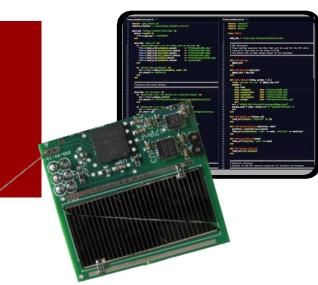
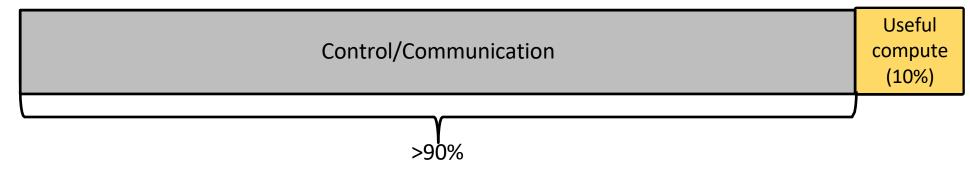
## Energy-minimal Computing Edge architectures for extreme efficiency





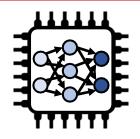
### Existing architectures are extremely inefficient

Instruction energy\* breakdown:



#### **Extreme Edge Computing Goal:**

increase energy-efficiency and preserve programmability



# Where does all the energy go in existing computer architectures?

#### Something is fundamentally wrong here:

Instruction energy\* breakdown:

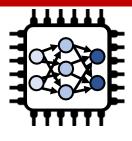
Fetch/Decode (40-50%)

Register file (20%)

Other control

(10%)

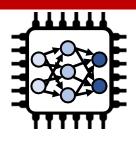
ASICs/Accelerators would improve this, but forfeit programmability



### Fundamental extreme edge trade-offs

### High **Programmability** Well-studied (GPUs, OOO, SIMD) Opportunity for the extreme edge High **Low Power** Performance

Well-studied (ASICs)



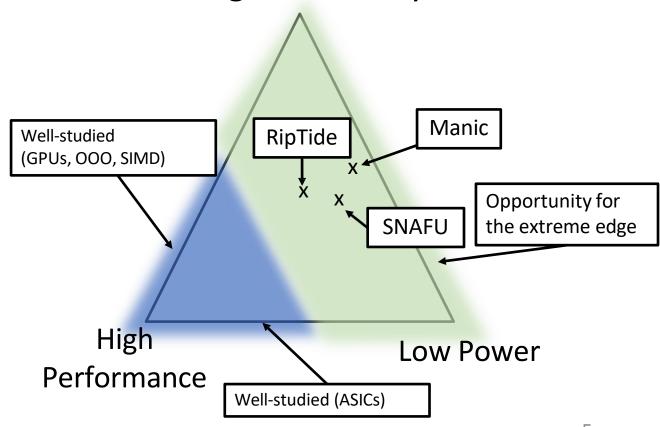
### Fundamental extreme edge trade-offs

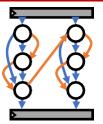
#### **Key Idea:**

Different architecture, different set of tradeoffs

Extreme edge applications demand **programmable & energy- minimal** architectures

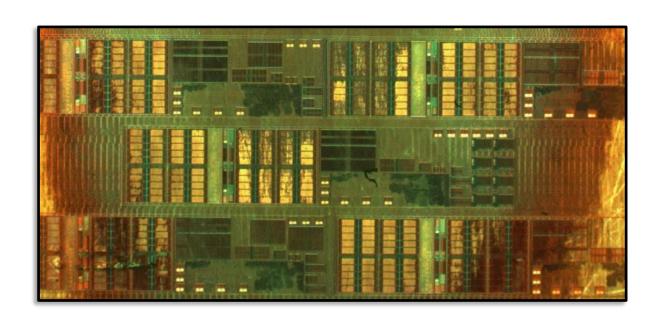




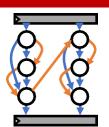


### MANIC: Extreme Edge Vector-dataflow processor

- Reduce instruction supply energy + VRF energy
- Maintain high-degree of programmability to support future kernels



Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	$\Box$	Û	Û
Vector- Dataflow	Ţ	Û	Û



#### Example Program

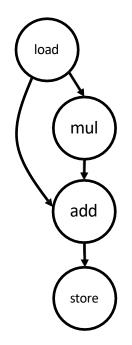
for i in 0...3:

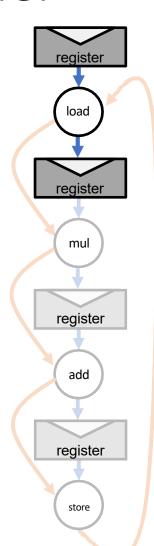
**b**oad r0, &a[i]

mul r1, r0, r0 add r2, r1, r0

store &b[i], r2

#### **Dataflow**



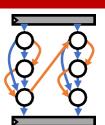


Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	⇧	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û

Dataflow

Control-flow

Related: MSP430, ARM M0



#### Example Program

for i in 0...3:

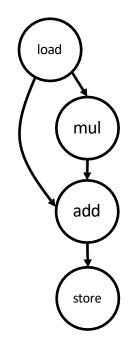
load r0, &a[i]

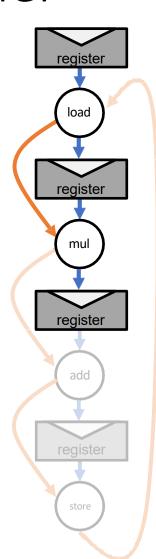
mul r1, r0, r0

add r2, r1, r0

store &b[i], r2

#### **Dataflow**

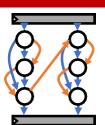




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û







#### Example Program

for i in 0...3:

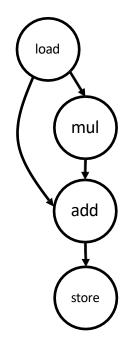
load r0, &a[i]

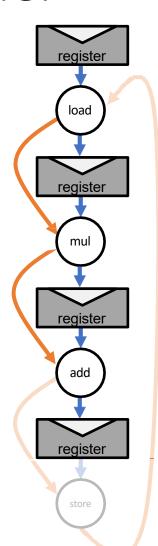
mul r1, r0, r0

**■**dd r2, r1, r0

store &b[i], r2

#### **Dataflow**

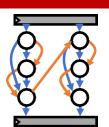




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û







#### Example Program

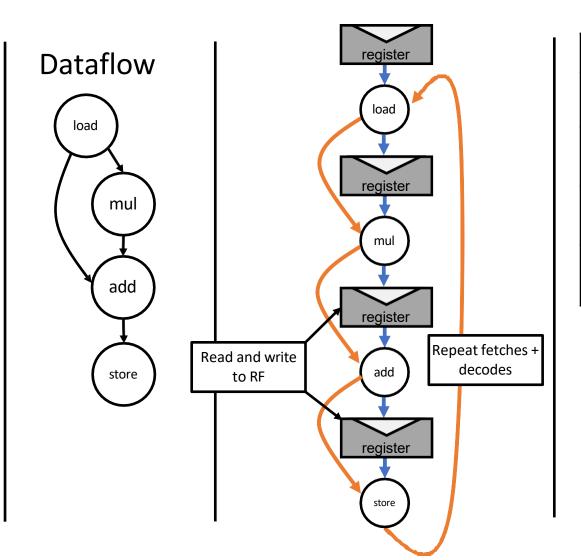
for i in 0...3:

load r0, &a[i]

mul r1, r0, r0

add r2, r1, r0

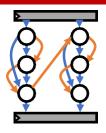
store &b[i], r2



Energy				
Model	Insns	RF Reads	RF Writes	
Scalar	Û	Û	Û	
Vector	Û	Û	Û	
Vector- Dataflow	Û	Û	Û	



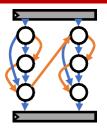




## Scalar execution is inefficient

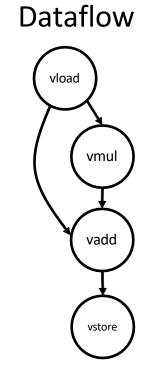
Energy wasted on instruction & data supply

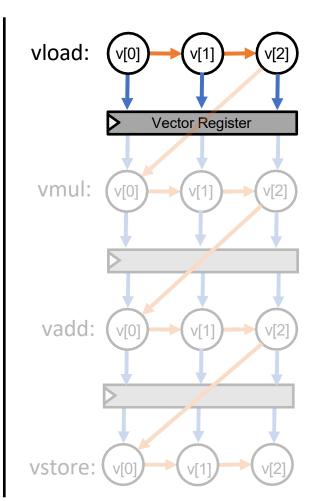
Memory	DCache access	ICache access	Compute + Control
--------	---------------	---------------	-------------------



#### **Example Program**

vload v0, &a
vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2

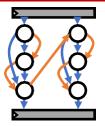




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	¢



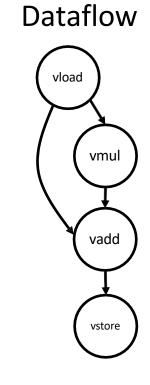


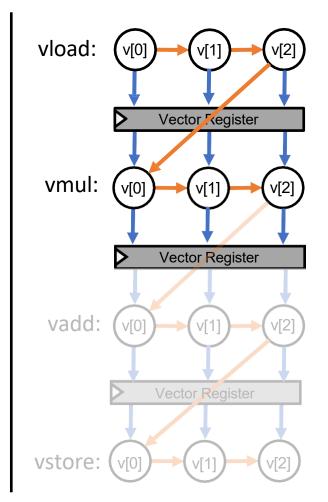


#### **Example Program**

vload v0, &a

vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2

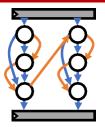




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û



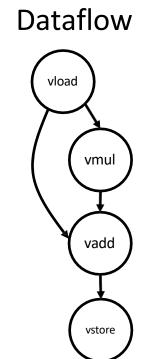


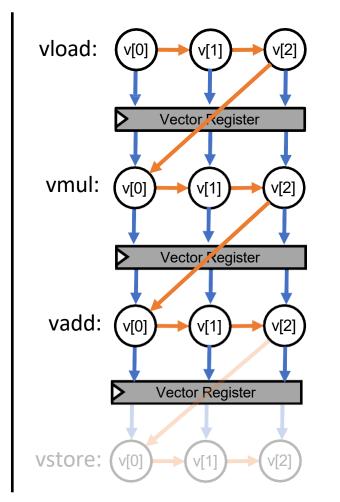


#### **Example Program**

vload v0, &a vmul v1, v0, v0

vstore &b, v2

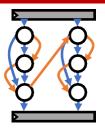




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û

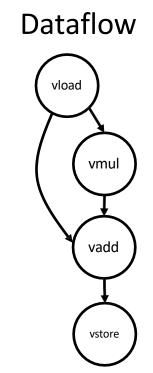


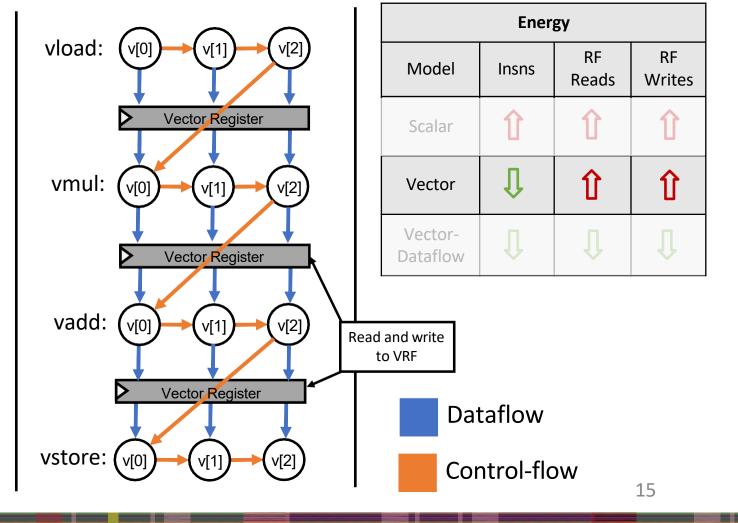


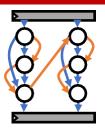


#### **Example Program**

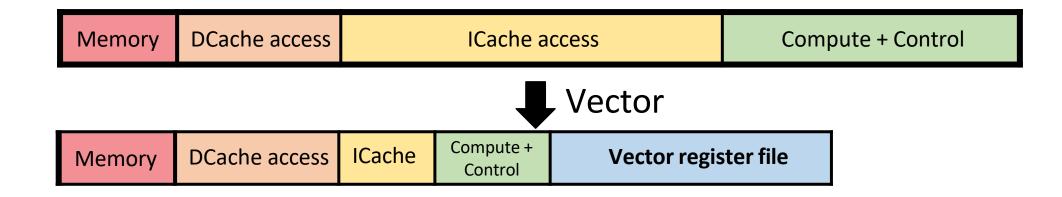
vload v0, &a
vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2

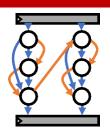






### Vector pays huge energy cost for VRF writes





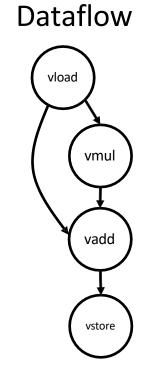
#### **Example Program**

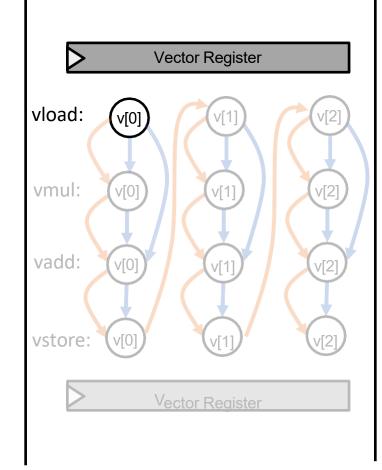
vload v0, &a

vmul v1, v0, v0

vadd v2, v1, v0

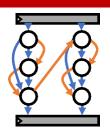
vstore &b, v2





Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û

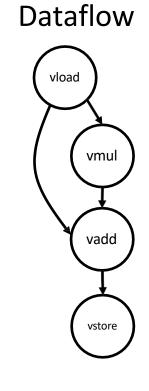


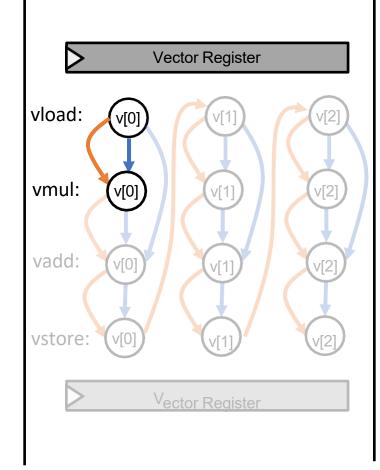


#### **Example Program**

vload v0, &a

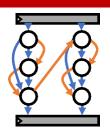
vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2





Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Ţ	Û	Û

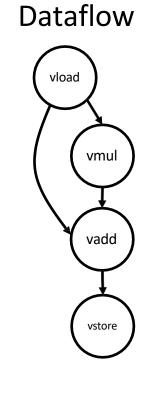


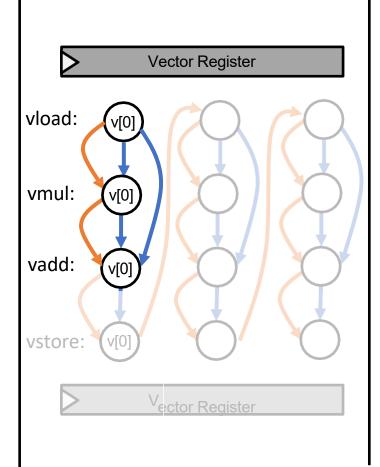


#### **Example Program**

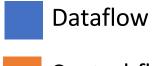
vload v0, &a vmul v1, v0, v0

vadd v2, v1, v0 vstore &b, v2

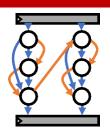




Energy			
Model	Insns	RF Reads	RF Writes
Scalar	Û	Û	Û
Vector	Û	Û	Û
Vector- Dataflow	Û	Û	Û



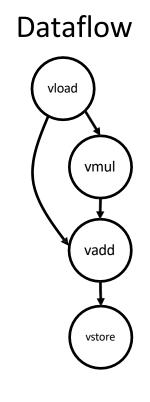


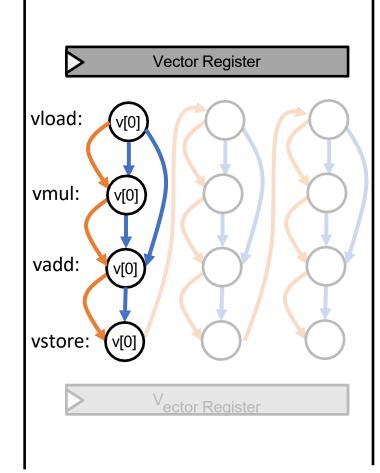


#### **Example Program**

vload v0, &a vmul v1, v0, v0 vadd v2, v1, v0

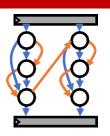
vstore &b, v2





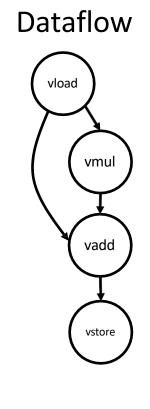
Energy				
Model	Insns	RF Reads	RF Writes	
Scalar	Û	Û	Û	
Vector	Û	Û	Û	
Vector- Dataflow	Û	Û	Û	

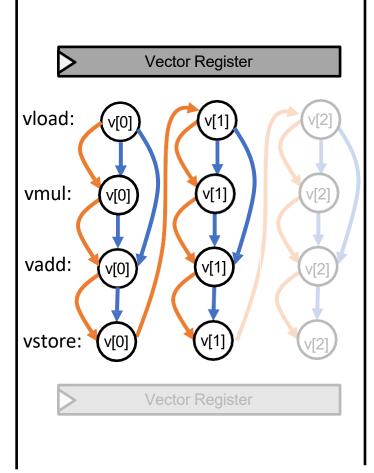




#### **Example Program**

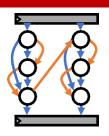
vload v0, &a
vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2





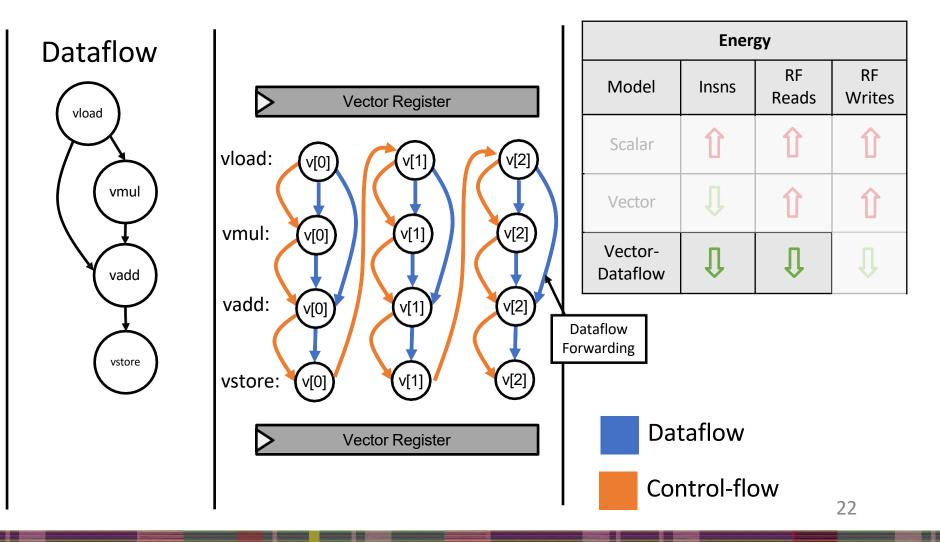
Energy				
Model	Insns	RF Reads	RF Writes	
Scalar	Û	Û	Û	
Vector	Û	Û	Û	
Vector- Dataflow	Ţ	Û	Û	

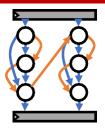




#### **Example Program**

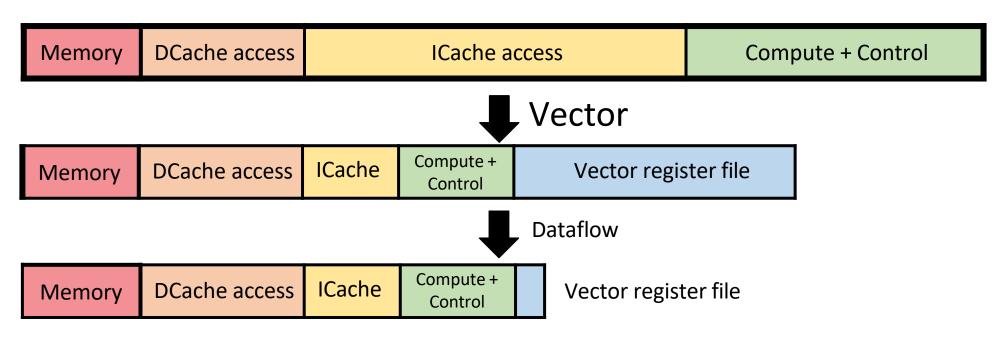
vload v0, &a
vmul v1, v0, v0
vadd v2, v1, v0
vstore &b, v2





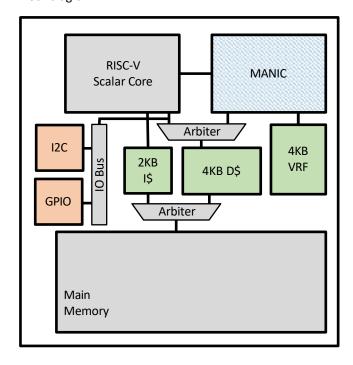
## Vector-dataflow reduces energy without costing programmability

- Vector-dataflow execution
  - Vector execution reduces instructions fetched
  - Dataflow execution eliminates VRF reads
- Software support to eliminate VRF writes



# MANIC is an energy-minimal computer architecture implementing vector-dataflow

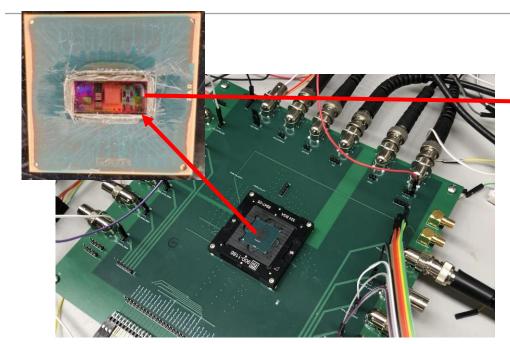
#### Block diagram



#### **Implementation Characteristics**

- Complete standalone system
- Scalar, Vector & MANIC designs
- Intel 22nm bulk FinFet (HVT)
- Embedded MRAM
- SRAM, logic, MRAM power isolated

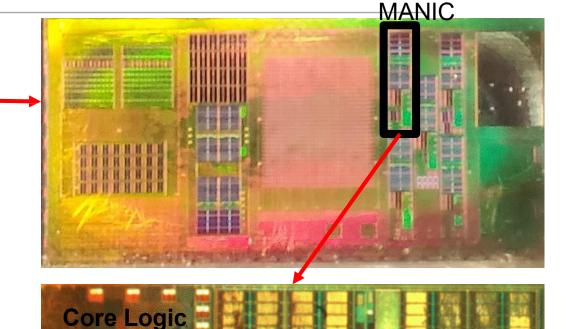
### Evaluating MANIC's efficiency in a silicon prototype



Intel 22nm FinFET 8 metal layers, MANIC + Vector + Scalar, 256kB MRAM + 64kB SRAM

**Evaluation Goals:** Energy characterization of first ever vector-dataflow chip.

**Key Result:** Power low & efficiency high enough to run on tiny solar panel indoors



**Operational Characteristics:** 

SRAM

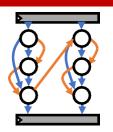
\$ + VRF

Frequency: 4-50MHz Voltage: 0.4-1.0V

Power: 19.1uW Efficiency: 256 GOPS/W

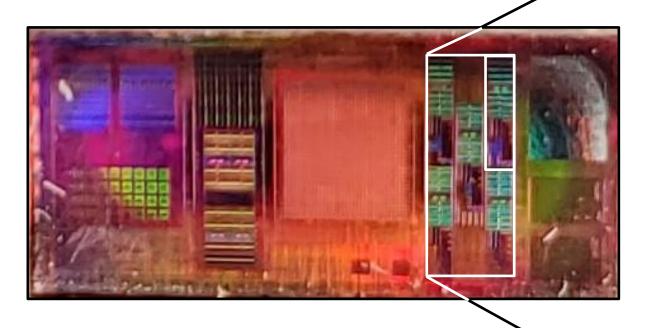
Carnegie Mellon University
Electrical & Computer Engineering

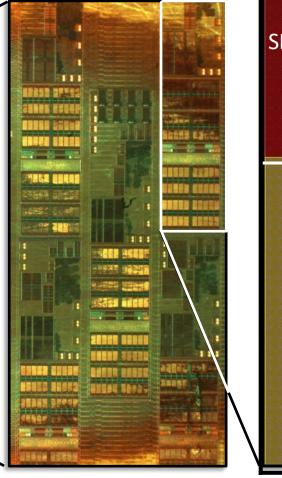
MRAM

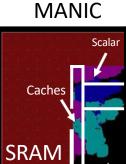


### More Pretty Chip Micrographs

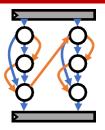
Courtesy CMU's Nanofabrication Laboratory & their electron microscope



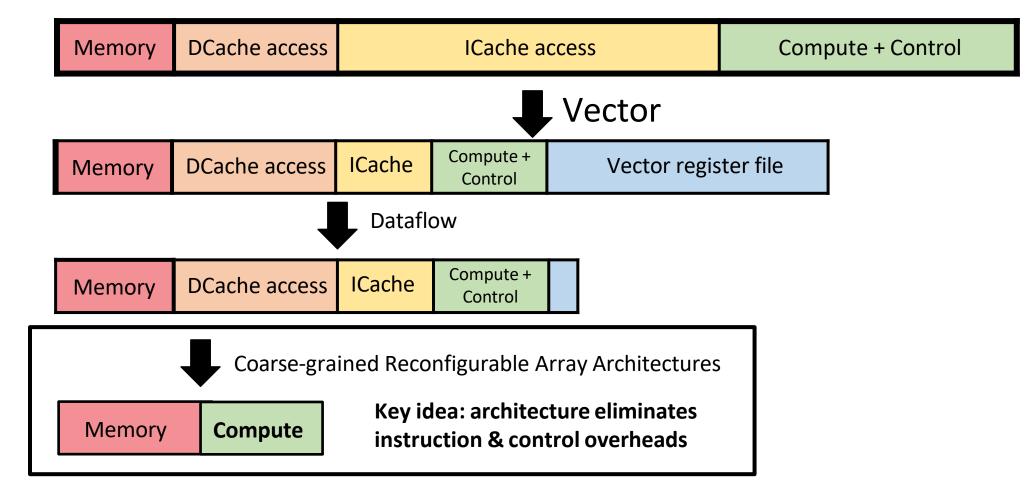


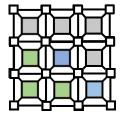


MANIC



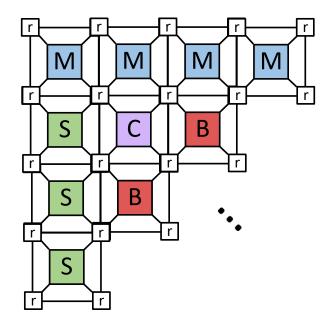
### Can we do even better? Let's eliminate all instruction control & caching costs!

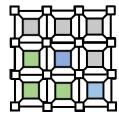




### **CGRA Overview**

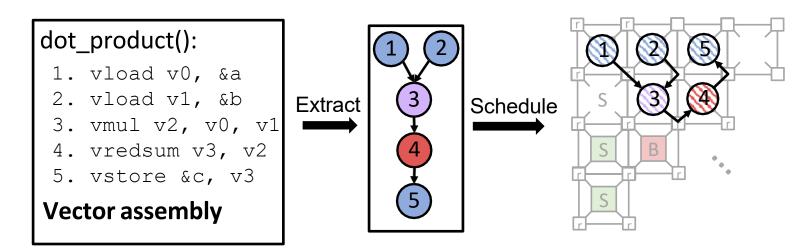
- Processing elements (PE) connected by Network-on-Chip (NoC)
  - Heterogenous PE capability
  - Connections configured by software compiler



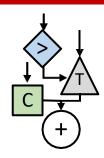


### **CGRA Overview**

- Collection of processing elements (PE) connected via NoC
  - Configure PE once, use many times: no instruction fetch/control costs
  - Data move directly PE to PE: no RF/VRF/Cache costs
  - Stream data through fabric: Reduced memory costs

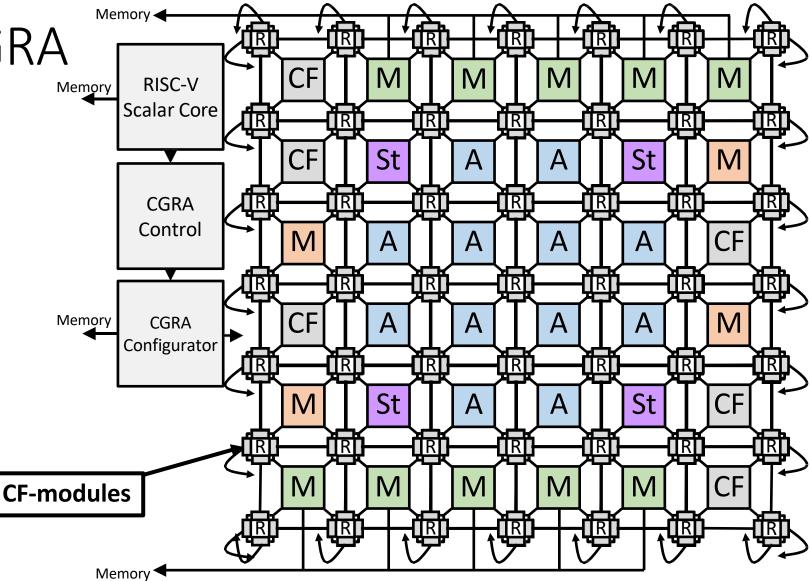


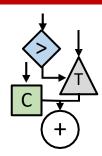
Nearly all energy for actually useful computation!



### RipTide CGRA

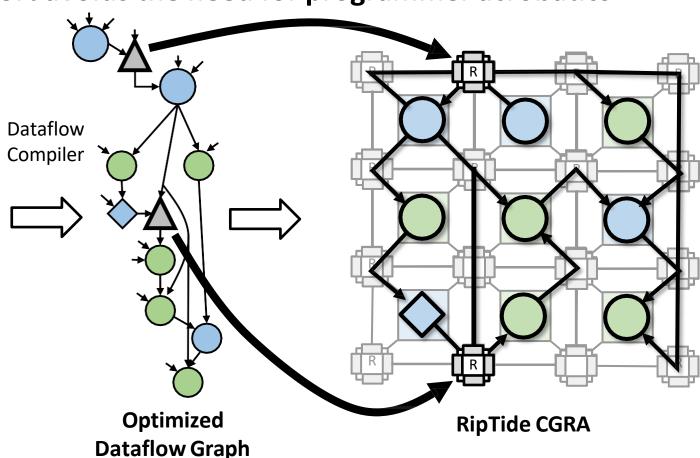
- Memory
- Multiplier
- St Stream
- A Arithmetic
- **CF** Control-flow





### CGRAs provide efficiency & programmability

**Dataflow** compiler support avoids the need for programmer acrobatics



### Dataflow Architecture

LABORATORY FOR COMPUTER SCIENCE



MASSACHUSETTS INSTITUTE OF TECHNOLOGY

The Varieties of Data Flow Computers

Computation Structures Group Memo 183-1 August 1979 Revised December 1979

Jack B. Dennis

### Dataflow Architecture: Dataflow Program Graphs

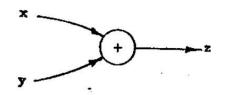


Fig. 2. Data flow actor.

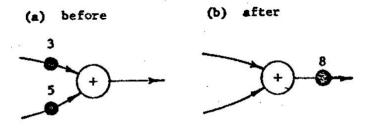


Fig. 3. Firing rule.

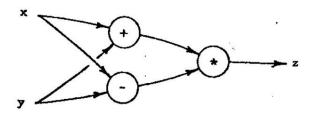
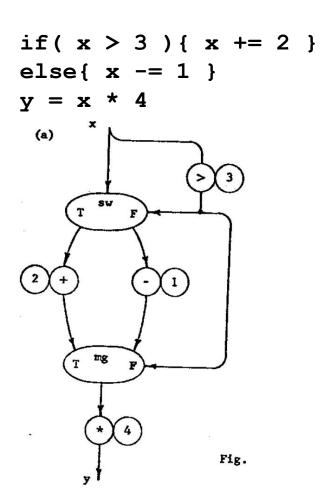
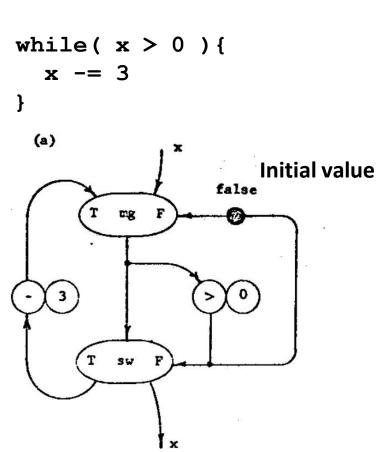


Fig. 4. Interconnection of operators.

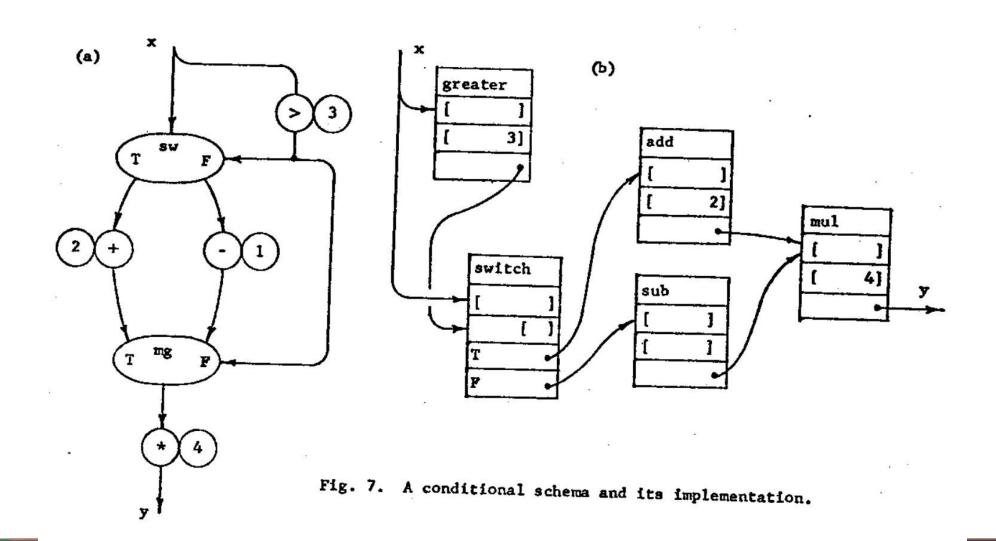


Conditional operations represented as dataflow

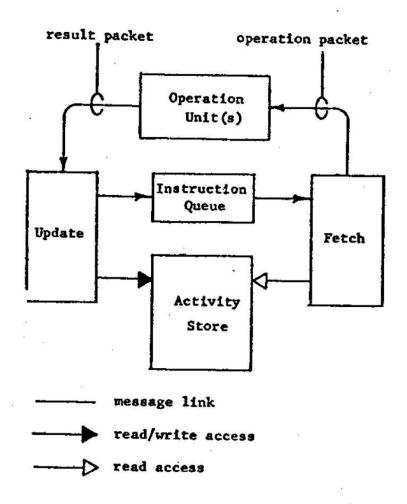


Loop represented as dataflow

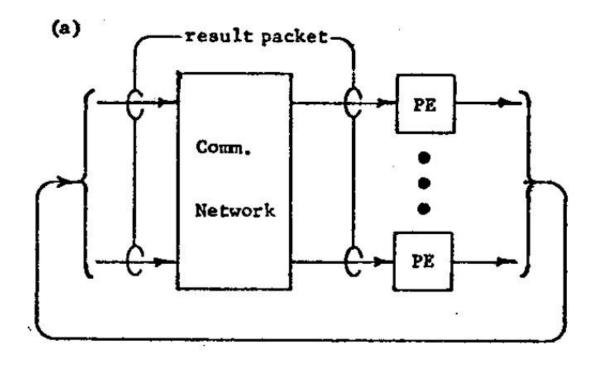
### Dataflow: "Activity Template" implementation



### Dataflow: Processing Element & Interconnect Arch.



**Processing Element Architecture** 



**Processing Element Interconnection Architecture** 

Question: what do we need to specify in this ISA?

### Dataflow: MIT Dataflow Architecture

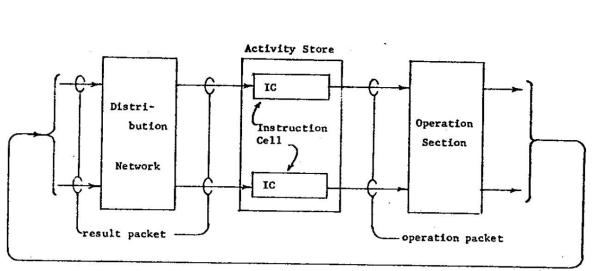


Fig. 14. MIT data flow processor.

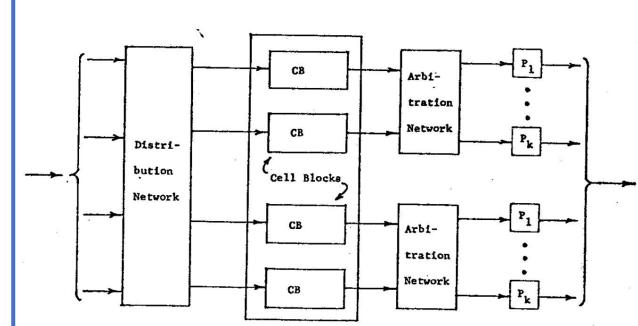
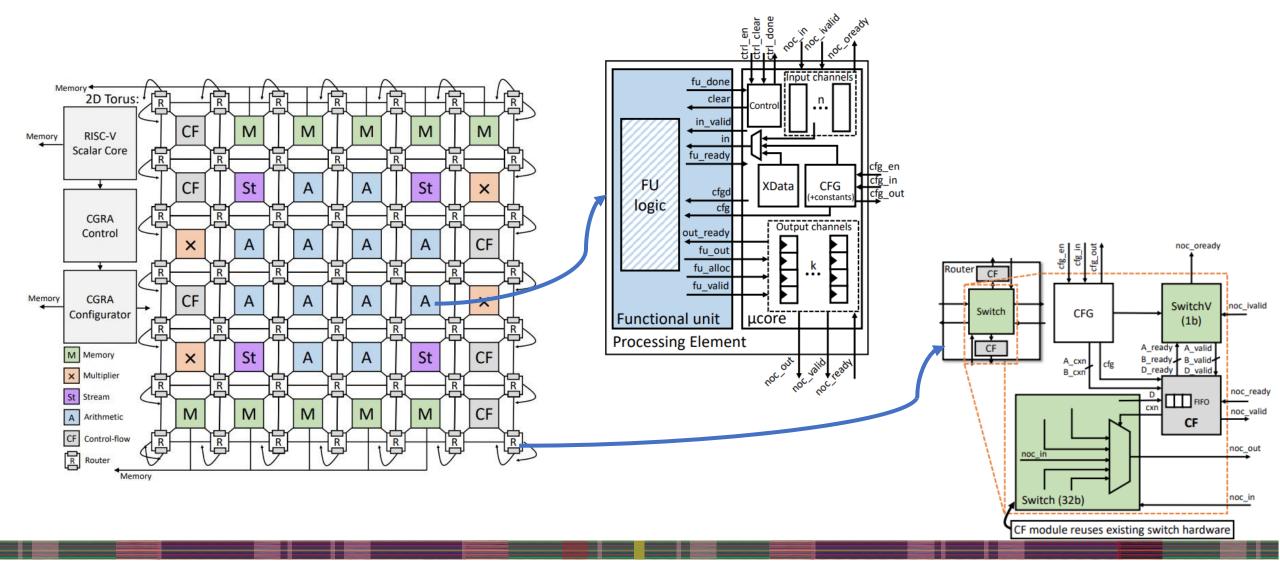


Fig. 15. Practical form of the MIT architecture.

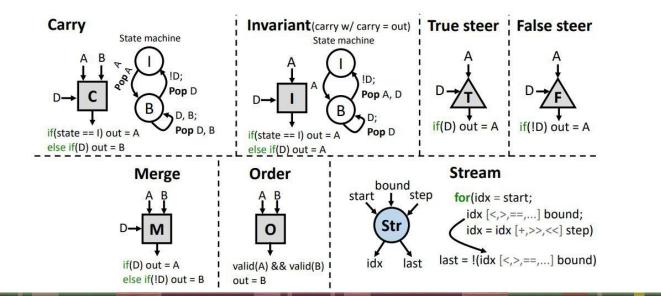
What is the main difference in this architecture versus the

# Dataflow: The Riptide Ordered Dataflow Machine



# Dataflow: The Riptide ISA

Operator(s)	Category	Symbol(s)	Semantics
Basic binary ops	Arithmetic	+, -, <<, !=, etc.	<i>a</i> op <i>b</i>
Multiply, clip	Multiplier	*, clip	$a  ext{ op } b$
Load	Memory	ld	1d base, idx(, dep)
Store	Memory	st	st base, idx, val(, dep)
Select	Control Flow	sel	cond? val0: val1
Steer, carry, invariant	Control Flow	(T   F), C, I	See Fig. 3
Merge, order	Synchronization	M, O	See Fig. 3
Stream	Stream	STR	See Fig. 3



What program constructs do the **carry** and **invariant** ISA ops support?

What does the **order** ISA op do?

What program construct(s) does the stream ISA op support?

# Background: What is a Dataflow Machine?

A Preliminary Architecture for a Basic Data-Flow Processor

Jack B. Dennis and David P. Misunas
Project MAC
Massachusetts Institute of Technology

An example of a program in the elementary data-flow language is shown in Figure 1 and represents the following simple computation:

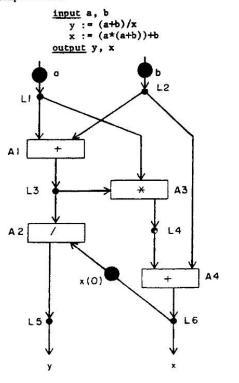
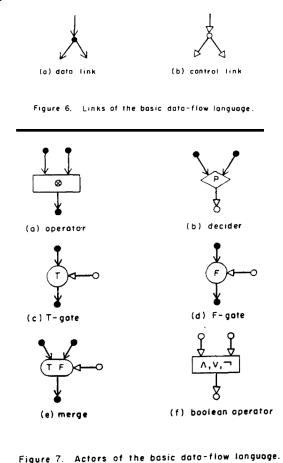
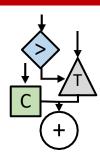


Figure 1. An elementary data - flow program



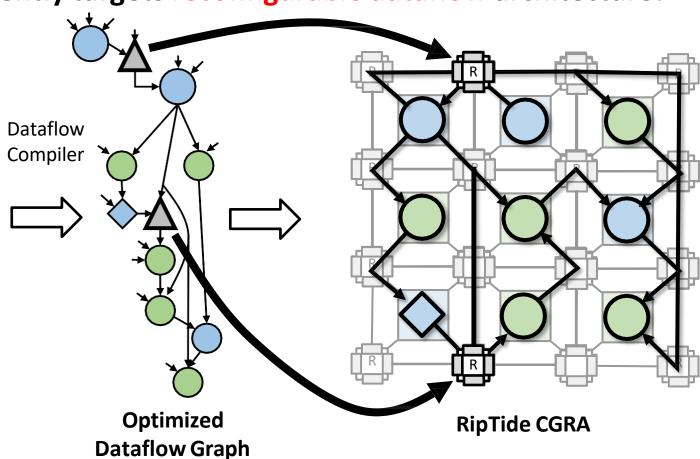
Operation Unit Operation Unit m-1 data operation packets packets Instruction Distribution Arbitration Network Memory Network Instruction Cell n-1

Figure 2. Organization of the elementary data-flow processor.

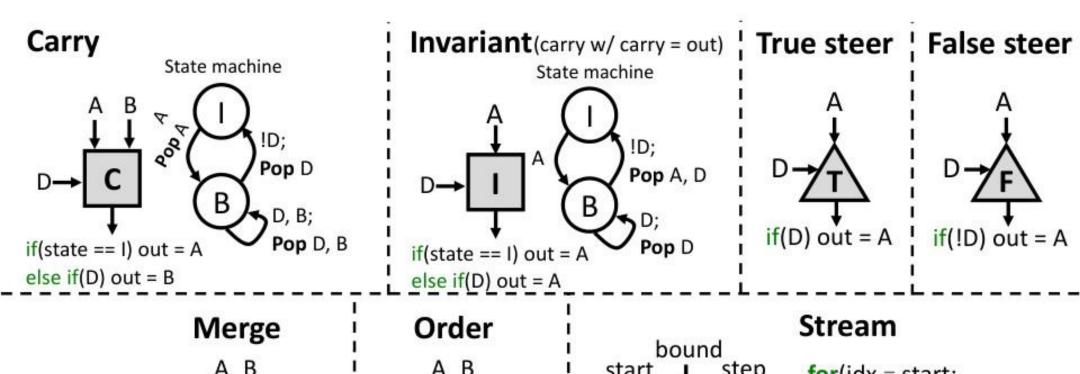


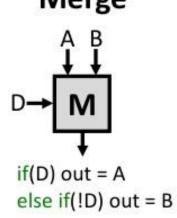
# Dataflow ISA matches CGRA Architecture

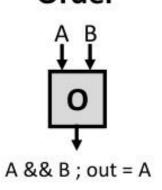
Dataflow compiler efficiently targets reconfigurable dataflow architecture!

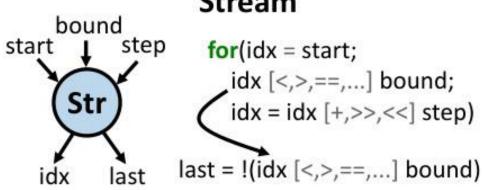


### Intermediate representation for dataflow compilation

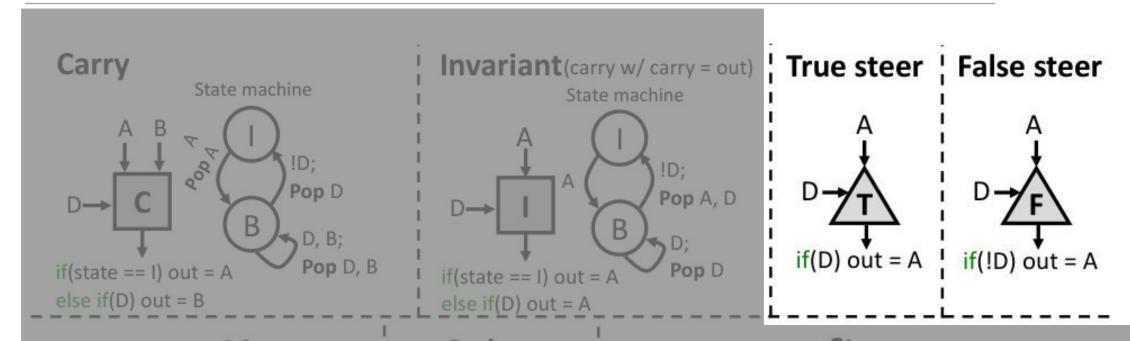






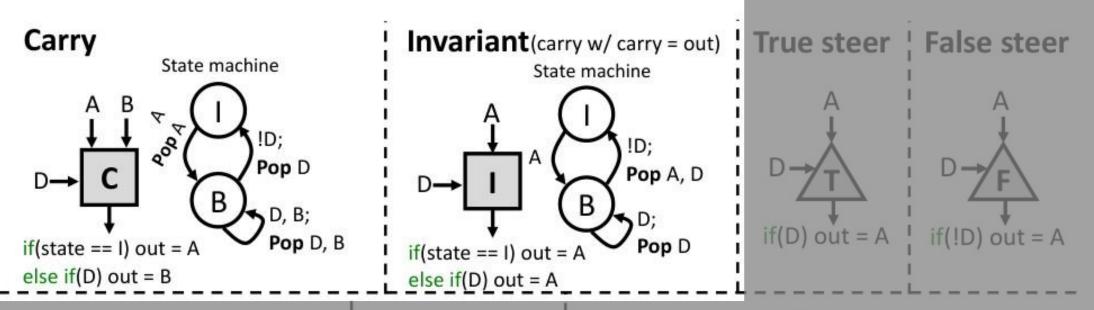


### Steering Sends Values Only Where They Are Needed



Steering instead of predication Avoids evaluating both sides of predication does of pr

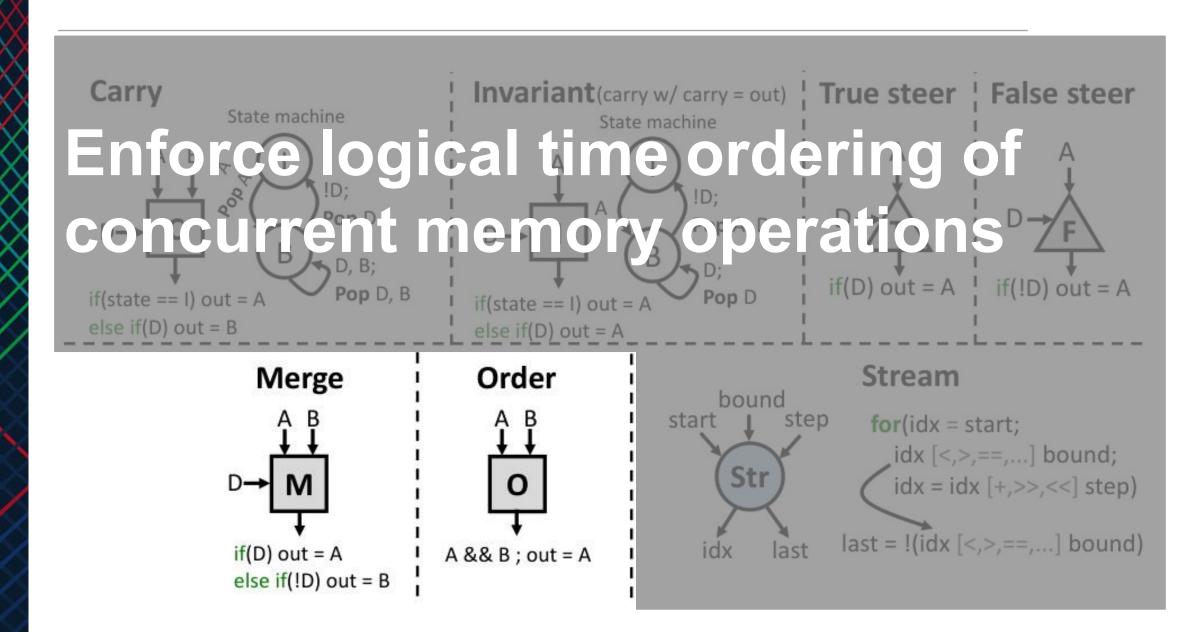
### Control-flow Operators Handle Loop-Carried Dependences



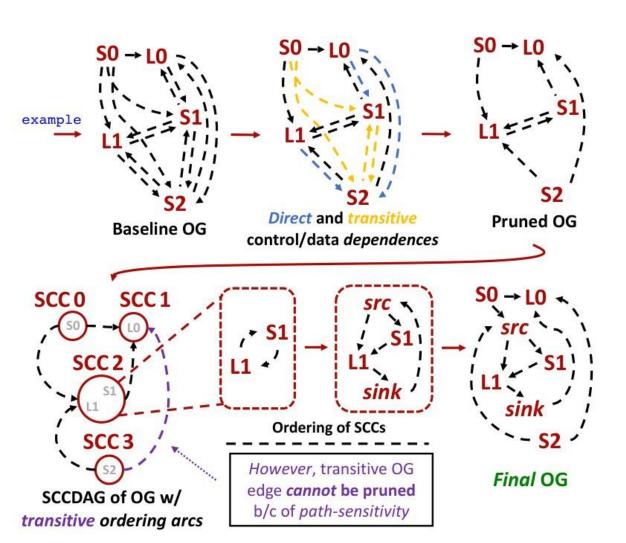
Generate fresh value tokens for loop carried dependences and [+,>>,<] step)

loop-invariant values idx last last =!(idx [<,>,==,...] bound)

### Memory Ordering Gates Maintain Memory Consistency



### Memory Ordering Reduction Analysis

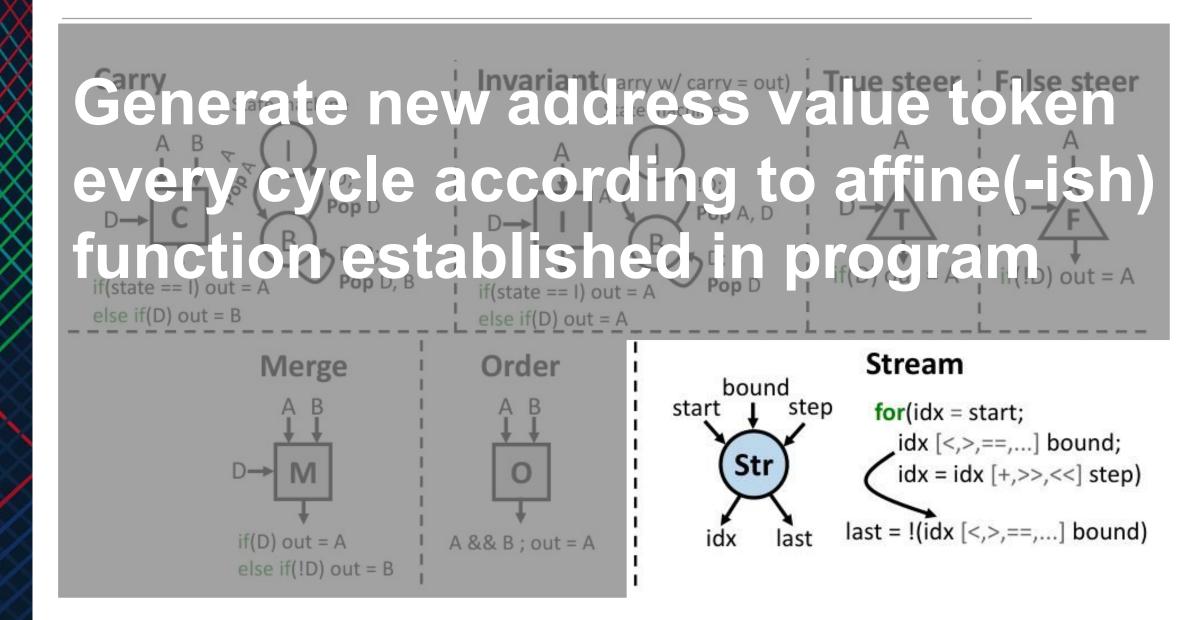


Existing memory ordering analyses rely on *transitive reduction* of ordering graph

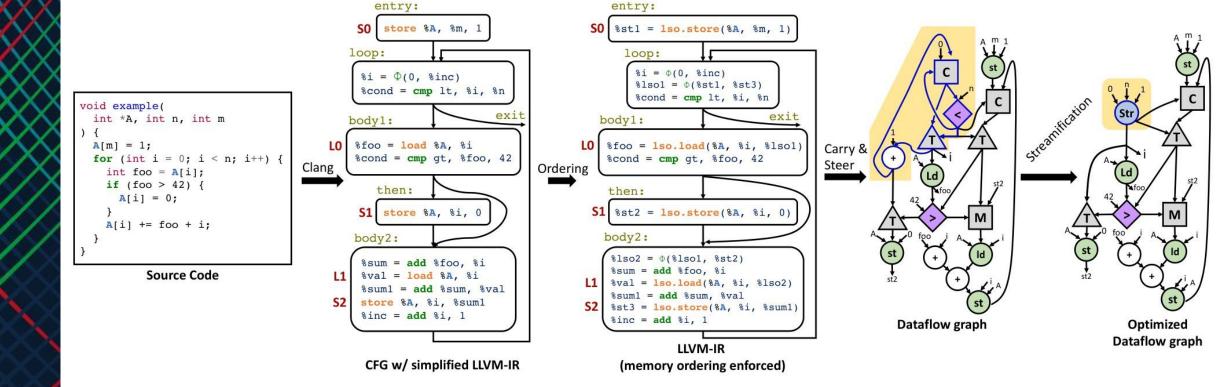
Dataflow ordering reduction requires *path* sensitivity or cuts required orders.

Carnegie Mellon University
Electrical & Computer Engineering

### Stream Gates Optimize Patterned Address Computation



### End-to-end compilation flow

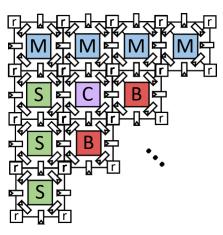


# Energy-Minimal Network-on-Chip Tricks: No Buffers & Control-Flow in NoC

#### ✓ Multi-hop, bufferless NoC

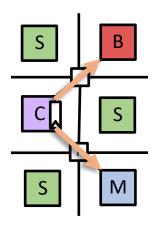
V.

Buffer on every link



Duplicates data in multiple buffers

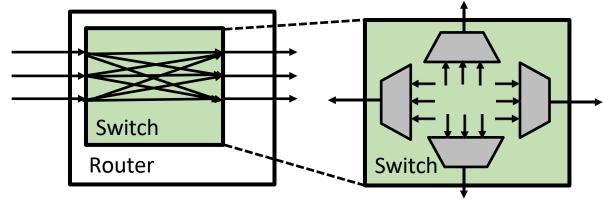
Buffers @ producer

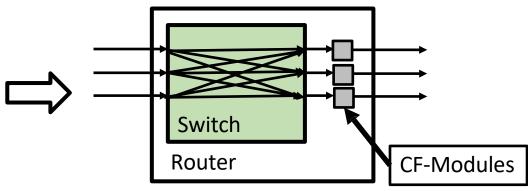


Allows broadcast to multiple consumers w/out duplication

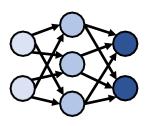
#### ✓ Control-flow in the NoC

- Control-flow operations are numerous, but simple
- → Wasteful to assign to PEs

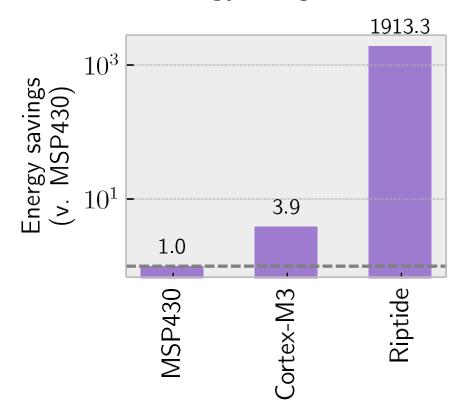




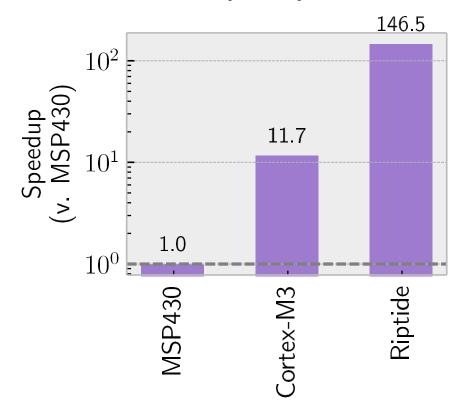
# Riptide vs. COTS Extreme Edge MCUs

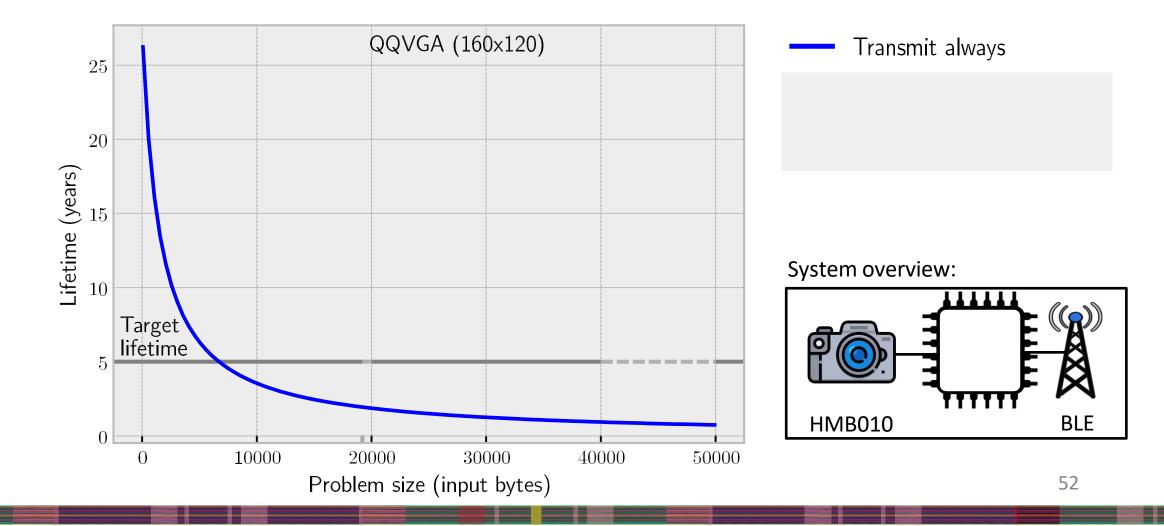


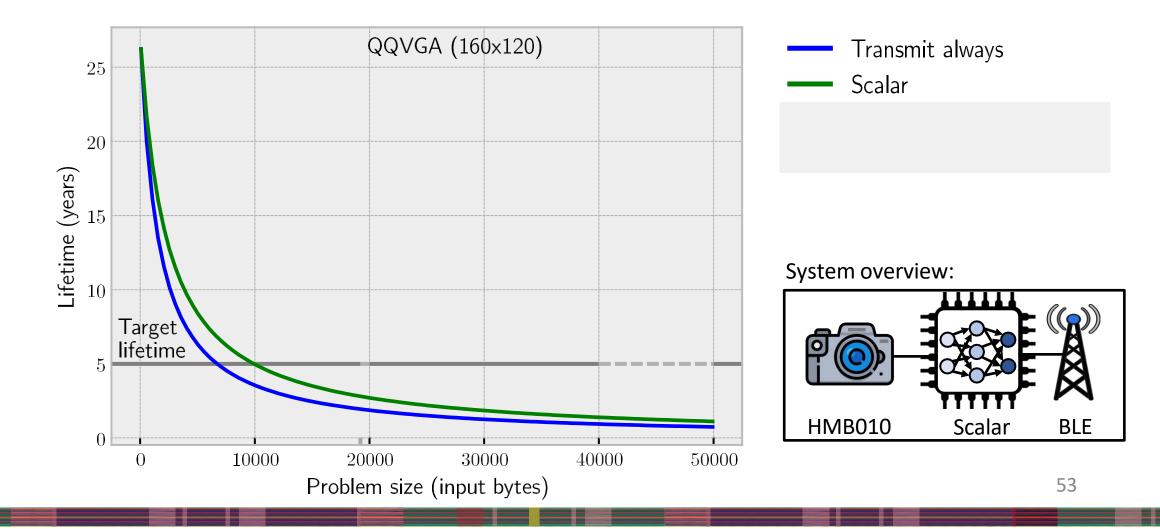
#### **Energy savings**

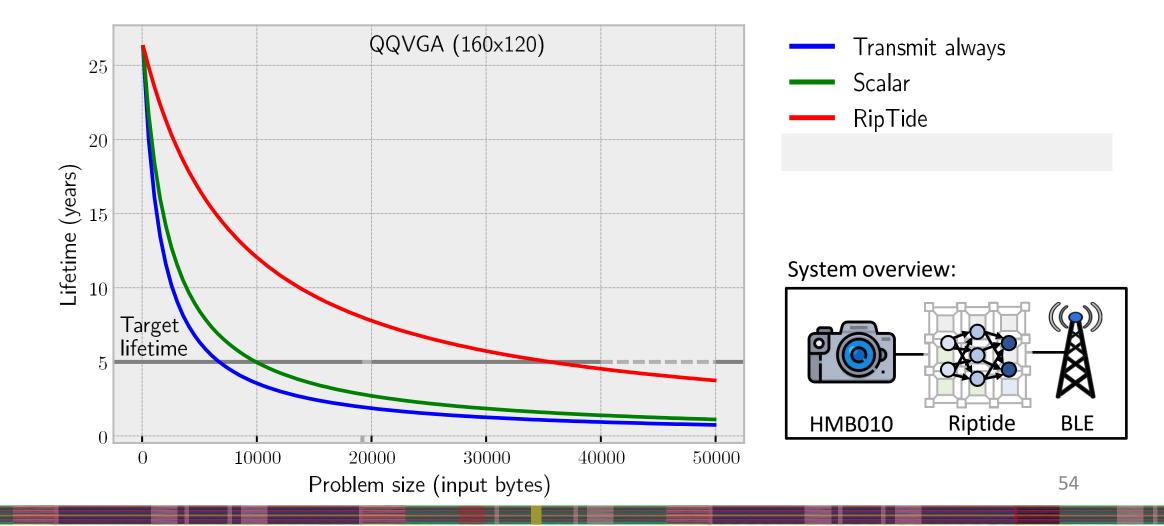


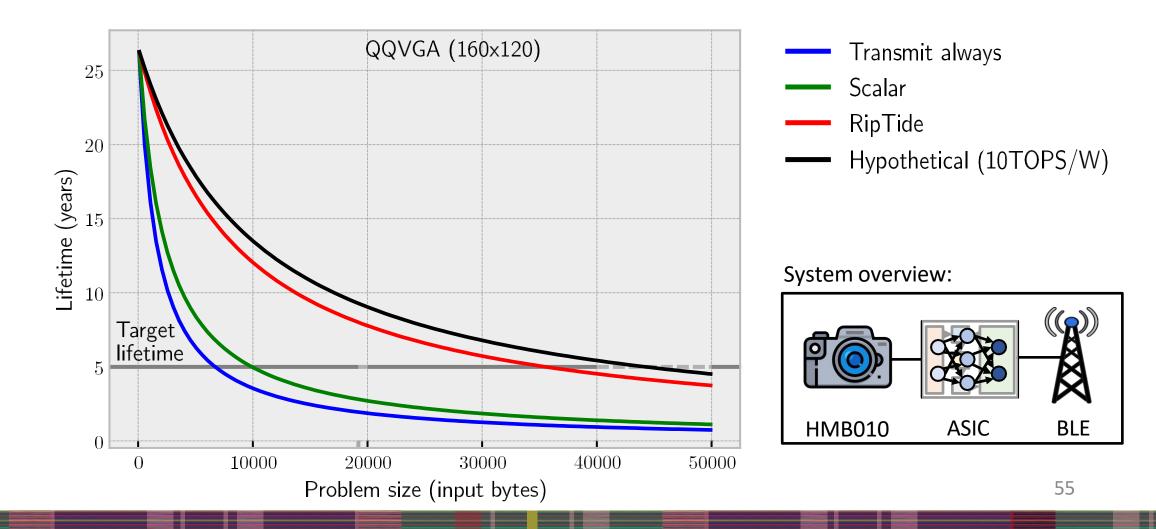
#### **Speedup**



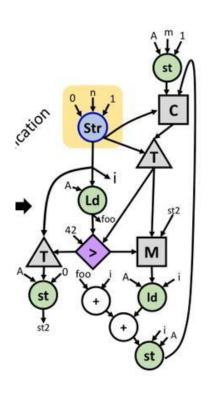








### That was "Ordered Dataflow"

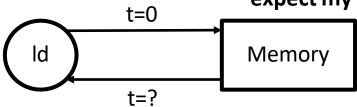


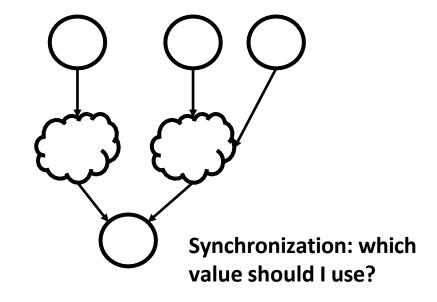
Axiom: Tokens proceed through the graph in the order of their generation

How do we ensure that tokens flow through the dataflow graph in order?

# What about allowing token reordering?

(Memory) Latency: when can I expect my value to come back?





"Tagged-token dataflow architectures"

Two issues: Latency & Synchronization

Latency: time between when operation is issued and when completes

Synchronization: need to assure data properly written before read

#### Monsoon: an Explicit Token-Store Architecture

Gregory M. Papadopoulos
Laboratory for Computer Science
Masachusetts Institute of Technology

David E. Culler Computer Science Division University of California, Berkeley

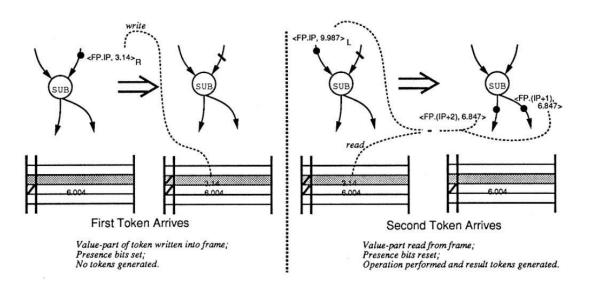


Figure 1: ETS Representation of an Executing Dataflow Program

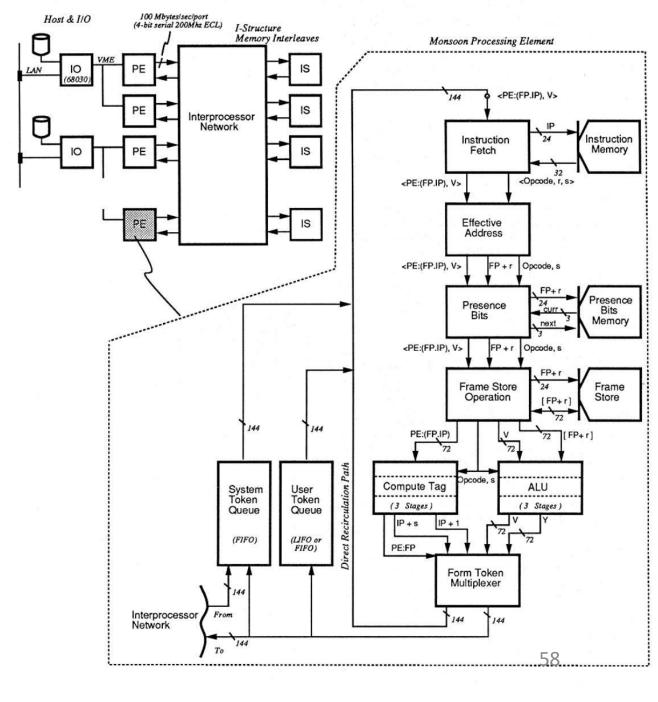


Figure 2: Monsoon Processing Element Pipeline

### Tagged-token Dataflow Architecture

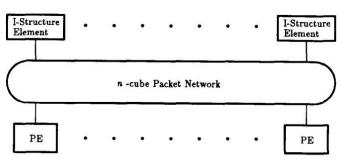


Fig. 17. Top-level view of TTDA.

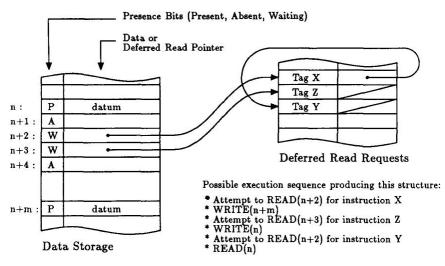
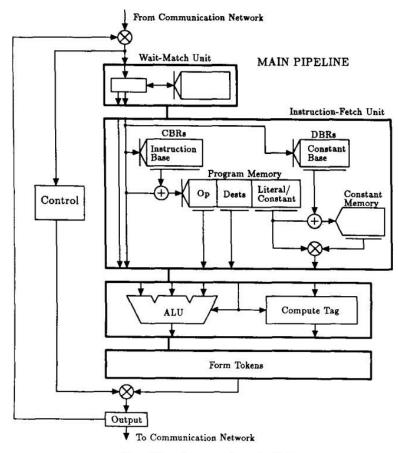


Fig. 7. I-structure memory.

I-structures: latency-tolerant memory
I-fetch send rd tok w/ addr+continuation
P: read & run; A/W: queue
I-store: send wr tok to populate table

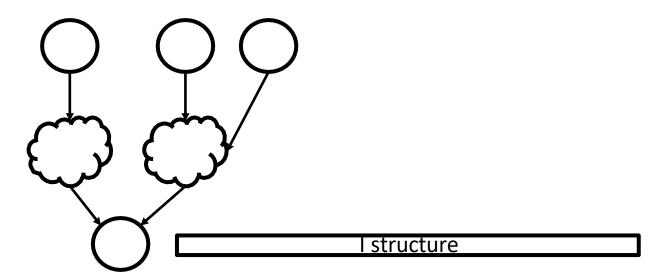
A/W are non-blocking (why?)

I-allocate: make storage for fetch/stores



Token matching: synchronization of out-of-order inputs, using i-structs

# Parallelism is a resource congestion problem

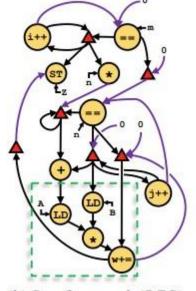


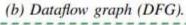
Synchronization: which value should I use? Many options accumulating over time

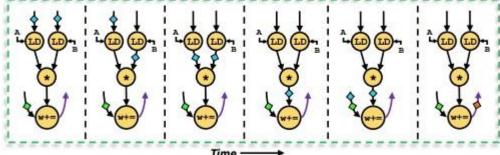
### Varieties of Dataflow Execution

Figure 3: A running example: dense matrix-vector multiplication. We zoom in on the innermost loop body (green box) to demonstrate dataflow execution.

(a) Pseudocode.







(c) Execution trace of innermost loop.

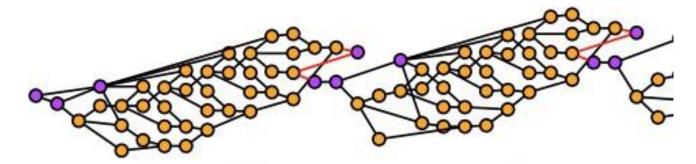
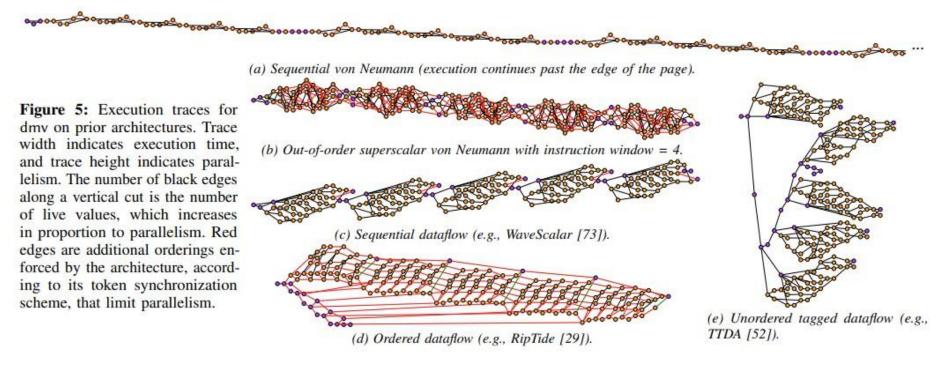


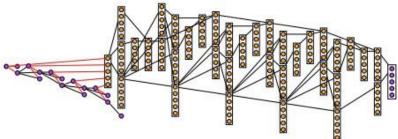
Figure 4: Partial execution trace of dmv.

### Varieties of Dataflow Execution



sense, vN's sequential ordering restricts execution to a "depth-first" traversal of a program's full dynamic execution graph (Fig. 1). The result is minimal parallelism, so that program execution takes a long time (the graph is wide) but live state is minimized (the graph is short).

Parallelism can be increased by having multiple vN execution streams, i.e., multithreading. Token synchronization now includes a token's thread id, in addition to its vN ordering. Multithreading punts all of the problems of parallelism.



(f) Data-parallel (e.g., vector, GPU).