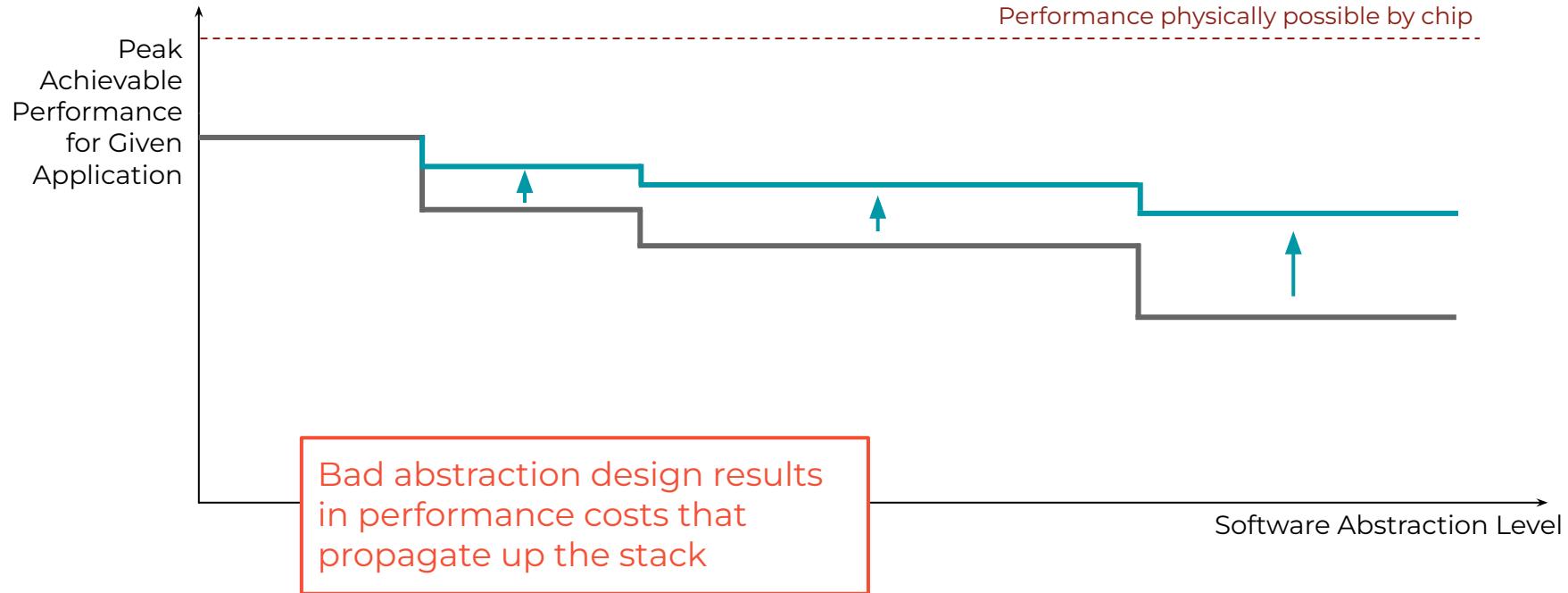
Figure 7. Potential speedup of matrix multiply in Python for four optimizations.



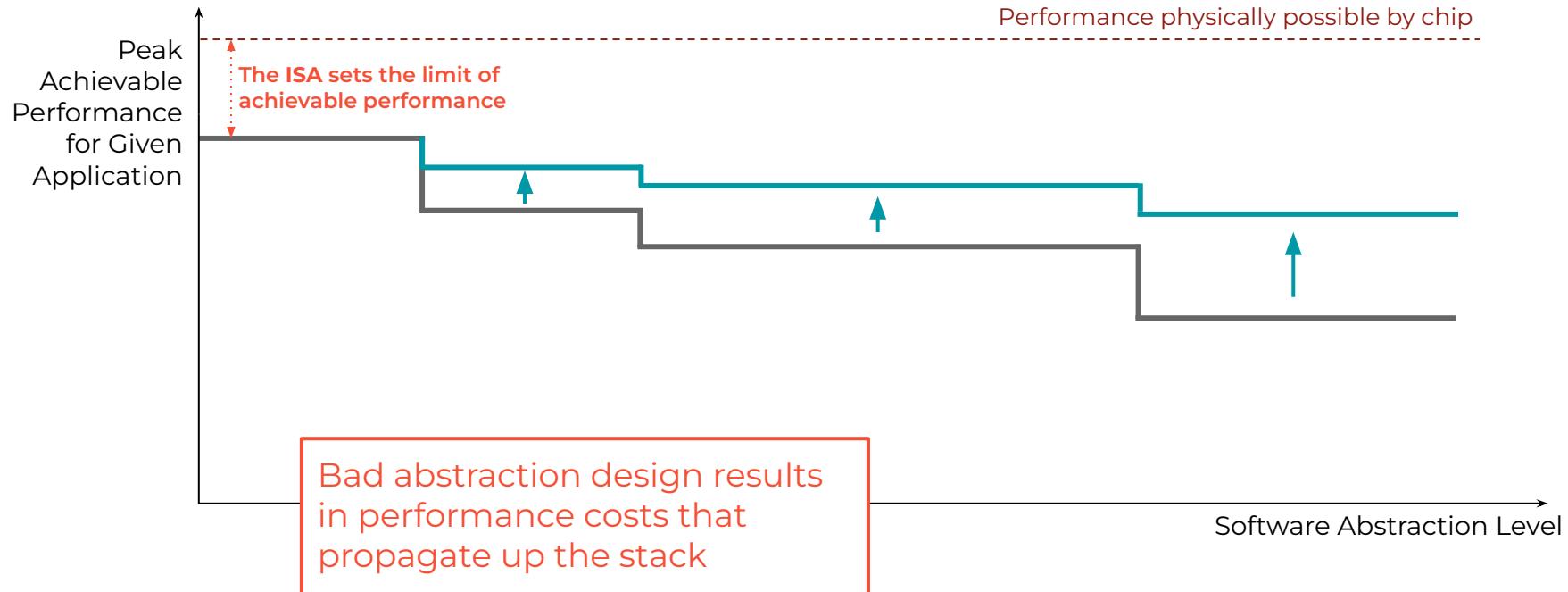
# The Abstraction Trade-off



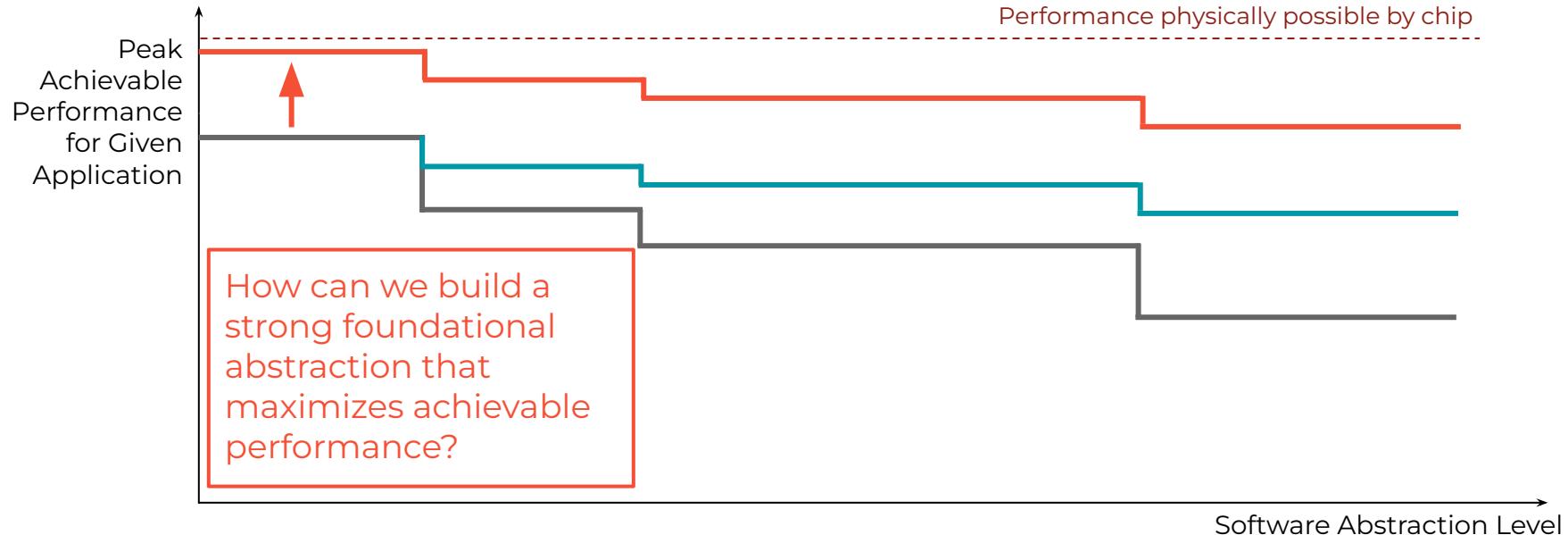
# The Abstraction Trade-off



# The Abstraction Trade-off



# The Abstraction Trade-off



ARCHITECTURE

# The Bedrock of Abstractions

# The First Golden Age

For computer architecture

**Complex Instruction Set Computer (CISC) was born in an era when programs were primarily written manually using assembly languages**

Complex instructions that performed multiple low-level operations (e.g. load + ALU + store) helped **raise the abstraction level** for assembly writers while also reducing instruction memory footprint

With the advent of high-level languages, the question changed:

- “What assembly language would programmers use?” became “What instructions would compilers generate?”



# CISC → RISC

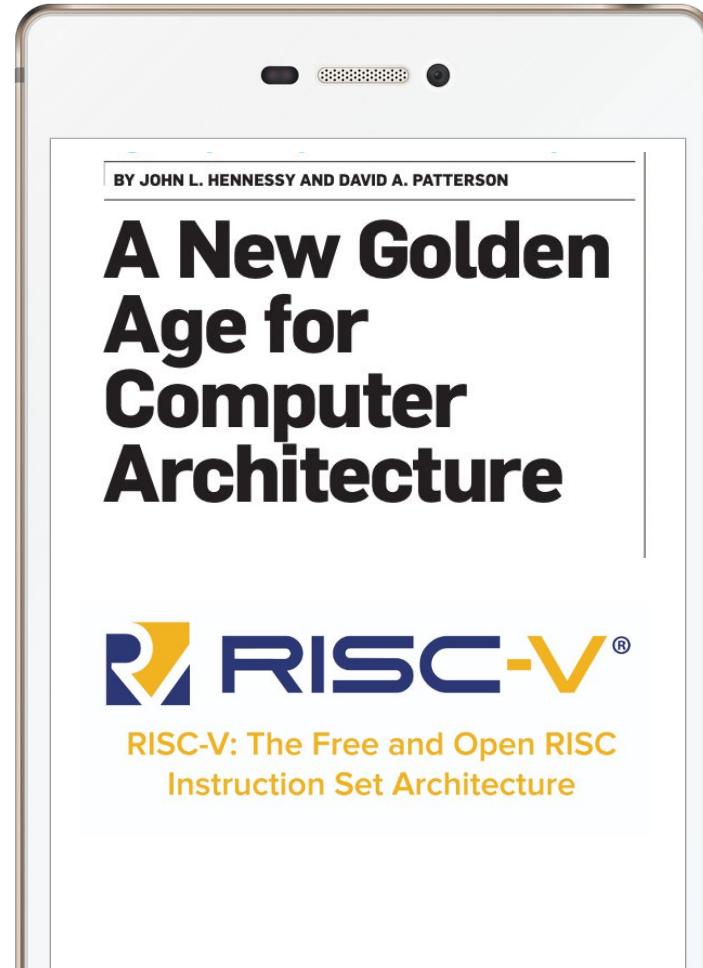
Computer architects began observing that using a Reduced Instruction Set Computer (RISC) improved performance and efficiency

Moving from CISC to RISC **lowers the abstraction level of the ISA**

- Moving complex control logic from hardware to software → more chip area available for compute and memory

As noted in the Hennessy & Patterson paper - experiments showed CISC ISAs using 75% of instructions per program, but executed 5–6 more cycles per instruction → **overall RISC was 4x faster**

- Software able to make better decisions when given lower-level control of the hardware



# RISC → VLIW?

**Even with RISC, CPU architecture hides a lot of complex control logic from software**

Caching, prefetching, out-of-order execution, branch prediction

- Area cost of this control hardware reduces available area for raw compute

## **Very Long Instruction Word (VLIW) to the rescue?**

The next natural step in the trend of shifting work and control from hardware to software

- Execution ordering, data prefetching, branch resolution
- Compiler responsible for exploiting instruction-level parallelism (ILP) available in a program

## **Itanium...**

Compiler struggled to achieve good ILP due to limited view of dynamic hardware behaviour

- Cache misses, dynamic mem dependencies, dynamic branches

# The Hardware Lottery

**In 2012, AlexNet was the breakthrough for Deep Neural Networks**

Convolutions and Backprop have existed since the 1980s

But CPUs were poorly suited for DNNs

- CPUs pay significant chip area for handling control-flow dominated programs
- DNNs are data-flow heavy workloads, with little-to-no dynamic control-flow

**Evolution of GPUs in the 2000s opened up an architecture better suited for dataflow-heavy workloads like DNNs**

Reduce dedicated control-flow logic = increased chip area dedicated to data-flow compute (i.e. fixed- and floating-point multipliers and adders)

**Watch Sara Hooker's GroqDay talk: [groq.link/gd2hwlottery](https://groq.link/gd2hwlottery)**

The image shows a smartphone with a white case and a black screen. On the screen, the title 'The Hardware Lottery' is displayed in a large, bold, black font. Below it, the author's name 'Sara Hooker' is in a smaller black font. Underneath that, it says 'Google Research, Brain Team' and provides an email address 'shooker@google.com'. At the bottom of the screen, the word 'Abstract' is visible. The background of the slide features a large orange bar at the bottom containing the call-to-action text.

**The Hardware Lottery**

Sara Hooker

Google Research, Brain Team  
shooker@google.com

**Abstract**

Hardware, systems and algorithms research communities have historically had different incentive structures and fluctuating motivation to engage with each other explicitly. This historical treatment is odd given that hardware and software have frequently determined which research ideas succeed (and fail). This essay introduces the term hardware lottery to describe when a research idea wins because it is suited to the available software and hardware and *not* because the idea is superior to alternative research directions. Examples from early computer science history illustrate how hardware lotteries can delay research progress by casting successful ideas as failures. These lessons are particularly salient given the advent of domain specialized hardware which make it increasingly costly to stray off of the beaten path of research ideas. This essay posits that the gains from progress in computing are likely to become even more uneven, with certain research directions moving into the fast-lane while progress on others is further obstructed.

# GPUs → GPGPUs?

**What started as a  
Domain-specific Accelerator  
(DSA) began to change to a more  
general-purpose processor**

Moving some control-flow  
responsibilities to software [e.g.  
execution order, prefetching]

- Didn't we try this with VLIW?  
What's changed?

Emergence of Software 2.0  
[see: Andrej Karpathy] → data-flow  
dominated workloads as better  
alternatives to solving problems



# A New Golden Age in Computing

## The end of CPU hegemony

Growing demand in dataflow dominated compute

Slowing of Moore's Law and Dennard Scaling

CPU "abstraction" no longer the only foundation for developing software

## Dawn of a new golden age in computing

Hennessy & Patterson: "A New Golden Age for Computer Architecture"

Lattner: "A New Golden Age of Compilers"

Karpathy: "Software 2.0" → A New Golden Age for Algorithms

## In the past, hardware (CPUs) has defined our software

The hardware-software abstraction contract has been reopened!

Opportunity now to achieve "software-defined hardware"

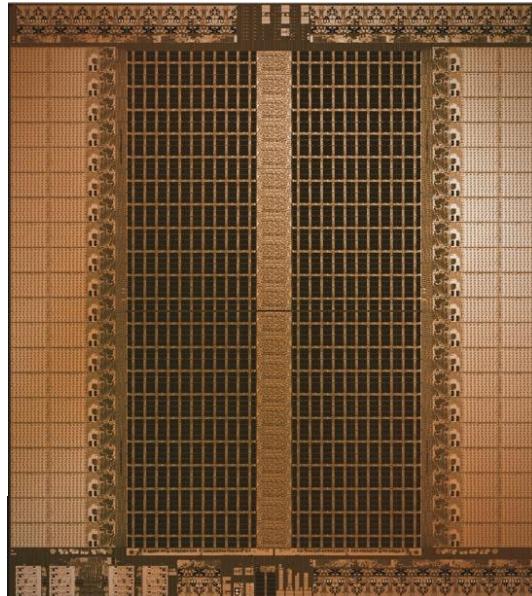
# Software-defined Hardware

# Groq is Radically Simplifying Compute

Less  
hardware  
control  
=

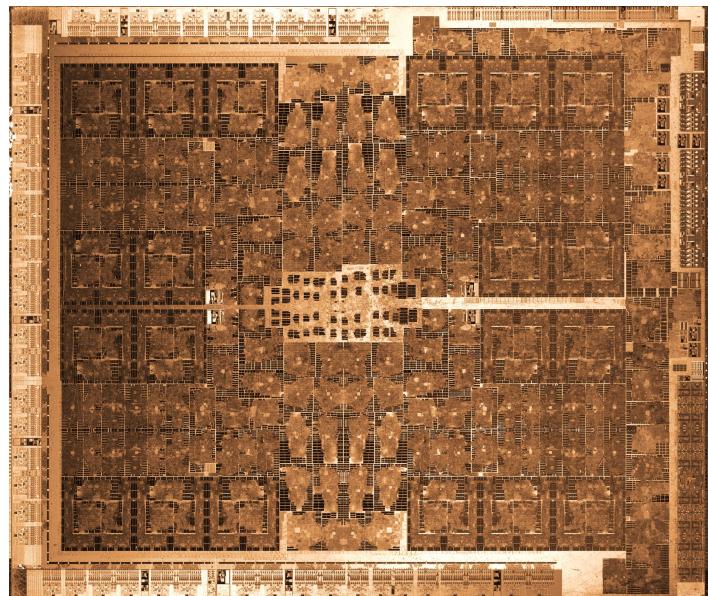
more  
compute &  
memory!

**GroqChip™ 1**

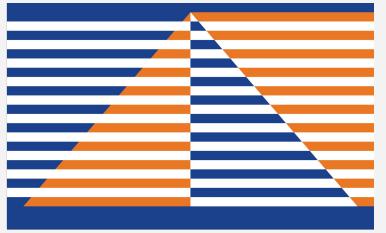


Simplified design enables  
**Compute Performance**

**Competitive Chip Example**



Complex design for processing data results in  
**Compute Costs**



2020

Appears in the *Proceedings of the International Symposium on Computer Architecture* (ISCA-47), June 1, 2020.  
Copyright 2020 IEEE Computer Society.

# Think Fast: A Tensor Streaming Processor (TSP) for Accelerating Deep Learning Workloads

Dennis Abts, Jonathan Ross, Jonathan Sparling, Mark Wong-VanHaren, Max Baker, Tom Hawkins, Andrew Bell, John Thompson, Temesghen Kahsai, Garrin Kimmell, Jennifer Whay, Rebekah Leslie-Hurd, Michael Bye, E.R. Creswick, Matthew Boyd, Mahitali Venigala, Evan Laforgo, Jon Purdy, Purushothaman Kamath, Dinesh Maheshwari, Michael Beidler, Geert Rossel, Oam Ahmid, Gagan Gagarin, Richard Czechalski, Ashutosh Rane, Sabil Parmar, Jeff Werner, Jim Sproch, Adrian Macias, Brian Kurtz

**Groq, Inc. · Mountain View, California**

**Abstract**—In this paper, we introduce the Timer Streaming Processor (TSP) architecture, a functionally-sparse microprocessor with memory bandwidth limited by vector and matrix deep learning operations. The TSP is designed to support a wide range of functionally sparse applications, such as those requiring deep learning operations. The TSP is built based on two key observations: (1) machine learning workloads exhibit a high degree of parallelism, and (2) the computation in hardware is fast, parallelized, and deterministic. The TSP is designed with productivity in mind, allowing users to quickly model, validate, and control their hardware configurations. The TSP is designed to exploit parallelism inherent in machine-learning applications, including instruction-level, memory-concurrent, and memory-parallel execution. The TSP is designed to be driven by elements at all active elements in the hardware (e.g., arbiters, and caches). Early ReNe50 image classification results demonstrate 20.4K GFLOPs per square millimeter, and 1.4 GFLOPs per square millimeter for 44 layers. Our first ASIC implementation of the TSP achieved 1.4 GFLOPs per square millimeter, and 0.7 GFLOPs per square millimeter for 44 layers. The TSP is designed to operate at a nominal silicon frequency of its 25.25  $\times$  25.14 mm chip area, yet a predicted power envelope of 1.1 mW/m<sup>2</sup> on machine-learning workloads.

The world is increasingly turning to computationally-intensive deep learning algorithms to solve important problems in science, transportation, security, and beyond. These workloads continue to grow both in size and complexity, presenting serious scalability, performance, and usability challenges for traditional CPU and GPU architectures. Keeping pace with this demand requires that we provision abundant on-chip memory and computation at near-zero power consumption. However, current execution, memory, and bandwidth approaches in microarchitectures make it difficult to reason about system runtime stalls. Furthermore, while microarchitectural enhancements such as caches, branch predictors, and prefetchers help tremendously in improving performance, they do not bound worst-case performance.

Over the last decade, data center operators have installed many-core systems as a stalwart fixture in warehouse-scale

The diagram illustrates the functional blocks of the Memory controller and their connections to the Image ALU Instruction pipeline. The Memory controller consists of four main functional units: Memory controller logic, Memory controller interface, Memory controller memory, and Memory controller cache. The Memory controller logic is connected to the Image ALU Instruction pipeline via a bidirectional bus. The Memory controller interface connects to the Memory controller memory and cache. The Memory controller memory and cache are also interconnected. The Memory controller interface is further connected to the Memory controller logic and the Memory controller cache.

Fig. 1. Conventional 2D mesh of cores (a) reorganized into a functionally

2022

# A Software-defined Tensor Streaming Multiprocessor for Large-scale Machine Learning

Dennis Abts	Garrison Kimmell	Andrew Ling	John Kim
Groq Inc.	Groq Inc.	Groq Inc.	KAIST/Groq Inc.
Matt Boyd	Andrew Bitar	Sahil Parmar	Ibrahim Ahmed
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Roberto DiCicco	David Han	John Thompson	Michael Bye
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Jennifer Hwang	Jeremy Fowers	Peter Lillian	Ashwin Murthy
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Elyas Mehtabuddin	Chetan Tekur	Thomas Sohmers	Kris Kang
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
	Stephen Maresh	Jonathan Ross	
	Groq Inc.	Groq Inc.	

## ABSTRACT

We describe our novel commercial software-defined approach for large-scale interconnection networks of tensor streaming processing (TSP) elements. The system architecture includes packaging, routing, and flow control of the interconnection network of TSPs. We describe the communication and synchronization primitives of a bandwidth-rich substrate for global communication. This scalable communication fabric provides the backbone for large-scale distributed systems.

CCS CONCEPTS

- Computer systems organization → Interconnection architectures.

N

the TSB consists of a variety of workloads, both memory and compute. We have shown that the TSB can support a wide range of streaming programmes, able to include global memory which is implemented as logically shared, but physically distributed DRAM on chip memory. Each TSB contributes 220 MB/bytes to the global memory capacity, with the TSB contributing up to 10% of the total system memory. The TSB also provides a highly programmable endpoint for each memory buffer of endpoints in the system. The TSB acts as both a tensor processor (endpoint) and network switch by moving tensors across the connections. We describe a novel software-controlled networking approach to avoid the need for a separate network interface card, controlled by dynamic contention for network links. We describe the topology, routing and flow control to characterize the performance of the network for a large-scale parallel image learning system, achieving up to 14.4 GFLOPs and more than 2 Terabytes of global memory accessible in less than 3 microseconds of end-to-end system latency.

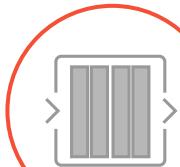
to make digital or hard copies of part or all of this work for personal or internal use is granted without fee provided that copies are not made or distributed for commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for third-party components of this work must be honored. Other uses, contact the owner/author(s).  
22, June 18–22, 2022, New York, NY, USA  
Copyright held by the owner/author(s).  
ISBN 978-1-4503-8610-2/22/06  
[doi.org/10.1145/3479456.3527405](https://doi.org/10.1145/3479456.3527405)

vidual processing elements work in unison to collectively execute the different "layers" of a deep learning network. It is this set of sub-tasks, expressed as individual PE programs, that are distributed among the computing elements and responsible for carrying out, or executing, the specifics of the machine learning model.

# GroqChip Scalable Architecture

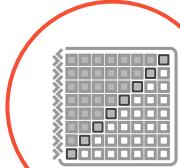
## SRAM Memory

Massive concurrency  
80 TB/s of BW  
Stride insensitive



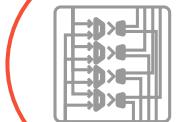
## Groq TruePoint™ Matrix

4x Engines  
320x320 fused dot product  
Integer and floating point



## Programmable Vector Units

5,120 Vector ALUs\* for  
high performance



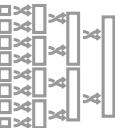
## Networking

480 GB/s bandwidth  
Extensible network scalability  
Multiple topologies



## Data Switch

Shift, Transpose, Permuter for  
improved data movement and  
data reshapes

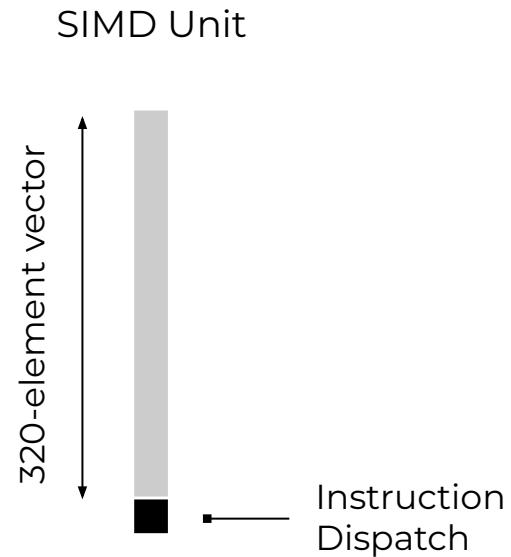


## Instruction Control

Multiple instruction queues for  
instruction parallelism



# GroqChip Building Blocks



# GroqChip Building Blocks

Build different types of specialized SIMD units



**MXM**  
Matrix-Vector /  
Matrix-Matrix Multiply



**VXM**  
Vector-Vector  
Operations



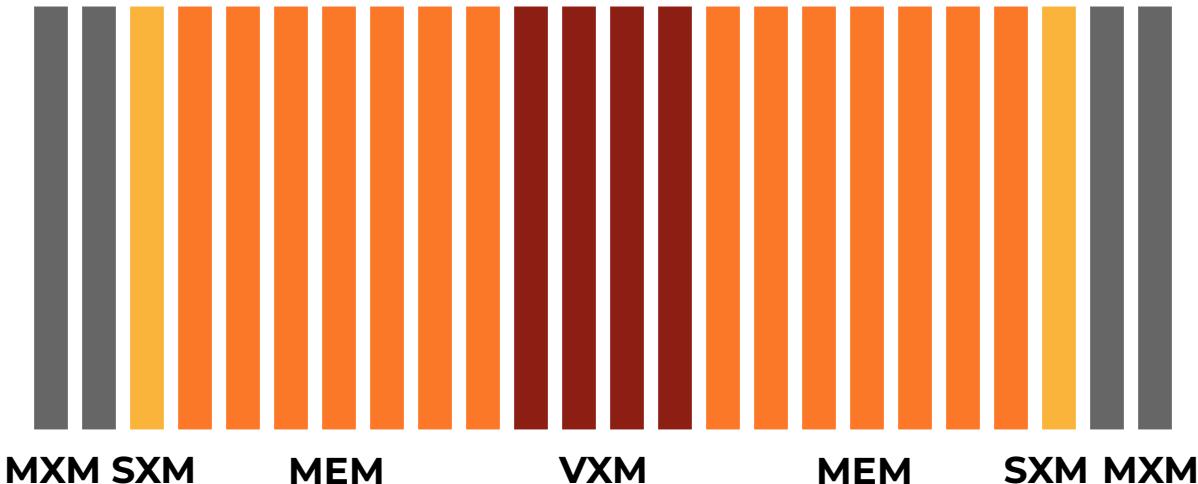
**SXM**  
Data Reshapes



**MEM**  
On-chip SRAM

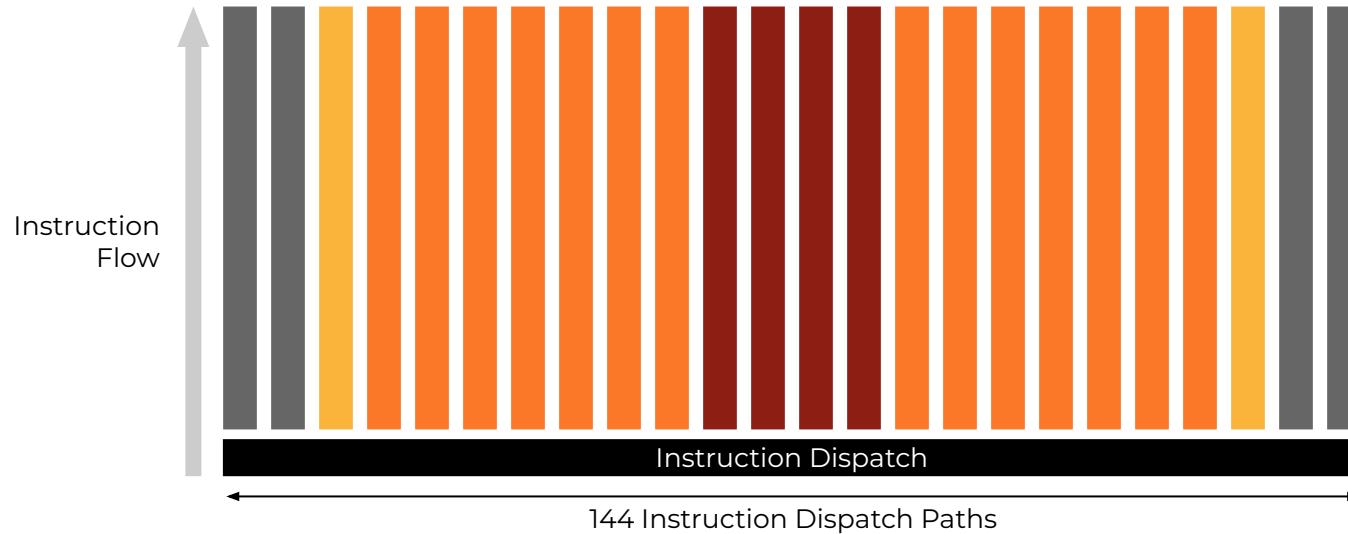
# GroqChip Building Blocks

Lay out SIMD units across chip area



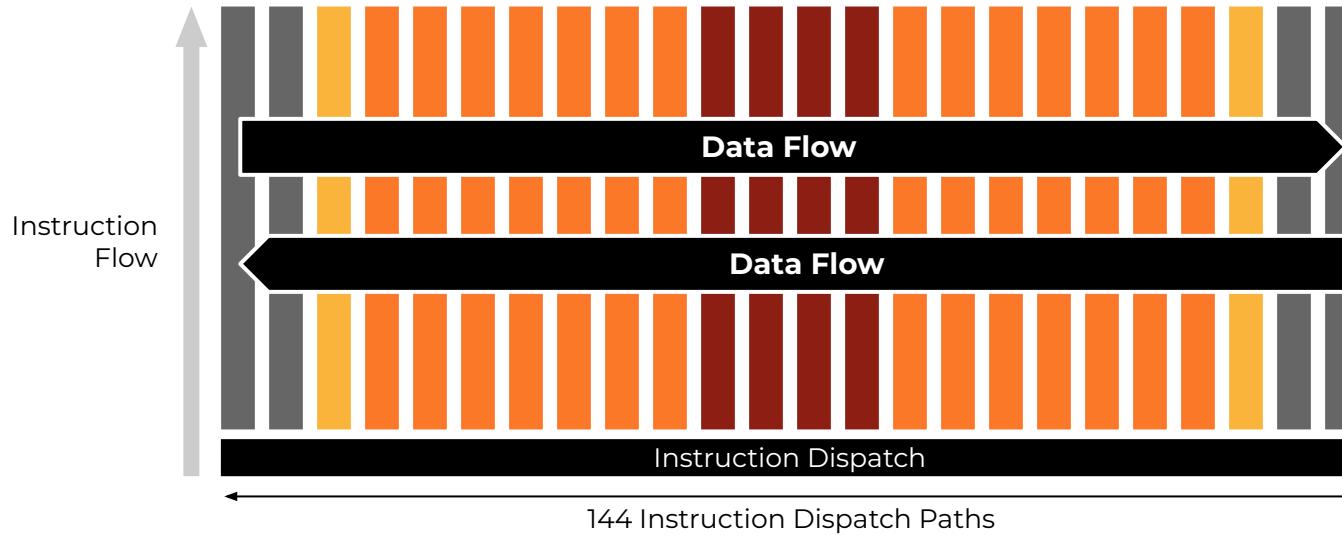
# GroqChip Building Blocks

Synchronized instruction dispatch across all SIMD units for lockstep execution



# GroqChip Building Blocks

High-bandwidth “Stream Registers” for passing data between units



# BUILDING A Compiler For GroqChip

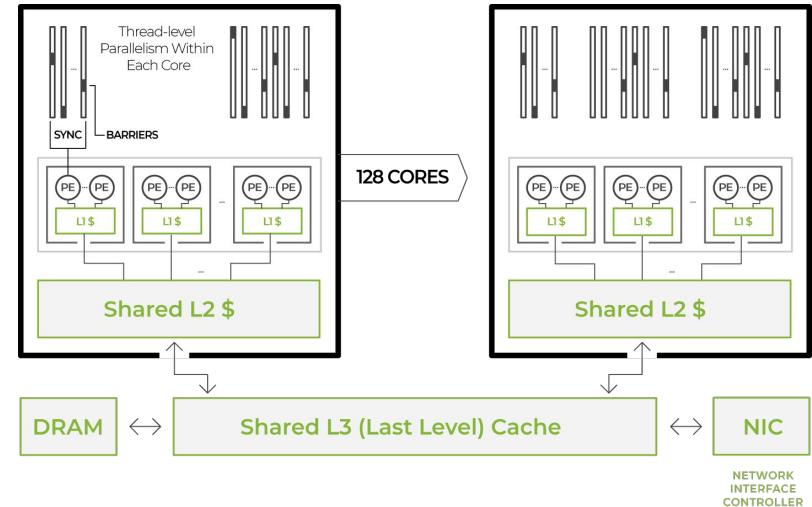
# Building Compilers For Spatial Architectures

**Challenges when building an automated compiler for spatial architectures**

## Unpredictable memory latency

- Non-deterministic memory (caching, DRAM) reduces compiler ability to hide memory latency

A100



# Building Compilers For Spatial Architectures

## Challenges when building an automated compiler for spatial architectures

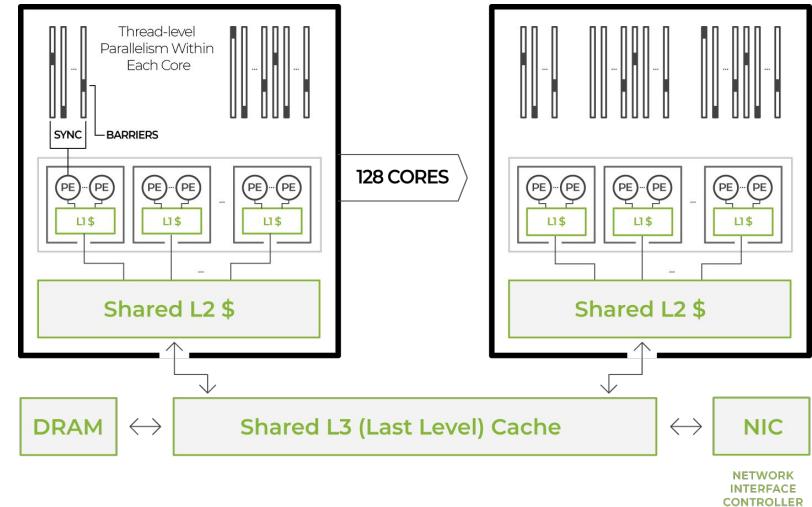
### Unpredictable memory latency

- Non-deterministic memory (caching, DRAM) reduces compiler ability to hide memory latency

### Thread synchronization

- Unsynchronized thread execution makes compile time resource allocation a challenge

A100



# Building Compilers For Spatial Architectures

## Challenges when building an automated compiler for spatial architectures

### Unpredictable memory latency

- Non-deterministic memory (caching, DRAM) reduces compiler ability to hide memory latency

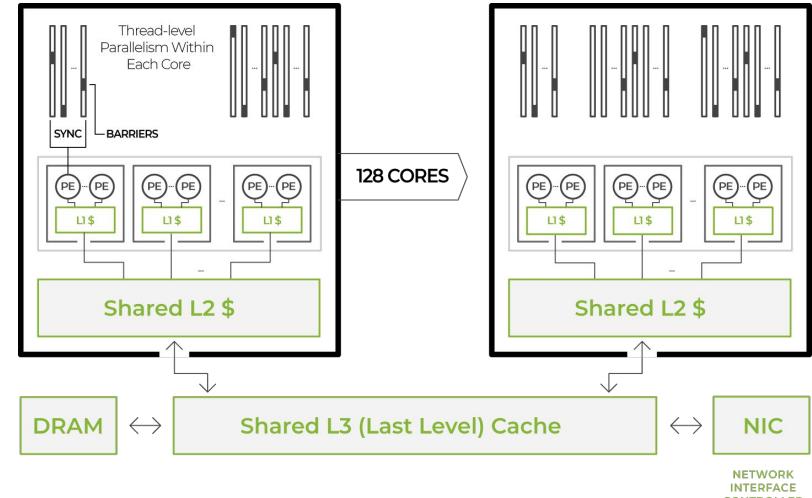
### Thread synchronization

- Unsynchronized thread execution makes compile time resource allocation a challenge

### Complex interconnect

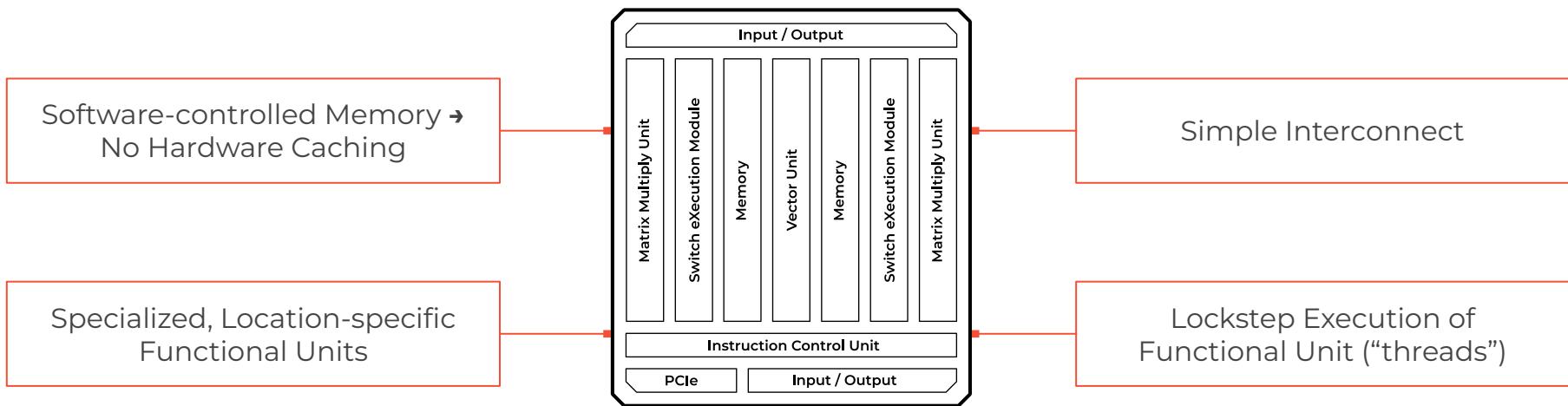
- Dynamic hardware (e.g. arbiters, queues, virtual channels) makes it difficult to resolve data movement at compile time
  - E.g. See Networks-on-Chip (NoC) in CMPs

A100



# Building a Compiler For GroqChip

## Groq Arch Abstraction



# An ISA That Empowers Software

**Software-controlled memory  
enabled by low-level abstraction  
exposed by architecture**

No dynamic hardware caching

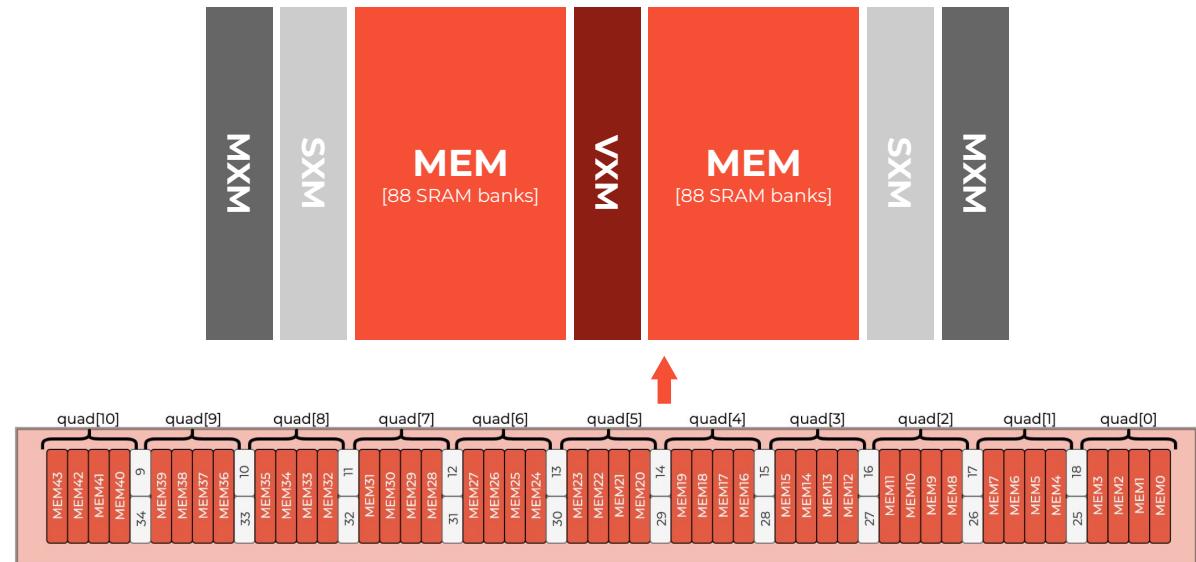
- Compiler aware of all data locations at any given point in time

Flat memory hierarchy (no L1, L2, L3, etc)

- Memory exposed to software as a set of physical banks that are directly addressed

Large on-chip memory capacity (220 MB) at high-bandwidth with (80TBps) reduces need to spill to non-deterministic DRAM

- Provides enough “scratchpad” memory to hide external memory accesses behind compute



# An ISA That Empowers Software

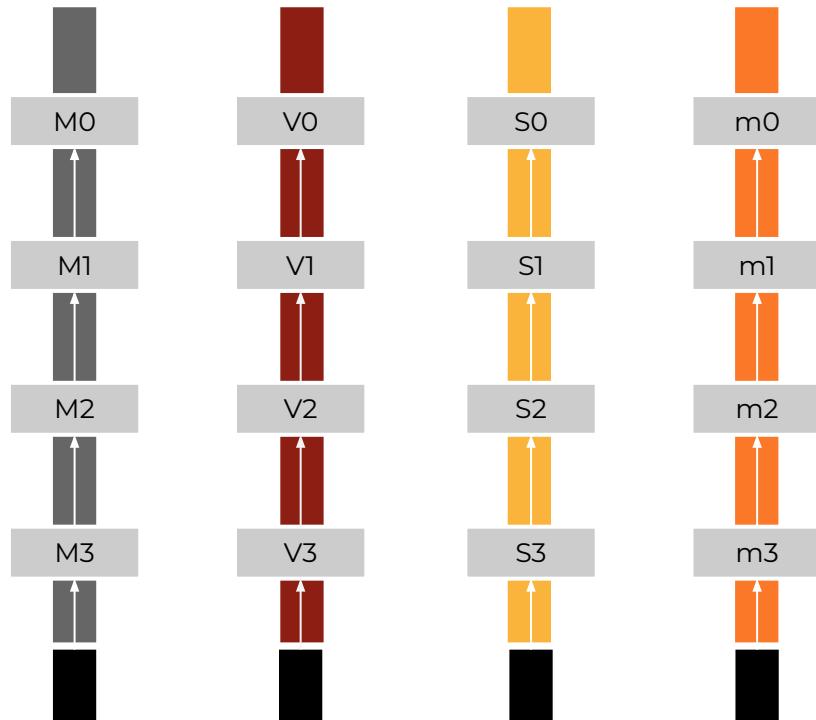
## Compiler empowered to perform cycle-accurate instruction scheduling

Functional units execute in lockstep

- One instruction issued per cycle at each dispatch path
  - Can be viewed as fully-pipelined 144-wide VLIW instructions

Little hardware control needed for managing instruction execution

- < 3% area overhead for instruction dispatch logic



# An ISA That Empowers Software

**Explicit time and space of instruction execution exposed by ISA to the compiler**

Each Functional Unit (FU) type provides its own low-level instruction set

Number of FUs of each type and relative positions on chip are exposed to software

Compiler can choose to leverage multiple FUs for more concurrency or more pipelining

Instruction Set	Function	Instruction
Low level	<b>ICU</b>	NOP $N$ Ifetch Sync Notify Config Repeat $n,d$
320-vector ops	<b>MEM</b>	Read $a,s$ Write $a,s$ Gather $s, map$ Scatter $s, map$ Countdown $d$ Step $a$ Iterations $n$
Explicit resource selection	<b>VXM</b>	unary operation binary operation type conversions ReLU Tanh Exp RSqrt
Explicit scheduling	<b>MXM</b>	LW IW ABC ACC
	<b>SXM</b>	Shift $up/down N$ Permute $map$ Distribute $map$ Rotate $stream$ Transpose $sg16$
	<b>C2C</b>	Deskew Send Receive

# An ISA That Empowers Software

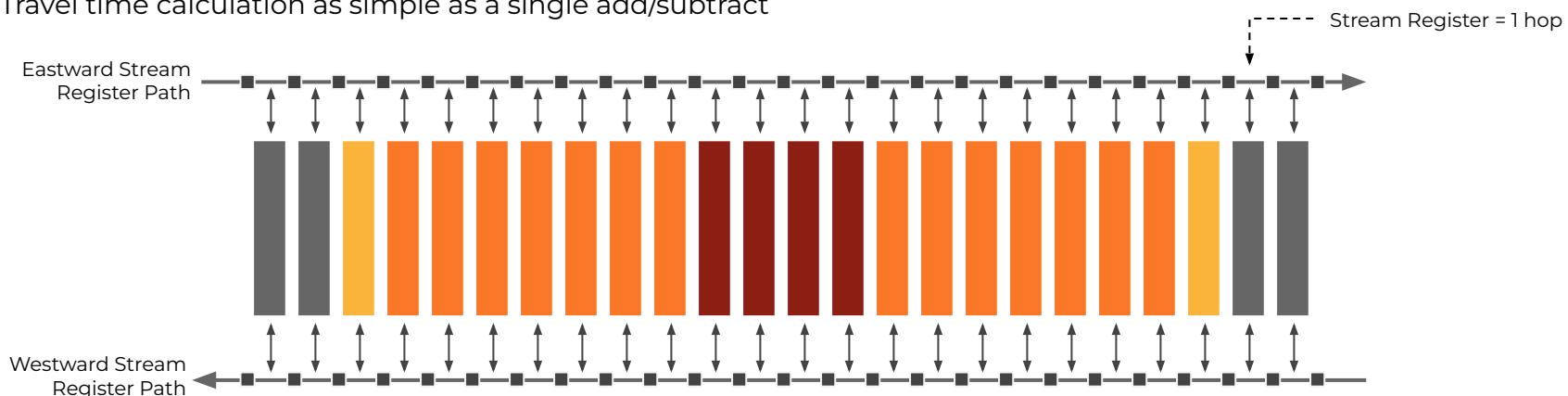
## Compiler can quickly reason about all data movement between FUs

Simple, one-dimensional interconnect for inter-FU communication

- Eastward and westward paths made up of arrays of “stream registers”
- Stream register = one-cycle hop

No arbiters / queues = software can easily reason about exact data movement without simulation

Travel time calculation as simple as a single add/subtract

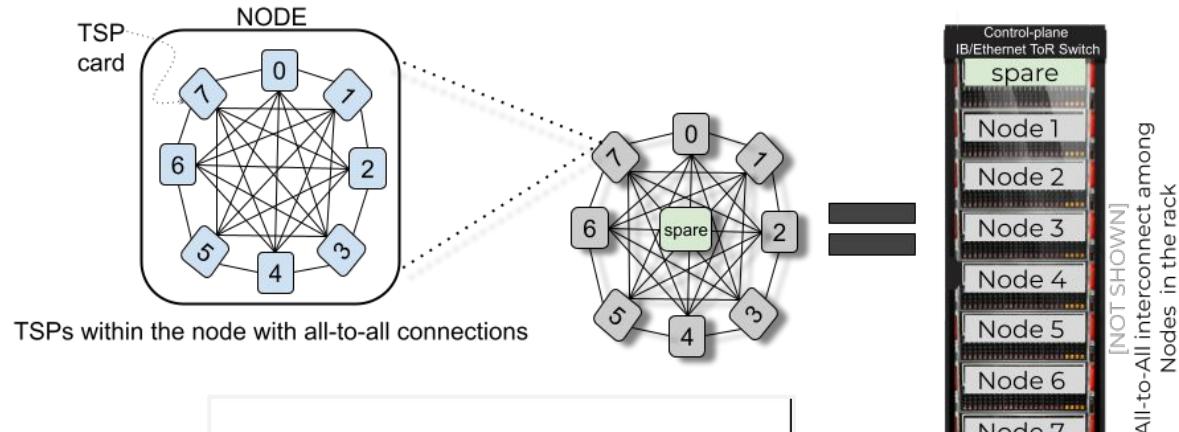


# An ISA That Empowers Software

## Lockstep execution extended to multiple GroqChip processors

C2C interface ensures synchronization across GroqChips  
Low-level abstraction allows compiler to view multiple GroqChip processors as a single device

- Can be scheduled as “one” “Software-scheduled network”



### A Software-defined Tensor Streaming Multiprocessor for Large-scale Machine Learning

Dennis Abts Groq Inc.	Garrin Kimmell Groq Inc.	Andrew Ling Groq Inc.	John Kim KAIST/Groq Inc.
Matt Boyd Groq Inc.	Andrew Bitar Groq Inc.	Sahil Parmar Groq Inc.	Ibrahim Ahmed Groq Inc.
Roberto DiCecco Groq Inc.	David Han Groq Inc.	John Thompson Groq Inc.	Michael Bye Groq Inc.
Jennifer Hwang Groq Inc.	Jeremy Powers Groq Inc.	Peter Lillian Groq Inc.	Ashwin Murthy Groq Inc.
Elyas Mehatabuddin Groq Inc.	Chetan Tekur Groq Inc.	Thomas Sohmers Groq Inc.	Kris Kang Groq Inc.
Stephen Maresh Groq Inc.	Stephen Maresh Groq Inc.	Jonathan Ross Groq Inc.	

#### ABSTRACT

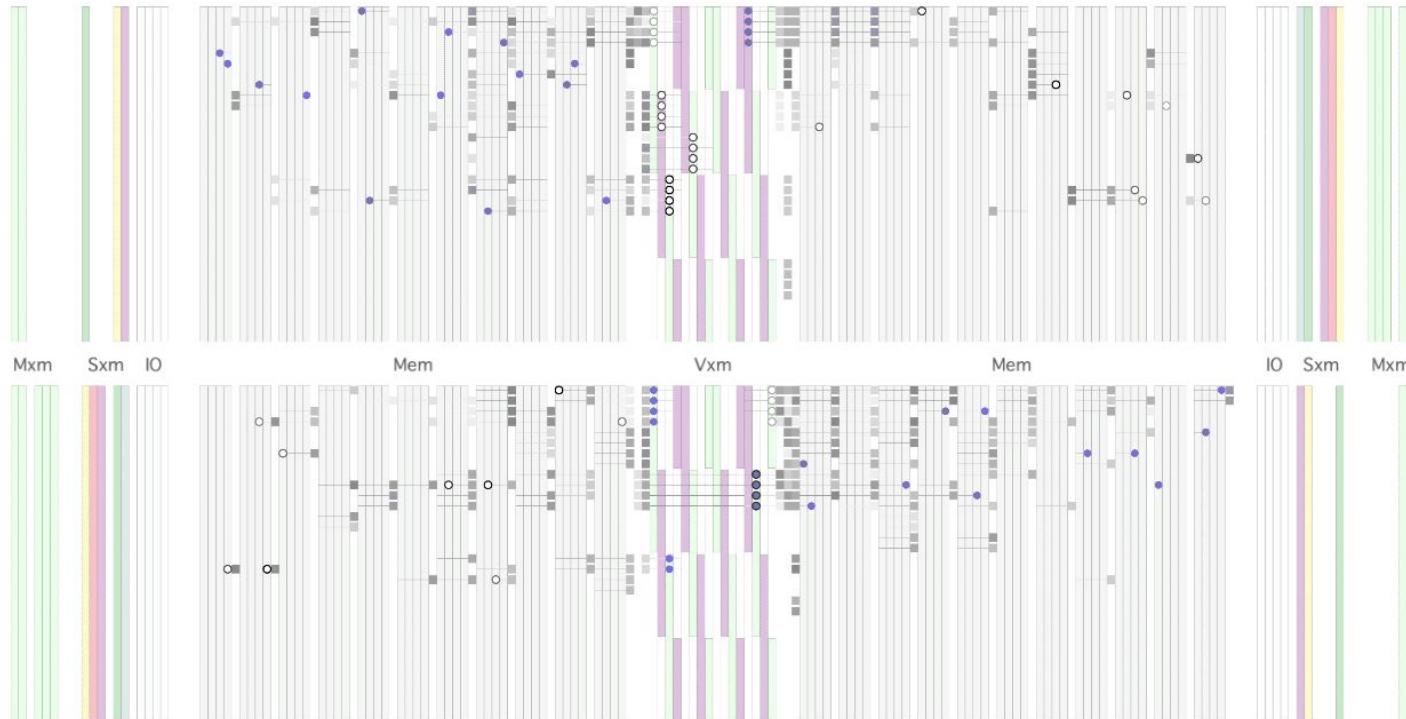
We describe our novel commercial software-defined approach for large-scale interconnection networking using Tensor Streaming Processing (TSP) components. The system architecture includes scheduling, routing, and flow control of the interconnection network of TSPs. We describe the communication and synchronization primitives of a bandwidth-aware rate controller for the TSPs. This is possible by using a hardware abstraction for the TSPs and a software abstraction for the network backbone for the TSPs. The results demonstrate systems based on a software-defined Dragonfly topology, ultimately

#### CCS CONCEPTS

• Computer systems organization → Interconnection architectures  
KEYWORDS  
Machine Learning, Tensor Streaming Processor, Dragonfly, Software Scheduling

# Power of Data Orchestration

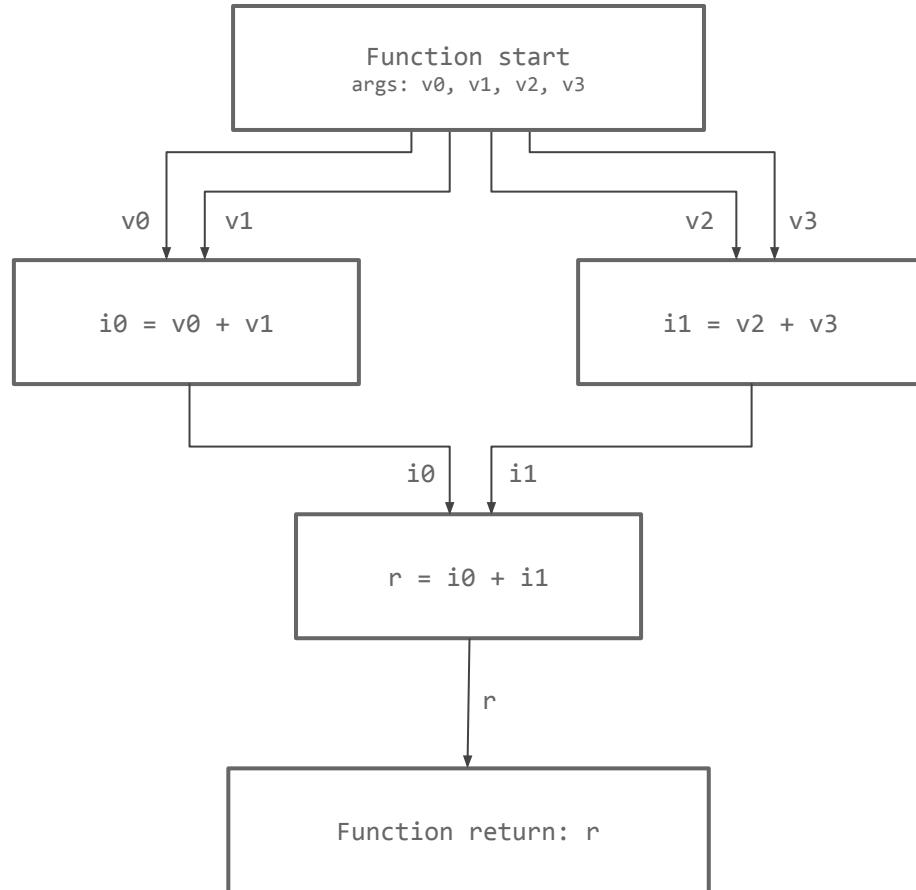
Given to Groq Compiler



# EXAMPLE

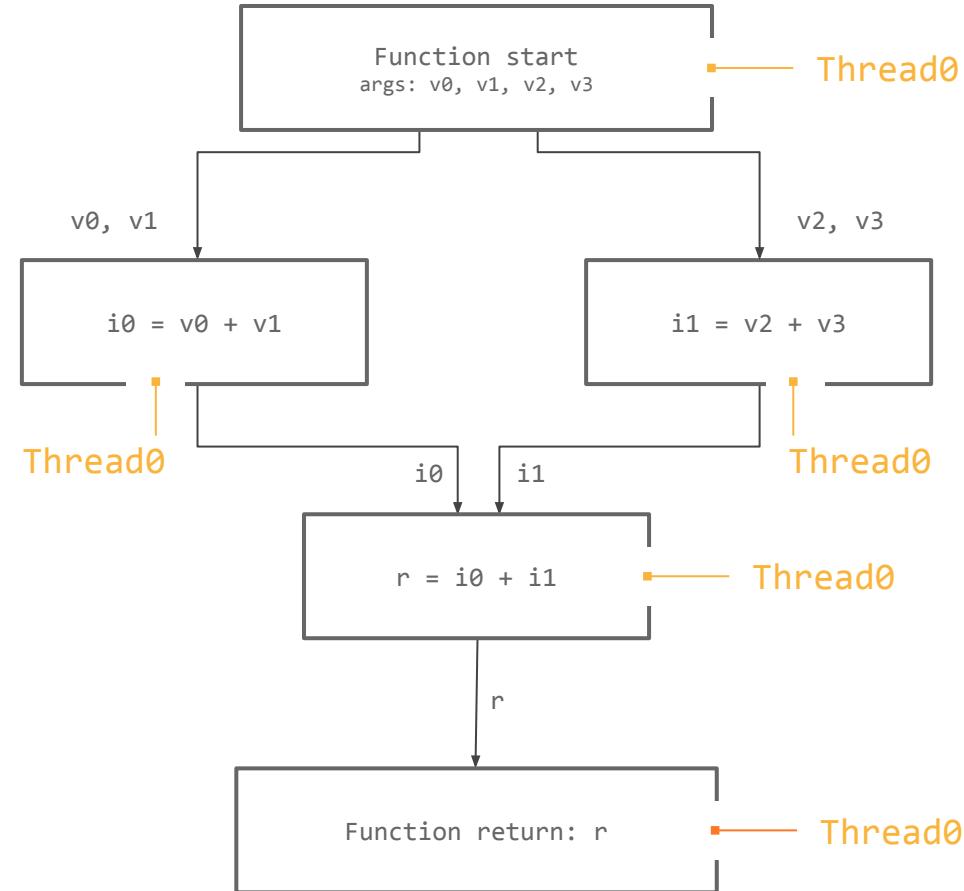
## Vectors

Add four vectors



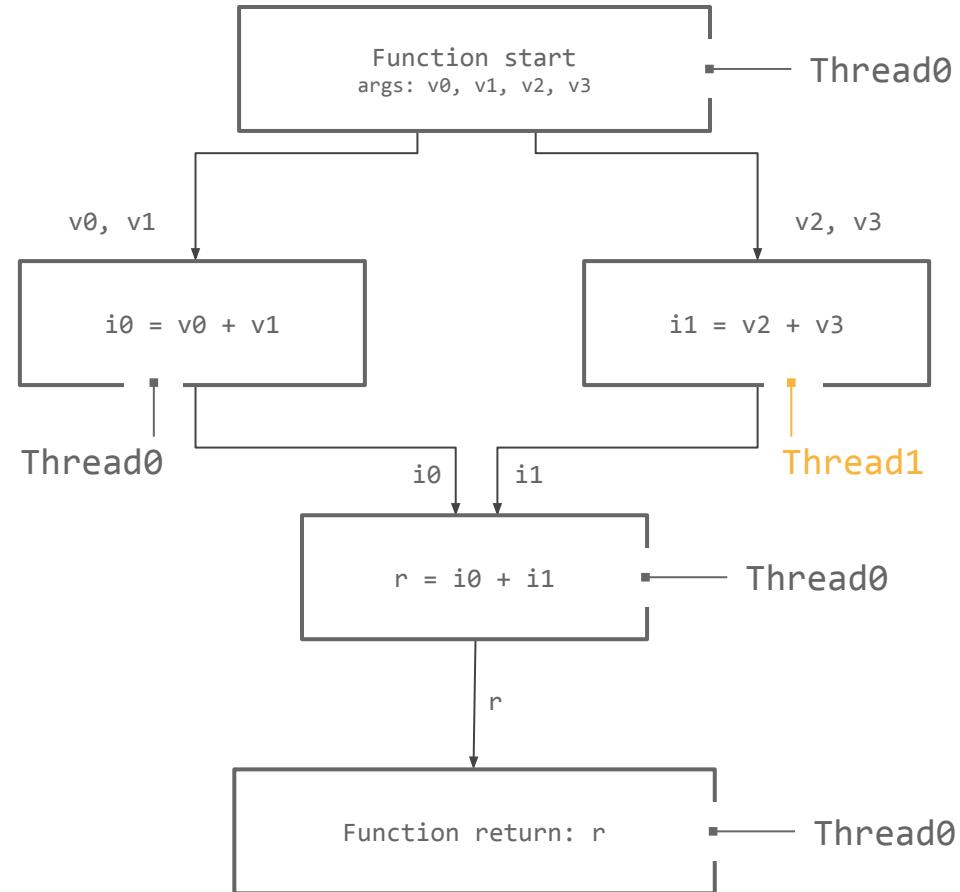
## EXAMPLE

# Singlethreaded



## EXAMPLE

# Multithreaded



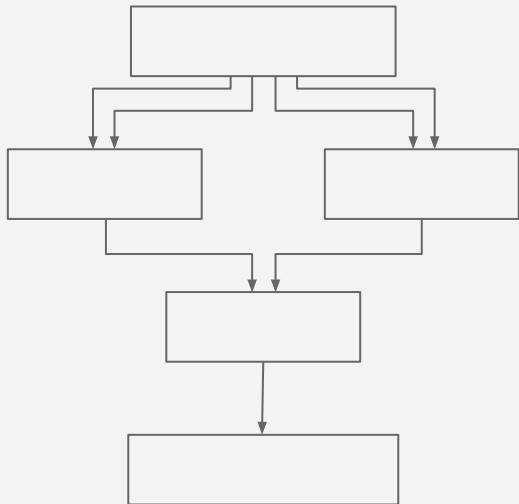
EXAMPLE

# Multithreaded

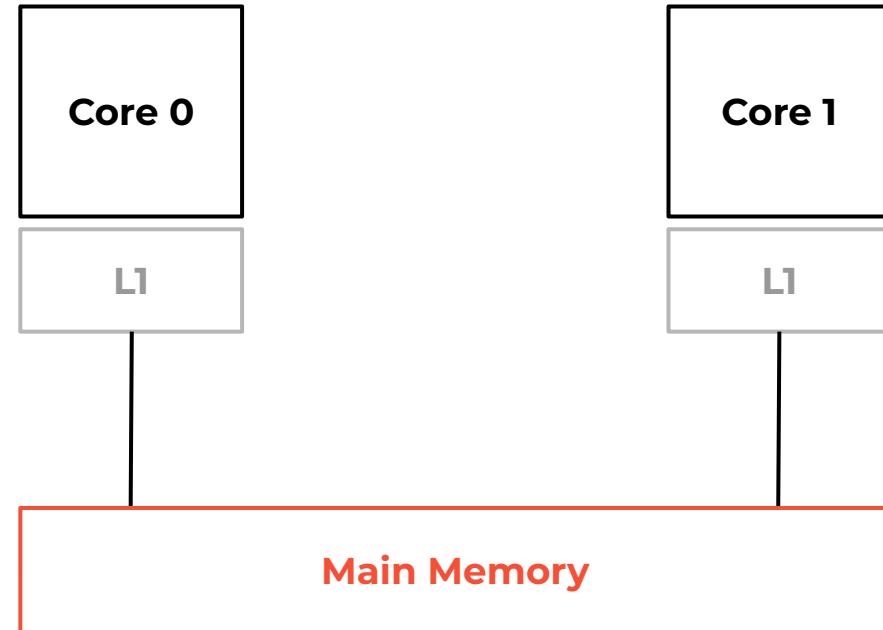
## Compiled program contents

- 1.** function addFour
  - a. dispatch Thread1(addTwo)
  - b. execute addTwo
  - c. join Thread1
  - d. execute addTwo
  - e. return
  
- 2.** function addTwo
  - a. load operand vectors
  - b. add operand vectors
  - c. store result
  - d. return

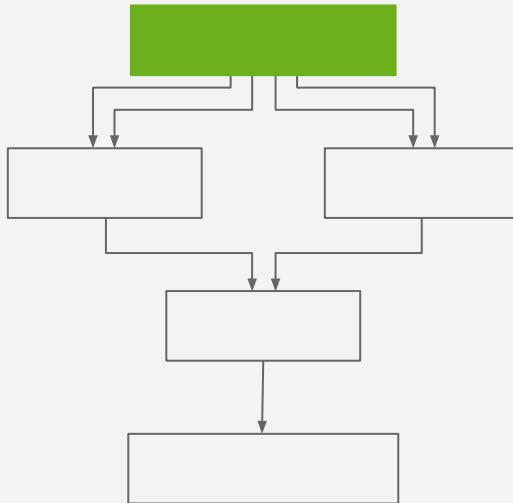
# On Multicore SIMD Arch



THEORETICAL MULTITHREADED SIMD PROCESSOR WITH CACHE

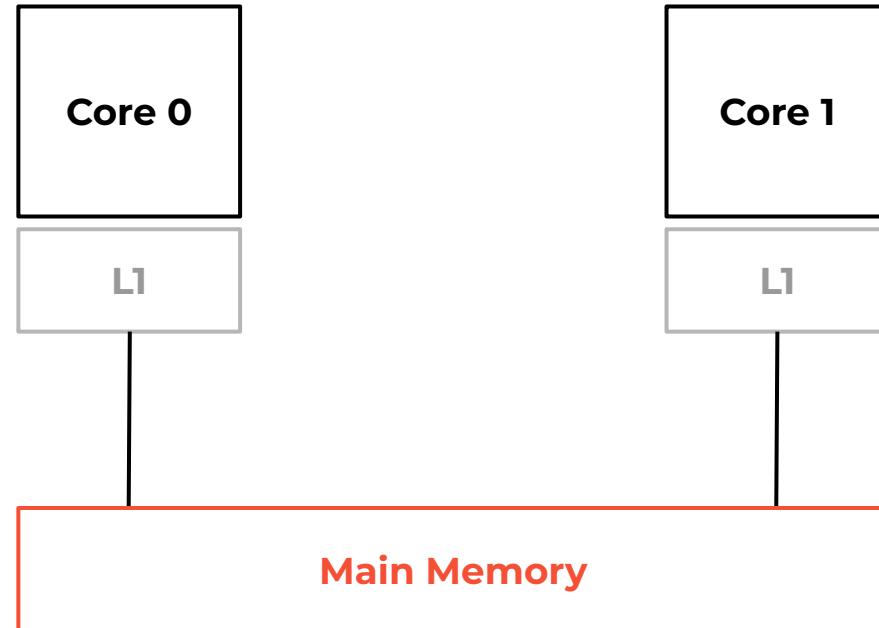


# On Multicore SIMD Arch

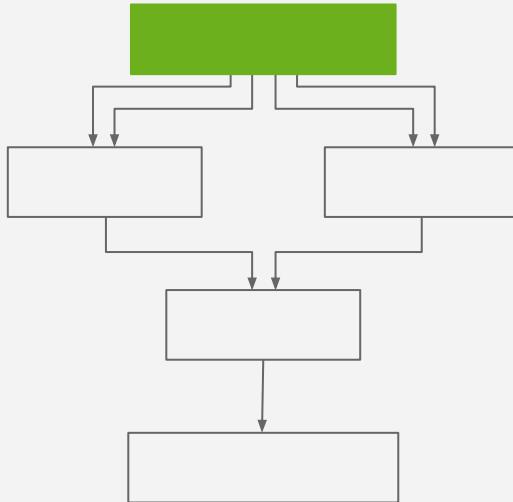


Time = 0

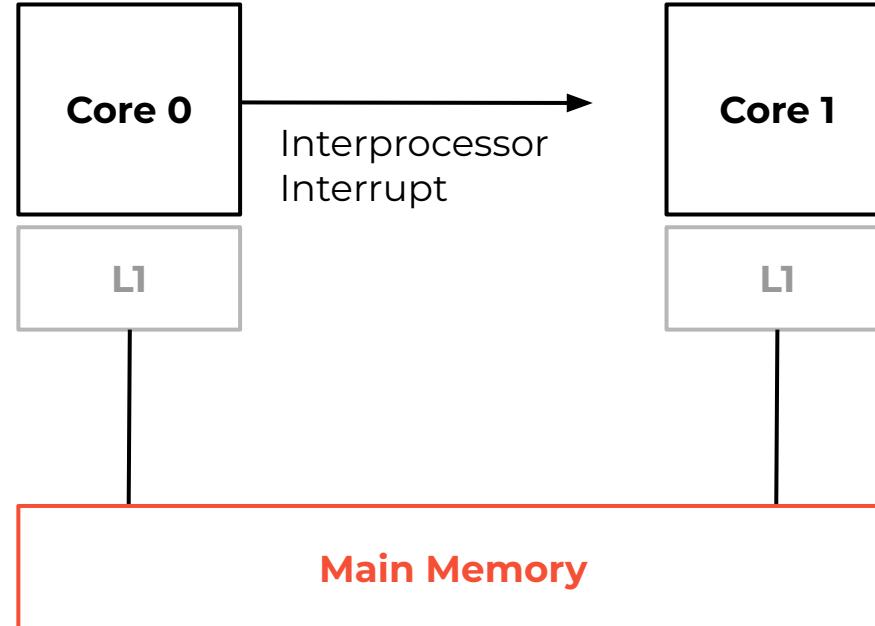
**Executing:**  
addFour



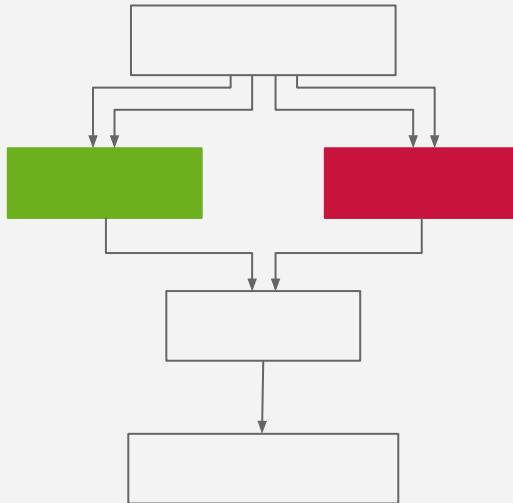
# On Multicore SIMD Arch



**Executing:**  
addFour

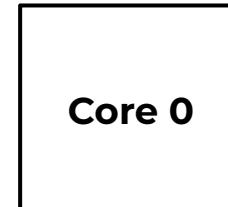


# On Multicore SIMD Arch

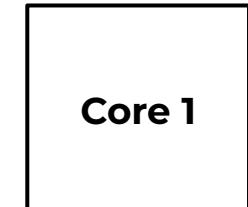


Time =  $t_{\text{dispatch}}$

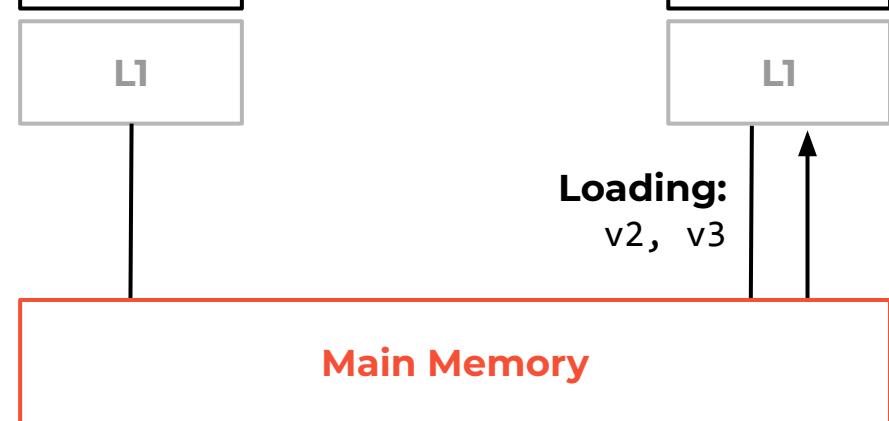
**Executing:**  
addTwo(v0, v1)



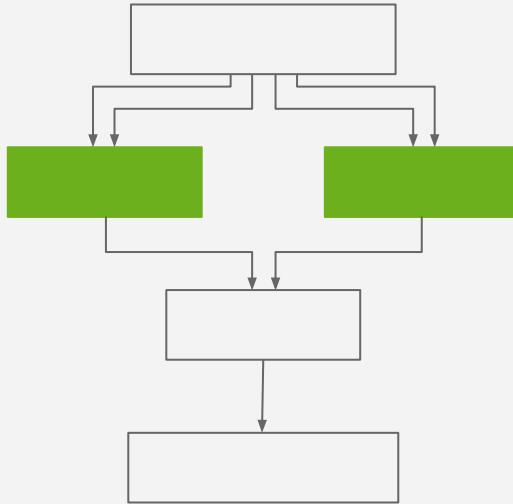
**Executing:**  
addTwo(v2, v3)



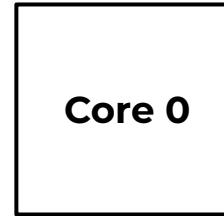
**Loading:**  
v2, v3



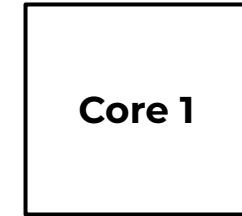
# On Multicore SIMD Arch



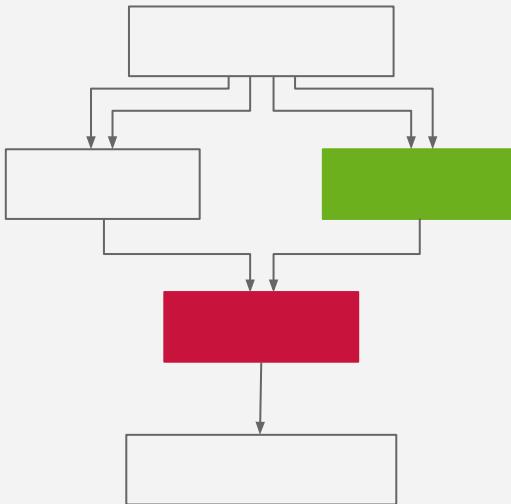
**Executing:**  
addTwo(v0, v1)



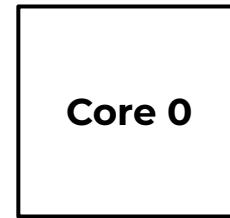
**Executing:**  
addTwo(v2, v3)



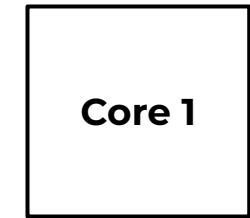
# On Multicore SIMD Arch



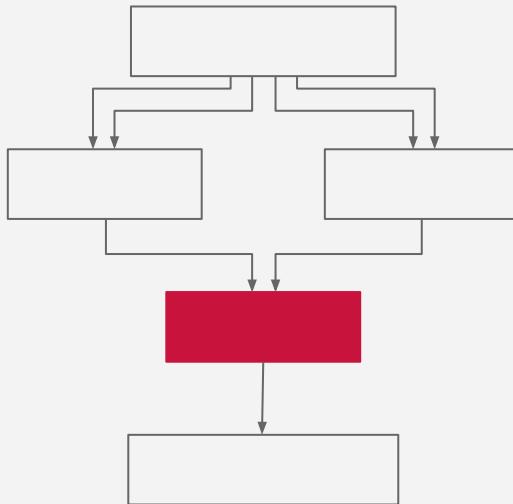
**Executing:**  
join Thread1



**Executing:**  
`addTwo(v2, v3)`

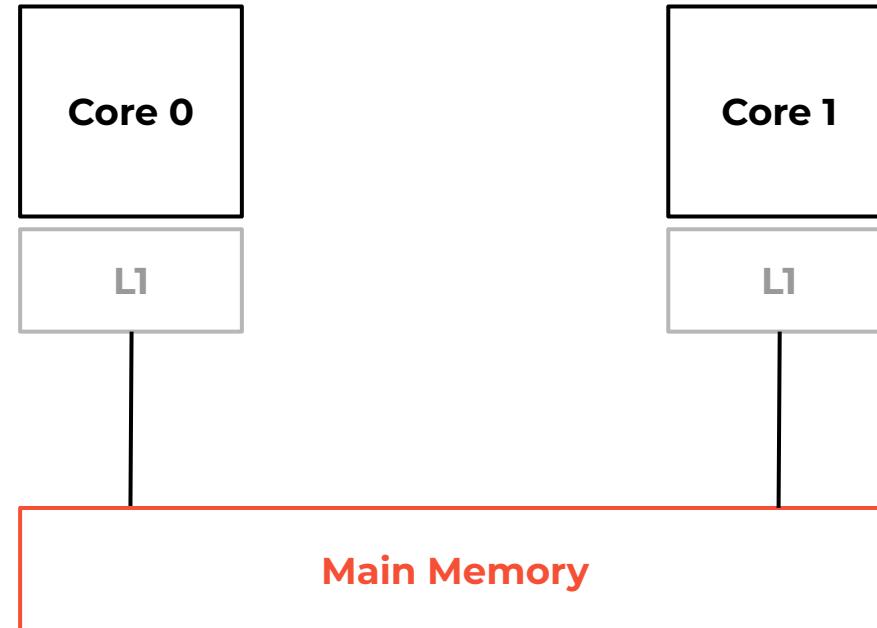


# On Multicore SIMD Arch

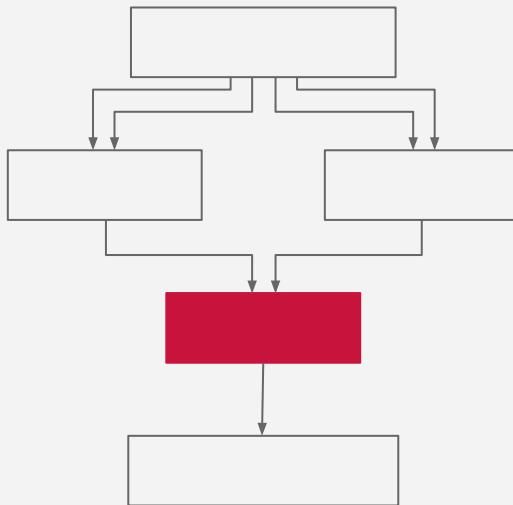


$$\text{Time} = t_{\text{dispatch}} + t_{\text{add}}$$

**Executing:**  
join Thread1

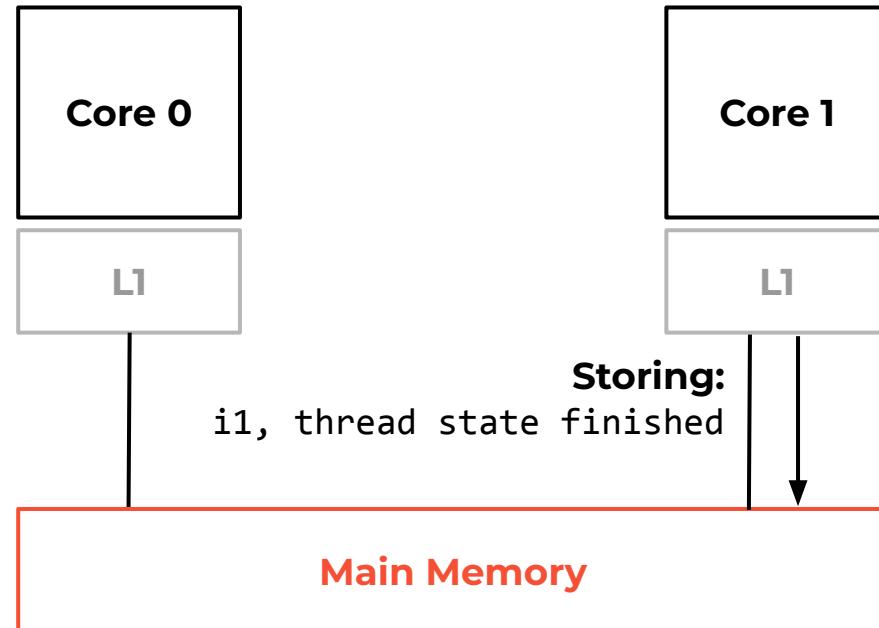


# On Multicore SIMD Arch

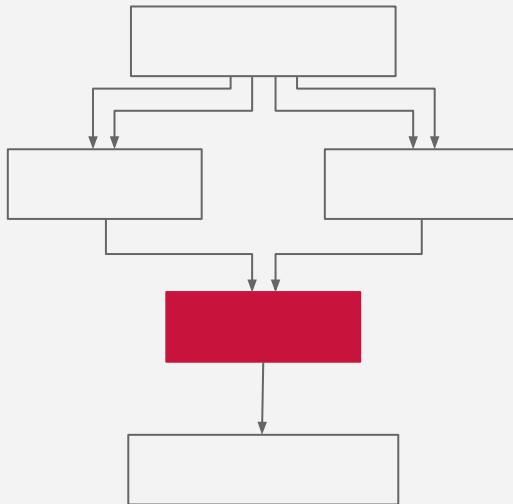


$$\text{Time} = t_{\text{dispatch}} + t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}}$$

**Executing:**  
join Thread1

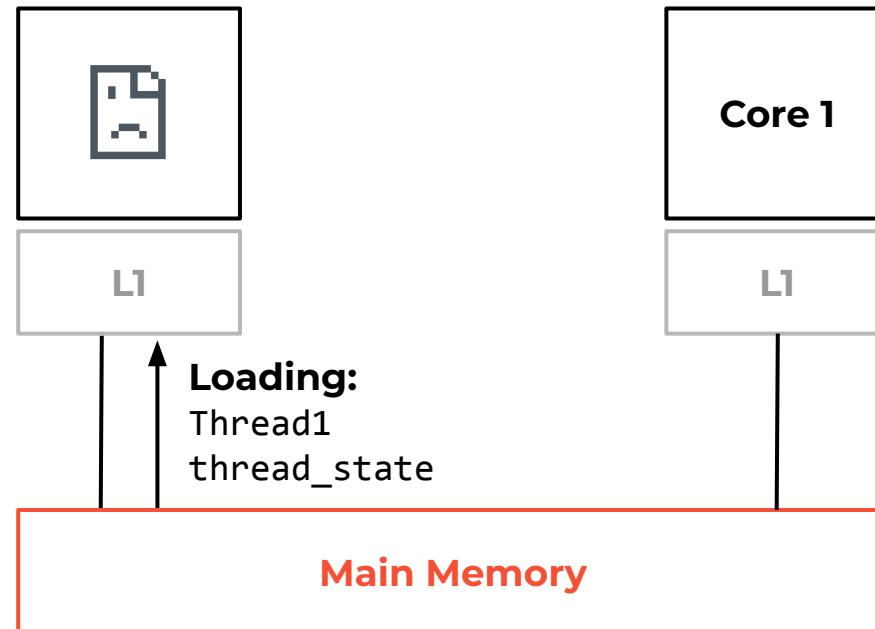


# On Multicore SIMD Arch

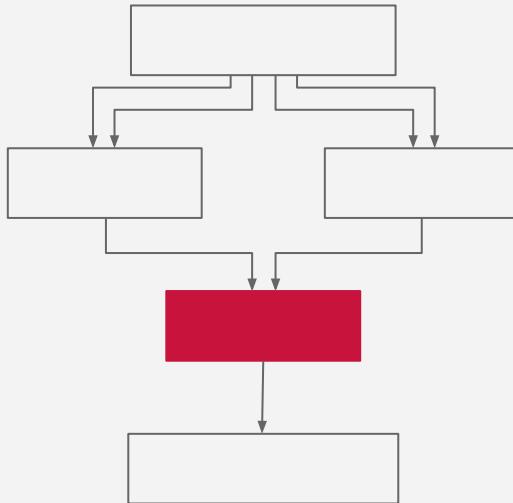


$$\text{Time} = t_{\text{dispatch}} + t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}}$$

**Executing:**  
cache miss on  
Thread1 thread\_state

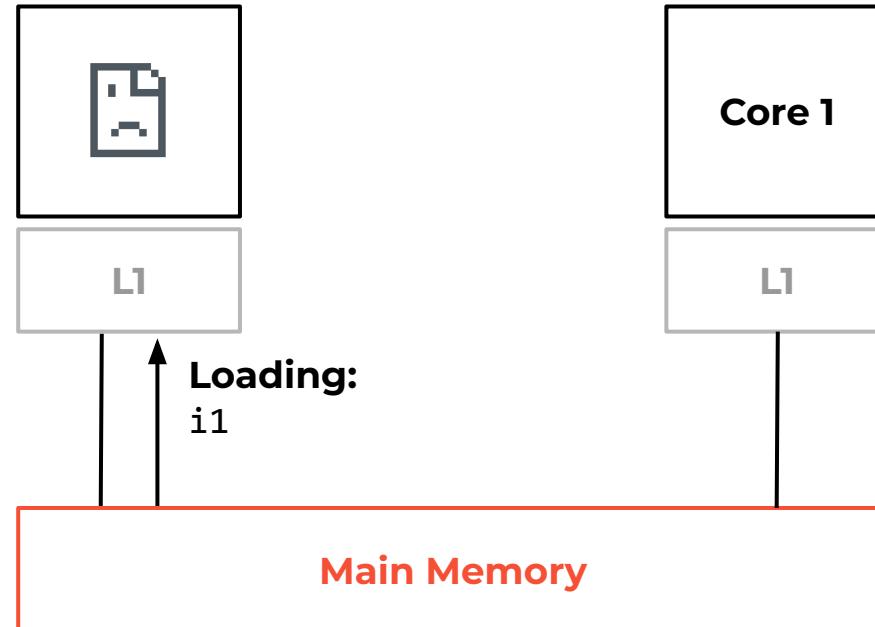


# On Multicore SIMD Arch

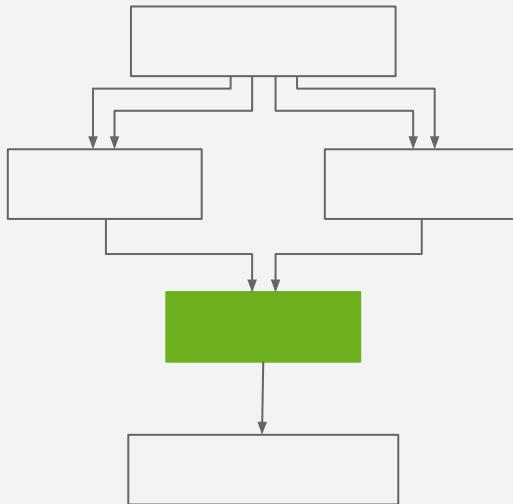


$$\text{Time} = t_{\text{dispatch}} + t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}} + t_{\text{load0}} + t_{\text{load1}}$$

**Executing:**  
cache miss on i1

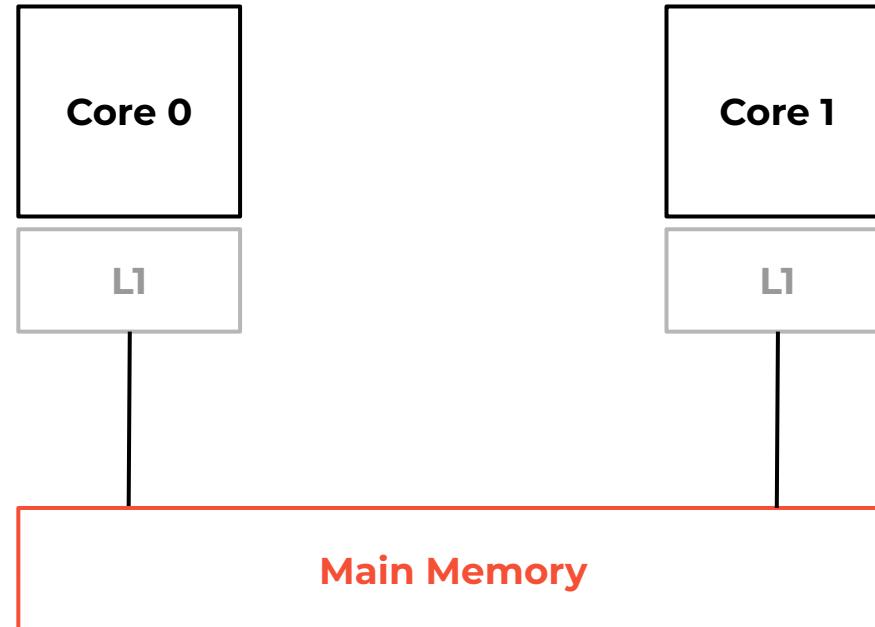


# On Multicore SIMD Arch

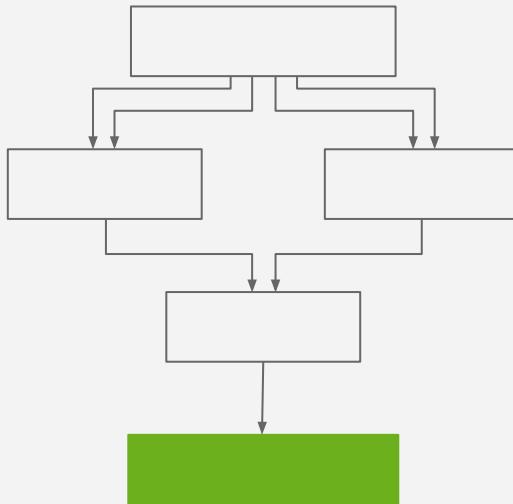


$$\text{Time} = t_{\text{dispatch}} + t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}} + t_{\text{load0}} + t_{\text{load1}}$$

**Executing:**  
addTwo(i0, i1)

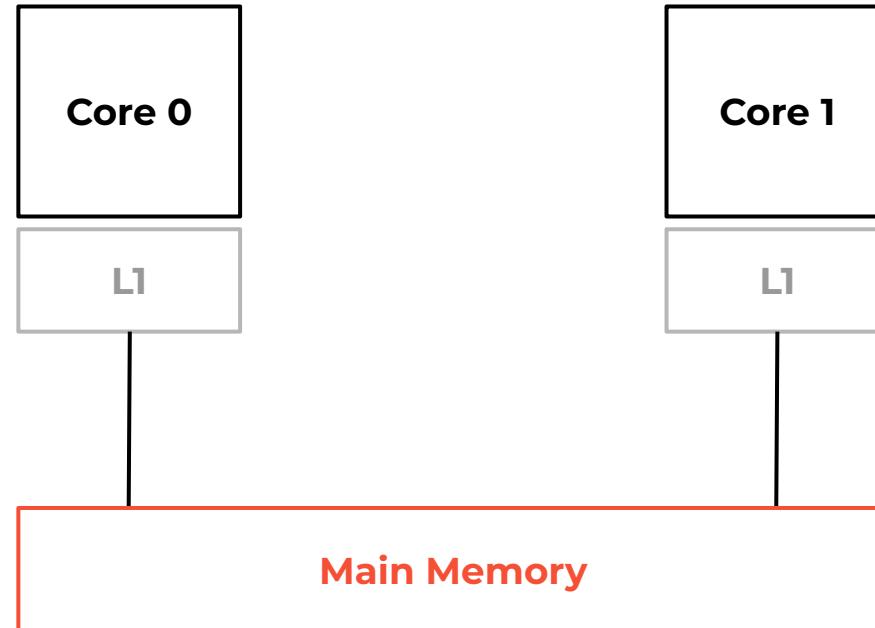


# On Multicore SIMD Arch



$$\text{Time} = t_{\text{dispatch}} + 2 * t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}} + t_{\text{load0}} + t_{\text{load1}}$$

**Executing:**  
return



# On Multicore SIMD Arch

**Code for synchronization must be present due to unresolvable times**

$$\text{Latency} = t_{\text{dispatch}} + 2 * t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}} + t_{\text{load0}} + t_{\text{load1}}$$

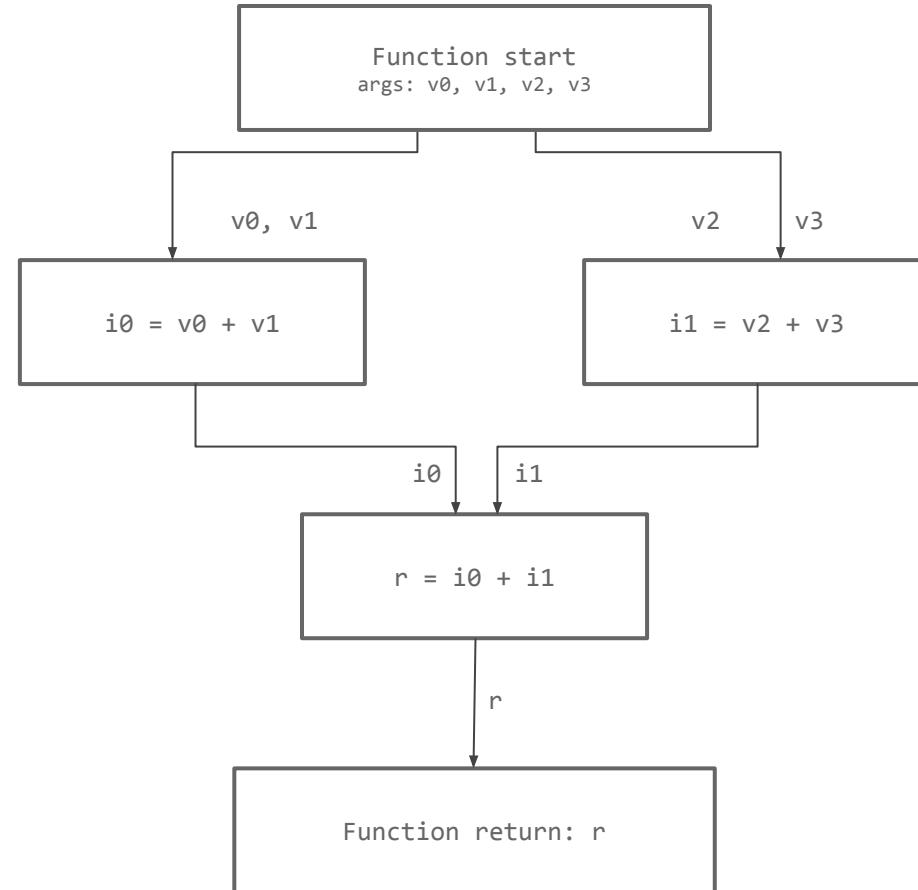
$t_{\text{dispatch}}$  unknown, implemented via interrupt

$t_{\text{store0}}$ ,  $t_{\text{store1}}$ ,  $t_{\text{load0}}$ ,  $t_{\text{load1}}$  unknown, will vary based on memory performance of system

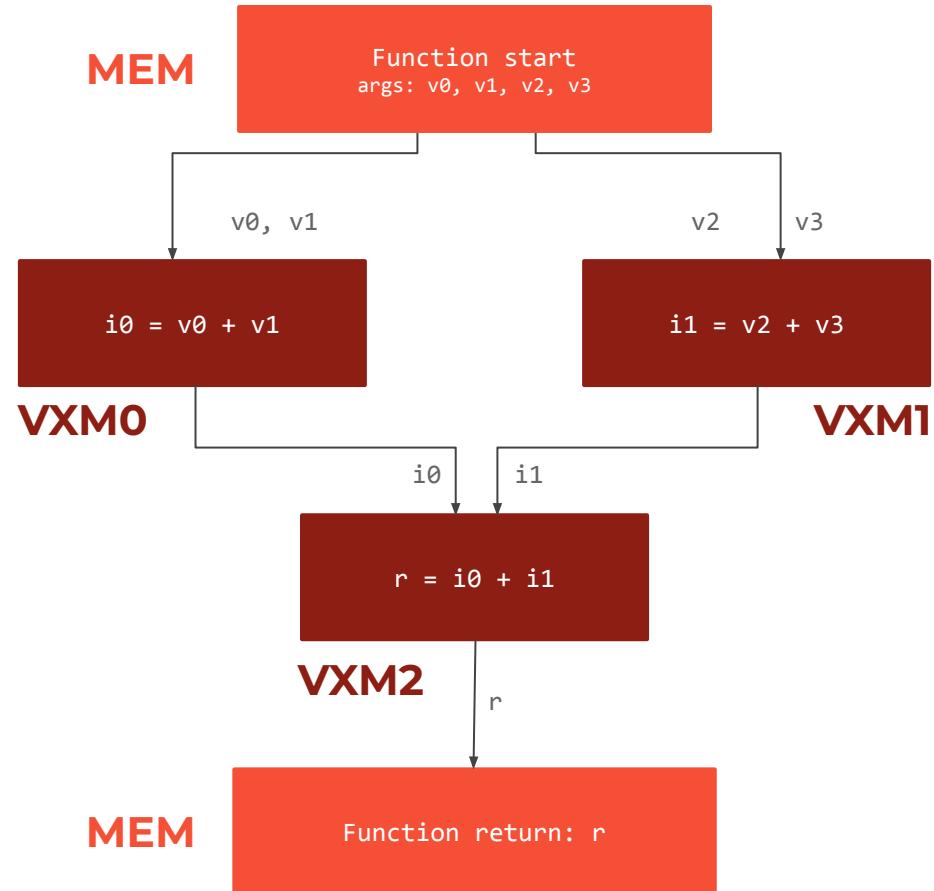
**Red** = unresolvable at compile time

**Green** = resolvable at compile time

# On GroqChip



# On GroqChip



# On GroqChip

## Compiled program contents

- 1.** MEM read v0, v1, v2, v3 → streams 0, 1, 2, 3 west
- 2.** VXM0 add streams 0, 1 → stream 0 west
- 3.** VXM1 add streams 2, 3 → stream 1 west
- 4.** VXM2 add streams 0, 1 → stream 0 east
- 5.** MEM store stream 0

# On GroqChip

## Compiled program contents: Instructions

MEM0 0 read v0 → streams 0 west store stream 0  $t_{read} + 2 * t_{travel} + 2 * t_{add}$  read v1 → streams 1 west

MEM1 0 read v1 → streams 1 west

MEM2 0 read v1 → streams 2 west

MEM3 0 read v1 → streams 3 west

VXM0 add streams 0, 1 -> stream 0 west  $t_{read} + t_{travel}$

VXM1 add streams 2, 3 -> stream 1 west  $t_{read} + t_{travel}$

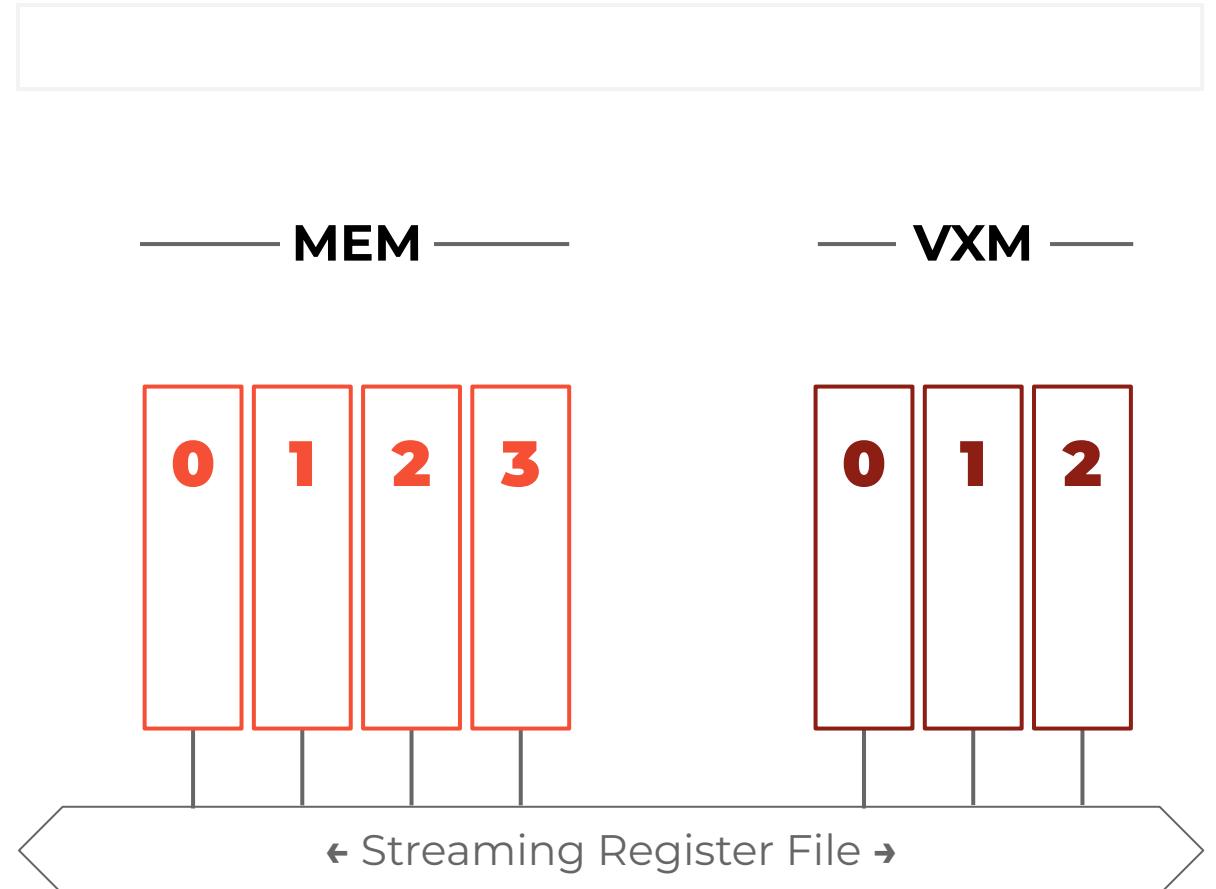
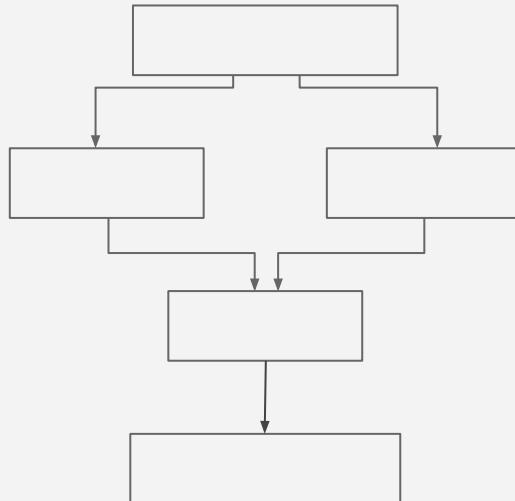
VXM2 add streams 0, 1 -> stream 0 east  $t_{read} + t_{travel} + t_{add}$

# On GroqChip

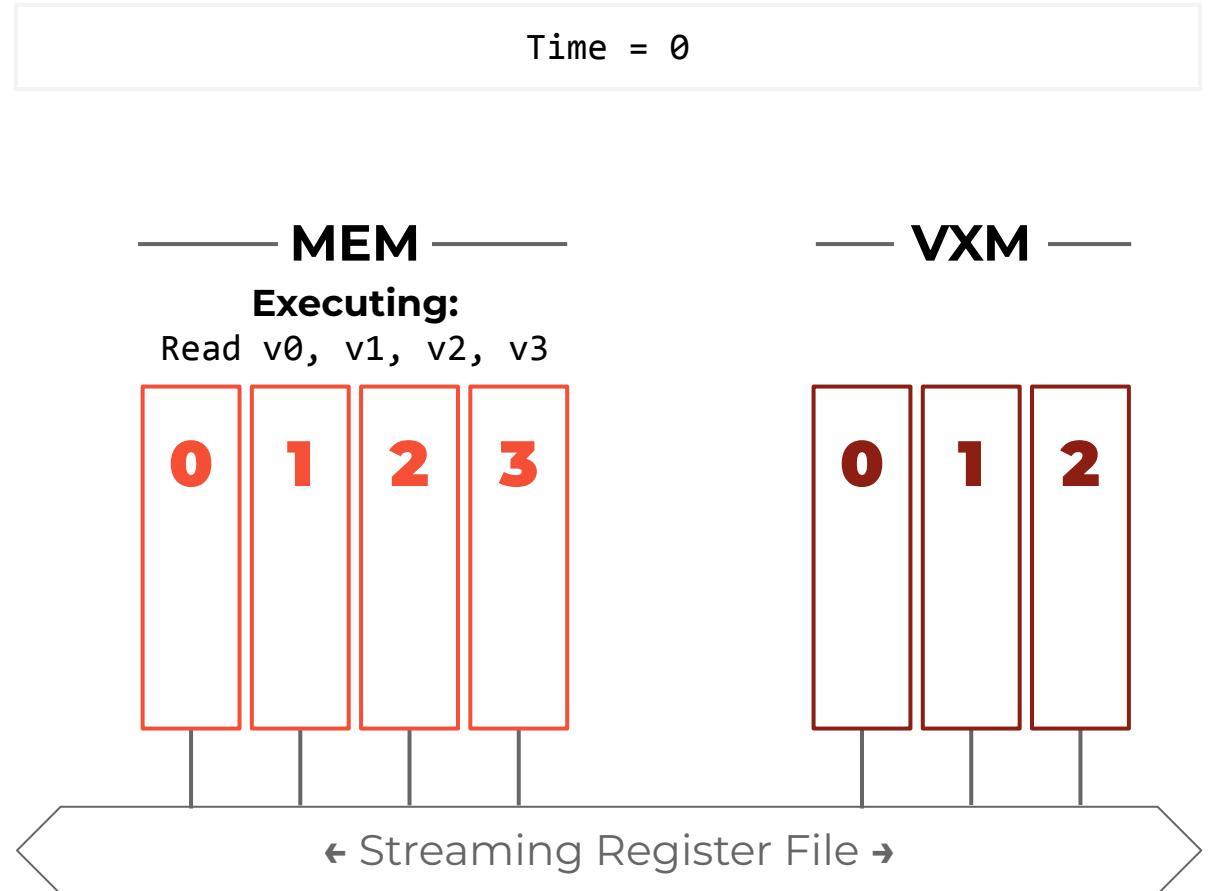
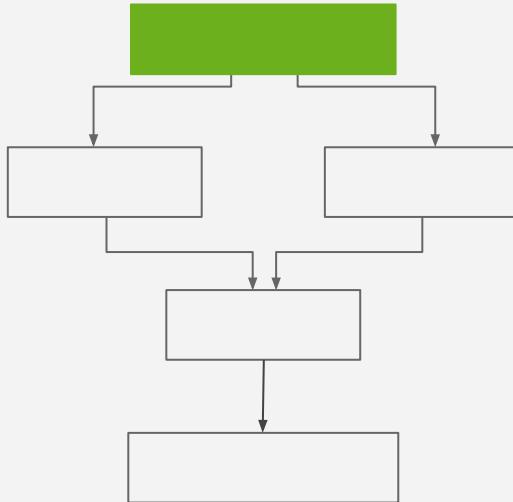
## Compiled program contents: Instructions

Block	Time	Instruction
MEM0	0	read v0 → streams 0 west
	$t_{read} + 2 * t_{travel} + 2 * t_{add}$	store stream 0
MEM1	0	read v1 → streams 1 west
MEM2	0	read v2 → streams 2 west
MEM3	0	read v3 → streams 3 west
VXM0	$t_{read} + t_{travel}$	add streams 0, 1 → stream 0 west
VXM1	$t_{read} + t_{travel}$	add streams 2, 3 → stream 1 west
VXM2	$t_{read} + t_{travel} + t_{add}$	add streams 0, 1 → stream 0 east

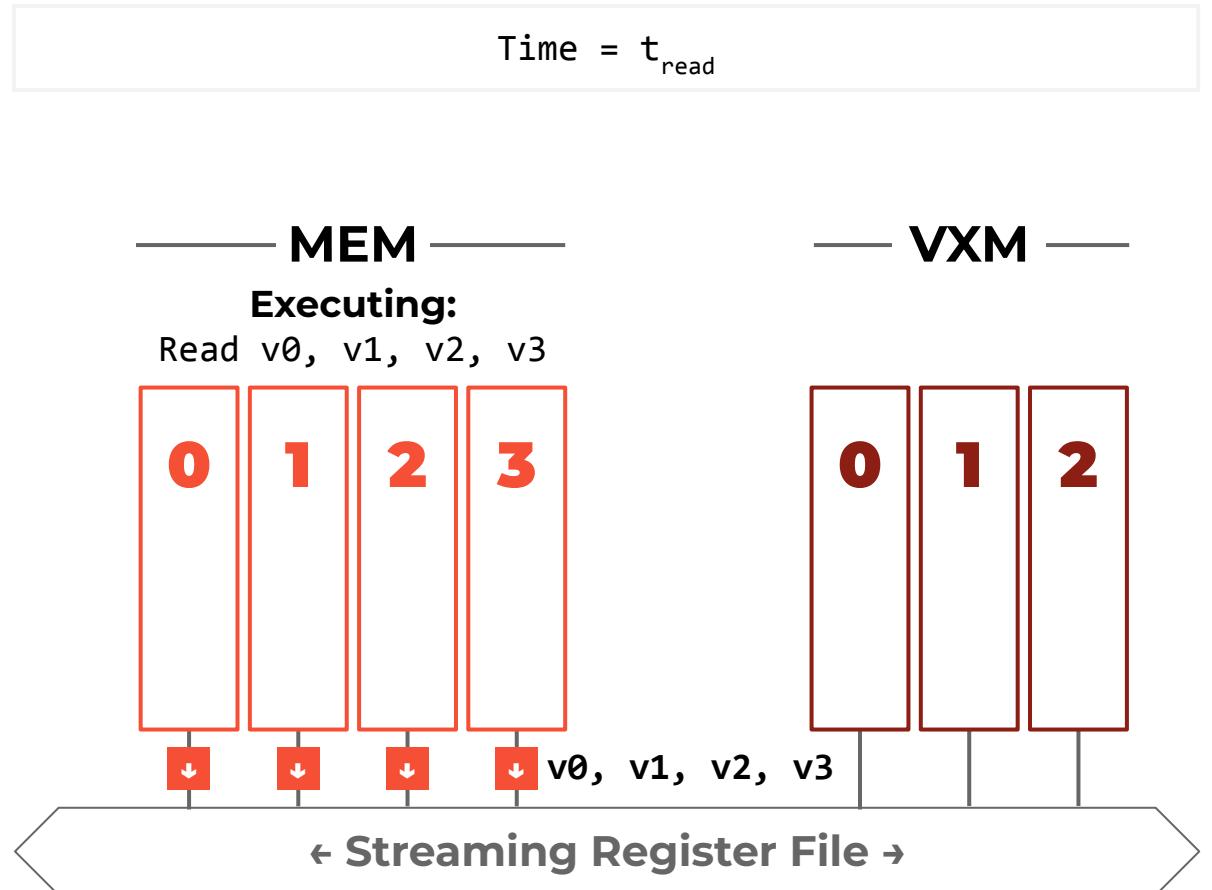
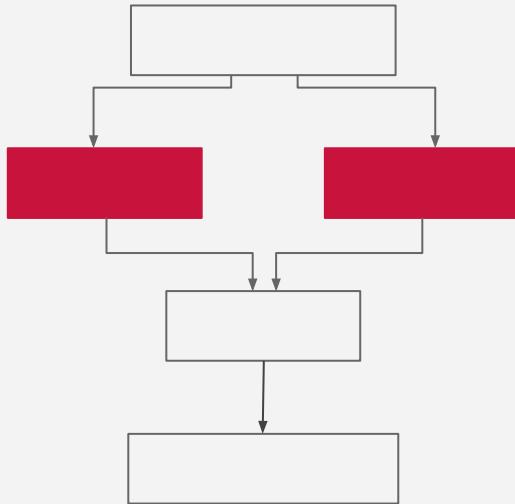
# On GroqChip



# On GroqChip

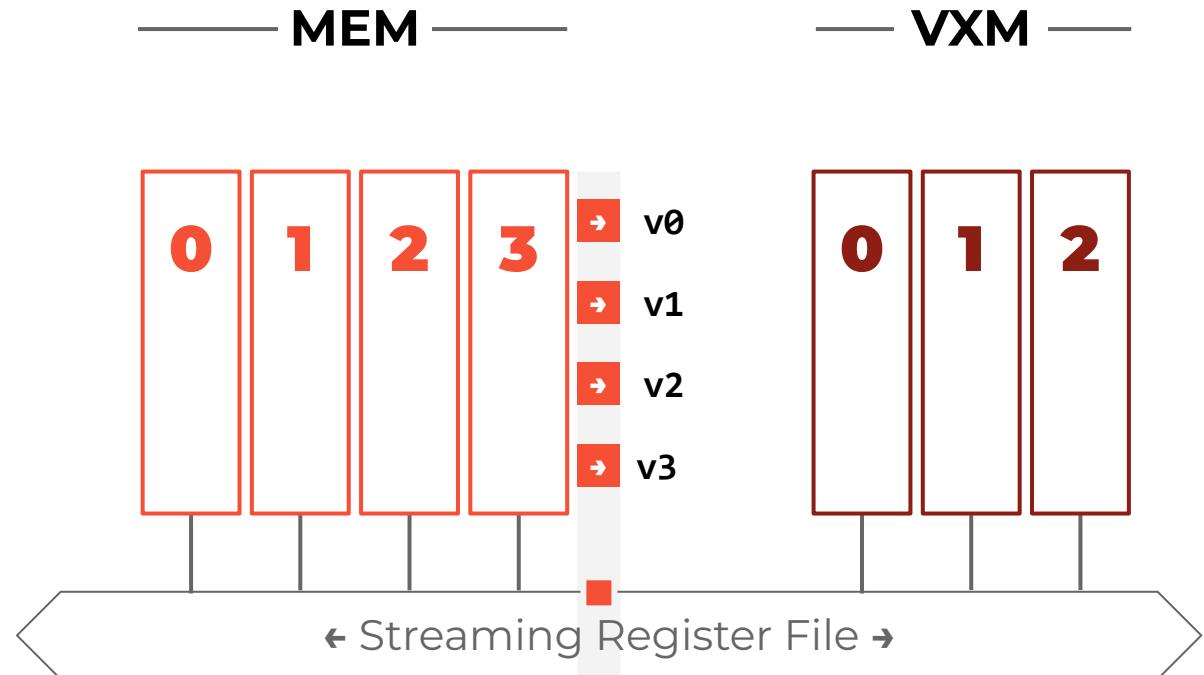
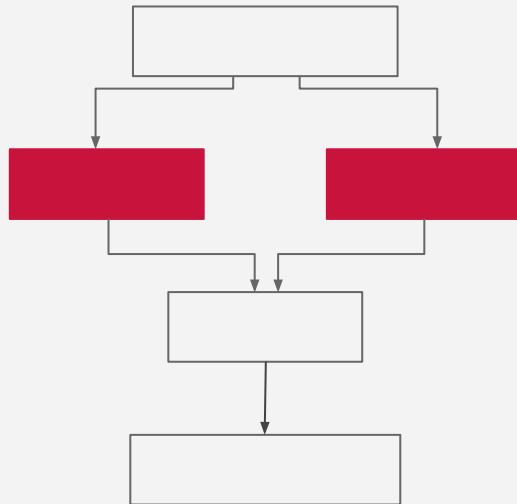


# On GroqChip

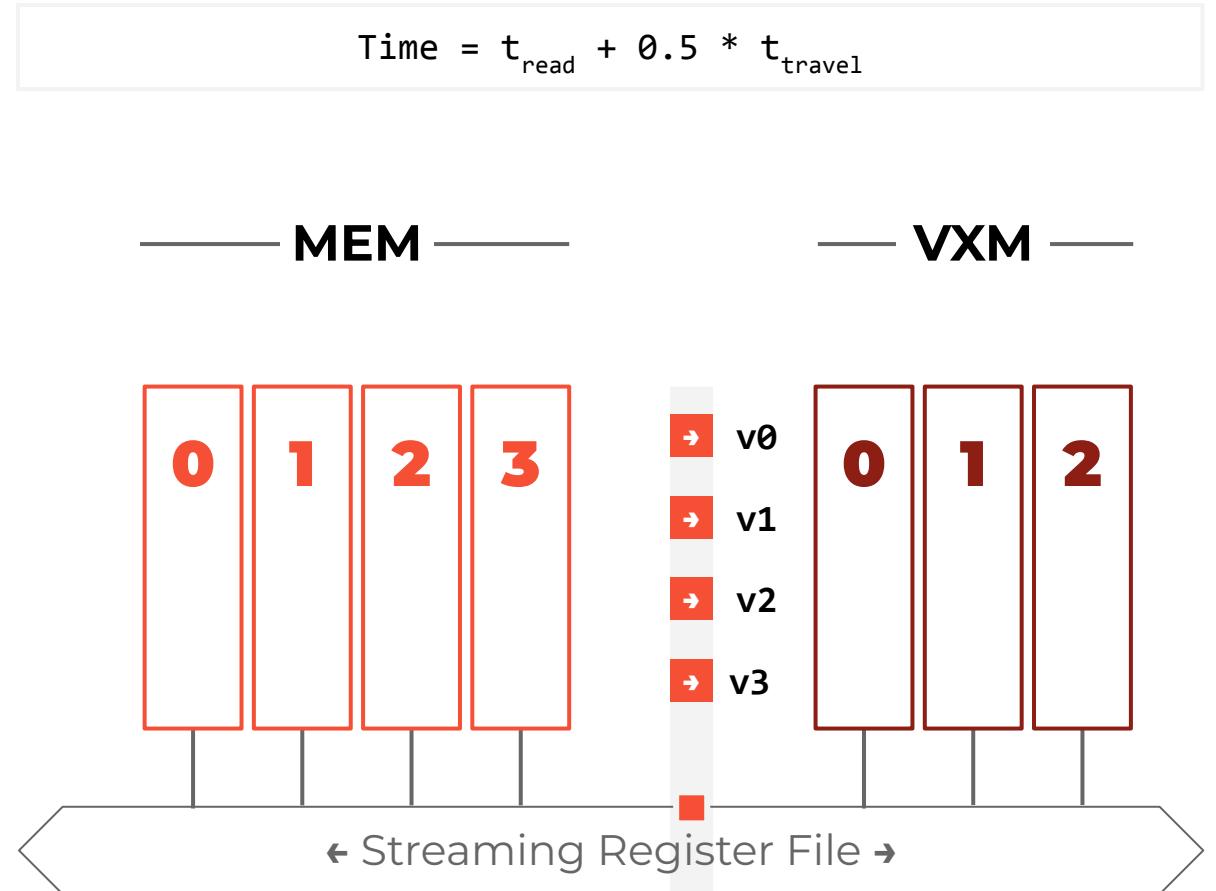
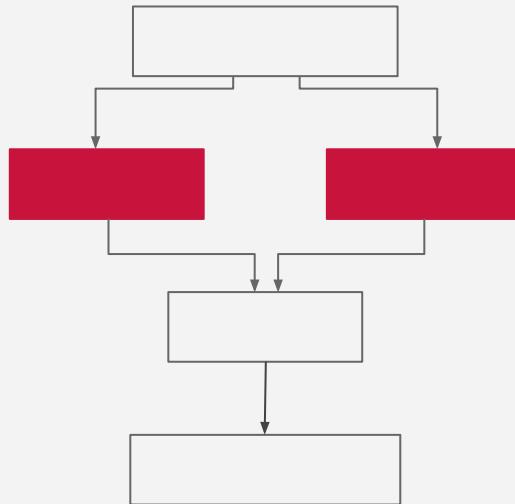


# On GroqChip

$$\text{Time} = t_{\text{read}} + 0.25 * t_{\text{travel}}$$

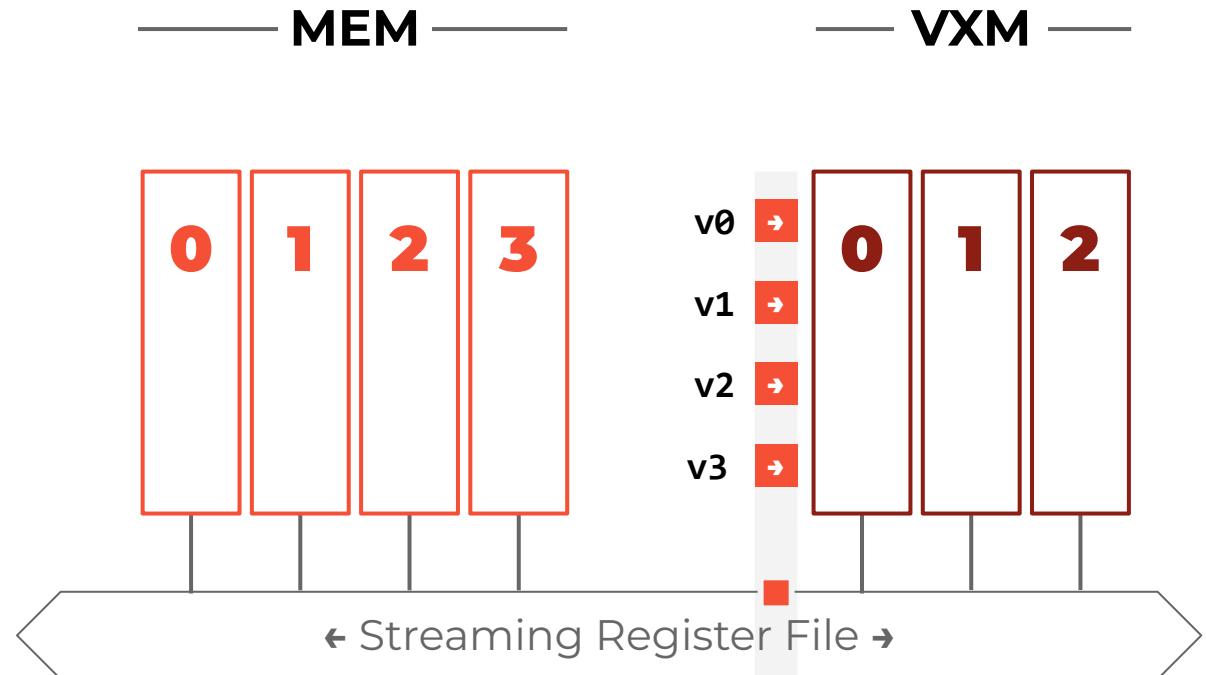
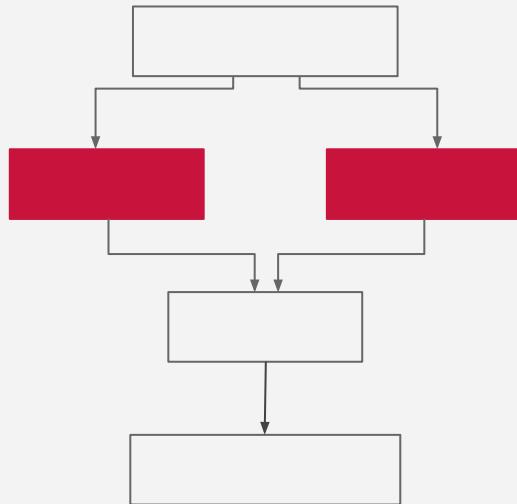


# On GroqChip

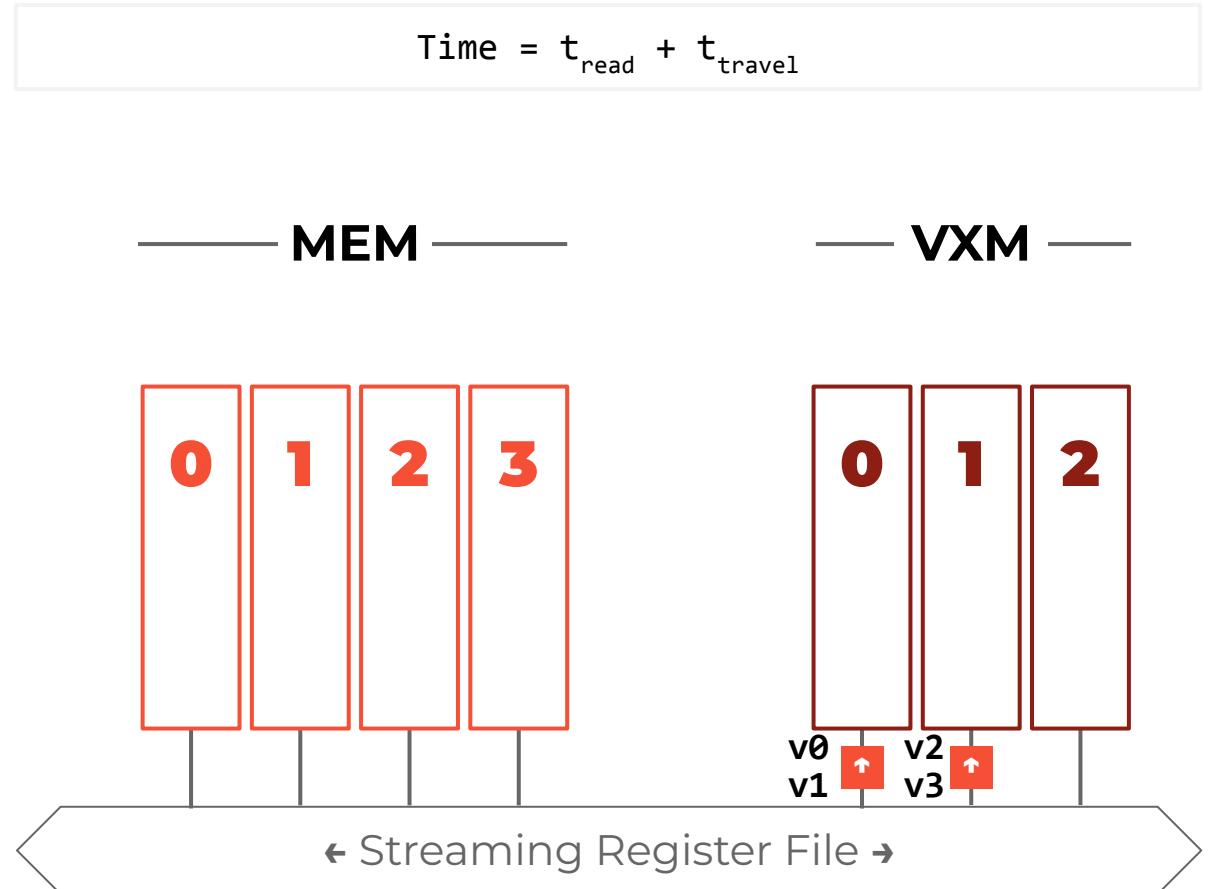
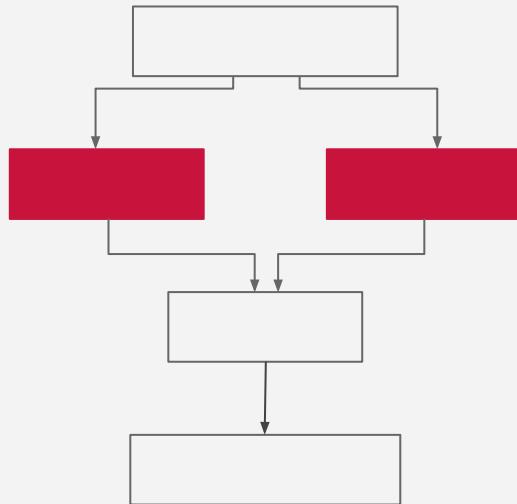


# On GroqChip

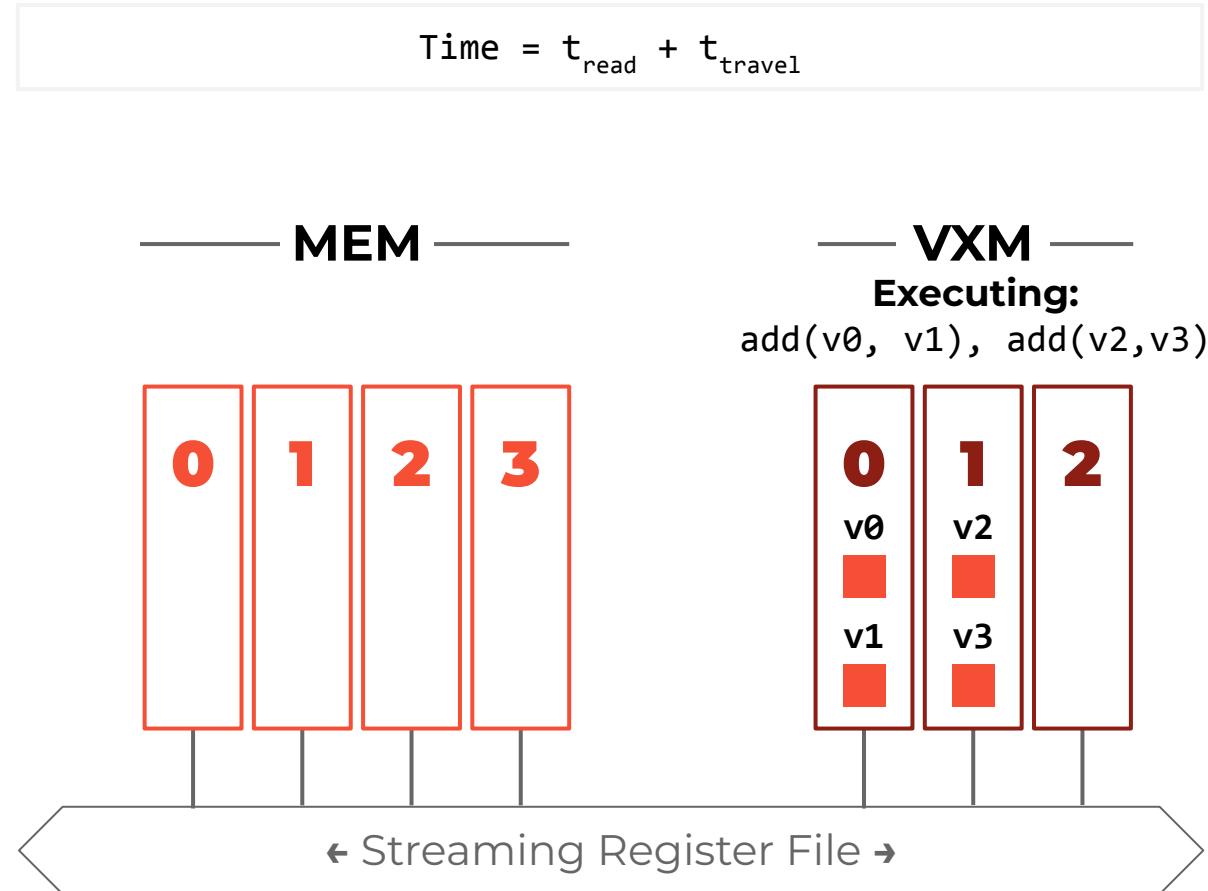
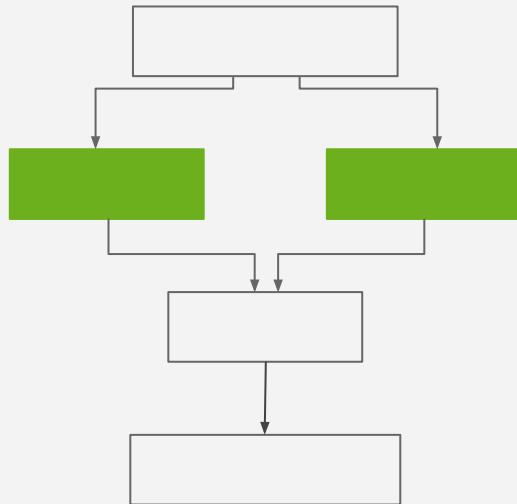
$$\text{Time} = t_{\text{read}} + 0.75 * t_{\text{travel}}$$



# On GroqChip

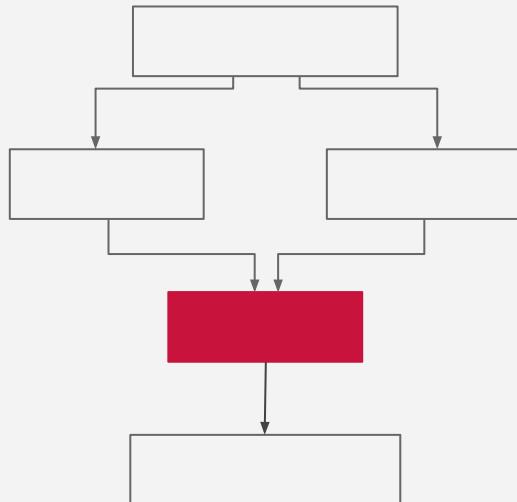


# On GroqChip



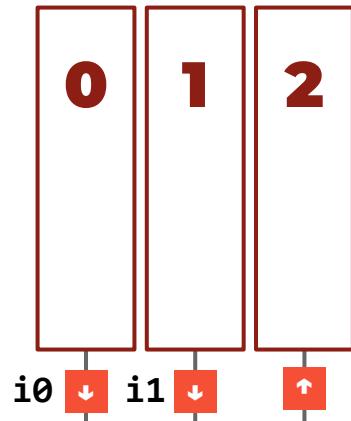
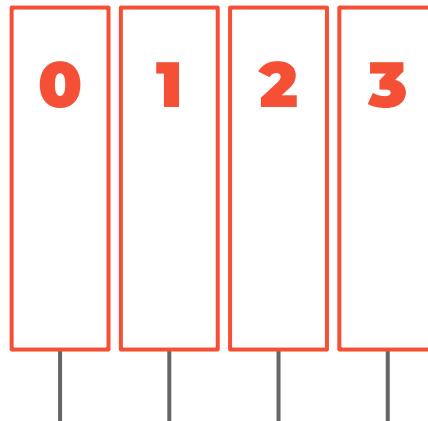
# On GroqChip

$$\text{Time} = t_{\text{read}} + t_{\text{travel}} + t_{\text{add}}$$



MEM

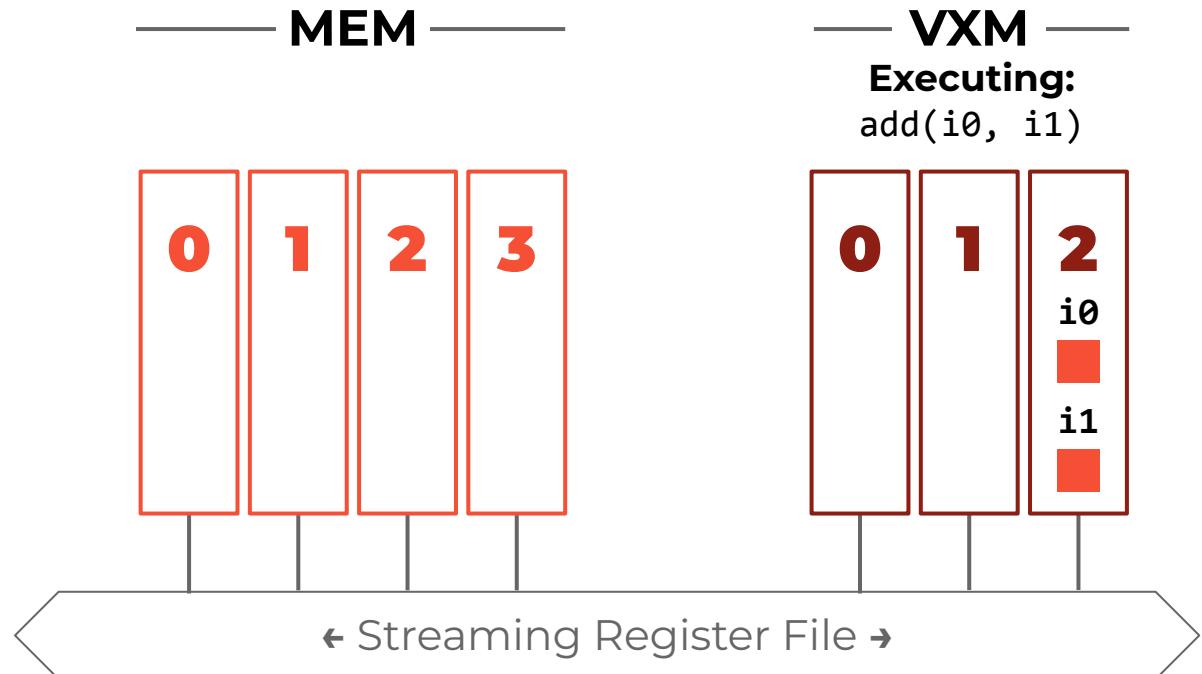
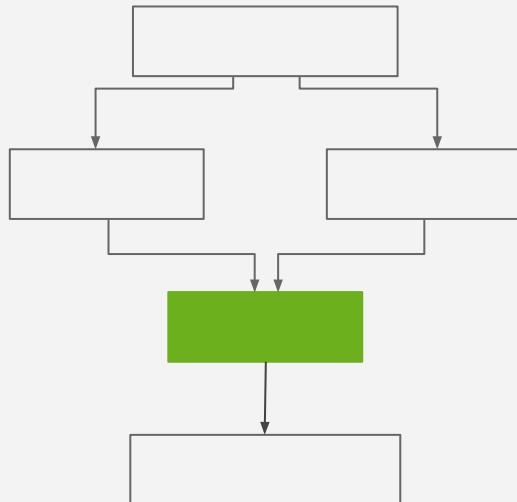
VXM



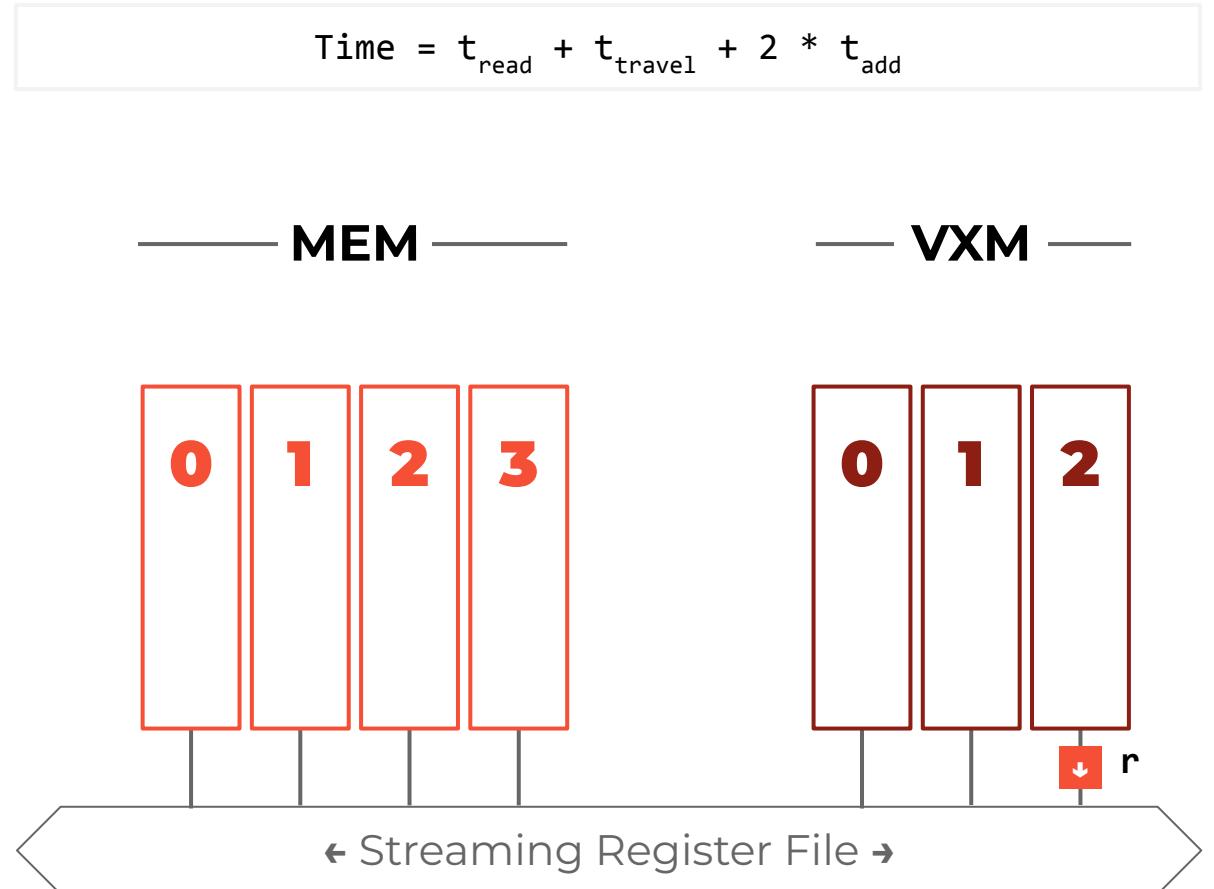
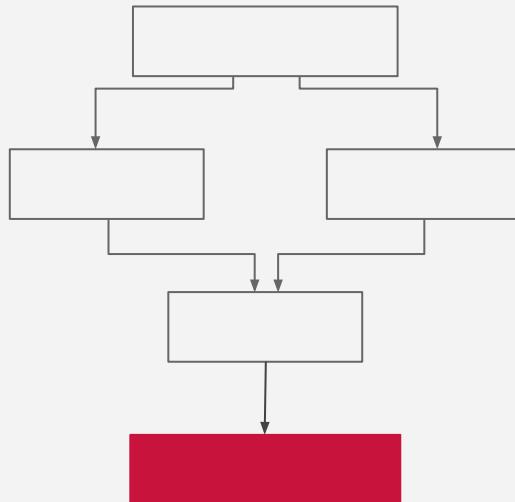
← Streaming Register File →

# On GroqChip

$$\text{Time} = t_{\text{read}} + t_{\text{travel}} + t_{\text{add}}$$

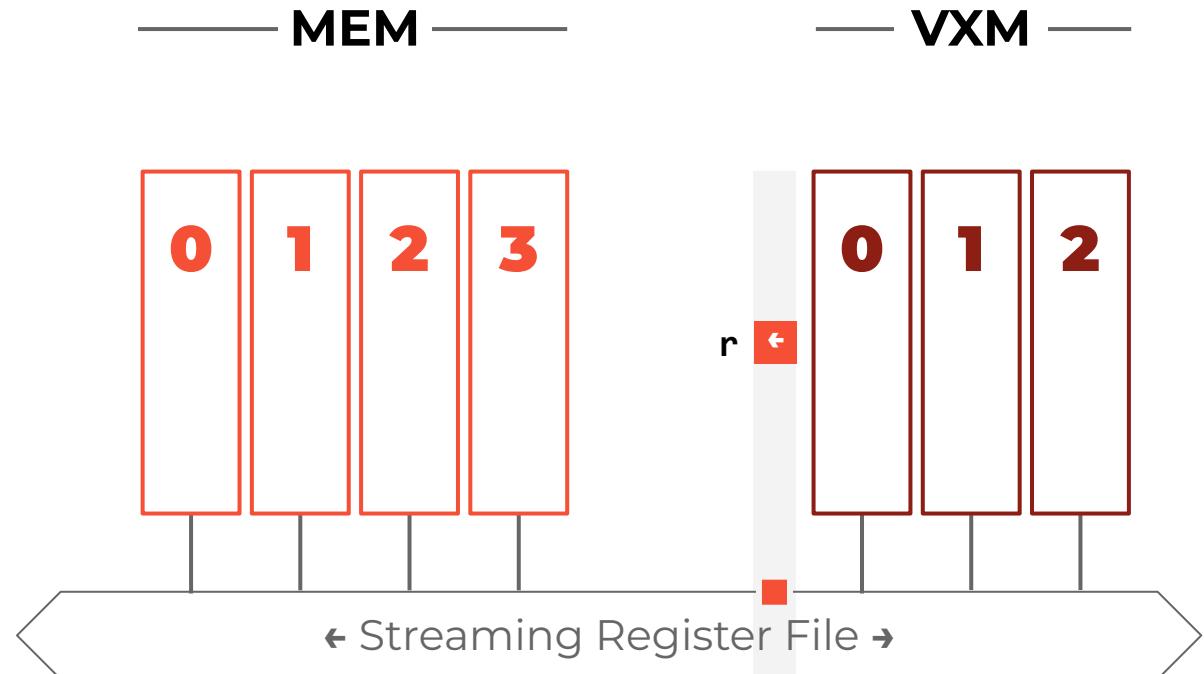
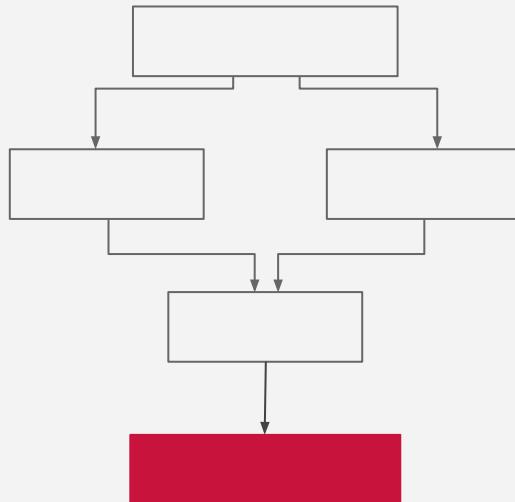


# On GroqChip

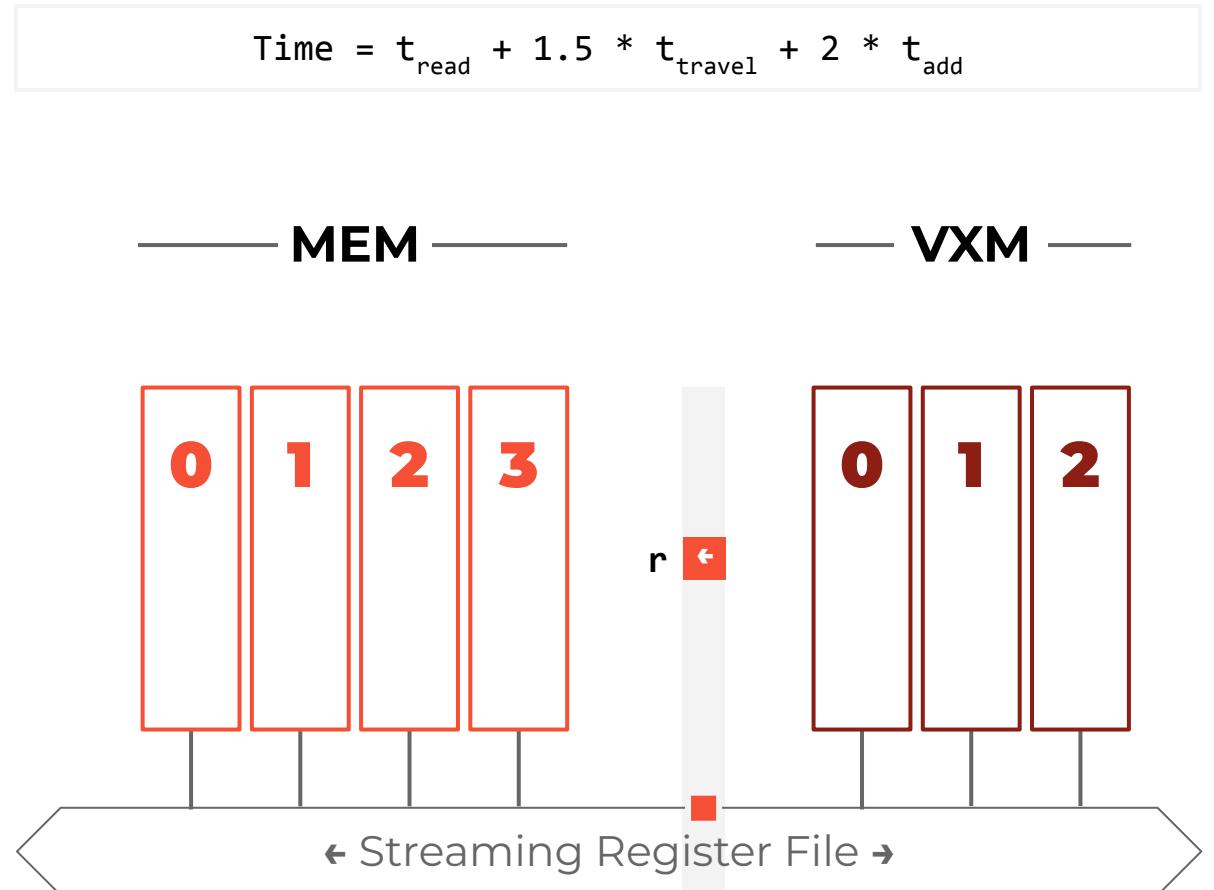
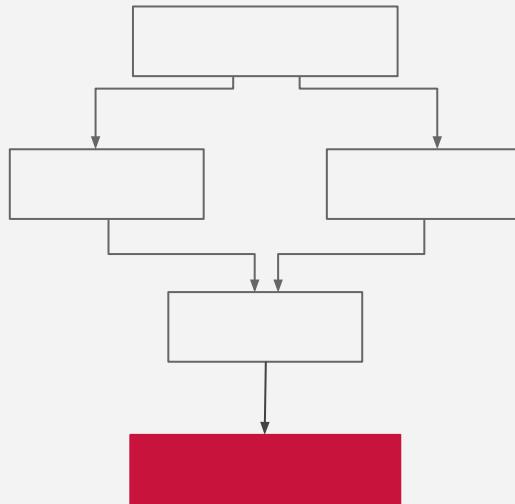


# On GroqChip

$$\text{Time} = t_{\text{read}} + 1.25 * t_{\text{travel}} + 2 * t_{\text{add}}$$

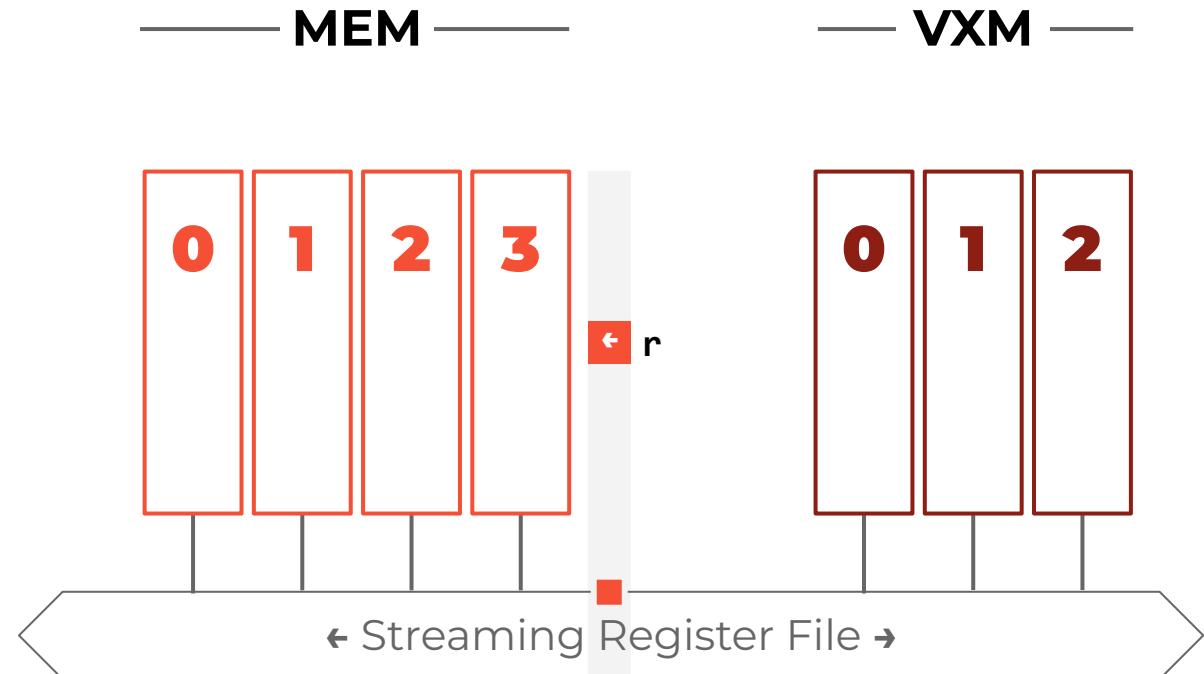
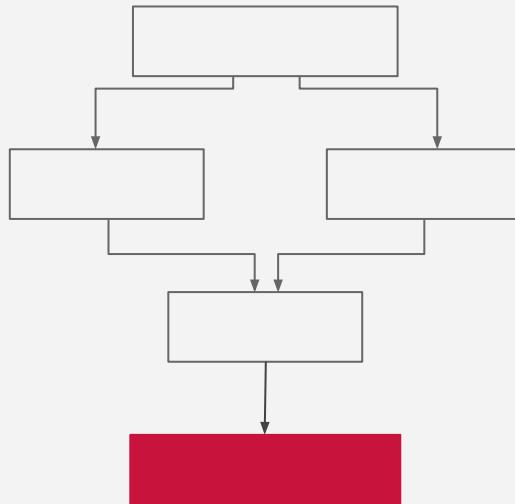


# On GroqChip

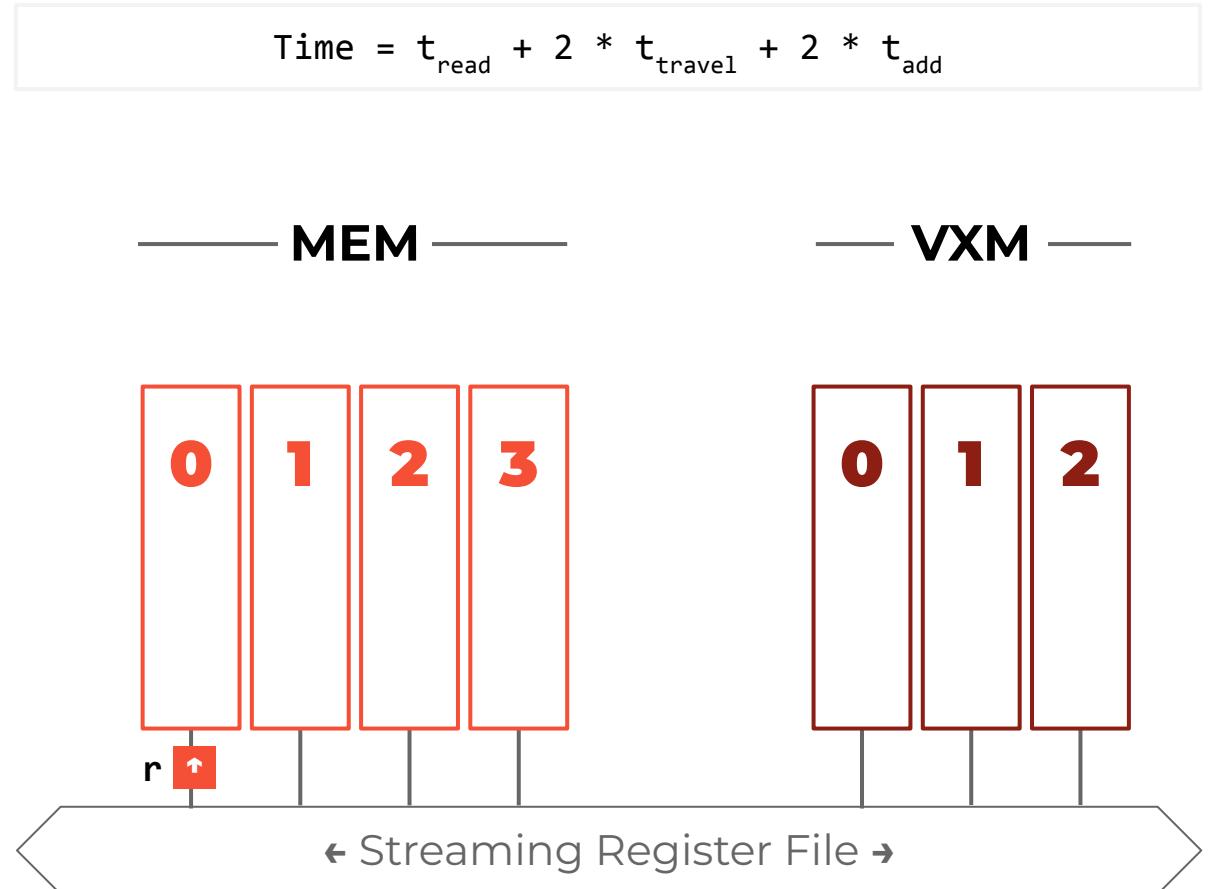
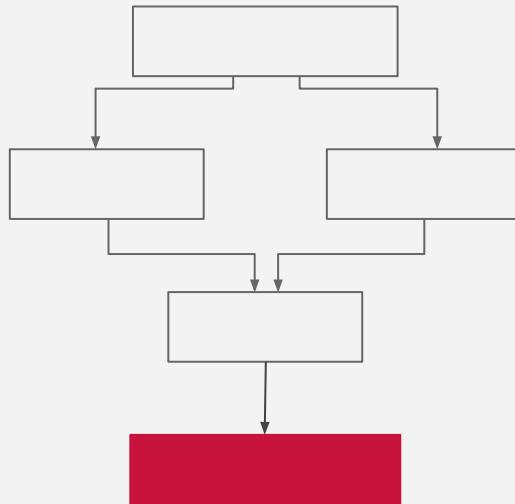


# On GroqChip

$$\text{Time} = t_{\text{read}} + 1.75 * t_{\text{travel}} + 2 * t_{\text{add}}$$

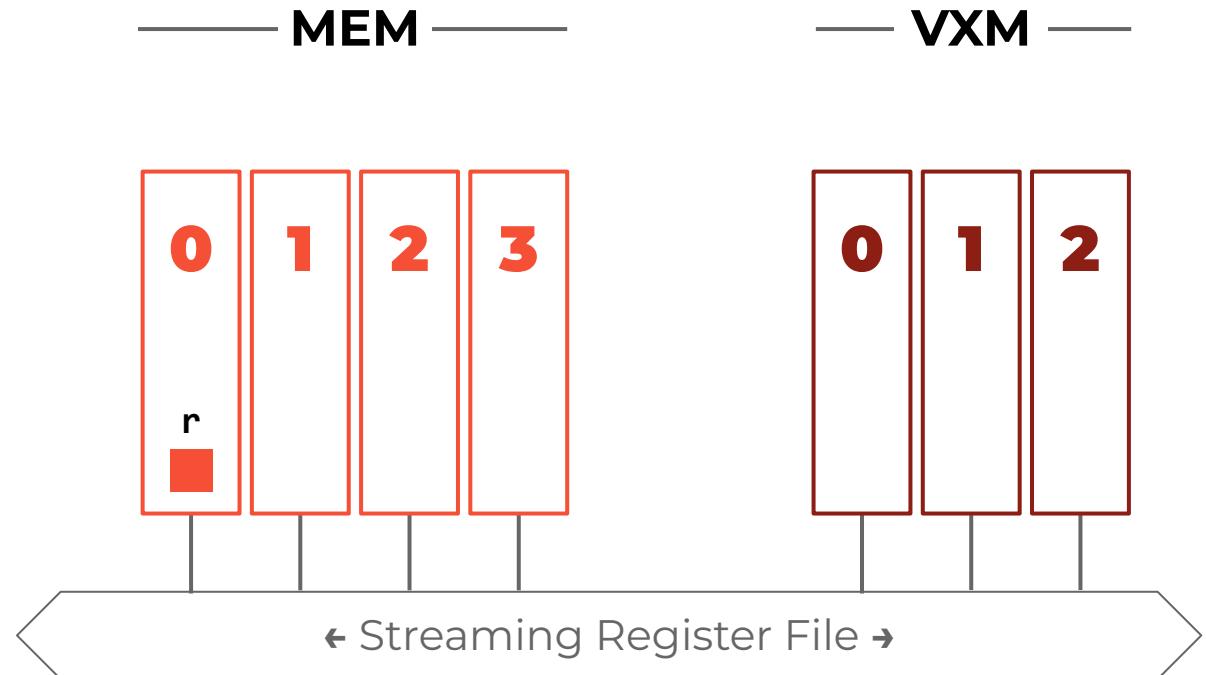
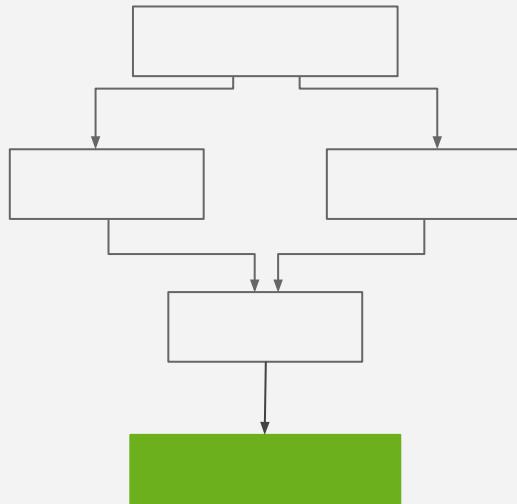


# On GroqChip



# On GroqChip

$$\text{Time} = t_{\text{read}} + 2 * t_{\text{travel}} + 2 * t_{\text{add}} + t_{\text{store}}$$



# Comparing Latencies

## GroqChip

$$\text{Latency} = t_{\text{read}} + 2 * t_{\text{travel}} + 2 * t_{\text{add}} + t_{\text{write}}$$

## Theoretical SIMD

$$\text{Latency} = t_{\text{dispatch}} + 2 * t_{\text{add}} + t_{\text{store0}} + t_{\text{store1}} + t_{\text{load0}} + t_{\text{load1}}$$

Unpredictable hardware causes pain for users

Predictable chips such as the GroqChip offers a holistic view of hardware

This holistic view enables effortless parallelism

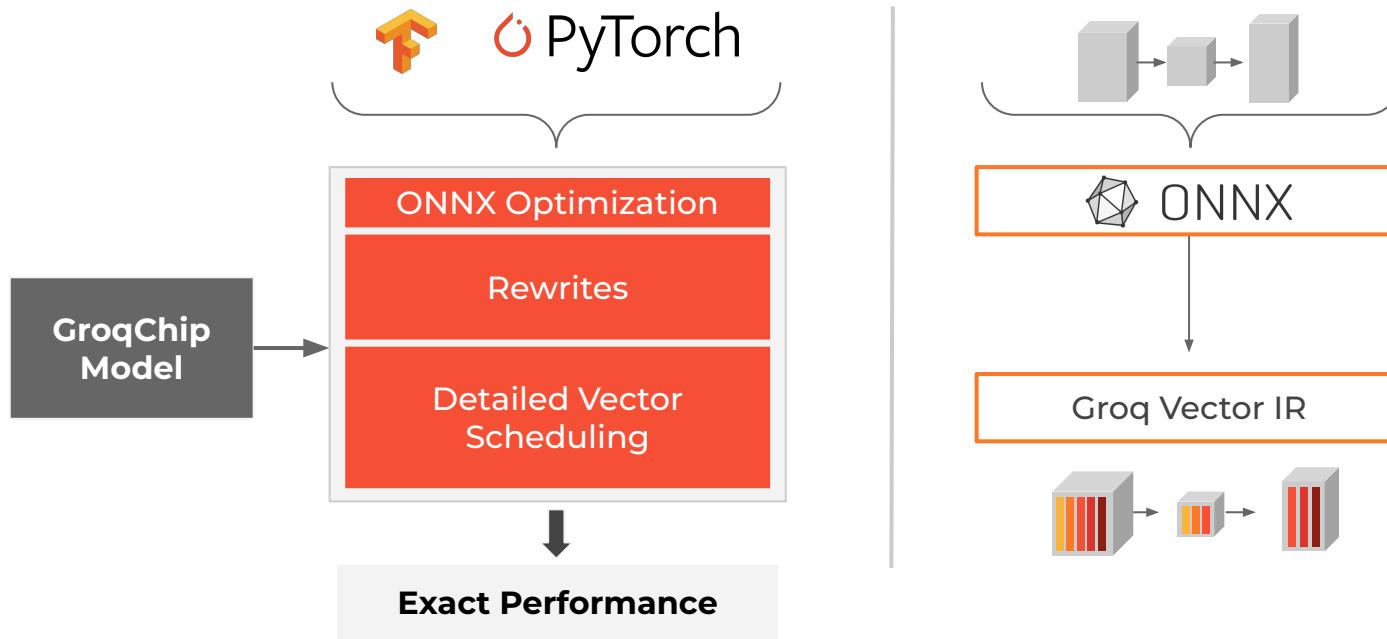
Red = unresolvable at compile time

Green = resolvable at compile time

# HARDWARE-SOFTWARE CO-DESIGN WITH The Groq Compiler

# Groq Compiler

Flow



**Kernel-less approach to HPC compilation**

# Building an Abstraction Model of GroqChip

**Architecture provides a natural fit for an abstraction model containing all information needed for compilation**

Specialized FUs with common interface → C++ templated polymorphism

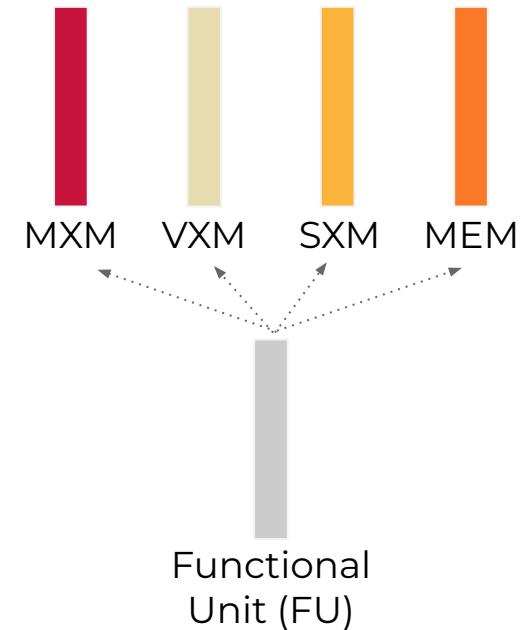
All instruction timing quickly resolvable at compile time via simple lookup table

All data movement quickly resolvable at compile time

- One dimensional interconnect abstracted as “timezone” indices

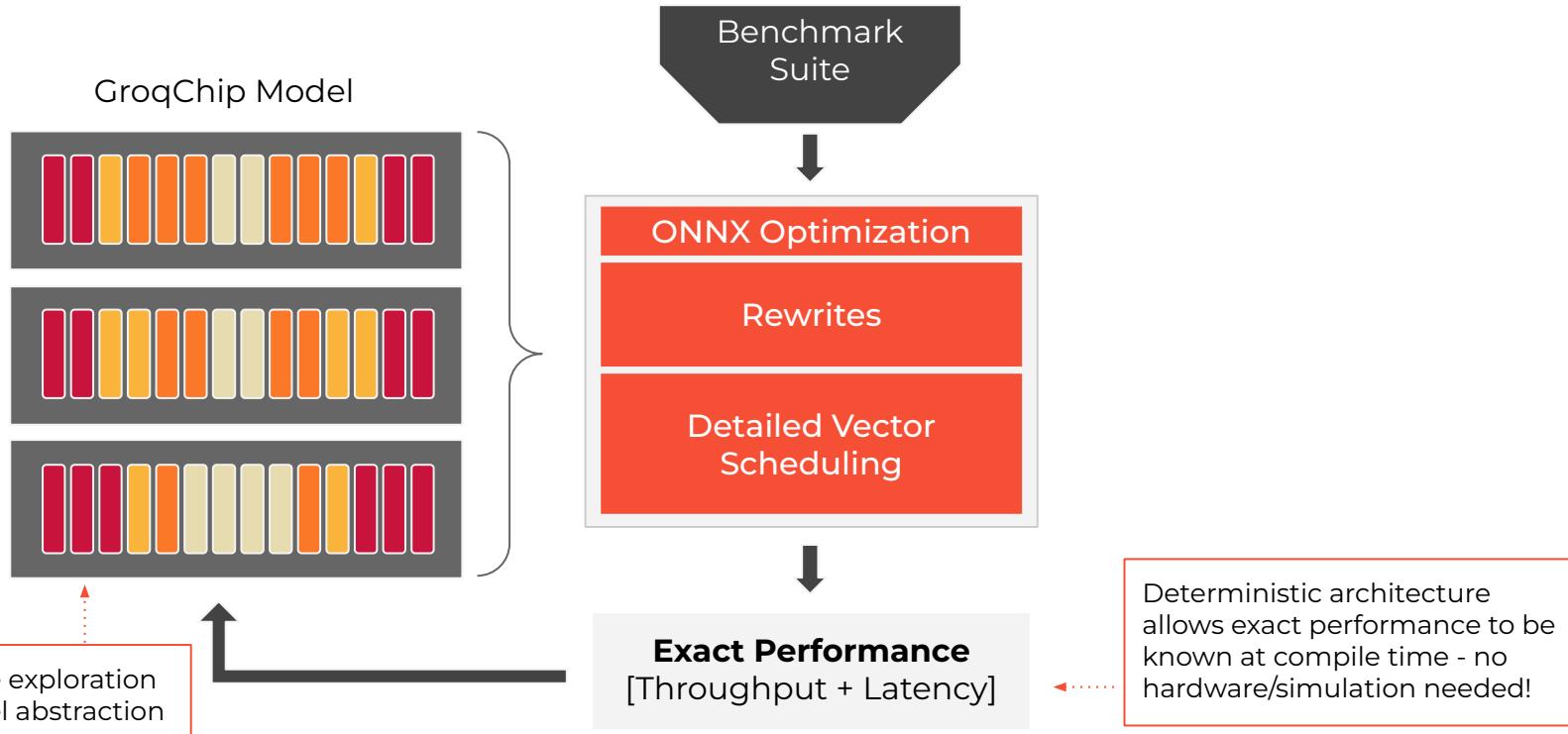
Efficient cycle-accurate resource allocation tracking via specialized data structures

- E.g. Bitvectors for fast range-allocation & lookup across time + space



# Hardware-software Co-design

Compiler-driven architecture exploration



# High-level Chip Models in FPGA CAD

**FPGA CAD tools have long been an example of hardware-software co-design via software abstractions of the chip**

FPGA CAD compiles HDL down to a bitstream configuring the chip (LUTs, DSPs, BRAM, routing)

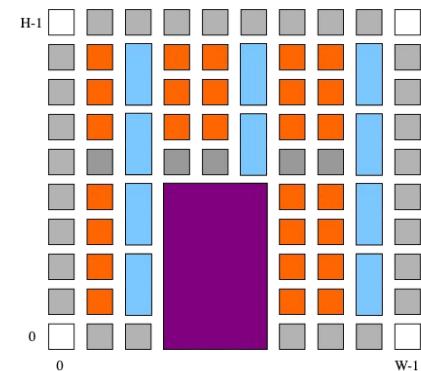
Metrics of the bitstream are statically determined  
(resource utilization, fmax)

Compilation requires a detailed, low-level chip model

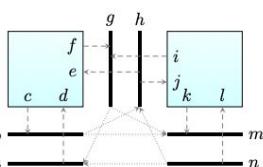
**Verilog-to-Routing (VTR) is an open-source FPGA CAD tool used for FPGA architecture exploration**

VTR can compile for any FPGA architecture that fits within its chip model framework

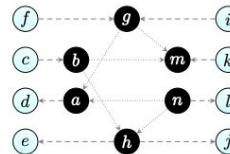
Originally introduced in VPR [Betz & Rose, FPGA'00]



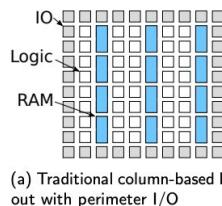
# High-level Chip Models in FPGA CAD



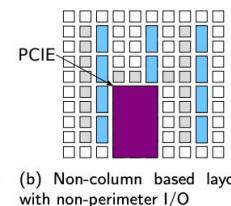
(a) Routing Architecture



(b) Corresponding Routing Resource Graph



(a) Traditional column-based layout with perimeter I/O



(b) Non-column-based layout with non-perimeter I/O

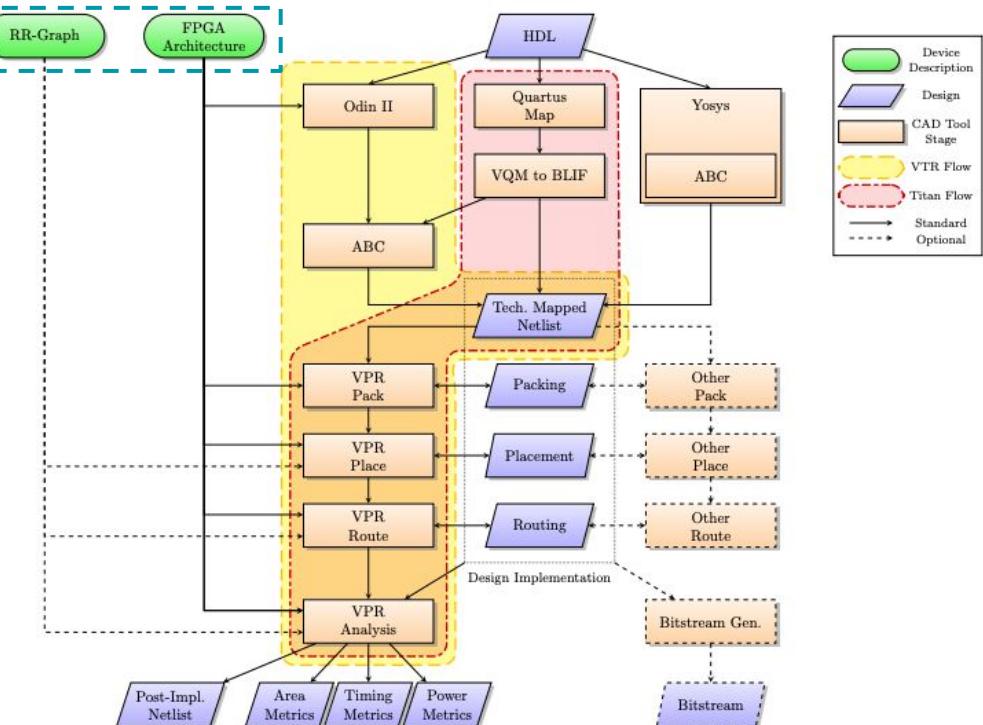


Fig. 1. VTR-based CAD Flow Variants

Murray, et al. VTR 8: High Performance CAD and Customizable FPGA Architecture Modelling

# Final Thoughts

## A New Golden Age of Computing

TSP, TPU, RISC-V, GPU, FPGA, CGRA

## The Contract Between Software and Hardware Has Been Reopened for Negotiation

What responsibilities are owned by the compiler vs. the hardware?

## Abstractions Are Key

If low-level control is given to software, can it be effectively used?

If low-level control is hidden from software, what performance may be lost?

What is the hardware cost?

Will the abstraction “leak”?

## Designing a Good Abstraction Requires Thinking Carefully About the Target Applications

Data-flow vs control-flow

Hardware-software co-design is key for developing end-to-end solutions that meet the challenging problems we face today



# Thank You

Contact me or  
the Groq team at  
[info@groq.com](mailto:info@groq.com)

**WE  
ARE  
HIRING**

FOLLOW US ON:



groq™

