



Accelerator Architectures for Machine Learning (AAML)

Lecture 6: Digital DNN Accelerator

Tsung Tai Yeh

Department of Computer Science
National Yang-Ming Chiao Tung University



Acknowledgements and Disclaimer

- Slides was developed in the reference with
Joel Emer, Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, ISCA 2019
tutorial
Efficient Processing of Deep Neural Network, Vivienne Sze, Yu-Hsin
Chen, Tien-Ju Yang, Joel Emer, Morgan and Claypool Publisher, 2020
Yakun Sophia Shao, EE290-2: Hardware for Machine Learning, UC
Berkeley, 2020
CS231n Convolutional Neural Networks for Visual Recognition,
Stanford University, 2020
CS224W: Machine Learning with Graphs, Stanford University, 2021



Outline

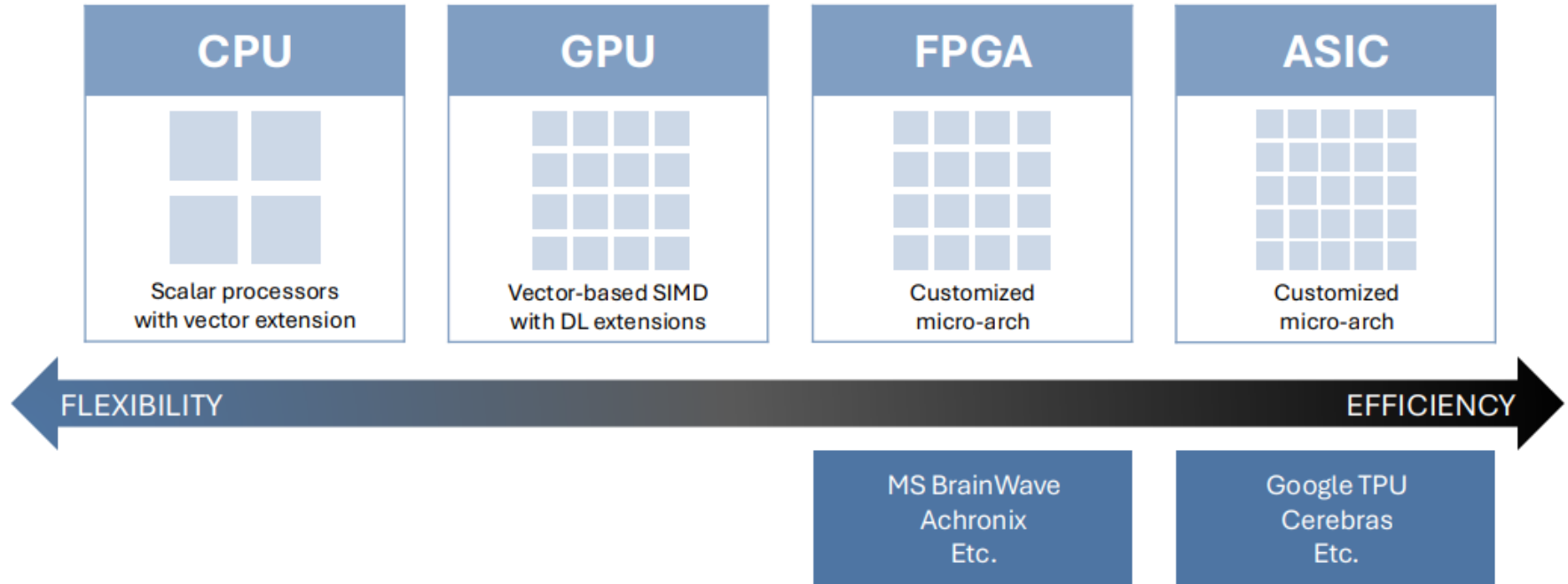
- Reconfigurable Deep Learning Accelerators
 - FPGA
 - SambaNova Reconfigurable Dataflow Unit (RDU)
 - Coarse grained reconfigurable array (CGRA)
 - GraphCore IPU
 - Wafer-scale AI chip -- Cerebras



Reconfigurable Deep Learning on FPGA



Spectrum of Architectures for Deep Learning

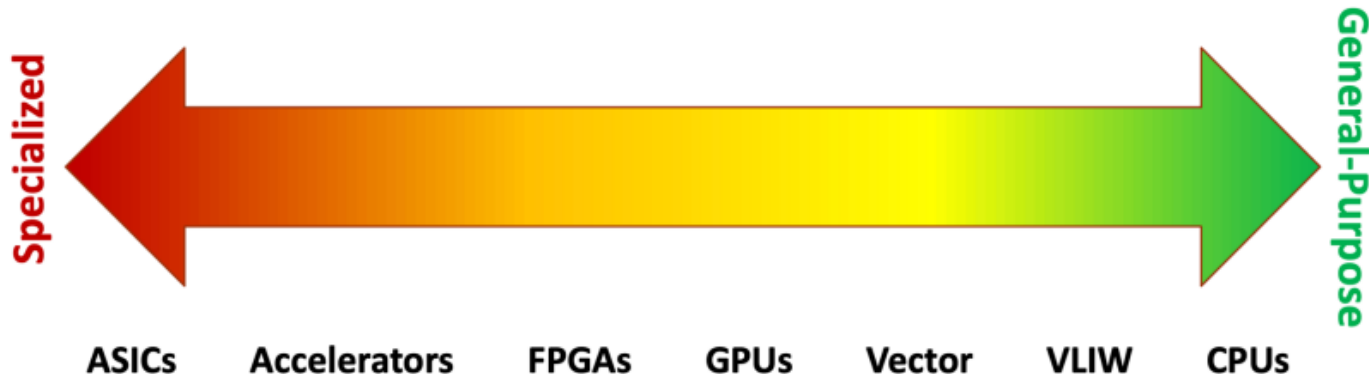




Why Reconfigurable Computing?

- AI accelerators improves 100X performance/energy compared to general-purpose processor
- But new hardware is sophisticated and expensive
 - Especially in cutting-edge manufacturing processes

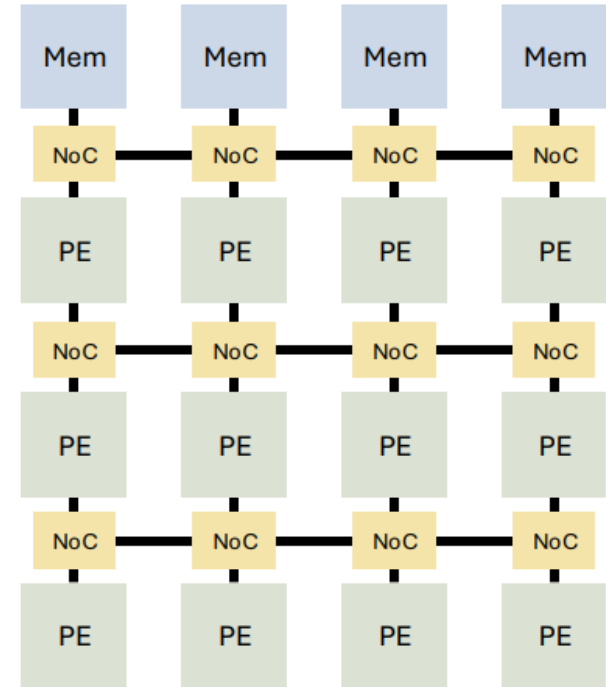
Can we bridge the gap btwn general-purpose & accelerators w/out custom hardware?





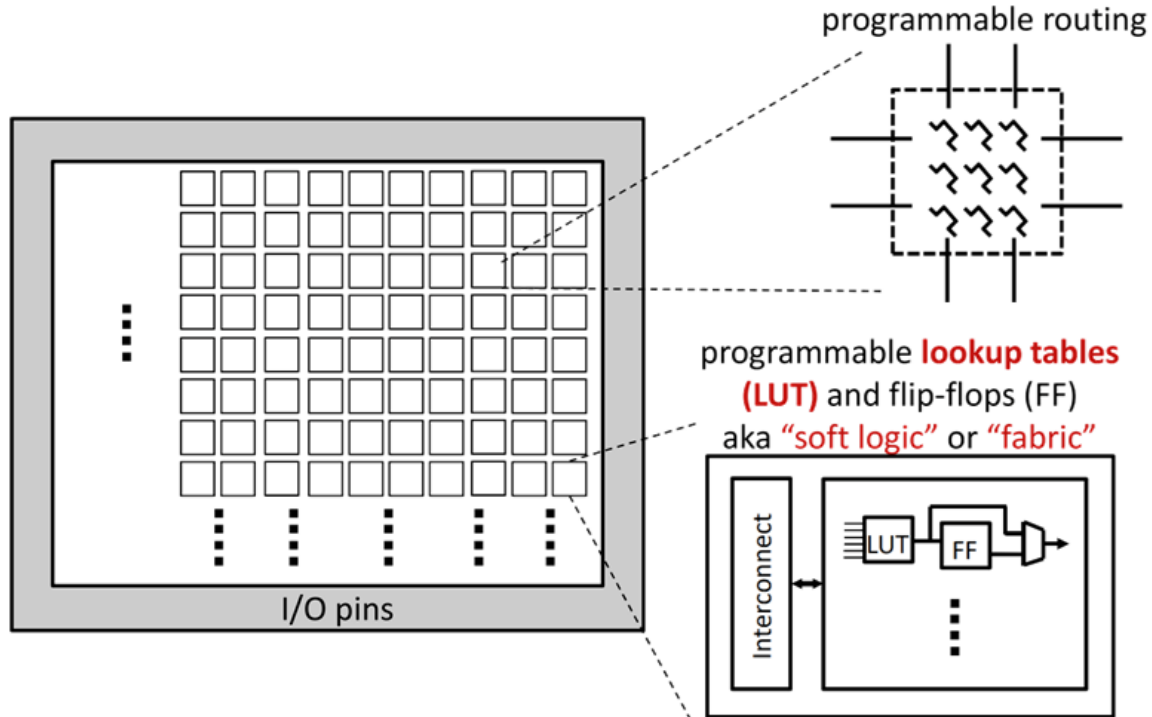
Reconfigurable Computing

- Basic idea
 - A spatial array of processing elements (PEs) & memories with a configurable network
 - Map your computation spatially onto the array
 - Goal: programmable with near ASIC efficiency





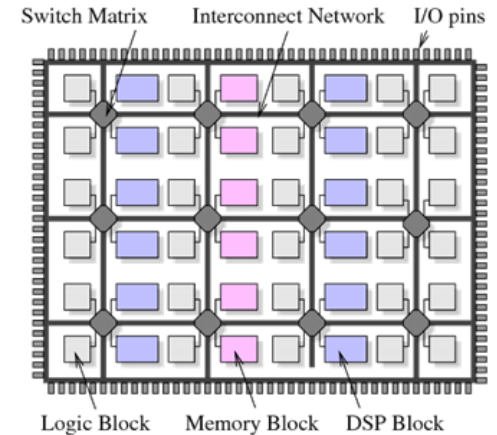
Basic FPGA Design



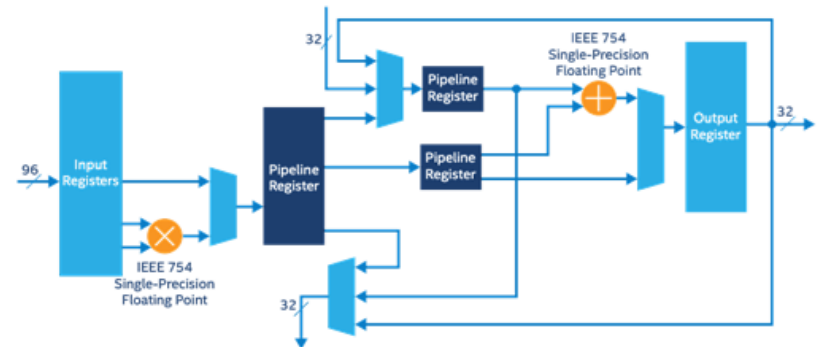


Modern FPGAs

- FPGAs are coarse-grain today
 - Hardened logic in LUTs
 - “DSP blocks” to implement wide add/mul efficiently
 - Dense memories distributed throughout fabric



Mode Name	Mathematical Function
Multiplication Mode	$X \cdot Y$
Adder or Subtract Mode	$(X + Y)$ or $(X - Y)$
Multiply-Add/Subtract	$(X \cdot Y) + Z$ or $(X \cdot Y) - Z$
Multiply Accumulate Mode	$(X \cdot Y) + \text{Acc}$ or $(X \cdot Y) - \text{Acc}$
Vector One Mode	$(X \cdot Y) + \text{Chain In}$



[Intel, Stratix 10]

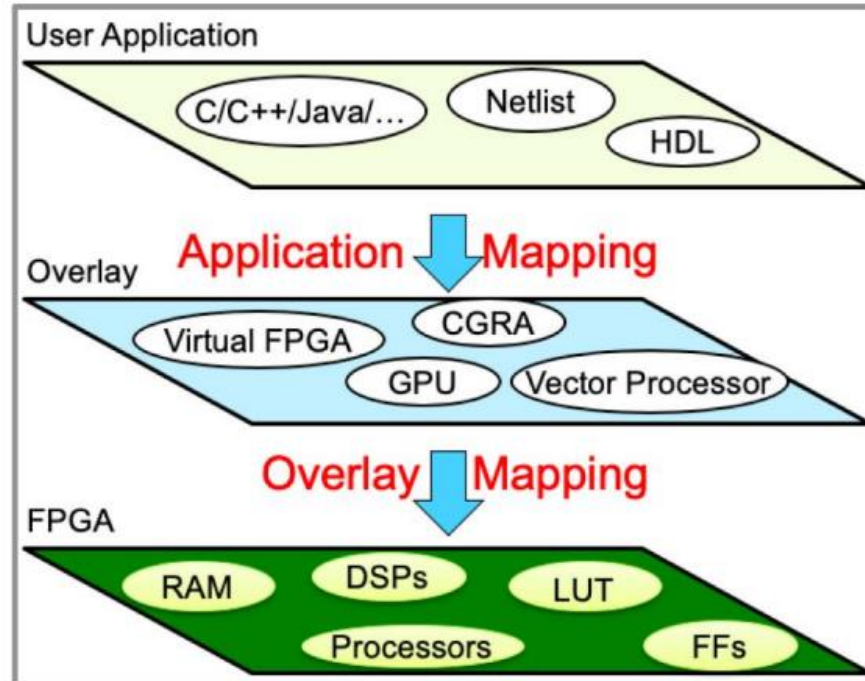
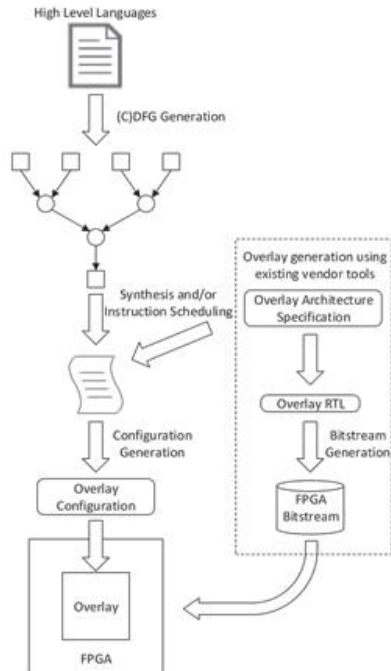


FPGA vs GPU

- Prefill stage
 - The latency of FPGAs in the prefill stage increases linearly, while the GPU ones almost remain constant as the model does not fully utilize GPUs
- Decode stage
 - FPGA-based accelerators are more efficient than GPUs



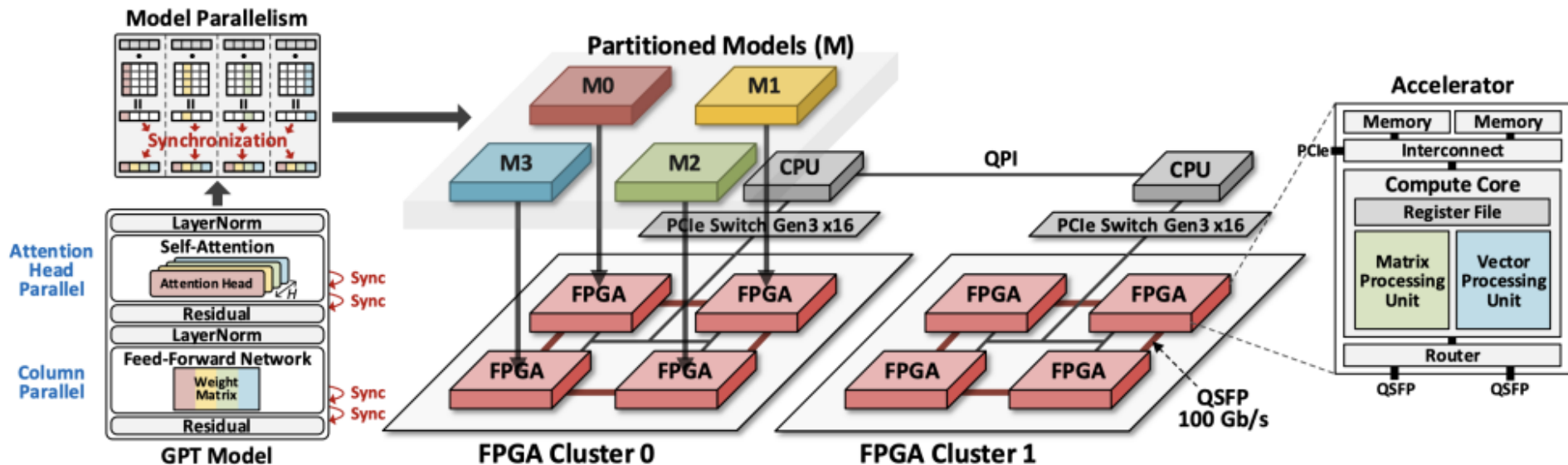
FPGA Overlay



[So, FPGAs for Software Programmers, 2016]



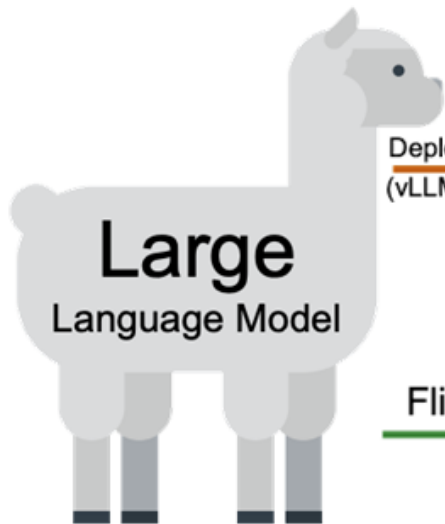
FPGA Overlay



[Hong, MICRO, 2022]



FlightLLM



Deploy with opt.
(vLLM, INT8, ...)

FlightLLM

GPU

V100
(12 nm)



Cost

Price Power
(~\$12K) (~220W)



Performance

~45 
token/s

1x throughput/price
1x throughput/power

FPGA

U280
(16 nm)



Cost

Price Power
(~\$8K) (~45W)



Performance

~55 
token/s

1.8x throughput/price
6.0x throughput/power



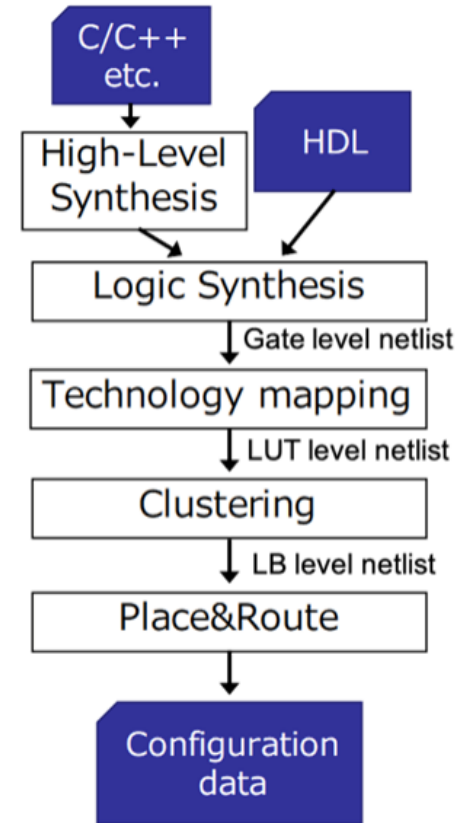
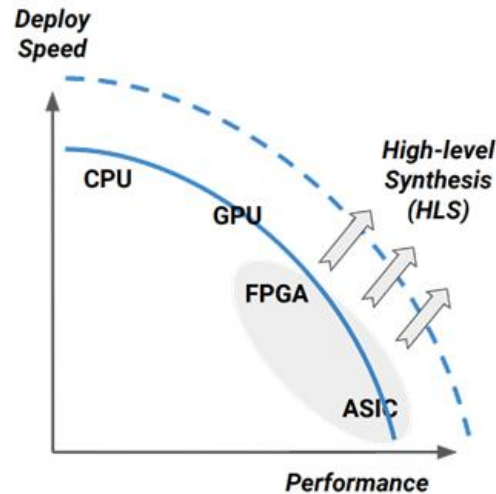
Challenges of FPGAs on LLM Inference

- Low computation efficiency
 - Hard to efficiently map sparse matrices onto DSP chains
- Underutilized memory bandwidth
 - Repeated off-chip memory accesses for each fine-grained kernel
- Huge instruction storage
 - Store instructions for all possible token lengths



High-Level Synthesis (HLS)

- Automated optimization and scheduling
- High portability against different PDK or PPA requirements
- Short design cycle





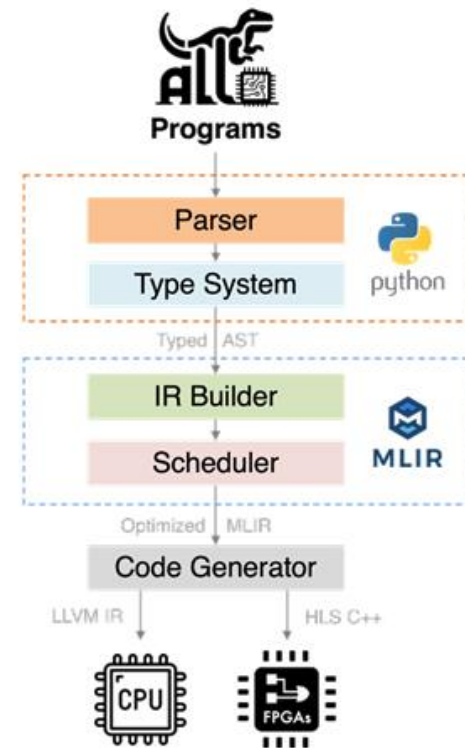
Challenges of HLS Accelerator Design

- **Time consuming**
 - Manual architecture and micro-architecture design, manual C/C++ code rewriting
- **Suboptimal**
 - Empirical parameter tuning, like parallel factors, buffer sizes, tiling sizes, etc..
- **Low flexibility**
 - Only support a small set of models



Accelerator Design Languages (ADLs)

- **Pythonic**
- **Maintainability**
 - Decoupled hardware customizations
- **Composability**
 - All the kernels, primitives, and schedules should be composable to form complex designs





FPGA on AI Accelerators Follows

- FPGA vendors doing what markets want
 - Future “FPGA” not sea-of-gates for RTL netlist
- Purposeful architectures for targeted applications
 - Make things easier/cheaper to do
 - Be very good at what it is intended to do
- Coping with architectural divergence
 - Soft-logic adds malleability to “architecture”
 - 2.5/3D integration allows specialization off a common denominator
 - Push reconvergence of abstraction up the stack



SambaNova Reconfigurable Dataflow Unit (RDU)



Plasticine Architecture

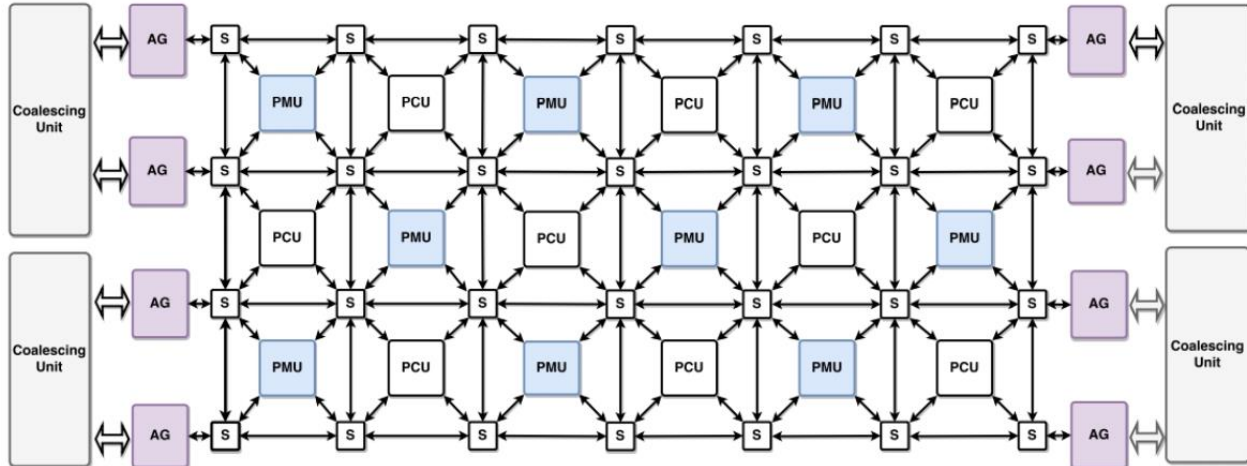
- **Plasticine architecture**

- A reconfigurable architecture for parallel patterns (Raghu, ISCA 2017)
- **Pattern Compute Unit (PCU)**
 - Reconfigurable pipeline with multiple stages of SIMD functional units (FUs)
- **Pattern Memory Unit (PMU)**
 - A banked scratchpad memory
- **The compiler**
 - Maps the computation of inner loops to PCUs
 - Most operands are transferred directly between FUs without scratchpad access or inter-PCU communication



Plasticine Architecture Overview

- Data access address calculation occurs while the PCU is working
- Each DRAM channel is accessed using several address generators (AG) on two sides of the chip
- Multiple AGs connect to an address coalescing unit for memory requests



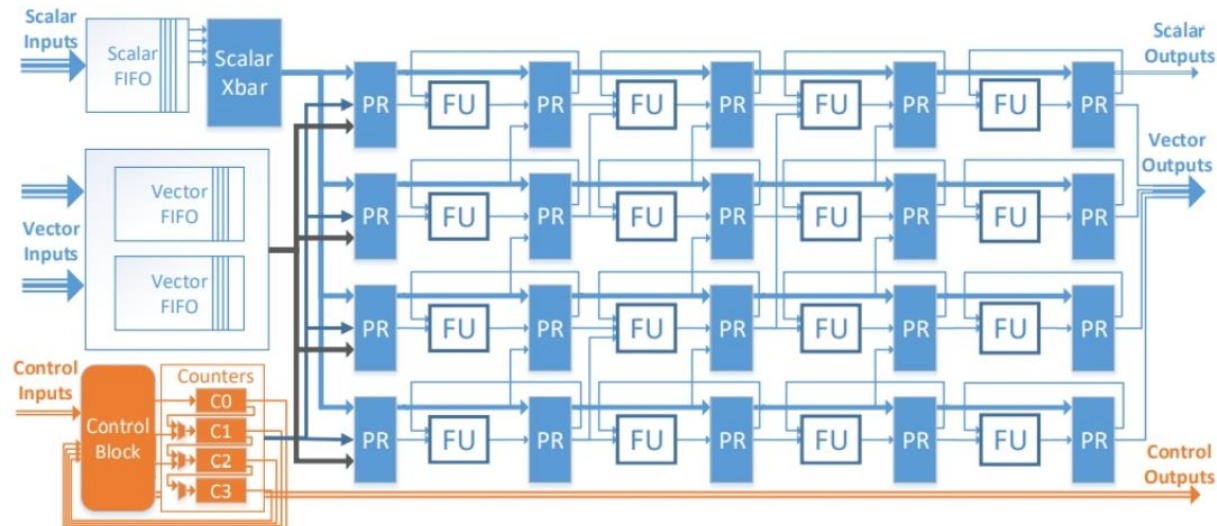


Plasticine PCU Architecture

- **Pattern Compute Unit (PCU)**

- Each stage's SIMD lane contains a FU and associated pipeline register (PR)

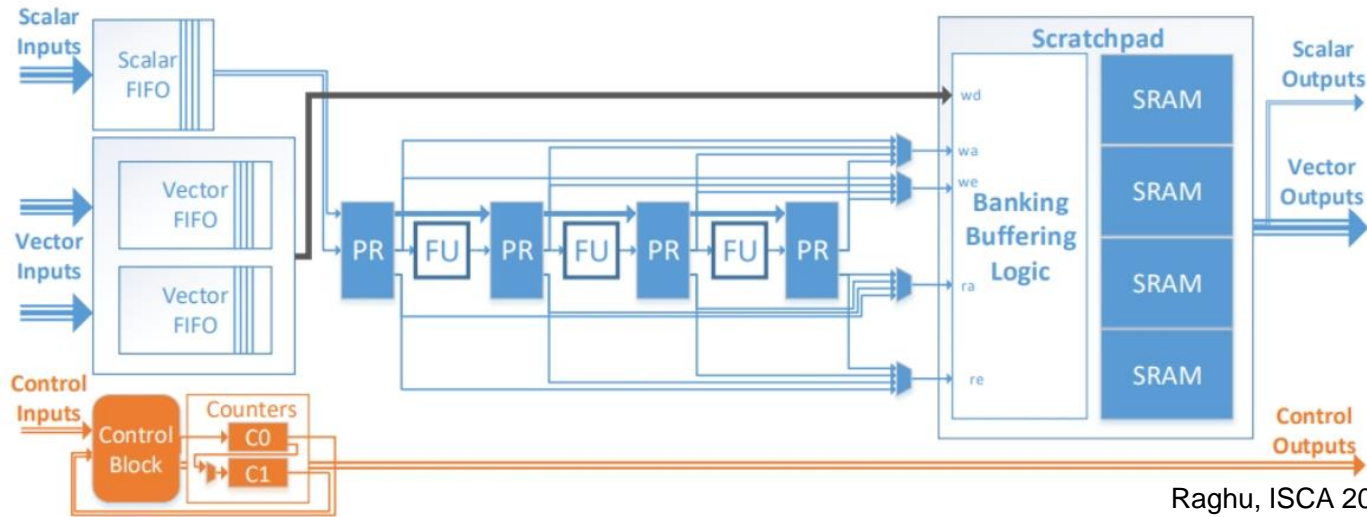
1. **Scalar**: uses to communicate single words of data
2. Each **vector** communicates one word per line in the PCU
3. **Control** signals at the start or end of execution of a PCU





Plasticine PMU Architecture

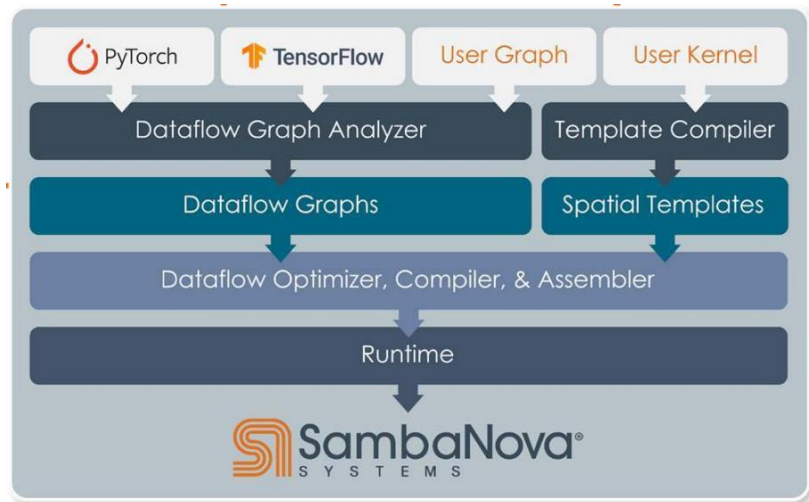
- **Pattern Memory Unit (PMU)**
 - Contains a scratchpad memory and address calculation
 - Calculates address only needs simple scalar math
 - Has simpler FUs than ones in PCUs





Reconfigurable Dataflow Unit (RDU)

- **SambaNova RDU**
 - **Pattern Compute Units**
 - BF16 with FP32 accumulation
 - Support FP32, Int32, Int16, Int8
 - **Pattern Memory Unit**
 - Memory transformation
 - **Dataflow optimization**
 - Tiling
 - Nested pipelining
 - Operator parallel streaming





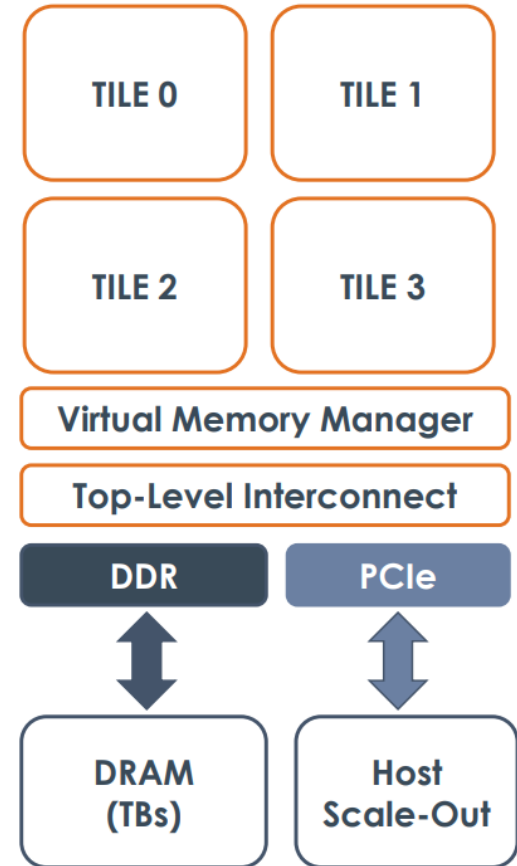
Dataflow Exploits Data Locality / Parallelism

- **Software-hardware co-design architecture**
 - Dataflow captures data locality and parallelism
 - Flexible time and space scheduling to achieve higher utilization
 - Flexible memory system and interconnect to sustain high compute throughput
 - Custom dataflow pipeline



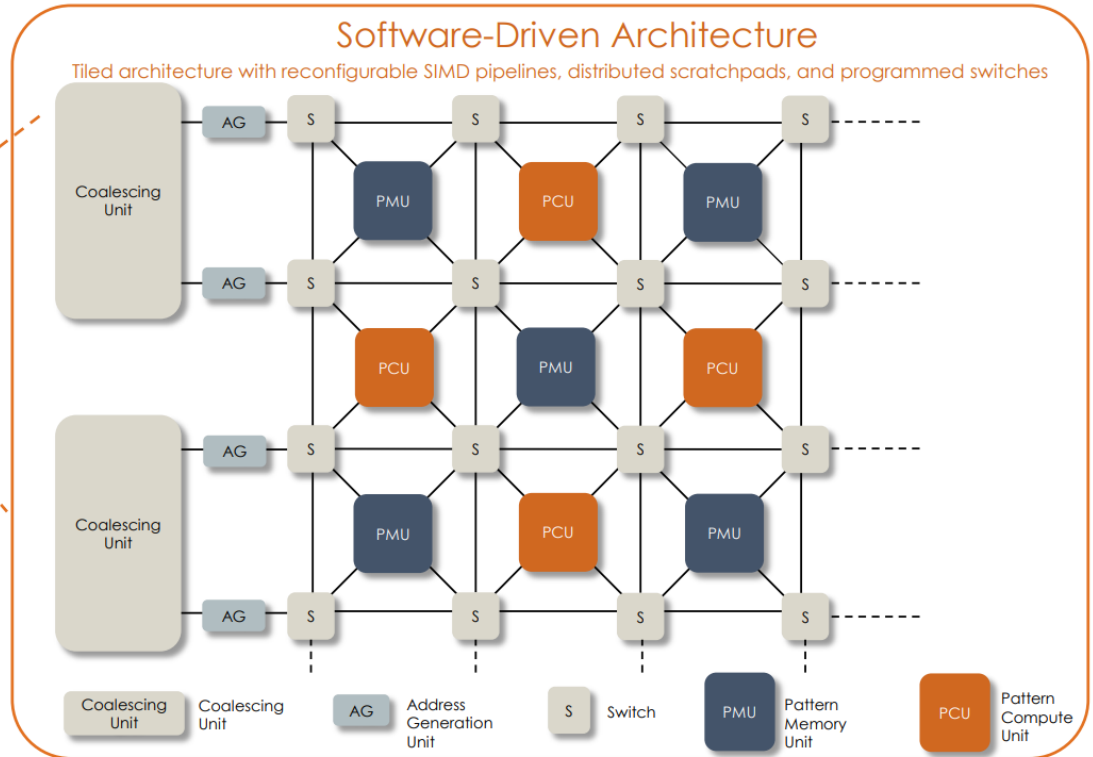
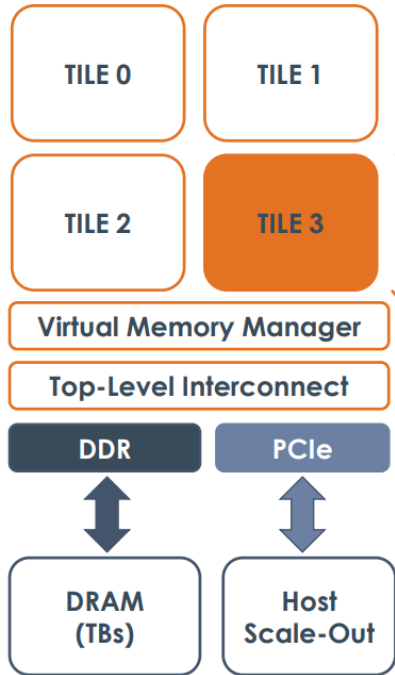
Chip and Architecture Overview

- **RDU Tile**
 - Compute and memory components
 - A programmable interconnect
- **Tile resource management**
 - Combine adjacent tiles to form a larger logical tile
 - Each tile controlled independently
 - Allow different applications on separate tiles concurrently
- **Memory access**
 - Direct access to TBs DDR4 off-chip memory
 - Memory-mapped access to host memory





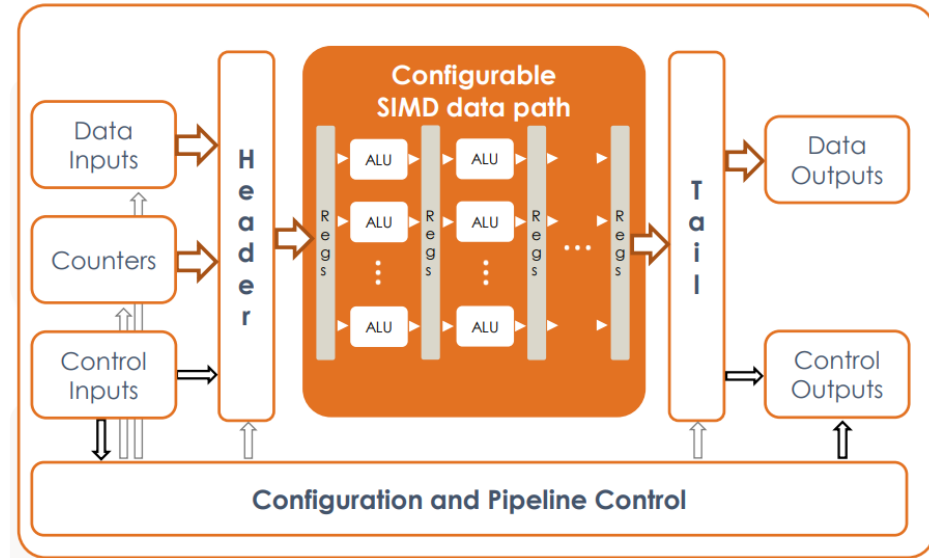
RDU Tile





Pattern Compute Unit (PCU)

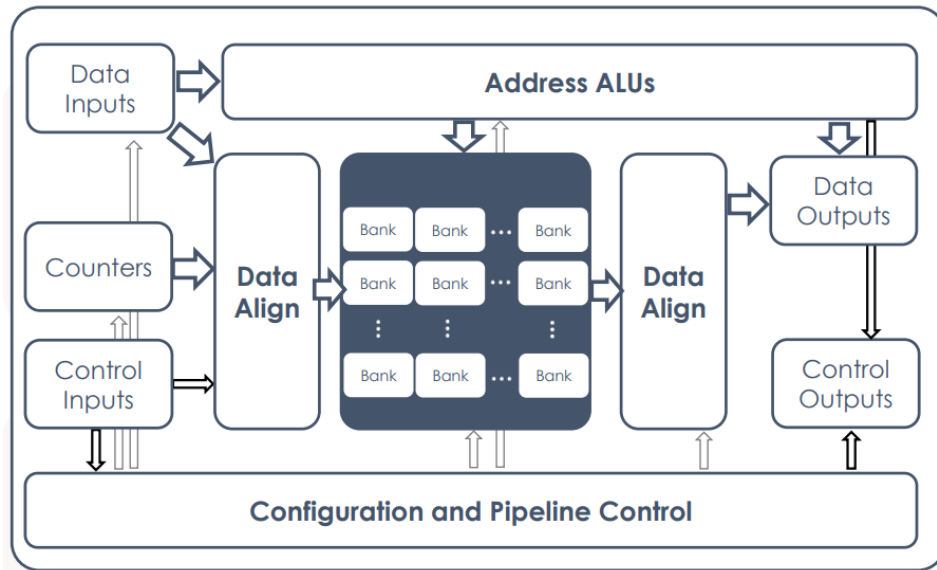
- **Pattern Compute Unit (PCU)**
 - Compute engine
- **Reconfigurable SIMD data path**
 - For dense and sparse tensor algebra in FP32, BF16, and integer data format
- **Programmable counters**
 - Program loop iterators
- **Tail unit**
 - Accelerates functions such as exp, sigmoid





Pattern Memory Unit (PMU)

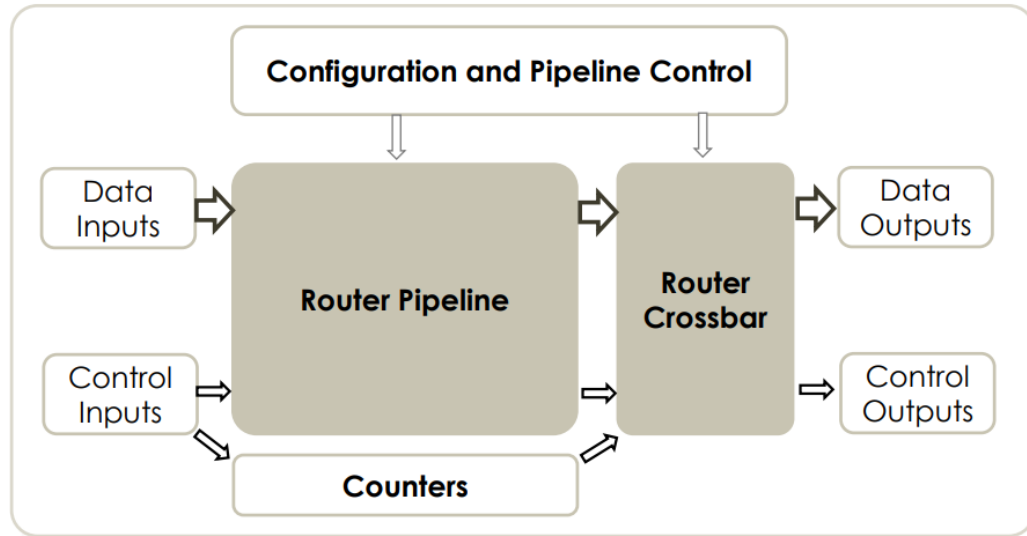
- **Pattern Memory Unit (PMU)**
 - **On-chip memory system**
 - **Banked SRAM arrays**
 - Write and read multiple high bandwidth SIMD data stream concurrently
 - **Address ALUs**
 - Address calculation for arbitrarily complex accesses
 - **Data align**
 - Tensor layout transformation





Switch and On-chip Interconnect

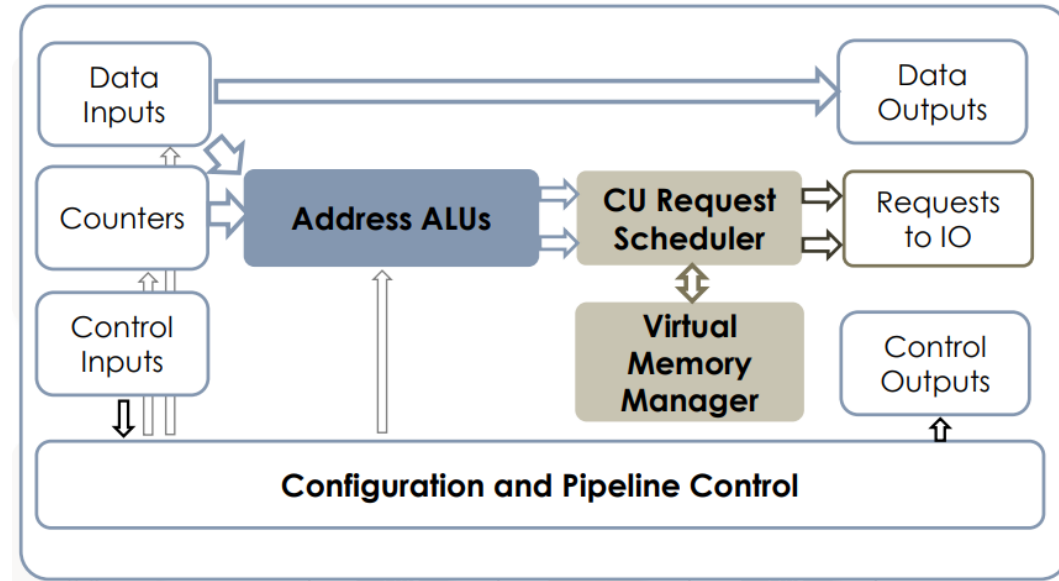
- **Switch**
 - Programmable packet-switched interconnect fabric
- **Independent data and control buses**
 - Suit different traffic classes
- **Programmable routing**
 - Flexible chip bandwidth allocation to concurrent stream
- **Programmable counters**
 - Outer loop iterators
 - On-chip metric collection





Interface to I/O Subsystem

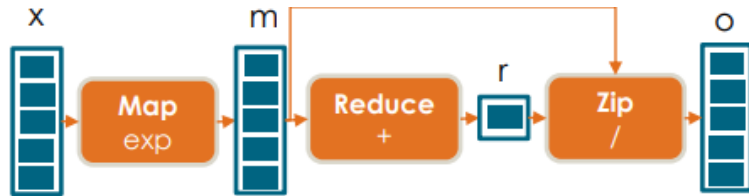
- **Address ALUs**
 - Address calculation for arbitrarily complex accesses
- **Coalescing Units**
 - Enable transparent access to memories across RDUs and host memory
- **Address space manager**
 - Programmable, variable length segments



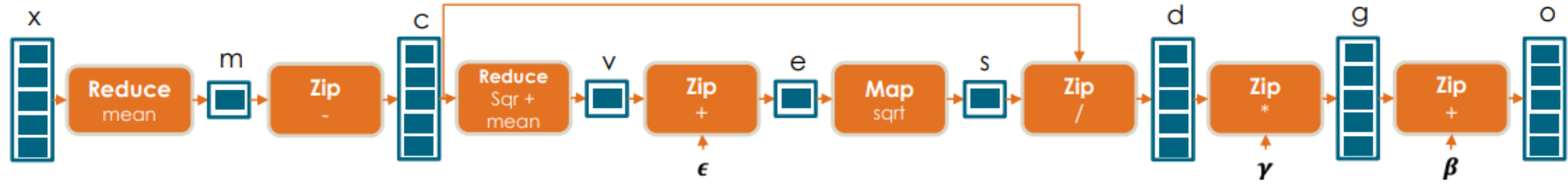


Operator Mapping

SOFTMAX:
$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$



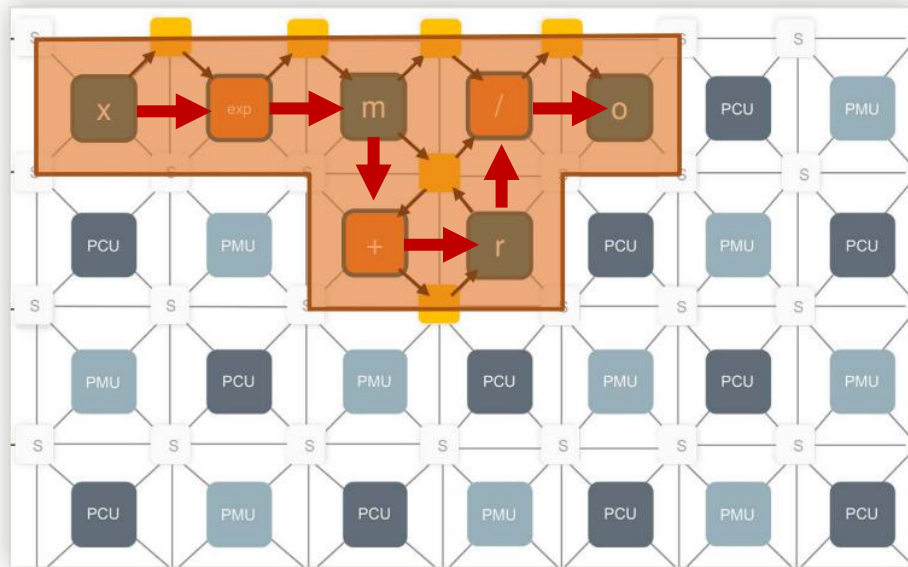
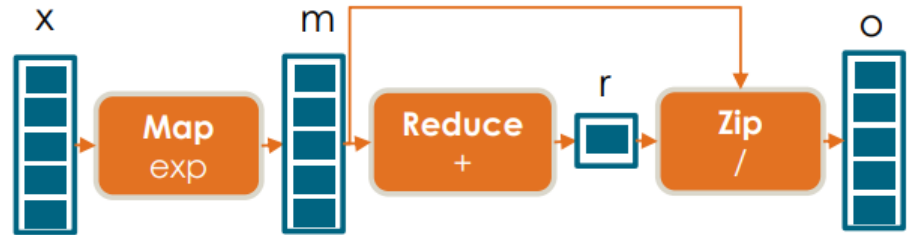
LAYERNORM:
$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$





Operator Mapping (Softmax)

SOFTMAX:
$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

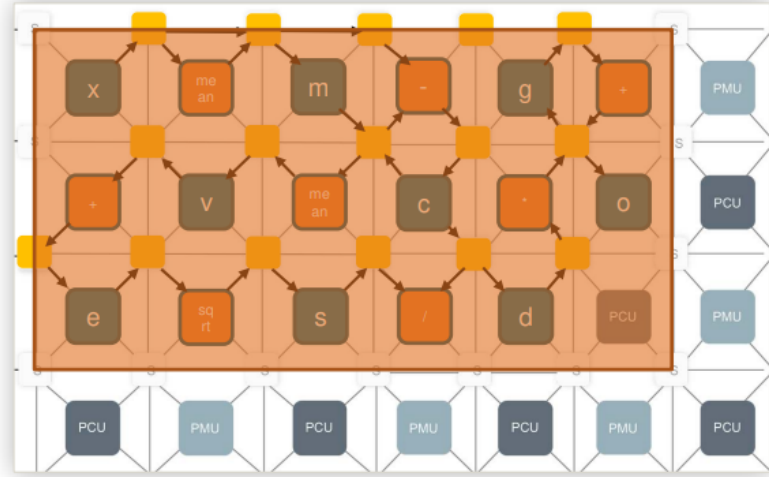
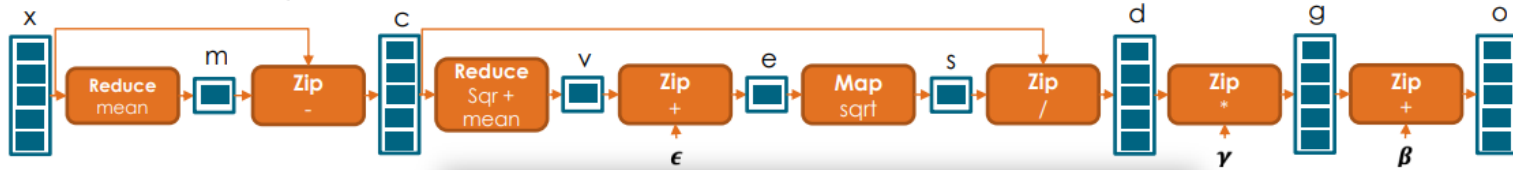




Pipelined in Space

LAYERNORM:

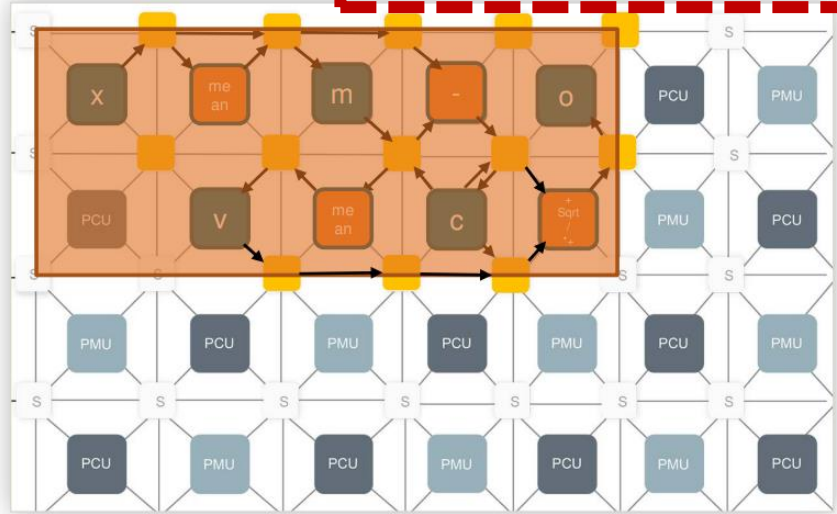
$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$





Pipelined in Space + Fused

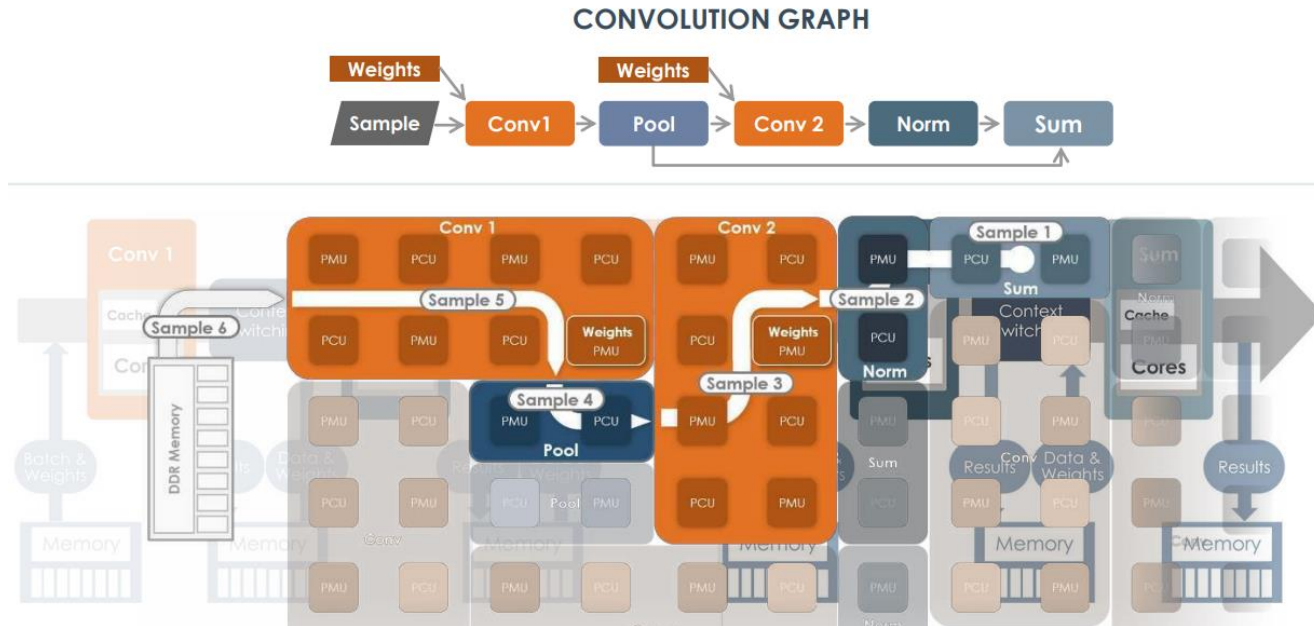
LAYERNORM:
$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$





Spatial Dataflow within an RDU

- **The dataflow removes**
 - Memory traffic and host communication overhead



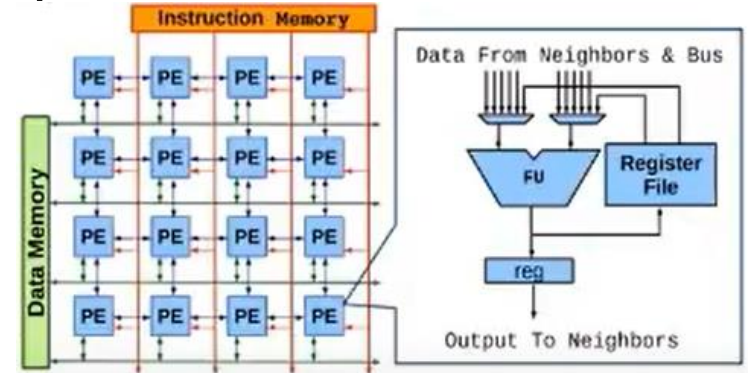


CGRA



Coarse grained reconfigurable array (CGRA)

- **Coarse grained reconfigurable array (CGRA)**
 - Multiple **processing elements (PEs)**
 - Each PE has ALU-like functional unit
 - **Array configurations vary by**
 - Array size
 - Functional units
 - Interconnection network
 - Register file architectures
 - CGRAs can achieve **power-efficiency** of several 10s of GOps/sec per Watt (why?)
 - Samsung SRP processor (embedded and multimedia apps)





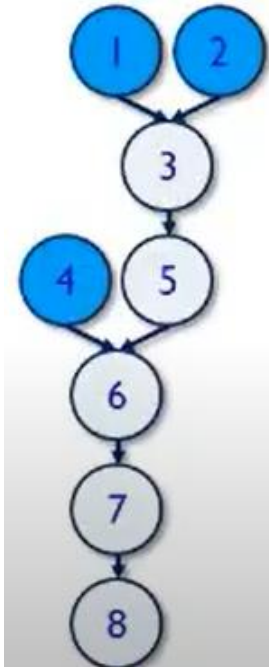
Key features of CGRA accelerators

- **Software-pipelining execution mapping**
 - **Accelerate loops with low parallelism**
 - Loops with loop-carried dependence, loops with high branch divergence
- **Avoid von-Neumann architecture bottleneck**
 - CGRAs are not subjected to dynamic fetch and decoding of instructions
 - CGRA instructions are in a **pre-decoded form** in the instruction memory
 - **PE transfers data directly** among each another
 - Without going through a centralized registers and memory



Loop execution on the CGRA

Data dependency graph



Mapping data
dependency
graph to CGRA

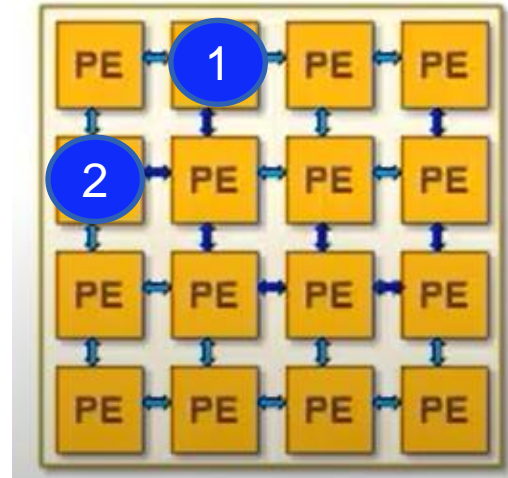


Loop:

$$t1 = (a[i]+b[i]-k)*c[i]$$

$$d[i] = \sim t1 \ \& \ 0xFFFF$$

Execution time: 1

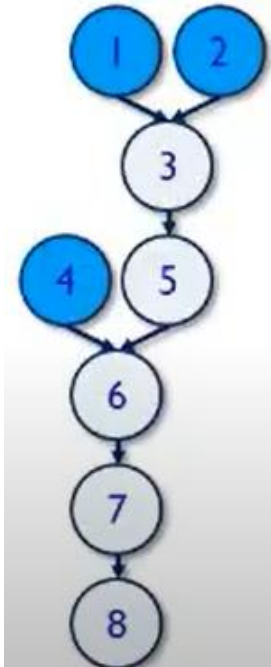




Loop execution on the CGRA

```
Loop:  
t1 = (a[i]+b[i]-k)*c[i]  
d[i] = ~t1 & 0xFFFF
```

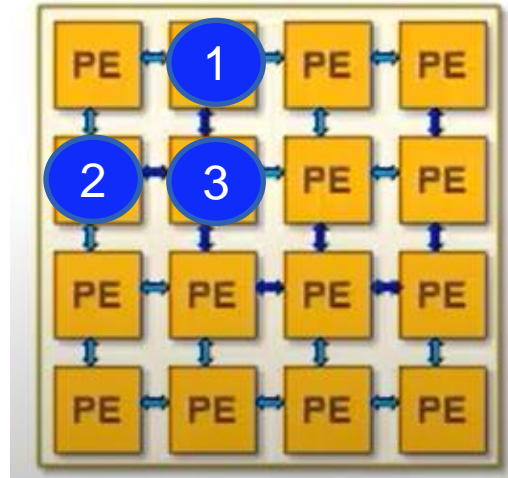
Data dependency graph



Mapping data
dependency
graph to CGRA



Execution time: 2





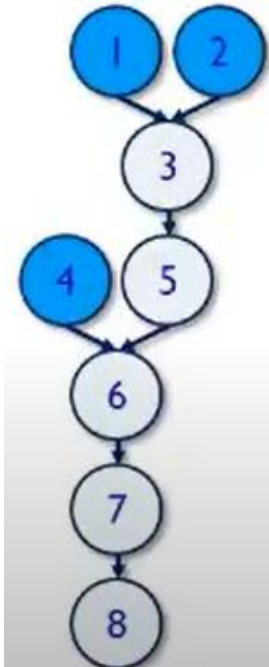
Loop execution on the CGRA

Loop:

$$t1 = (a[i]+b[i]-k)*c[i]$$

$$d[i] = \sim t1 \ \& \ 0xFFFF$$

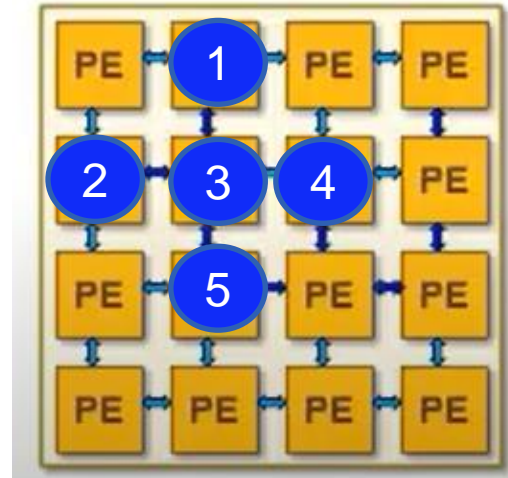
Data dependency graph



Mapping data
dependency
graph to CGRA



Execution time: 3

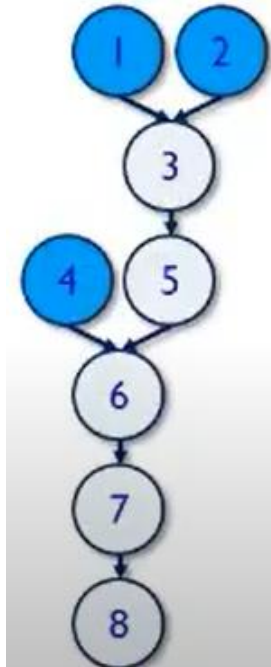




Loop execution on the CGRA

```
Loop:  
t1 = (a[i]+b[i]-k)*c[i]  
d[i] = ~t1 & 0xFFFF
```

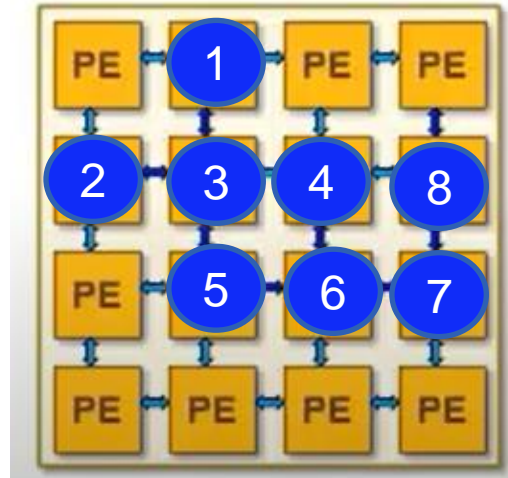
Data dependency graph



Mapping data
dependency
graph to CGRA



Execution time: 6





Takeaway Questions

- What are hardware components used by RDU ?
 - (A) Pattern computer unit (PCU)
 - (B) Pattern memory unit (PMU)
 - (C) Interconnect network router
- What are features of CGRAs ?
 - (A) Customized PEs
 - (B) Software-pipelining execution mapping
 - (C) Reconfigurable dataflow

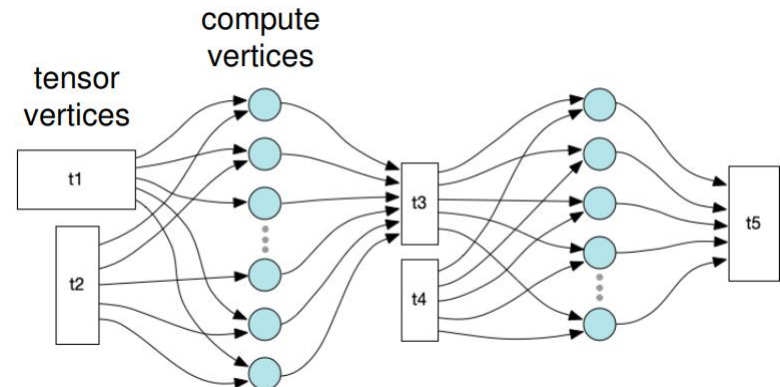


GraphCore IPU



GraphCore IPU

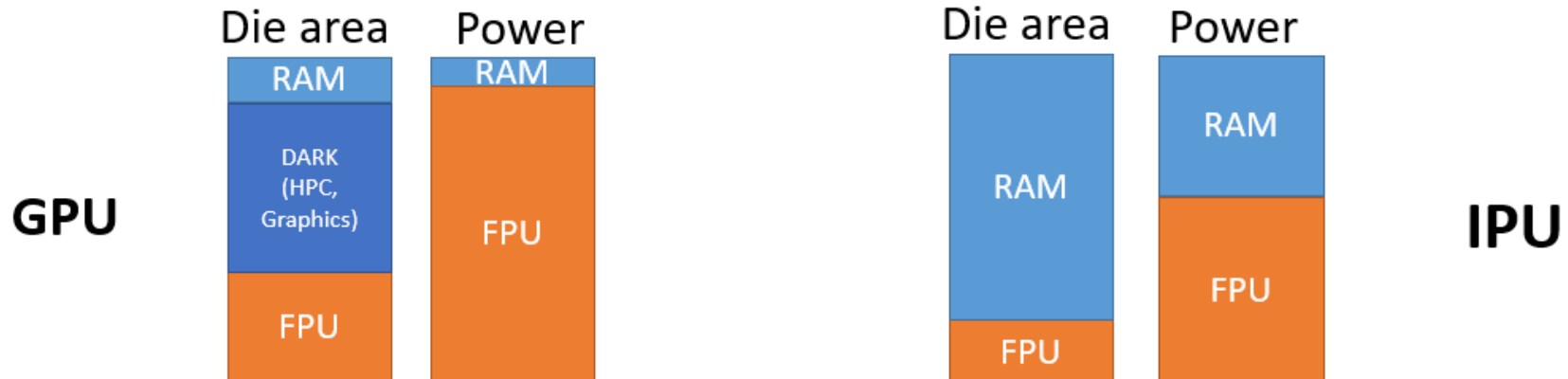
- **GraphCore Intelligent Processing Units (IPUs)**
 - Unlike GPU that is dedicated to accelerate large dense matrix
 - IPU supports **dynamic sparse training** and unstructured computation such as path tracing in 3D computer graphics
 - Multiple **tile processors**
 - **Poplar programming model**
 - Dedicated compiler (PopC)
 - Mapping compute graph to tile processors
 - Compute kernels (Codelets)





Graphcore IPU Approach

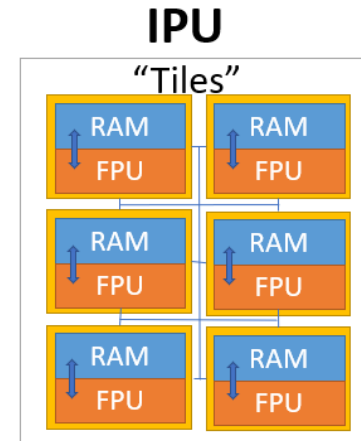
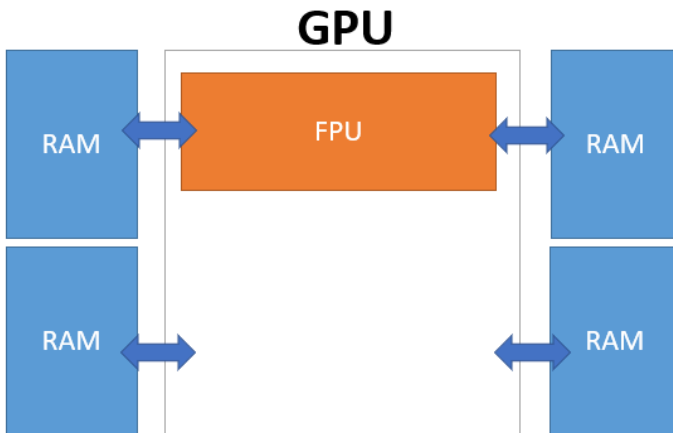
- Post-Dennard, **the silicon is power-limited**
 - We can put more logic on the die than we can power (dark silicon)
- **IPU architecture approach**
 - Replace dark silicon logic with on-chip RAM that has lower power density





Graphcore IPU approach

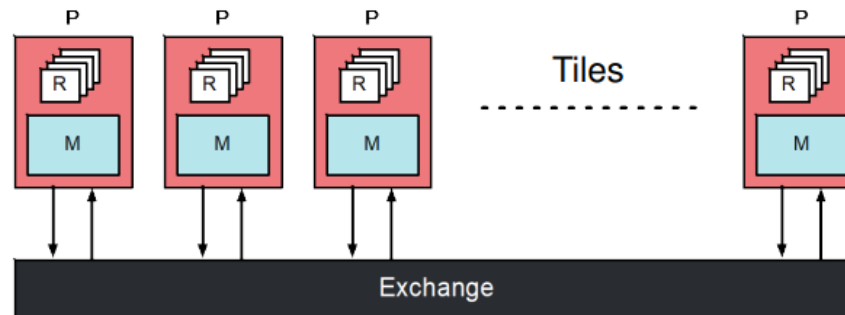
- **GPU approach**
 - Shared memory model with caches and memory hierarchy to reduce latency
- **IPU approach**
 - Move as much memory as possible into the chip local to the logic





Graphcore IPU Abstraction

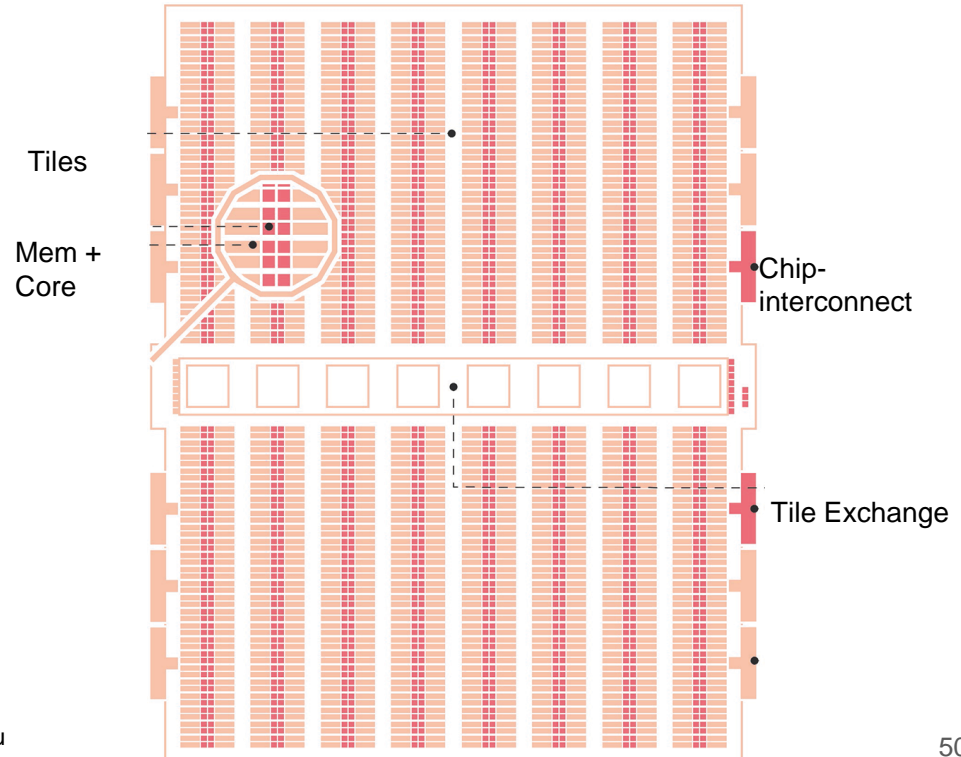
- Tile processors
 - Each tile is a multi-threaded processor and has its local memory
 - Tiles communicate through all-to-all, stateless exchange
- A tensor vertex
 - Can be distributed over many tiles





Distributed memory architecture

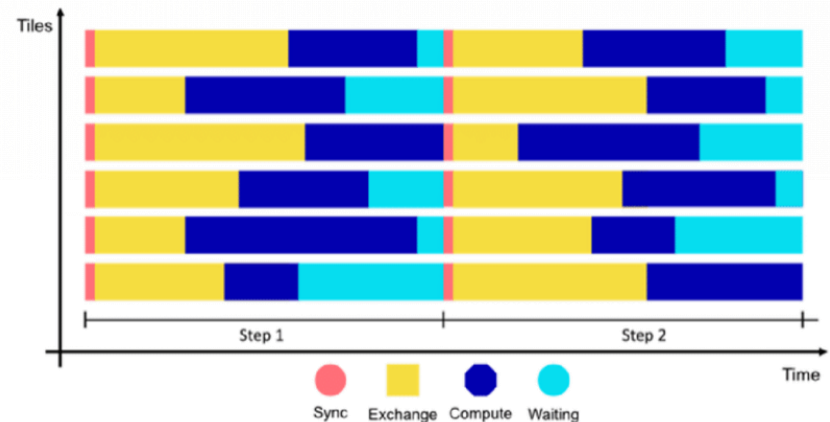
- 1472 tiles with 6 threads sharing 624 KiB of local SRAM
- Total of 896 MiB and 250Tflop/s in 8832 worker threads
- 7.8TB/s exchange between tiles
- Tiles have no shared memory or caches





Execution Model

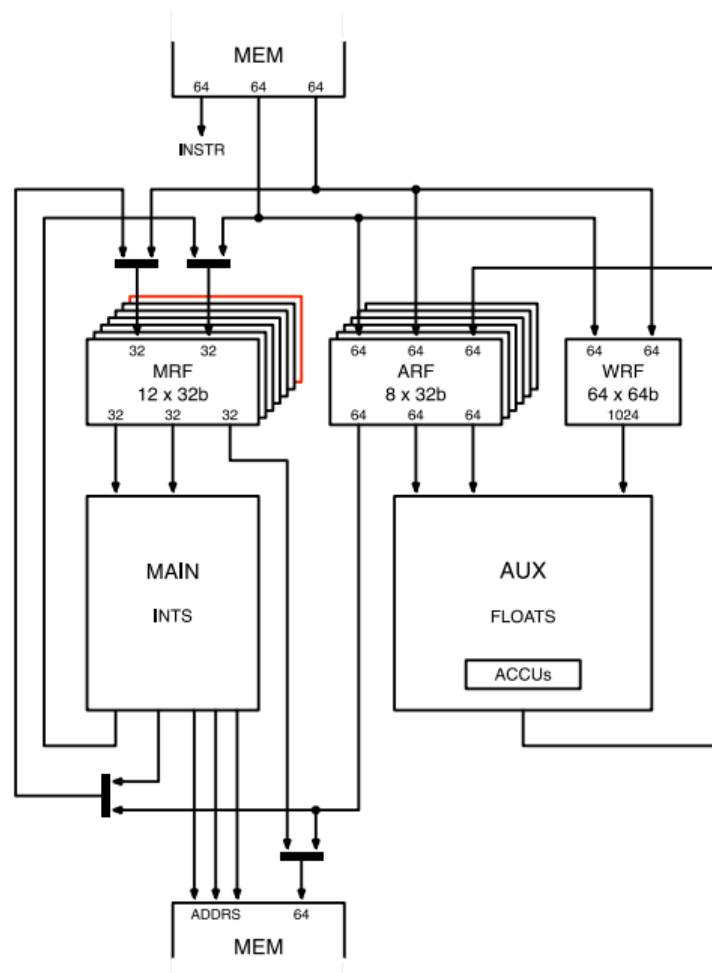
- Tile workers execute instructions independently in parallel (**MIMD**),
- Wait for sync, followed by all-to-all data exchange phase (hardware implementation of **Bulk Synchronous Parallel (BSP)** Model)
- No concurrency hazards (races, deadlocks etc.)
- Compiler faces hard job of scheduling and *load-balancing* compute-chunks on tile workers





Tile Processor

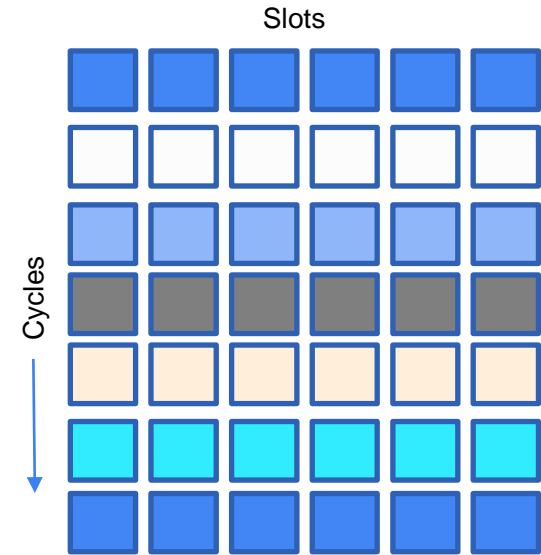
- 32b instructions, single or dual issue
- **Two execution paths:**
 - **MAIN:**
 - Control flow, integer/address arithmetic, load/store to/from either path
 - **AUX:**
 - Floating-point arithmetic for tensor operations + special instructions like log, tanh, PRNG etc.





Tile Processor

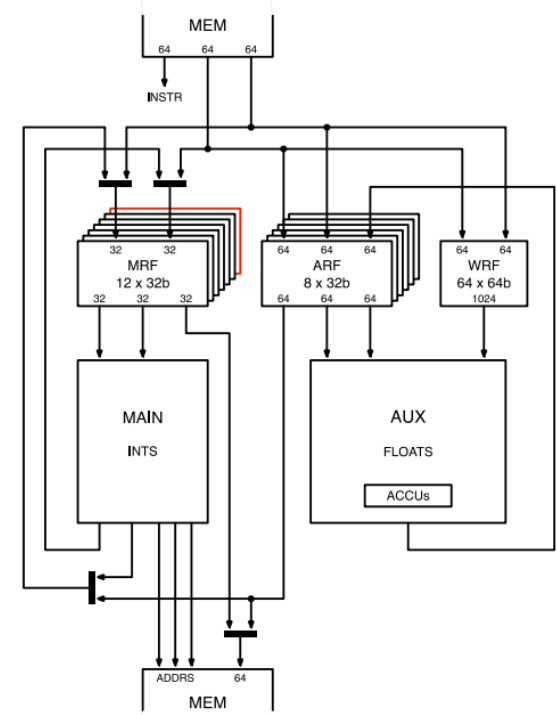
- Fine-grained multithreading that switches between 6 threads on every cycle in round-robin fashion
 - Issued worker programs run in a slot at 1/6 of the clock, so they can't see the pipeline, i.e., mem access, branches etc. all appear to take one cycle per instruction
 - This makes worker execution simple for the compiler to predict for easier load balancing





N + 1 barrel threading

- 7 program contexts
- 6 round-robin pipeline slots
- **The supervisor program**
 - A fragment of the control program
 - Orchestrating the update of vertices
 - Execute in all slots not yielded to workers
 - Dispatch workers by **RUN** instruction
- **A worker program** is a codelet updating a vertex
 - Execute in 1 slot at 1/6 of clock
 - Returns its slot to the supervisor by **EXIT** instruction





Sparse Load/Store

- **Large on-die SRAM memory**
 - 896 MiB on-die SRAM at 47TB/s (data-side)
 - Access arbitrarily-structured data which fits on chip
- **Ld/St instructions**
 - Support sparse gather in parallel with arithmetic at full speed via compact pointer lists
 - 16b absolute offsets to a base
 - 4b cumulative delta offsets to a base



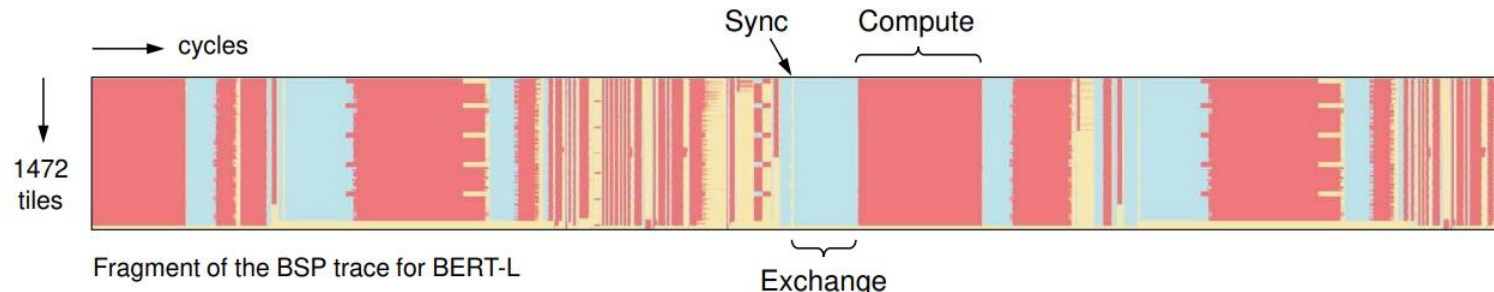
Global Program Order

- **Tile processors**

- Execute asynchronously until they need to exchange data
- Each tile executes a list of atomic codelets in one compute phase

- **Bulk Synchronous Parallel**

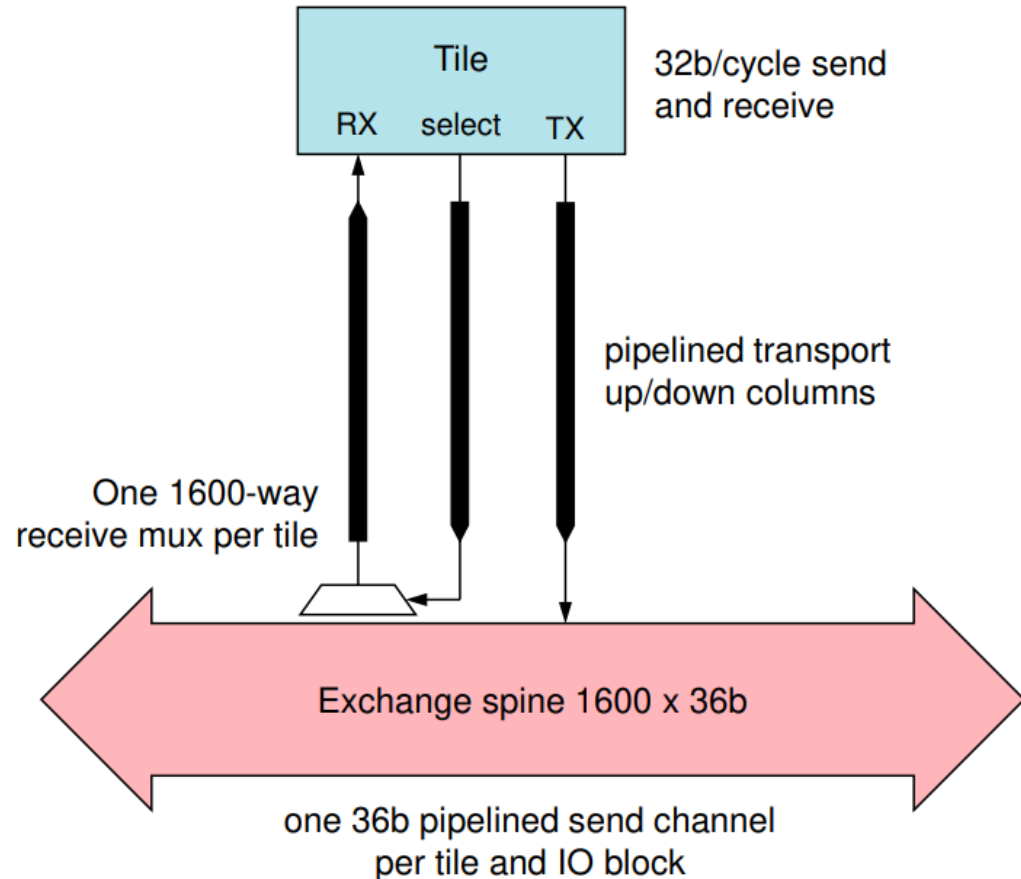
- Repeat {Sync; Exchange; Compute}
- Hardware global sync. In ~150 cycles on chip, 15 ns/hop between chips





Exchange Mechanics

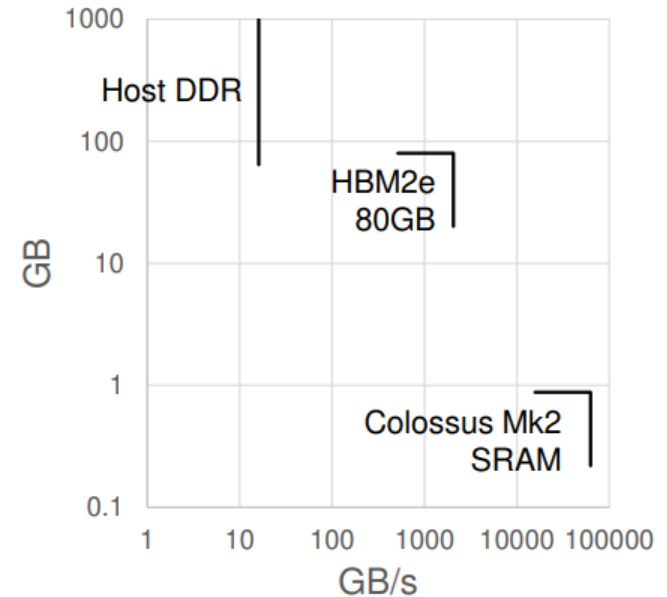
- **IPU POPLAR compiler**
 - Schedule transmit, receive and select at precise cycles from sync
 - Knowing all pipeline delays
- **Data movement**
 - At full bandwidth
 - No queues, arbiters, or packet overheads





Why no HBM Memory ?

- **Memory bandwidth** limits how fast AI can complete
- **GPU and TPU**
 - Solve for bandwidth and capacity using HBM
 - HBM is expensive, capacity-limited, and adds 100W+ to the processor thermal envelope
- **IPU**
 - Solves for bandwidth with SRAM, and for capacity with DDR





IPU hardware helps software

- **Simple mechanisms allow software evolution**
 - Native graph abstraction
 - Codelet-level parallelism
 - Pipeline-oblivious threads
 - BSP removes concurrency hazards
 - Stateless all-to-all exchange
 - Cacheless, uniform, near/far memory



Takeaway Questions

- Why GraphCore IPU employ large SRAM instead of HBM?
 - (A) Achieve high bandwidth
 - (B) Large memory capacity
 - (C) Save silicon area

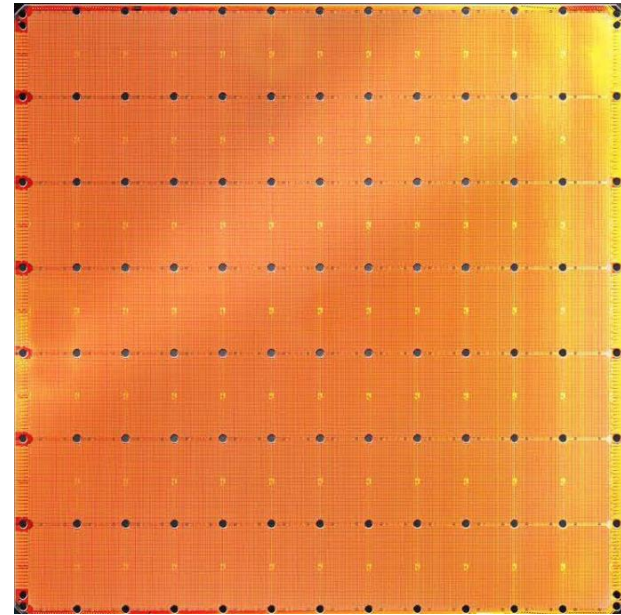


Wafer-scale AI chip -- Cerebras



Largest AI chip

- 46,225 mm² silicon
- 1.2 trillion transistors
- 400,000 AI optimized cores
- 18 Gigabytes of on-chip memory
- 9 Pbyte/s memory bandwidth
- 100 Pbit/s fabric bandwidth
- TSMC 16 nm process



Cerebras WSE

21.1 Billion
Transistors
815 mm² silicon



GPU



Why big chips ?

- Big chips process data more quickly
 - Cluster scale performance on a single chip
 - GB of fast memory 1 clock cycle from core
 - On-chip interconnect orders of magnitude faster than off-chip
 - Model-parallel, linear performance scaling
 - Training at scale, with any batch size, at full utilization



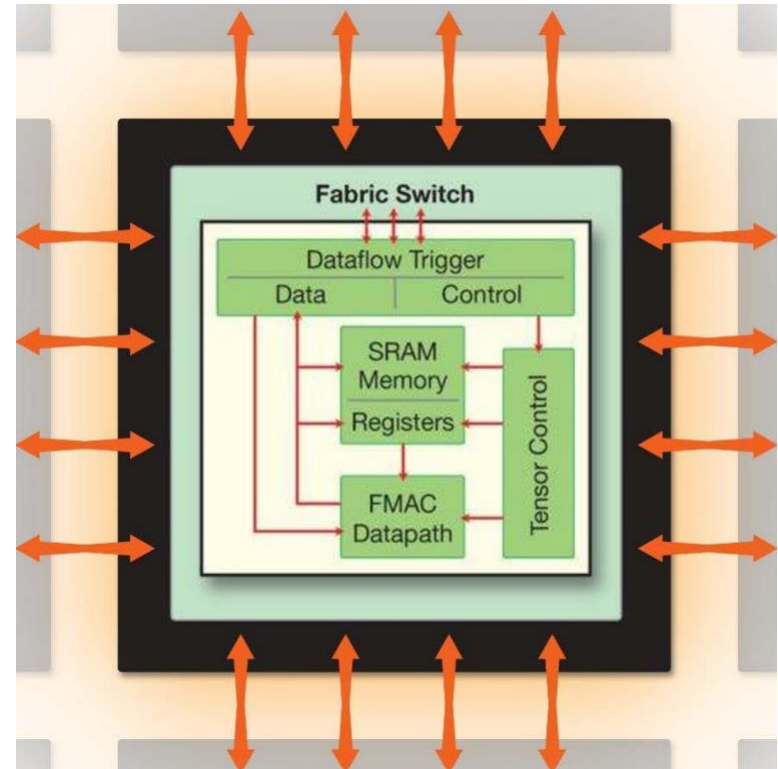
Cerebras Architecture

- Core optimized for neural network primitives
- **Flexible, programmable core**
 - NN models are evolving
- **Designed for sparse compute**
 - Workloads contain fine-grained sparsity (where are these sparsity from ?)
- **Local memory**
 - reusing weight & activations
- **Fast interconnect**
 - Layer-to-layer with high bandwidth and low latency



Cerebras programmable core

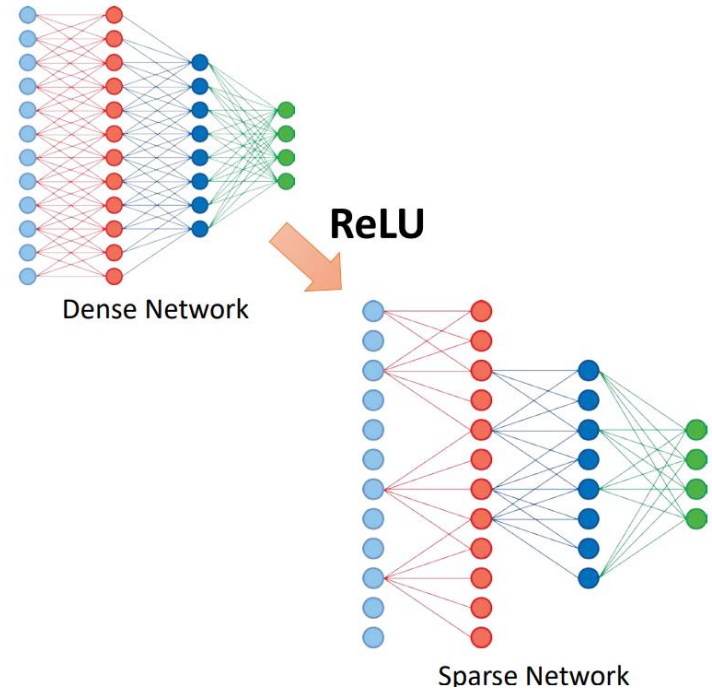
- Flexible cores optimized for **tensor operations**
 - General ops for control processing
e.g. arithmetic, logical, LD/ST, branch
 - Optimized tensor ops for data processing
 - **Tensor operands**
e.g. $\text{fmac } [Z] = [Z], [W], a$
3D 3D 2D





Sparse compute engine

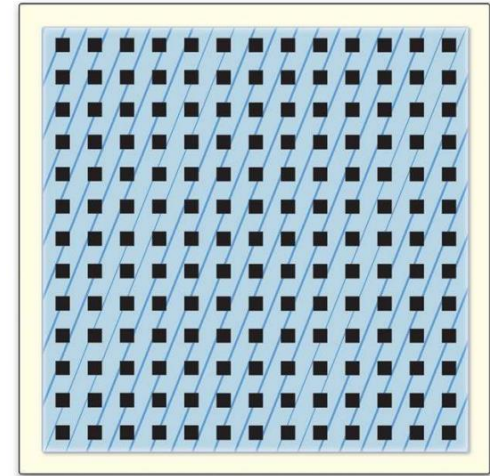
- **Nonlinear activations** naturally create **fine-grained sparsity**
- **Dataflow scheduling in hardware**
 - Triggered by data
 - Filters out sparse zero data
 - Skips unnecessary processing
- **Fine-grained execution datapaths**
 - Small cores with independent instructions
 - Efficiently processes dynamic, non-uniform work





Cerebras memory architecture

- **Traditional memory designs**
 - Centralized shared memory is slow & far away
 - Requires high data reuse (caching)
 - Local weights and activations are local -> low data reuse
- **Cerebras memory architecture**
 - All memory is fully distributed along compute
 - Datapath has full performance from memory



Memory uniformly distributed across cores

■ Core ■ Memory



High-bandwidth low-latency interconnect

- **2D mesh topology** effective for local communication
 - High bandwidth and low latency for local communication
 - All HW-based communication avoids SW overhead
 - Small single-word message



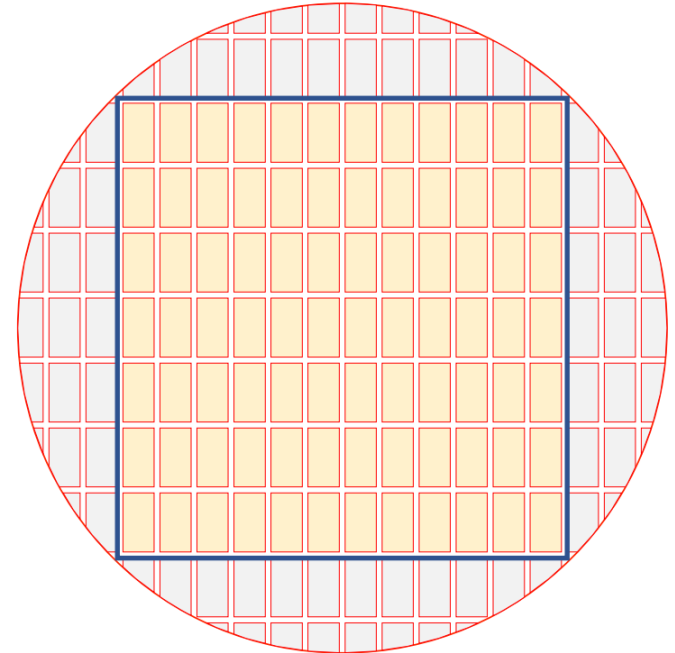
Challenges of wafer scale

- Building a 46,225 mm², 1.2 trillion transistor chip
- **Challenges include**
 - Cross-die connectivity
 - Yield
 - Thermal expansion
 - Package assembly
 - Power and cooling



Challenge 1: cross die connectivity

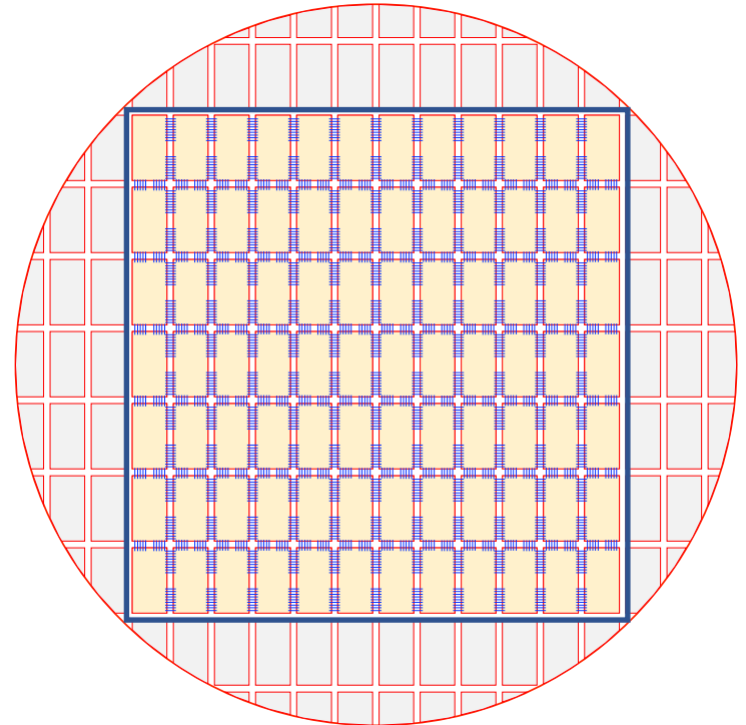
- Standard fabrication process requires die to be independent
- Scribe line separates each die
- Scribe line used as mechanical barrier for die cutting for test structures





Cross-die wires

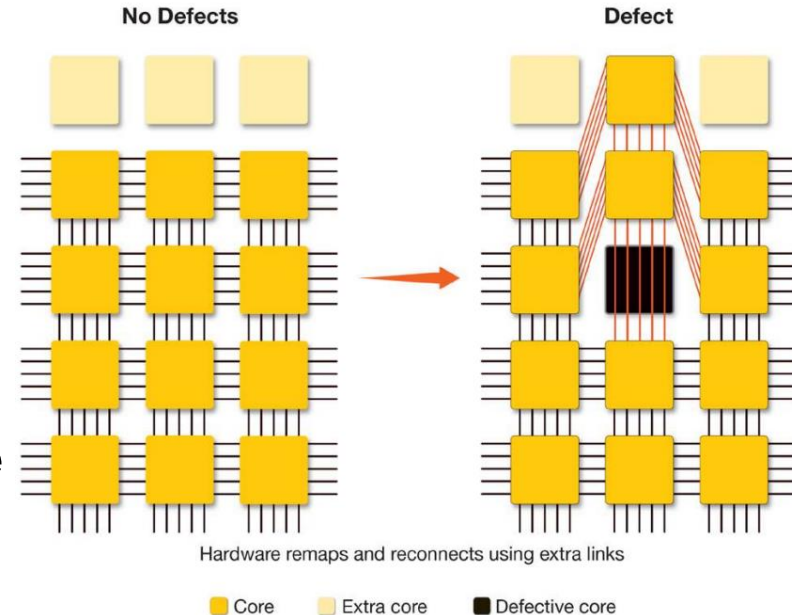
- Add wires across scribe line with TSMC
- Extend 2D mesh across die
- Same connectivity between cores and across scribe lines create a homogeneous array
- Short wires enable ultra high bandwidth with low latency





Challenges II: Yield

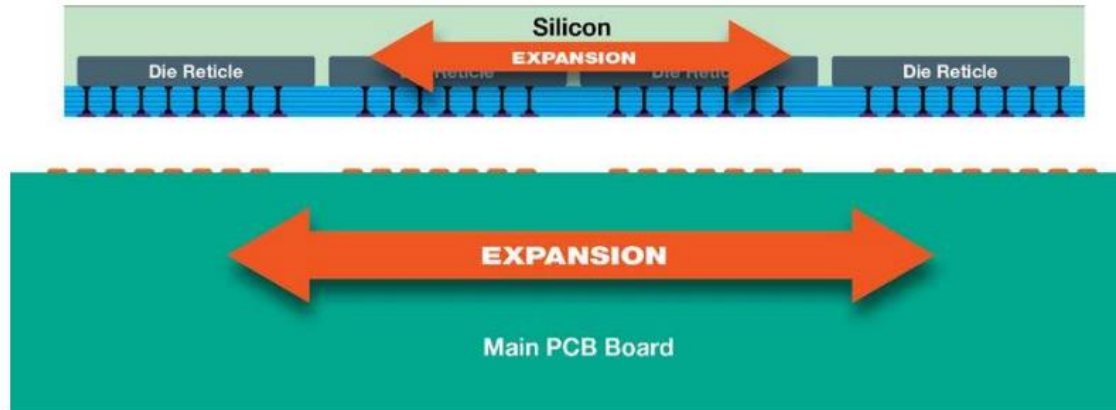
- Impossible to yield full wafer with zero defects
 - Silicon and process defects are inevitable even in mature process
- **Redundant cores**
 - Uniform small cores
 - Redundant cores and fabric links
 - **Redundant cores replace defective cores**
 - **Extra links** reconnect fabric to restore logical 2D mesh





Challenge III: Thermal expansion in package

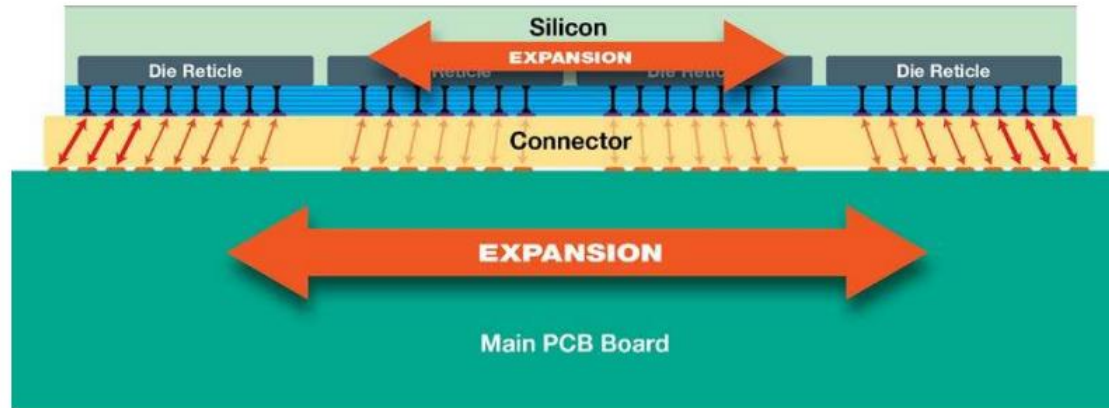
- **Silicon and PCB expand at different rates** under temperature
- Size of wafer would result in too much mechanical stress using traditional package technology





Connecting wafer to PCB

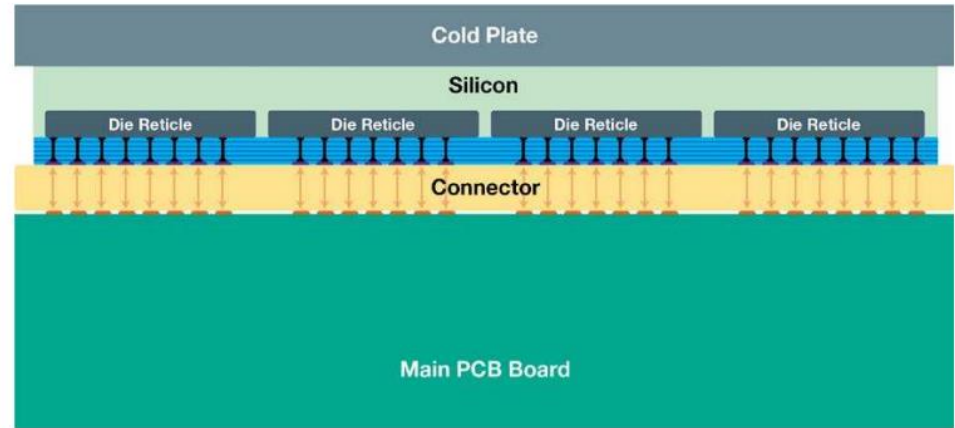
- Developed custom connector to connect wafer to PCB
- Connector absorbs the variation while maintaining connectivity





Challenge IV: Package assembly

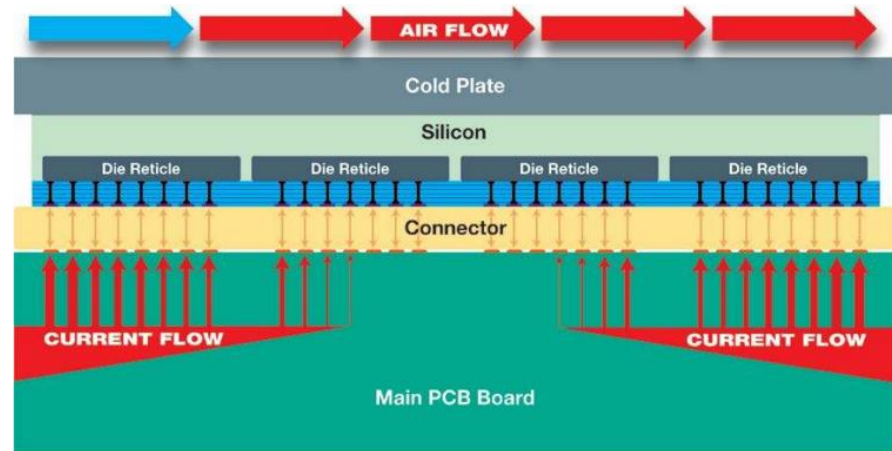
- Package includes
 - PCB
 - Connector
 - Wafer
 - Cold plate
- All components require precise alignment
- Developed custom machines and process





Challenge V: Power and cooling

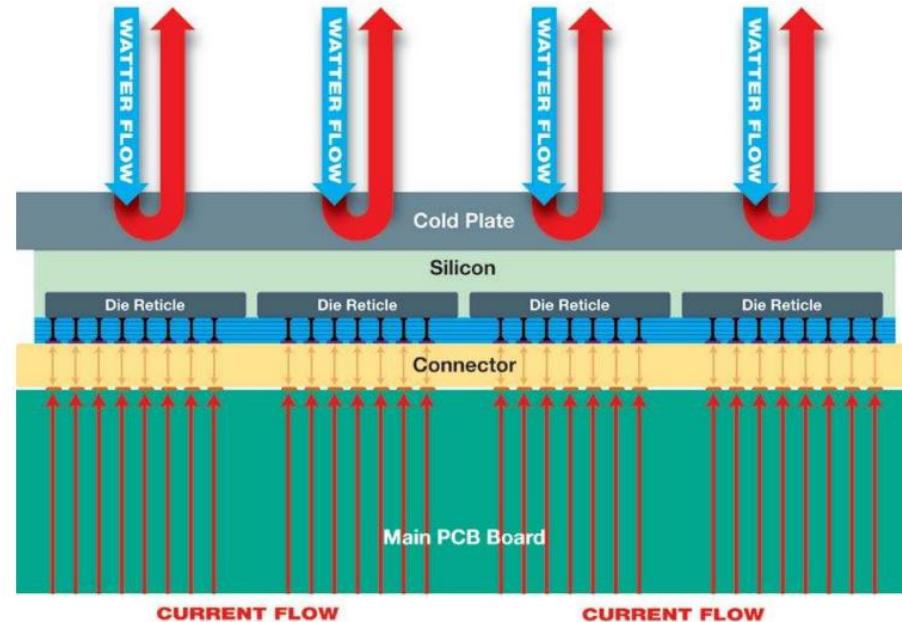
- Concentrated high density exceeds traditional power & cooling capacities
- **Power delivery**
 - Current density too high for power plane distribution in PCB
- **Heat removal**
 - Heat density too high for direct air cooling





Using the 3rd dimension

- **Power delivery**
 - Current flow distributed in 3rd dimension perpendicular to water
- **Heat removal**
 - Water carries heat from wafer through cold plate





Takeaway Questions

- What are challenges to build a large chip for NN applications ?
 - (A) Power and cooling
 - (B) Fault tolerance for defected dies
 - (C) Package assembly
- How does Cerebras tackle the DNN sparsity ?
 - (A) Customized sparse core
 - (B) Data-driven dataflow scheduling
 - (C) Filters out sparse zero data