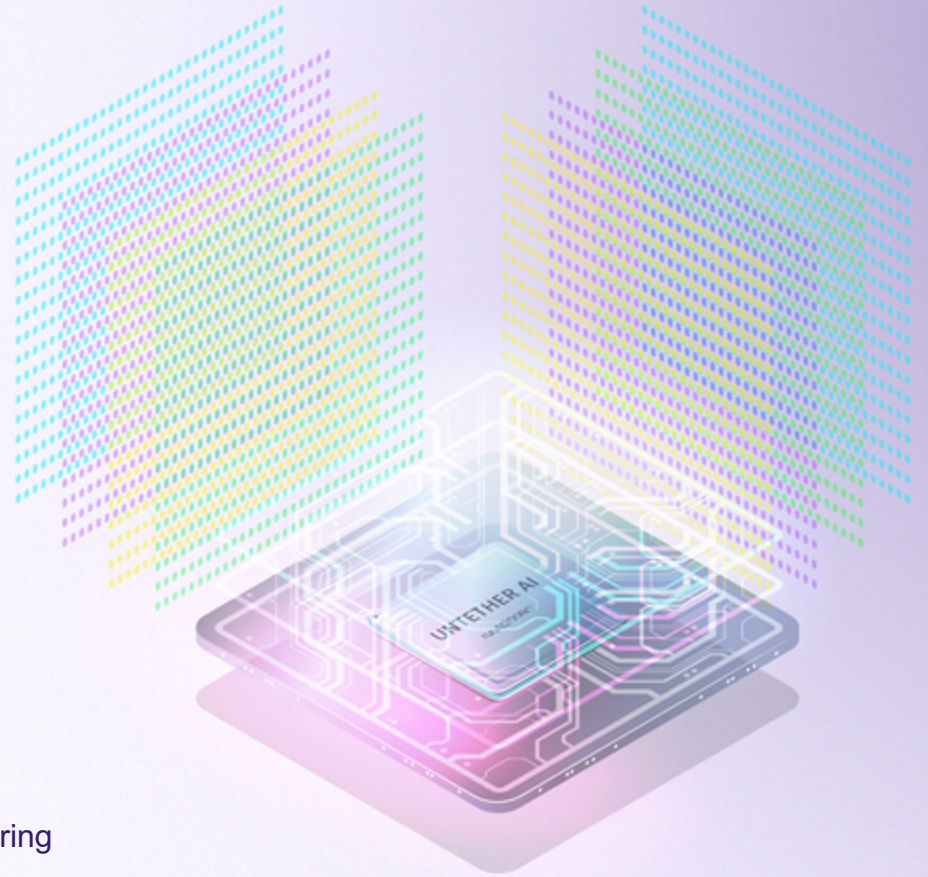


# UNTETHER AI

## Boqueria

Robert Beachler – VP of Product/Hardware Engineering

Dr. Martin Snelgrove – Co-founder and CTO



# A Brief History of the Current AI Summer

2012

2014

2016

2018

2020

2022

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky  
University of Toronto  
kriz@cs.utoronto.ca

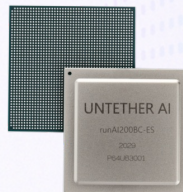
Ilya Sutskever  
University of Toronto  
ilya@cs.utoronto.ca

Geoffrey E. Hinton  
University of Toronto  
hinton@cs.utoronto.ca

Deepmind  
acquired  
by Google

AlphaGo beats  
Lee Sedol

Untether AI founded  
in Toronto



runAI200 introduced  
First at-memory inference accelerator  
500 INT8 TOPs  
200MB SRAM  
8 TOPs/W  
TSMC 16nm



Boqueria introduced  
2 PetaFlops FP8  
238MB SRAM  
30 TFLOPs/W  
TSMC 7nm

# AI Inference Presents 3 Key Challenges to Chip Makers



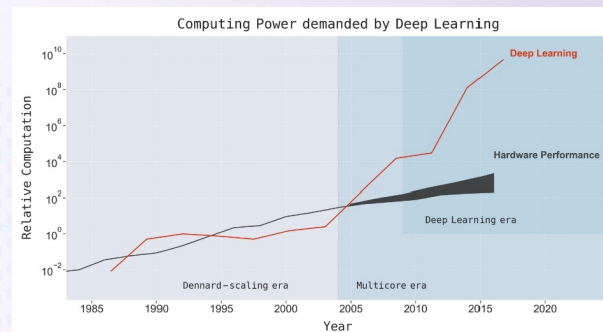
Increasing computational and power requirements



Scalability and flexibility for changing NN landscape

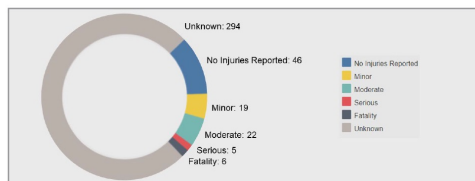


Accuracy loss costs \$millions and risks lives

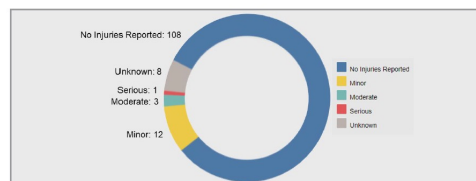


The Computational Limits of Deep Learning  
Neil C. Thompson<sup>1</sup>, Kristjan Greenewald<sup>2</sup>, Keeheon Lee<sup>3</sup>, Gabriel F. Manson

Level 2 ADAS Highest Injury Severity



ADS Highest Injury Severity



NHSTA report, June 2022 for July 2021 to May 2022

Model Category	Model Name	Model Size (Mparams)
Recommendation	Less complex	70,000
	More complex	>100,000
Computer Vision	ResNetXt101-32x4-48	44
	RegNetY	700
	FBNetV3 based model	28.6
Video Understanding	ResNetXt3D based	58
NLP	XLRM-R	558

First-Generation Inference Accelerator Deployment at Facebook

# Architecting an AI Inference Accelerator



Power-efficient throughput is required to meet NN compute demand

- Data movement is the costliest part of inference – 90% of energy consumption
- Data movement is different between training and inference
- Optimizing compute architecture to minimizing distance travelled results in inference-specific AI accelerators



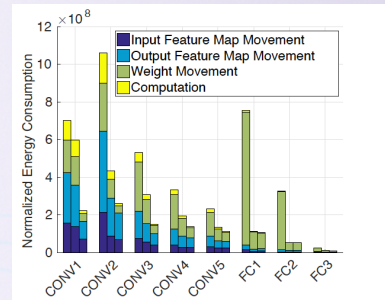
Proper level of granularity to create a scalable compute architecture

- Right balance between coarse-grained and fine-grained approach
- Don't over-fit for a particular application/NN

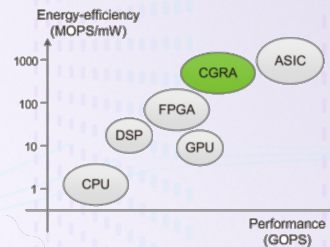


Utilize the most efficient datatype for a given application and accuracy requirements

- A mixture of datatypes provides the best results



Designing Energy-Efficient Convolutional Neural Networks using Energy-Aware Pruning  
Tien-Ju Yang, Yu-Hsin Chen, Vivienne Sze, Massachusetts Institute of Technology



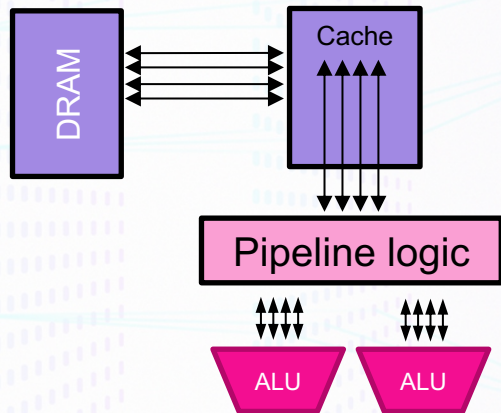
A Survey of Coarse-Grained Reconfigurable Architecture and Design: Taxonomy, Challenges, and Applications  
LEIBO LIU, JIANFENG ZHU, ZHAOSHI LI, et. AI.

Datatype	F1
FP32	1.000
BF16	1.001
FP8	0.996

BERT-Base SQuAD1.1 accuracy, no retraining, Untether AI

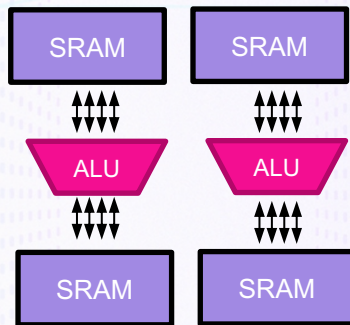
# At-Memory Compute Is the Sweet Spot for AI Acceleration

Near Memory/  
Von Neumann Architectures



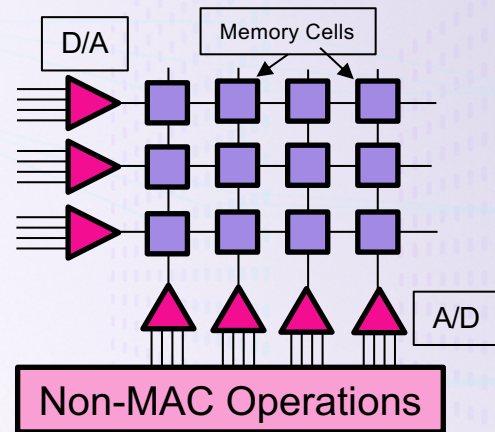
- Long, narrow busses
- Deep/shared cache

At-Memory Computation



- Short, massively parallel direct connections
- Dedicated, optimized memory for efficiency and bandwidth

In-Memory Computation



- Multi-value memory cell
- Analog techniques used for multiply-accumulate
- A/D and D/A support circuitry
- Digital processors for non-MAC operations

# Boqueria : A 2 PFLOPs, 30 TFLOPs/W At-Memory Inference Accelerator with 1,458 RISC-V Cores

## 729 Dual RISC-V memory banks provide unmatched performance

- 2,015 FP8 TFLOPs, 1,008 BF16 TFLOPs\*
- 1.35 GHz, TSMC 7nm

## At-memory compute provides energy efficiency and massive bandwidth

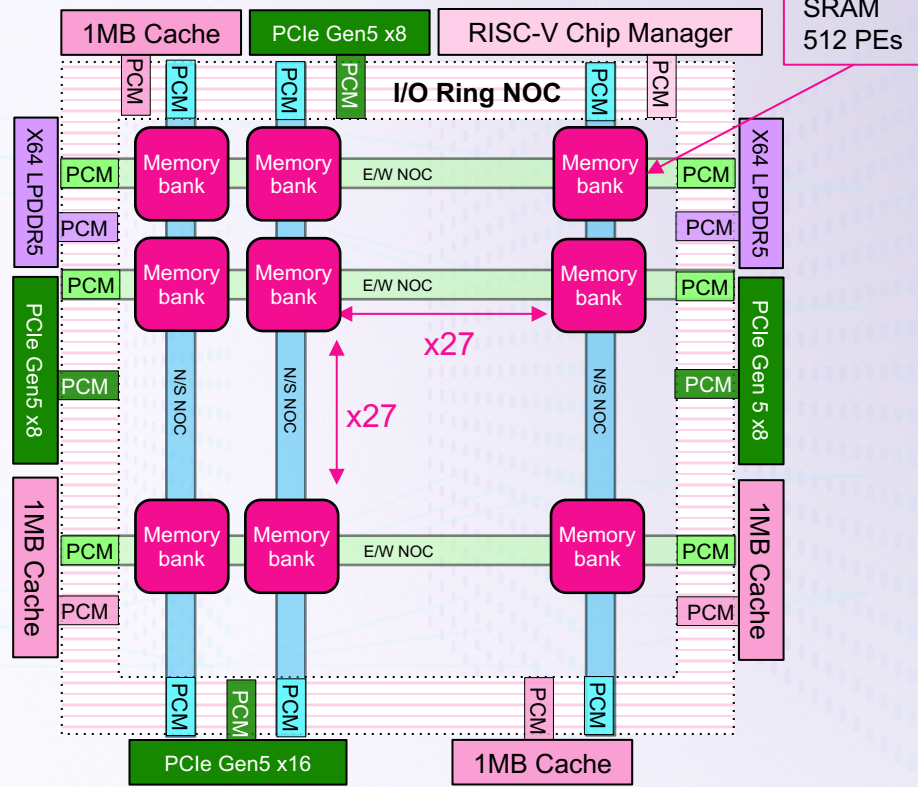
- 30 TFLOPs/W
- 238MB on-chip SRAM
- ~1 PB/s SRAM bandwidth

## Scalability

- External memory support
  - 32GB LPDDR5 across two x64 ports
  - 819 Gb/s DRAM bandwidth
- PCI-Express Host and chip-to-chip connectivity

## Accuracy

- Multiplicity of datatypes - INT4 to BF16





# The Memory Bank



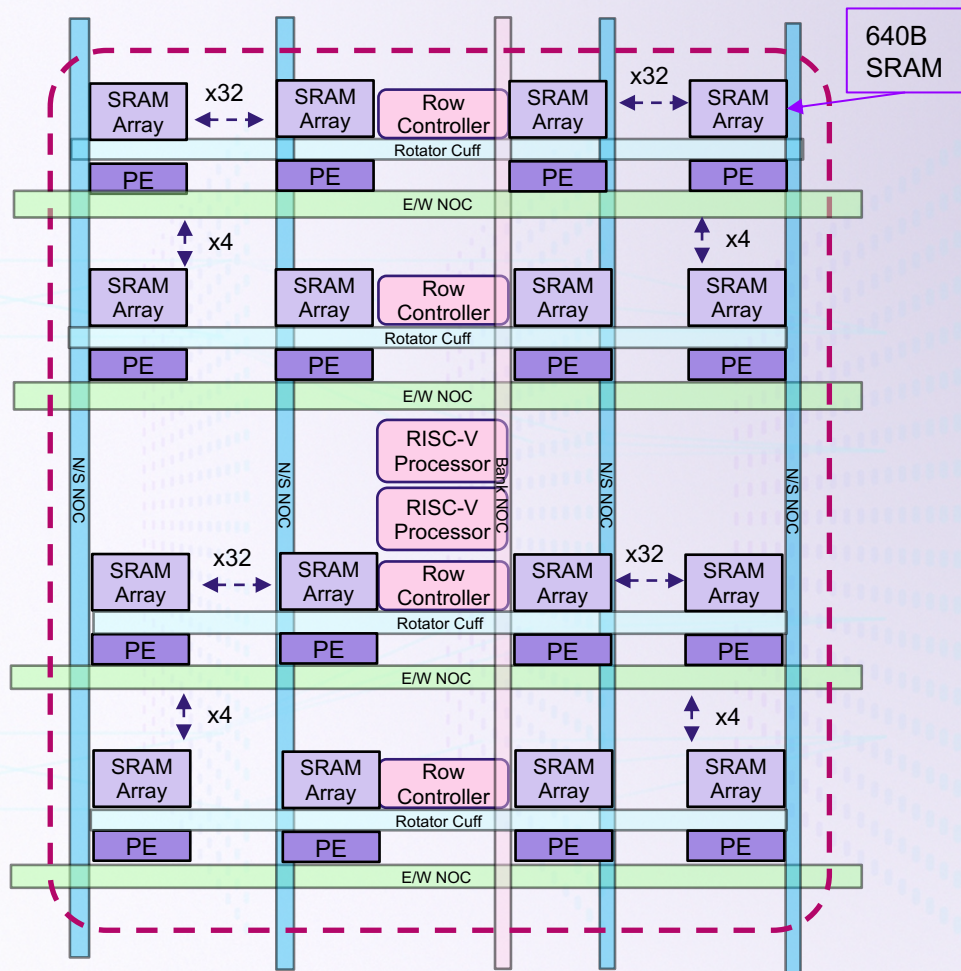
## Dual multi-threaded RISC-V for programming flexibility

- Each RISC-V manages 4 row controllers
- Row controllers operate independently
  - Command/control 64 SIMD PEs for GEMV calculations, scalable to GEMM



## Extreme connectivity – custom pipelined communication

- Rotator cuff moves activations between nearest neighbor PEs, conserving energy
- 8 E/W NOCs, each with 7GB/s bandwidth in each direction (56GB/s total in each direction)
- 1 N/S NOCs, with 70GB/s bandwidth in each direction
- Bank NOC allows communication between RISC-V Processors





# At-Memory Compute – Putting Compute Where the Data resides



## Low-Power SRAM Array

- Coefficient and data storage
- 0.4V datapath operation
- Custom drivers to minimize power when reading/writing to memory



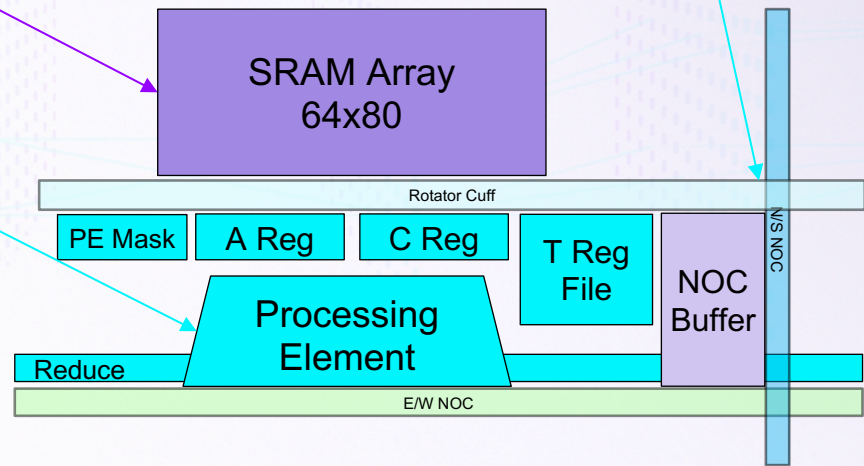
## Energy efficient data transfer

- Rotator cuff accesses activations from nearest neighbors
- E/W and N/S pipelined NOCs for spanning rows/columns of memory banks



## Processing Element

- INT4, INT8, FP8p/r, BF16 Support
- Zero detect for power savings
- Structured sparsity support
- Dedicated reduce circuitry for functions such as SoftMax and LayerNorm







# Untether AI FP8 – Designed for Inference Acceleration

**AI Inference requires precision – but full floating point is expensive to achieve needed accuracy targets**

**Untether AI researched various datatype and found FP8 provided the best balance of precision, throughput, and energy efficiency**

- But requires both range (FP8r) and precision (FP8p)
- Each additional mantissa doubles the precision, and reduces mean squared error by  $(1/2)^2$

**FP8 is 2x times more energy/die size efficient than if designed for native INT8**

UAI FP8 precision



UAI FP8 range



Training FP8 precision



Training FP8 range



# Accuracy Results – Relative Degradation

## ResNet-50 ImageNet accuracy

Datatype	Top-1 Accuracy
FP32	1.000
BF16	1.000
INT8	0.992
FP8	0.991

Accuracy degradation between FP8 and INT8 is negligible

FP8 quadruples throughput

BF16 for ultimate accuracy

FP8 to quadruple throughput

## BERT-Base SQuAD1.1 accuracy

Datatype	Exact Match	F1
FP32	1.000	1.000
BF16	1.000	1.001
FP8	0.988	.996



# Multi-threaded Custom RISC-V Processor – Adapted to AI Inference

## Standard RV32EMC Instruction Set

- Embedded
- Multiplication/division
- Compressed instructions
- ***UAI added 20+ custom instructions specific to at-memory compute and inference acceleration***

## Each Processor

- 6KB Memory
- 32-bit ALU
- 32-bit multiplier
- x16 register file with 4-way context switching to enable 4 threads

## Example Custom Instructions

Instruction	Description
<code>pe_move</code>	Copy PE register to PE register
<code>pe_rcv</code>	Receives packets from socket
<code>pe_send</code>	Send packets to socket
<code>pe_load</code>	Load PE registers from CRAM
<code>pe_store</code>	Store PE registers to CRAM
<code>pe_macc</code>	Perform MACC
<code>pe_gemv</code>	Perform GEMV
<code>pe_multimove</code>	Copy into multiple PE registers
<code>pe_broadcast</code>	Broadcast and copy to PE register
<code>pe_rotate</code>	Rotate for count cycles
<code>pe_convert</code>	Convert d and output to PE register
<code>pe_set_fmt</code>	Set number format
<code>pe_row_reduce</code>	Run row reduce function
<code>pe_reduce</code>	Run PE reduce function
<code>pe_k_stack</code>	Push and pop k stack
<code>rwc_nop</code>	nop for count cycles
<code>rwc_sleep</code>	Sleep with wakeup mask
<code>eq_set_cfg</code>	Set row extension configuration
<code>eq_save</code>	Save cuff registers to EQ
<code>rwc_magic</code>	Magic function to output RWC commands

# High-bandwidth I/O and Connectivity

## High-Bandwidth I/O NOC

- 141GB/s in clockwise direction
- 141GB/s in counter-clockwise direction

## Efficient, high throughput data NOCs

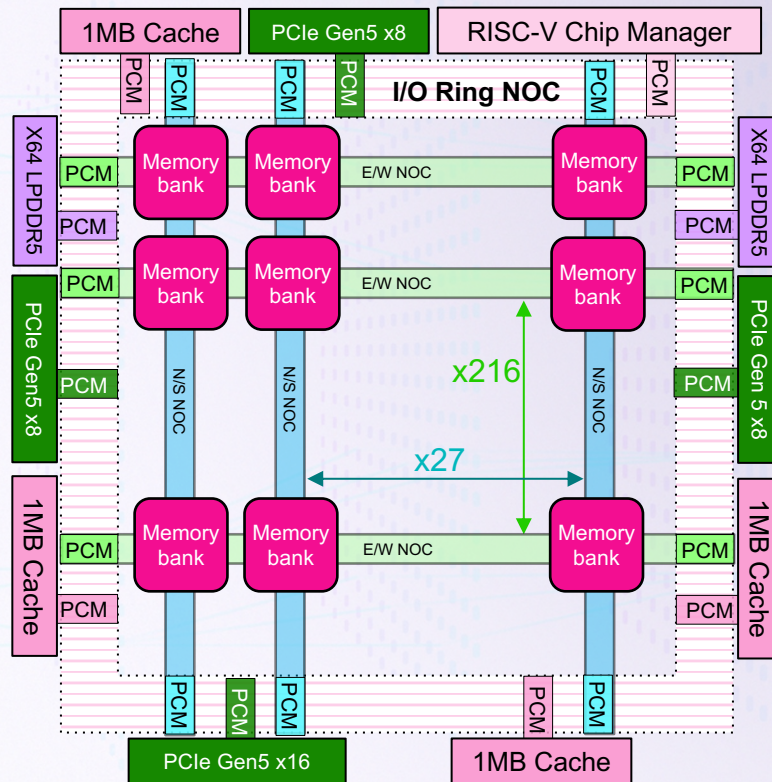
- 1.5 TB/s East and West throughput
- 1.9 TB/s North and South throughput

## Extremely scalable

- X16 PCIe Gen5 for host connectivity – 63GB/s
- 3 ports of PCIe Gen5 x8 for chip-to-chip and card-to-card connectivity – each 31.5GB/s

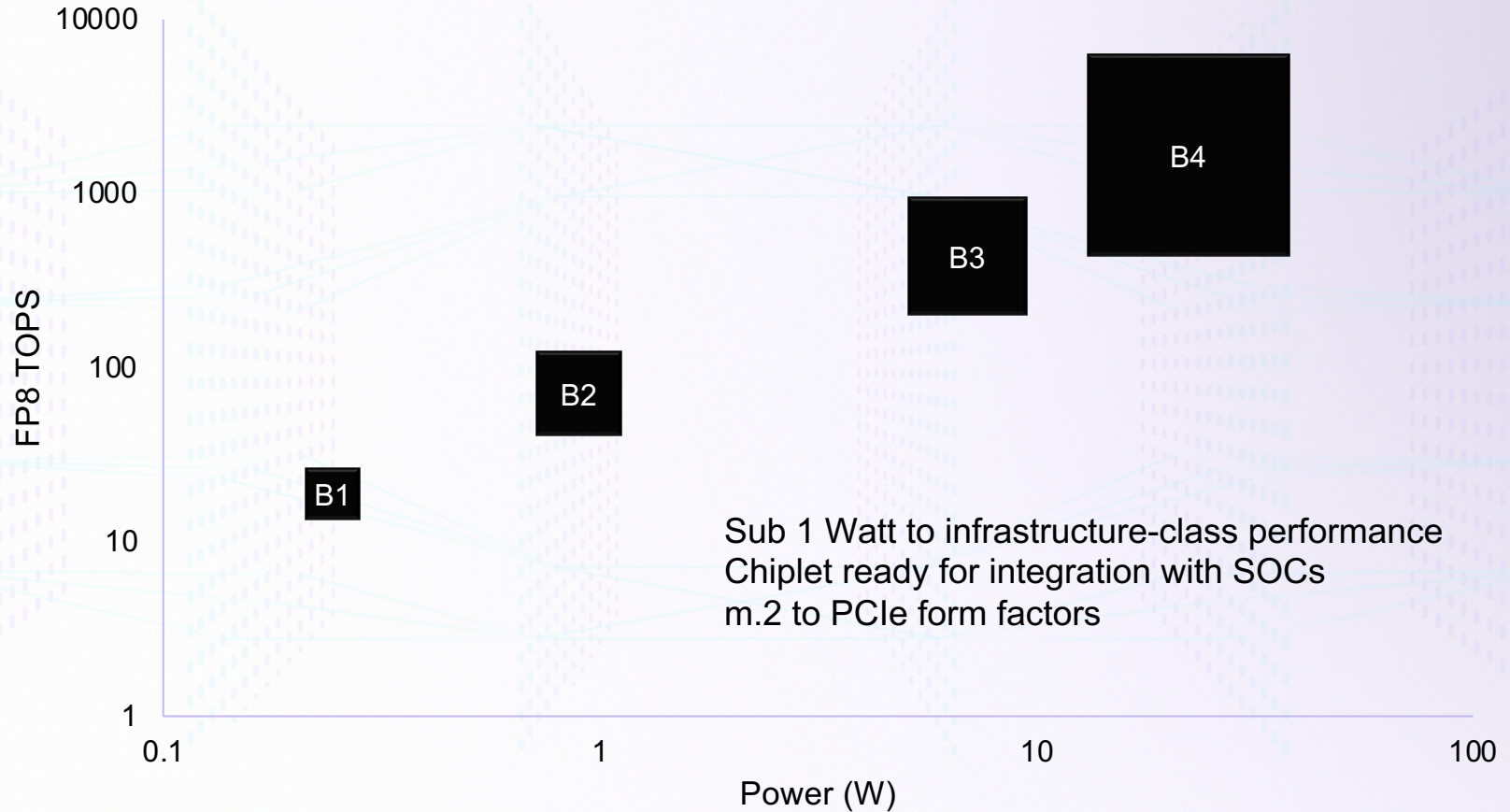
## Multi-level Memory Architecture

- 238MB of At-memory SRAM for efficient compute
  - ~1 PB/s bandwidth
- 4MB of scratchpad for data manipulation
- 32 GB of external LPDDR5 with >100GB/s bandwidth





# The Benefits of a Scalable Architecture



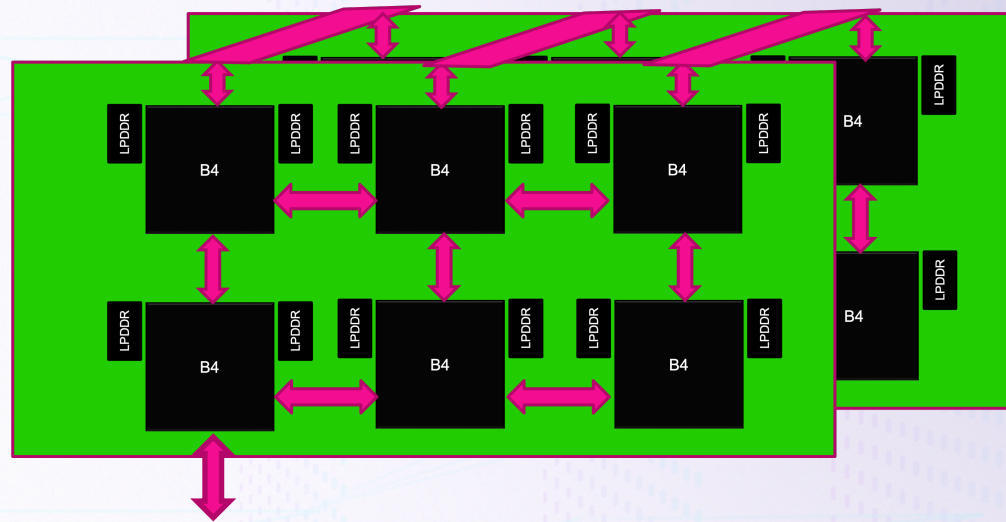
# Scalability For Large Language Models

## Maximum performance: 6-Chip PCIe Gen5 Card

- 1.4GB SRAM
- 12 PetaFLOPs of FP8 compute
- 192GB LPDDR5 DRAM

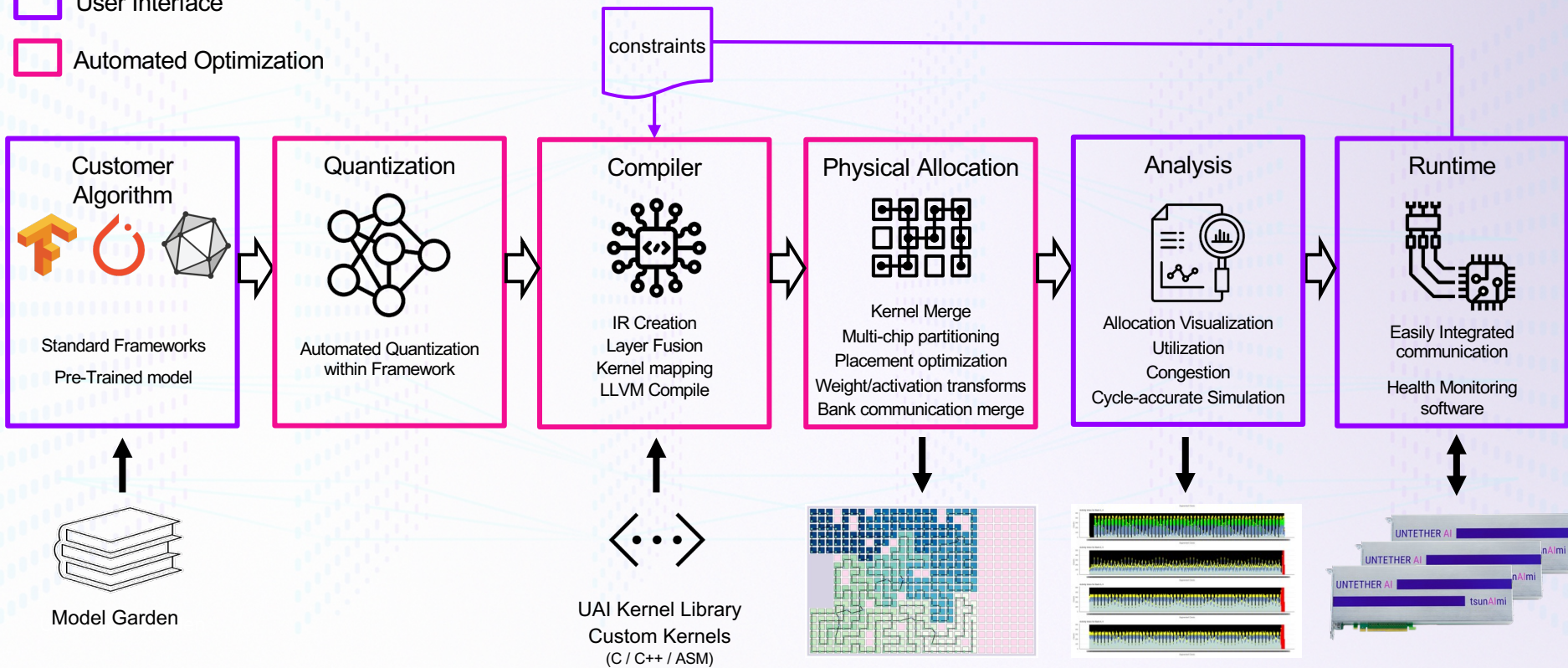
## Scalability

- PCIe GEN5 x8 chip-to-chip and card-to-card communication



# imAIgine SDK Tool Flow

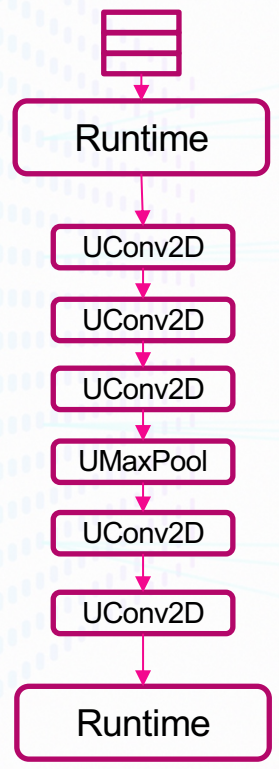
- User Interface
- Automated Optimization



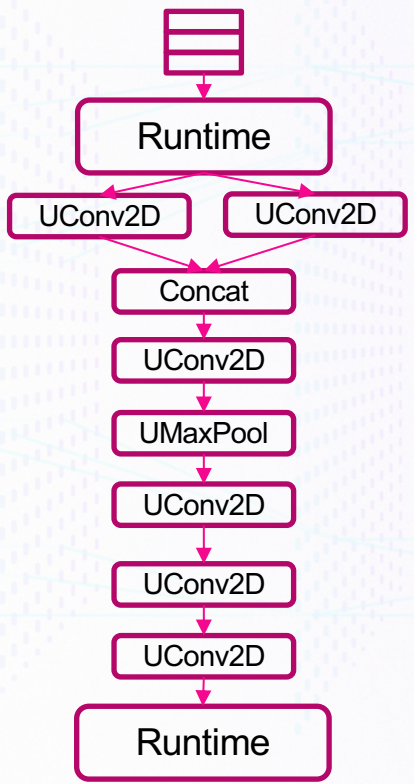


# imAIne SDK: Spatial Compilation Optimizations

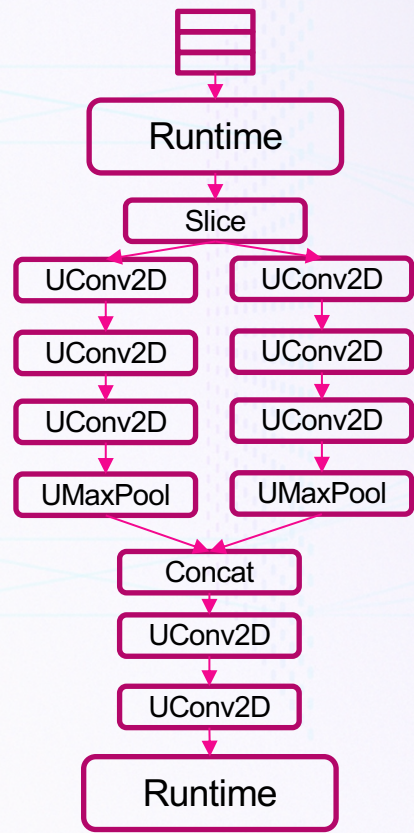
Pipelining



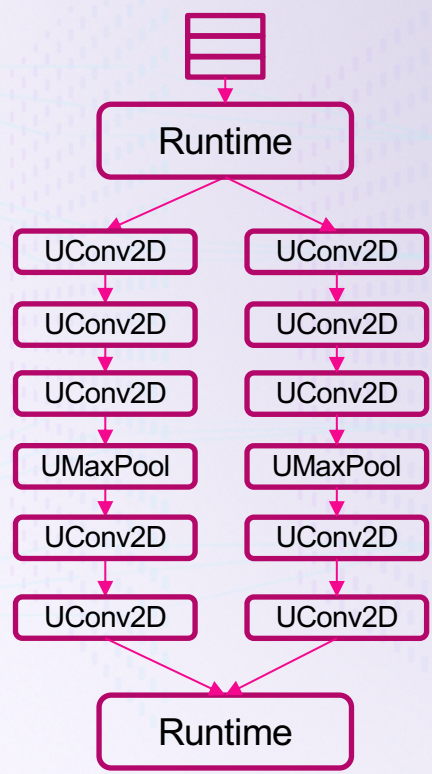
Layer Replication



Sub-Graph Replication



Multi-Instance



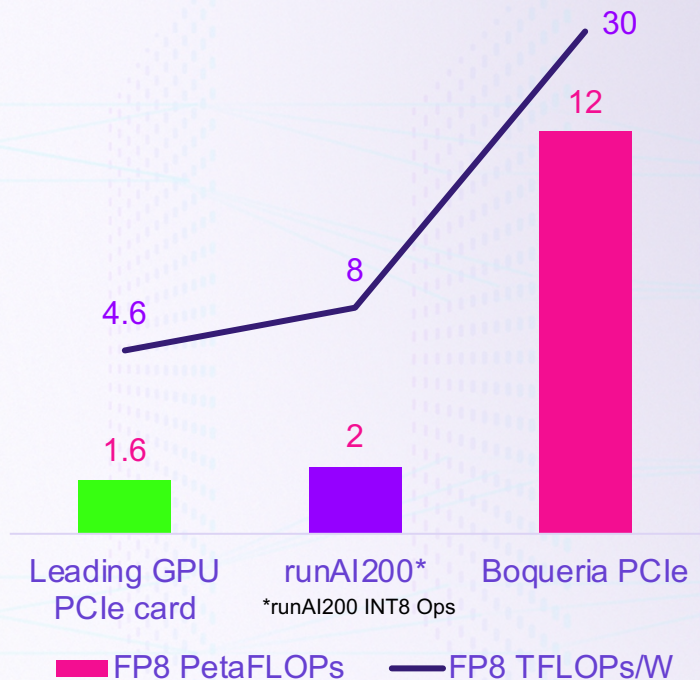




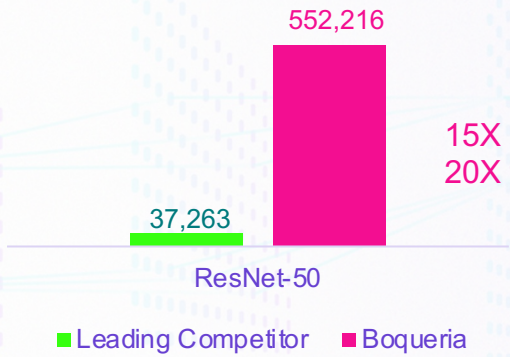
# The Efficiency of At-memory Compute

## Boqueria compared to leading GPU

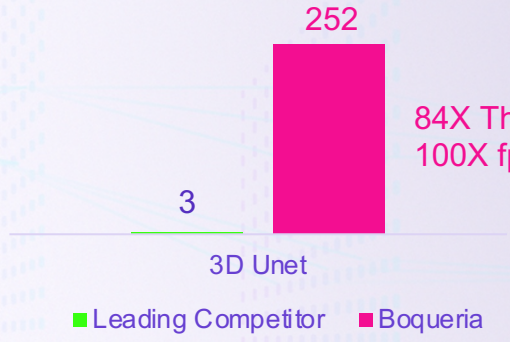
- 5x greater performance (PFLOPs/card)
- 7X greater efficiency (TFLOPs/W)



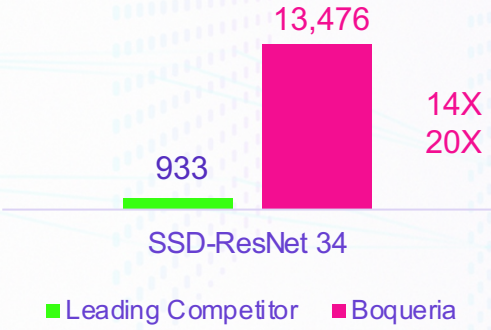
# Energy Efficiency Translates to Throughput



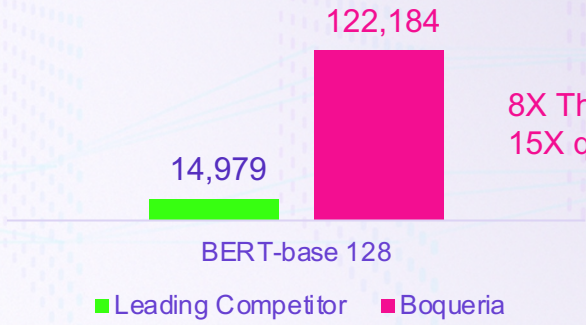
15X Throughput  
20X fps/W



84X Throughput  
100X fps/W



14X Throughput  
20X fps/W



8X Throughput  
15X qps/W

Competitor information based on MLPerf 2.0 best performance normalized to single PCIe-card and TDP except BERT-base 128 based on published benchmark by vendor and TDP  
Boqueria information based on kernel code cycle counts and simulations

# Boqueria – A 2 PFLOPs, 30 TFLOPs/W At-Memory Inference Accelerator



## Unrivalled throughput and efficiency

- Up to 100X efficiency, 80X throughput for a variety of neural network models

## Scalable from small to large models

- Spatial architecture and chip/card interconnects enables optimal performance/power



## Flexible to adapt to latest NN architectures

- Equally adept at vision and NLP networks



## Efficient and accurate datatypes

- FP8p/r for accuracy and throughput, BF16 for utmost accuracy