# RISC-V Near-Memory Processing Accelerators

## Peter Hsu, Ph.D.

Peter Hsu & Associates, S.L.
Barcelona, Spain

Email:  peter.hsu@phaa.eu
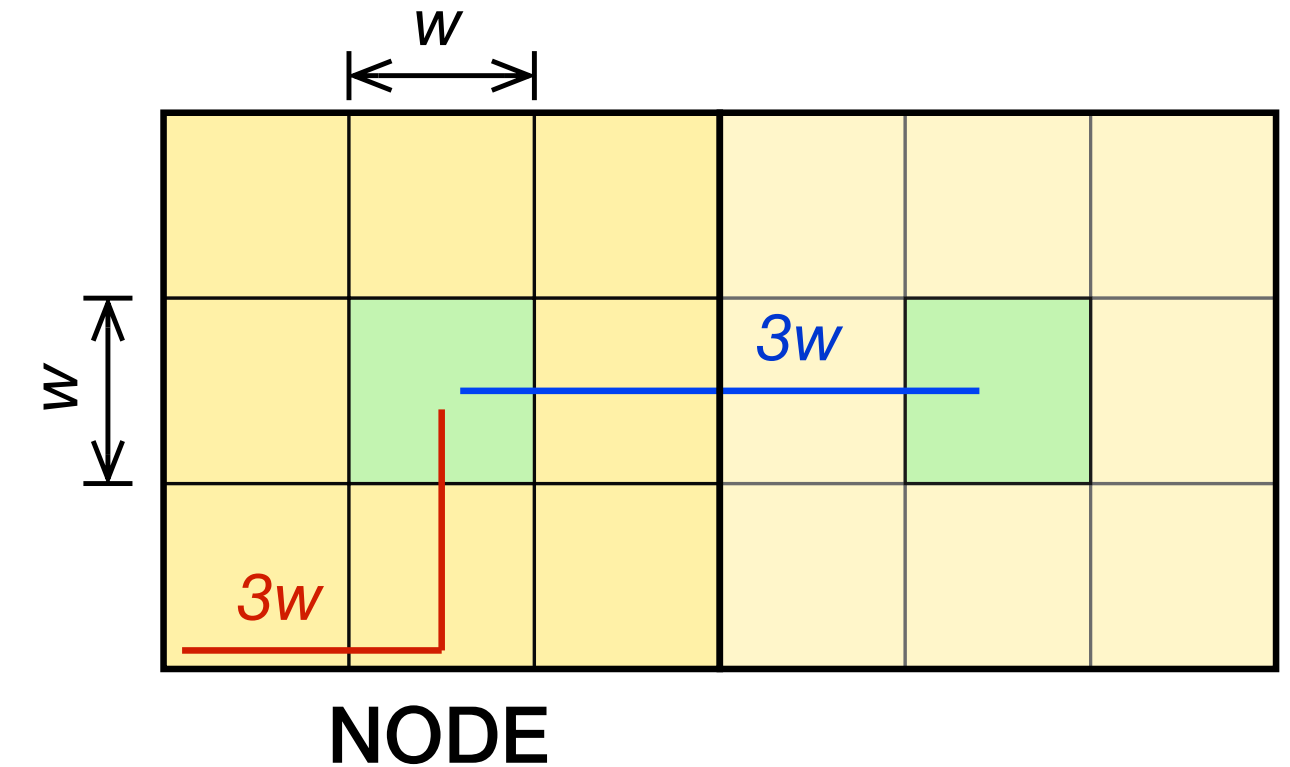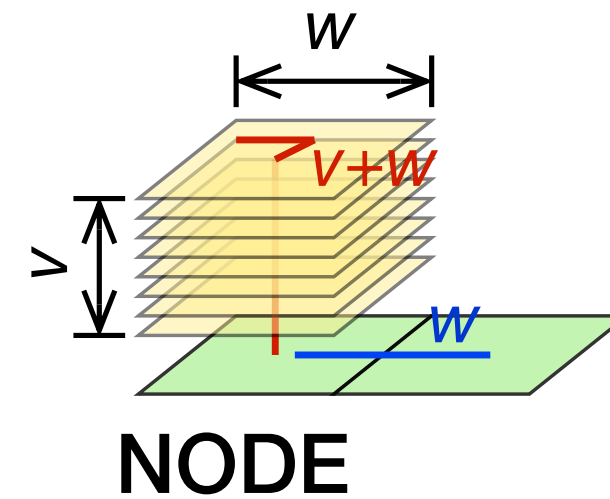LinkedIn:  www.linkedin.com/in/peter-hsu-122a315

# Introduction

- RISC-V opens many new opportunities for innovation and enables smaller organizations to design computer hardware

- But no matter how creative, a new chip design using advanced process technology faces enormous development cost

- We propose a "mix and match" paradigm combining state-of-the-art memory technology with mature ASIC logic to reduce development cost of near-memory computing architecture accelerators

- Machine learning and supercomputing accelerators examples

# Why Stacked Memory?

## Technical
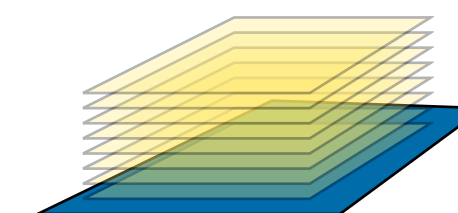
- Shorter SRAM access path, node communication distance

$w$

$v+w$

$v$

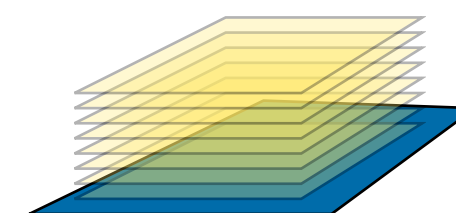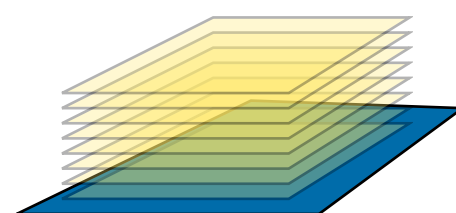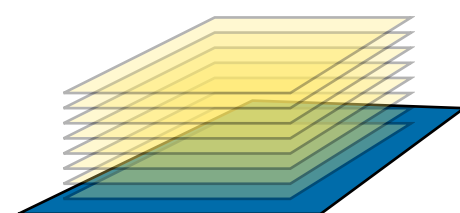$w$

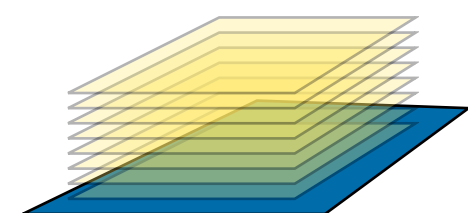**NODE**

$w$

$w$

$3w$

$3w$

**NODE**

## Business

- Reduce development cost

Save cost by building custom logic on mature ASIC technology

Reuse stacked memory for different designs

Multiple accelerators specialize for different applications

# Agenda

1. Economics & technology

2. Accelerator architecture

3. Examples

4. Programming model & evaluation
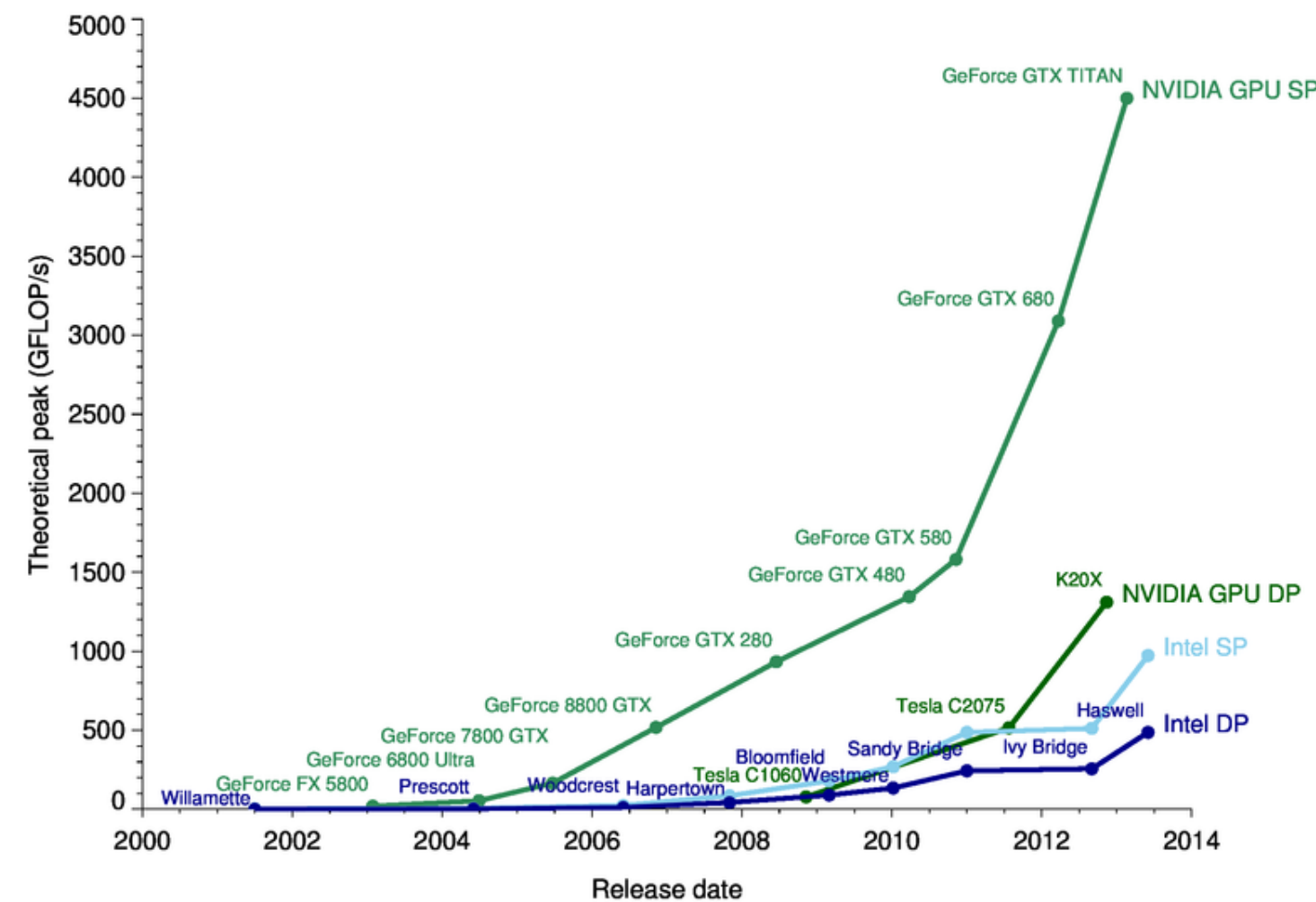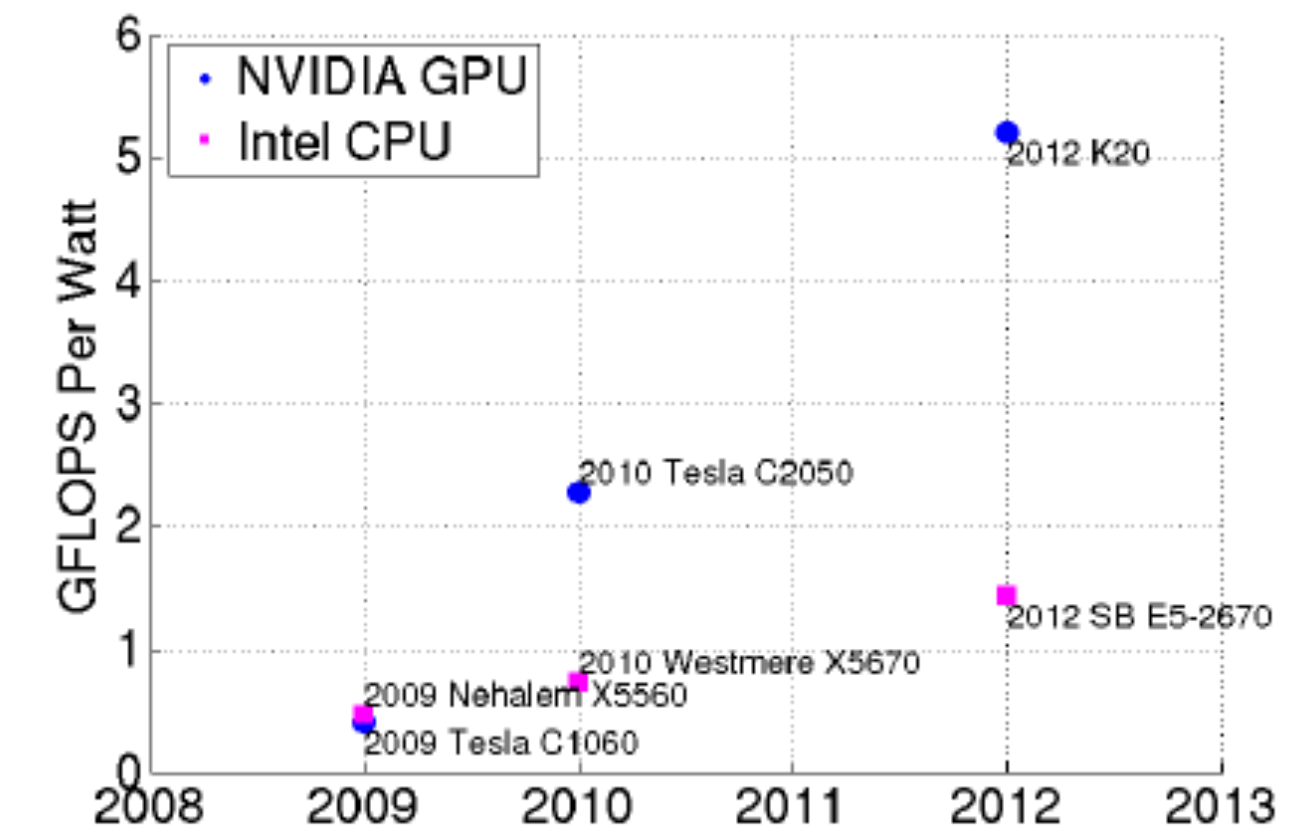
5. Memory technologies

6. Conclusion

Many interesting accelerator architectures are possible

We propose a technical solution for near-memory computing accelerators that is economically feasible for smaller organizations with limited resources

More details in paper "In-Memory Accelerators Using Stacked Memory" (PDF)

# Why Specialize?

- Performance

- Energy





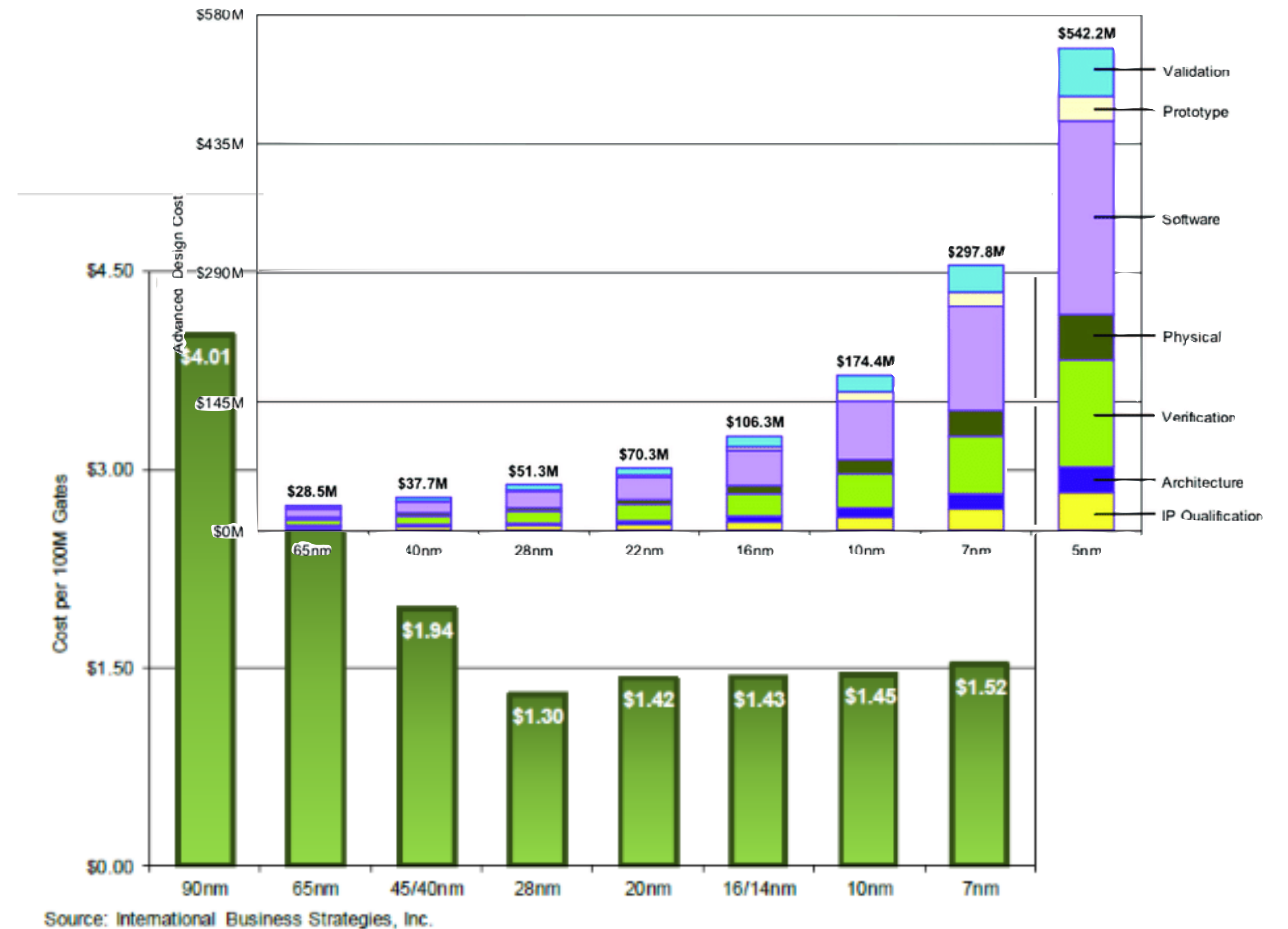| Hardware | CPUs | | | GPUs | | | |
|---|---|---|---|---|---|---|---|
| Manufacturer | Intel | Intel | Intel | NVIDIA | Sony, IBM, Toshiba | NVIDIA | NVIDIA |
| Model | Q9450 | Q9450 | Q9450 | 7900 GTX | PlayStation 3 | 8800 GTX | GTX 280 |
| # cores used | 1 | 4 | 4 | 4x96 | 2+6 | 4x128 | 4x240 |
| Implementation | MATLAB | MATLAB | SSE2 | Cg | Cell SDK | CUDA | CUDA |
| Year | 2008 | 2008 | 2008 | 2006 | 2007 | 2007 | 2008 |
| **Performance / Cost** | | | | | | | |
| Full System Cost (approx.) | $1,500** | $2,700** | $1,000 | $3,000* | $400 | $3,000* | $3,000* |
| Relative Speedup | 1x | 4x | 80x | 544x | 222x | 1544x | 2712x |
| Relative Perf. / $ | 1x | 2x | 120x | 272x | 833x | 772x | 1356x |

# Economic Considerations

Cost of silicon device

- Recurring cost (RE)

- Non-recurring cost (NRE)

5nm chip $500M NRE

- 50M devices ➙ $10/device

- 100K ➙ $5000/device



Source: International Business Strategies, Inc.

# Problem

- RISC-V enables hardware innovation by smaller organizations

  - No architecture licensing fee, no limits on customizing ISA

  - Open source software ecosystem with compilers, OS…

- But creating commercially competitive accelerator is challenging

  - Leading edge semiconductor design is extremely expensive

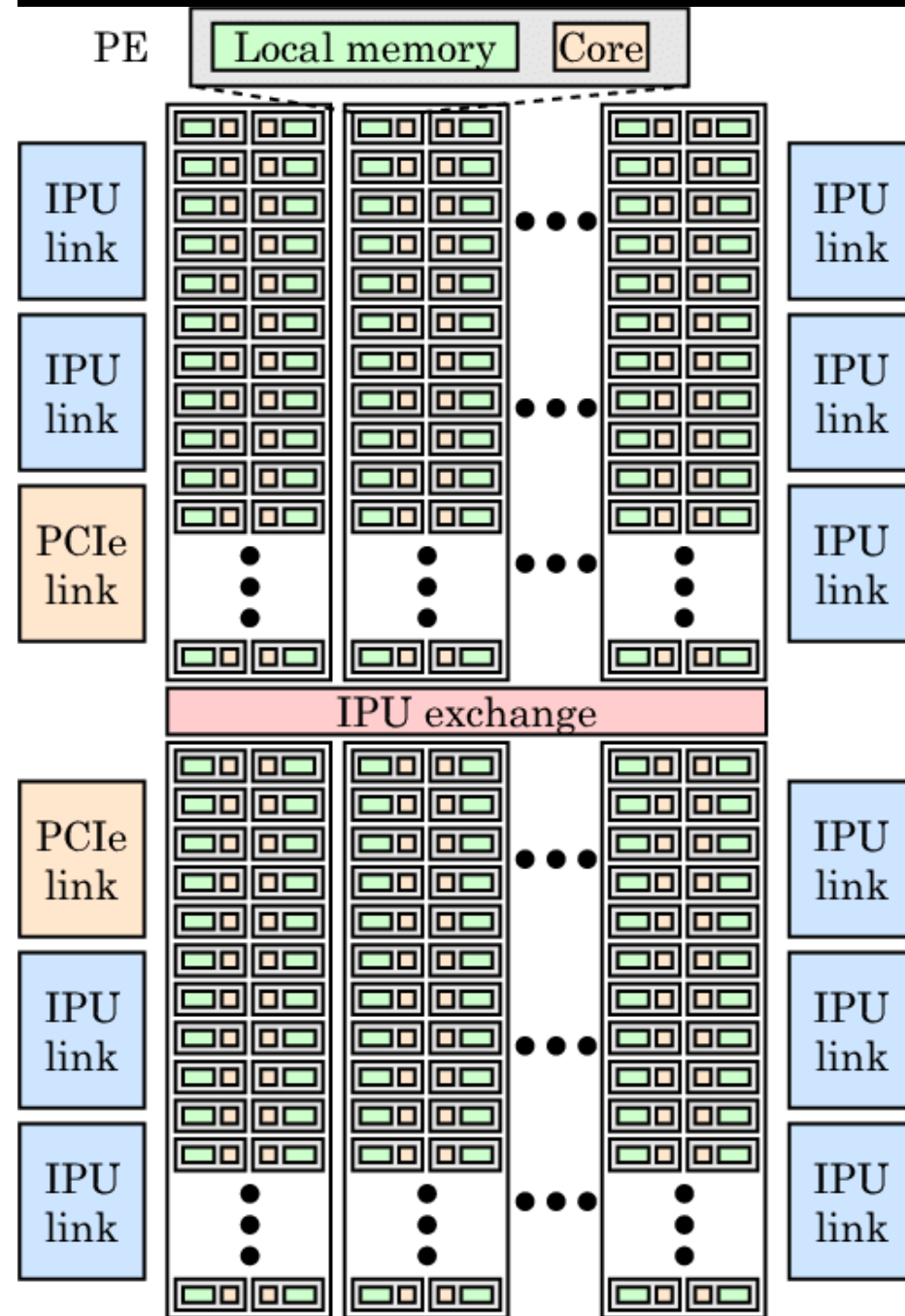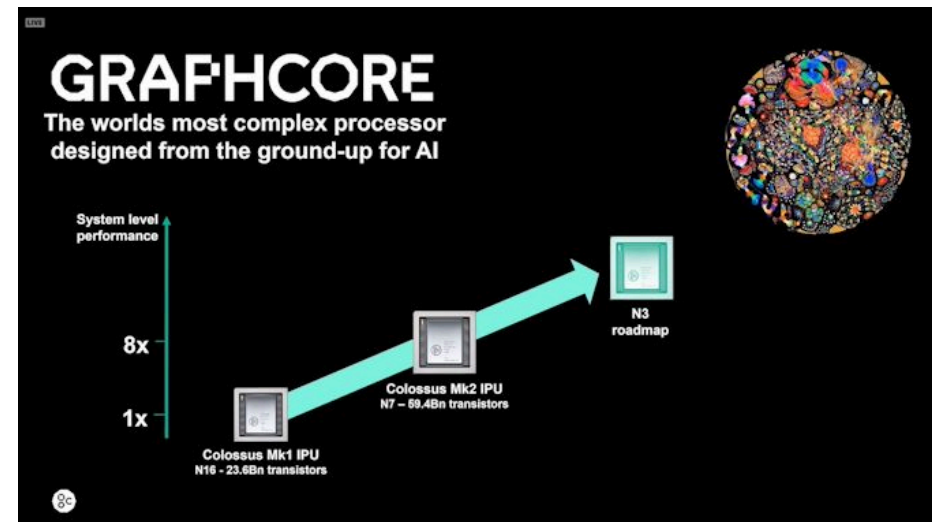  - Novel ideas take time to gain momentum and volume in market

# Solution

"End of Moore's Law" ➝ rapid advances in packaging technology

- Multiple chips on substrate (Open Chiplet Initiative)

- Chip/wafer stacking (TSMC 3DFabric™)

- Example:  GPU die with stacked HBM memories on substrate
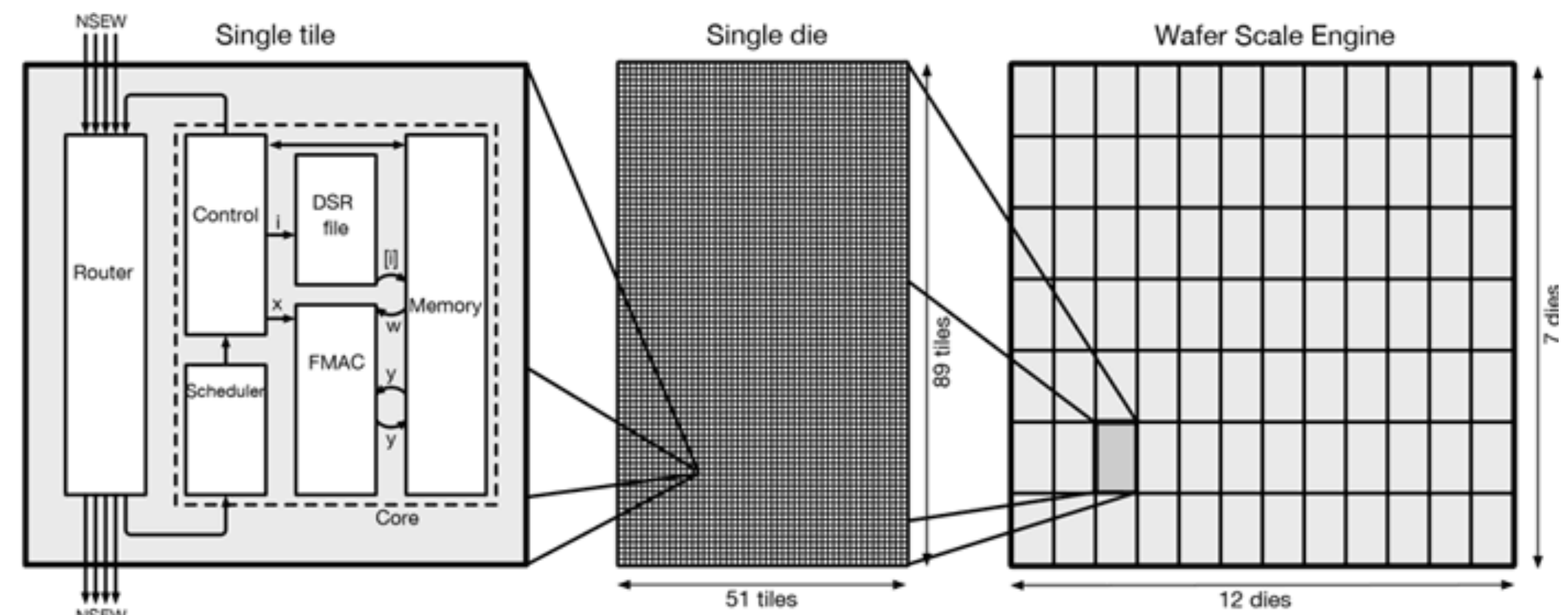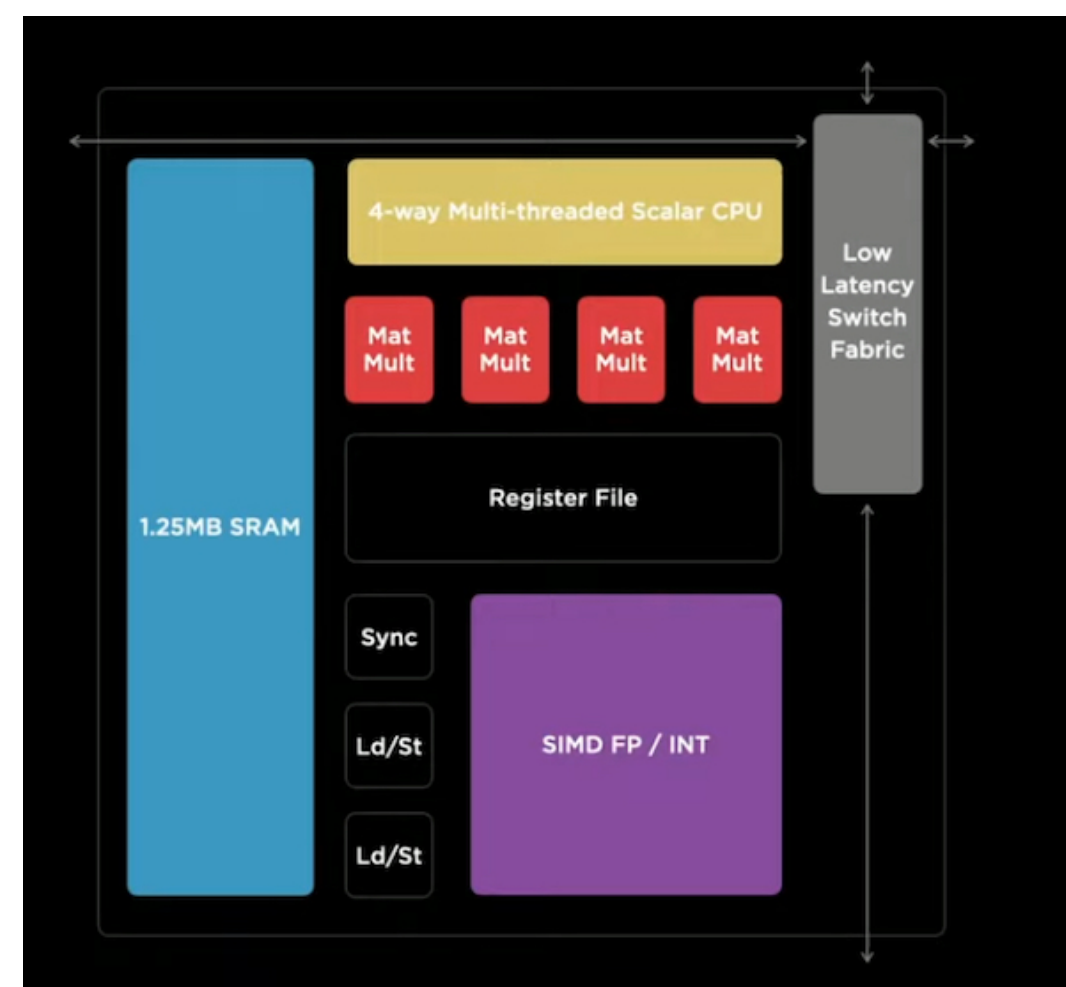
Mix and match advanced memory and mature logic technology

- Multiple accelerators share design cost of stacked memory

# Some Recent Accelerators

# In-, At-, Near-Memory Computing

**Global In-Memory Computing Market Is Expected to Reach USD 41.53 Billion by 2028 : Fior Markets**

**Large Amounts of Local Memory**



Compute-in-Memory IC

# Memory Chip Different from ASIC

FLASH, DDR, LPDDR, even HBM are high volume chips because

- Same chip used in multiple products

- Many chips in a single product

Memory cell density/flexibility tradeoff

- Dedicated memory fab → lowest cost, inflexible, standard parts

- Embedded logic fab → medium cost, customizable like ASIC

# Mix and Match Paradigm
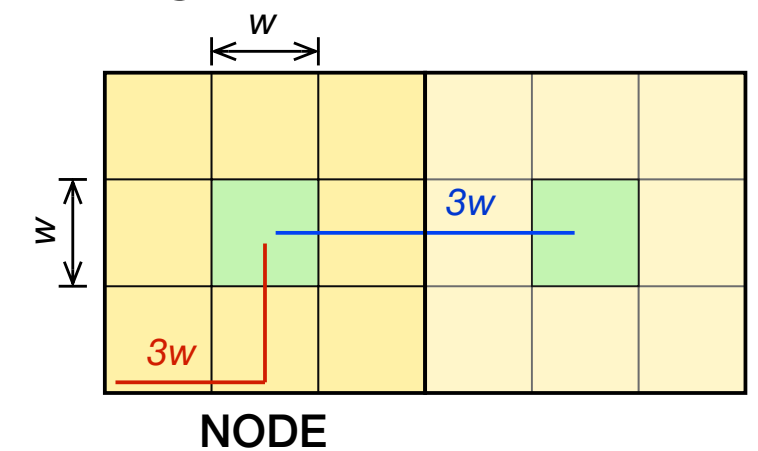
Monolithic in-memory computing chips use a lot of area for memory



- Expensive advanced logic process not fully utilized

- Long wires over memory area between logic islands waste energy

We propose separating memory and logic onto stacked dies



- Differently optimized process for memory and logic gates

- Contiguous logic reduces wire length, improves energy efficiency

# Wafer Stacking

- Same size dies, yield = one big die

- Low cost, good for memory chips

- Yield problem for logic die





Fig. 1 Process sequence for this experiment

# Chip Stacking

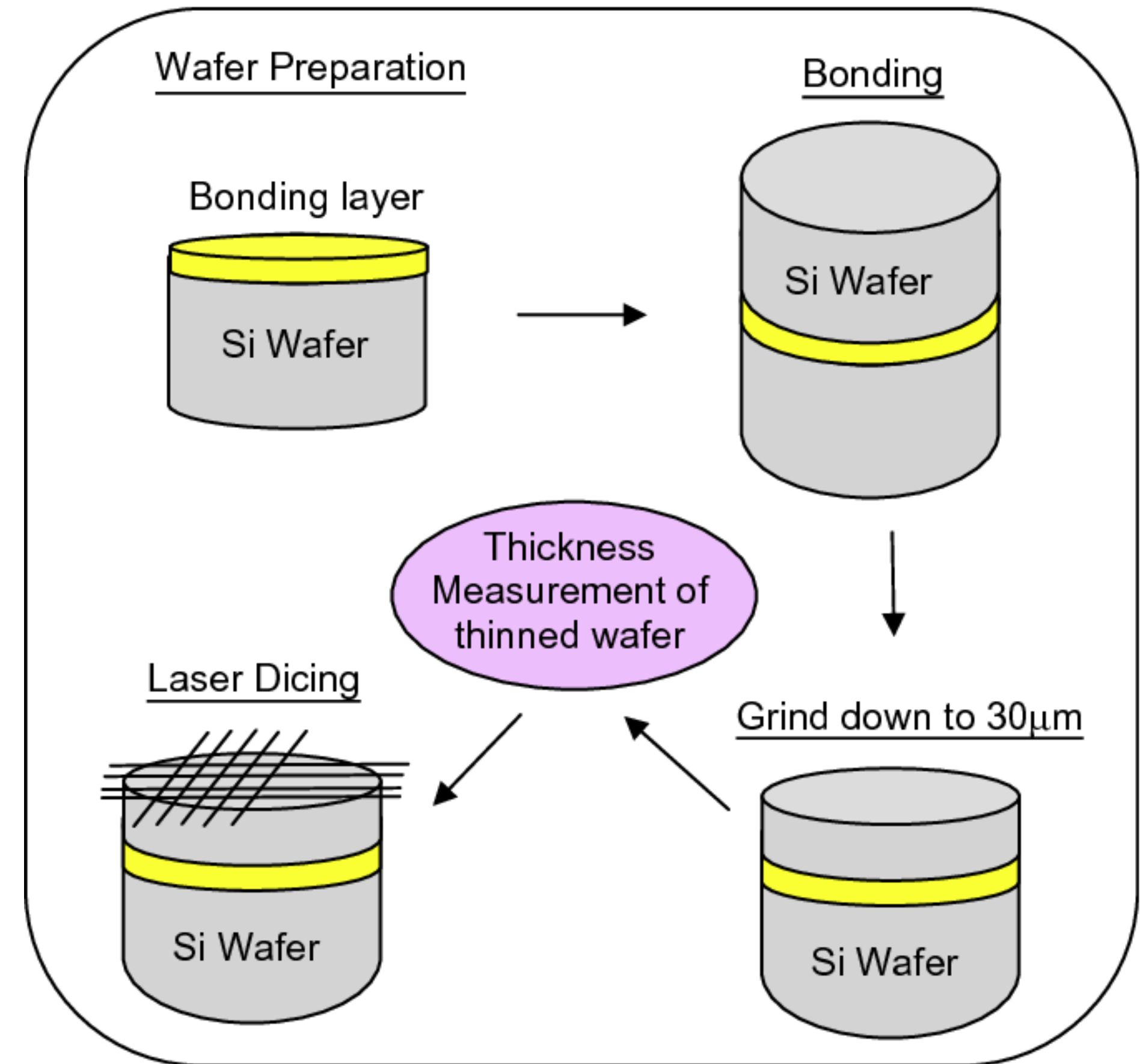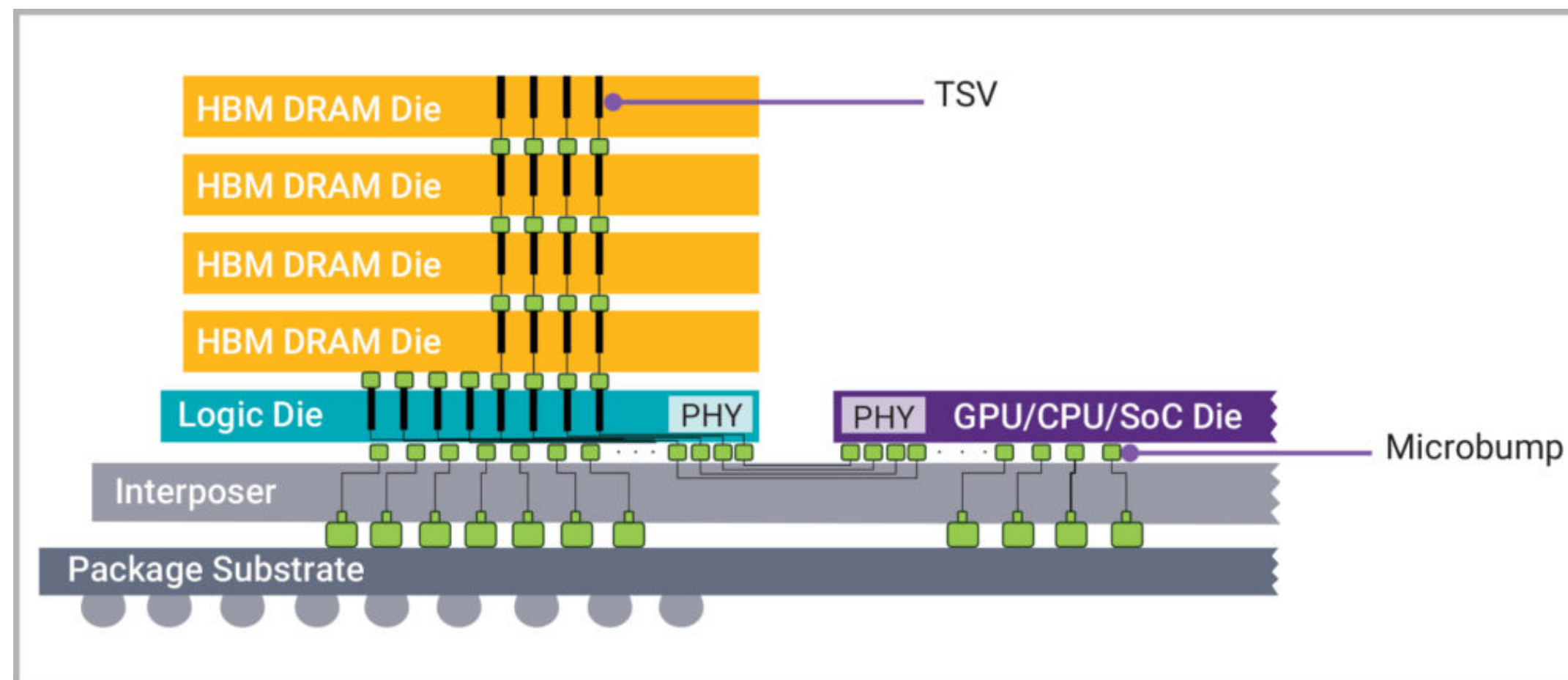- Accommodate different size dies, different processes

- Test before stacking (KGD) → large logic chip is possible

- Short vertical interconnect using TSV → good performance, power

- Higher cost

# Memory Stack

TSMC 7nm SRAM optimized process $\approx$ 2 MB/mm$^2$

- Chose 66 mm$^2$ stack area for reasonable yield

- 1 GB in 8 layers = 531 mm$^2$ silicon area

- 4 TB/s bandwidth using 10% TSV area overhead

Multiple stacks on large logic chip

- More capacity, bandwidth



Compute-in-Memory IC



SRAM STACK

LOGIC LAYER

0.28 mm

5.76 mm

11.52 mm

# Logic Chip

Mature technology

- GlobalFoundries 22FDX, TSMC 22ULP…

- Low CAD tools cost, many IP available

66 mm$^2$ stack ➞ 128 logic tiles each 0.52 mm$^2$

- 1.8 MGE (million gate equivalent) per tile
  ≈ 12 IEEE 64-bit FP multiply-add units

- 230 MGE per stack (1500 DP FPMAC)



Compute-in-Memory IC



MEMORY TILE

512 KB    TSV    512 KB

0.72 mm

0.72 mm

# Platform Architecture

8 SRAM stacks on 672 mm$^2$ logic chip (24×28 mm)

- 1024 tiles

- 1 GHz logic in 22 nm

- 2D torus NoC, 32-bit links

Tile is a complete computer

- Shared memory multicore architecture (SMP)

- Eight 32-bit memory accesses per cycle (32 GB/sec)

# Manufacturing Ecosystem

- TSMC has been developing stacked memory technology for some time

- CAD tools, TSV PHY already available

- Mixing mature logic chip is a business decision

## GLink-3D Application: SRAM on Top of Processor

- SRAM is separated from integrated chip and located on top of Processor
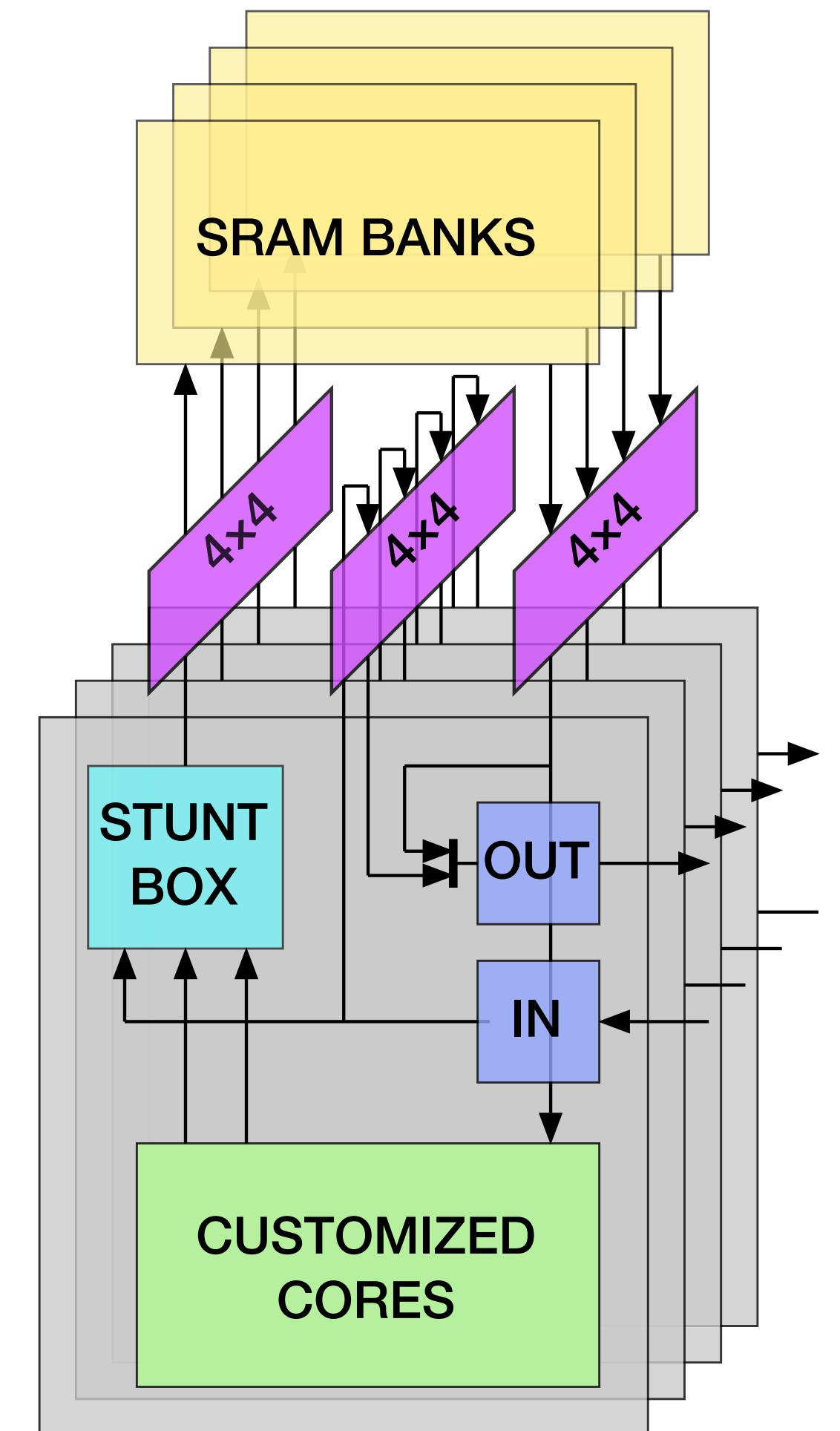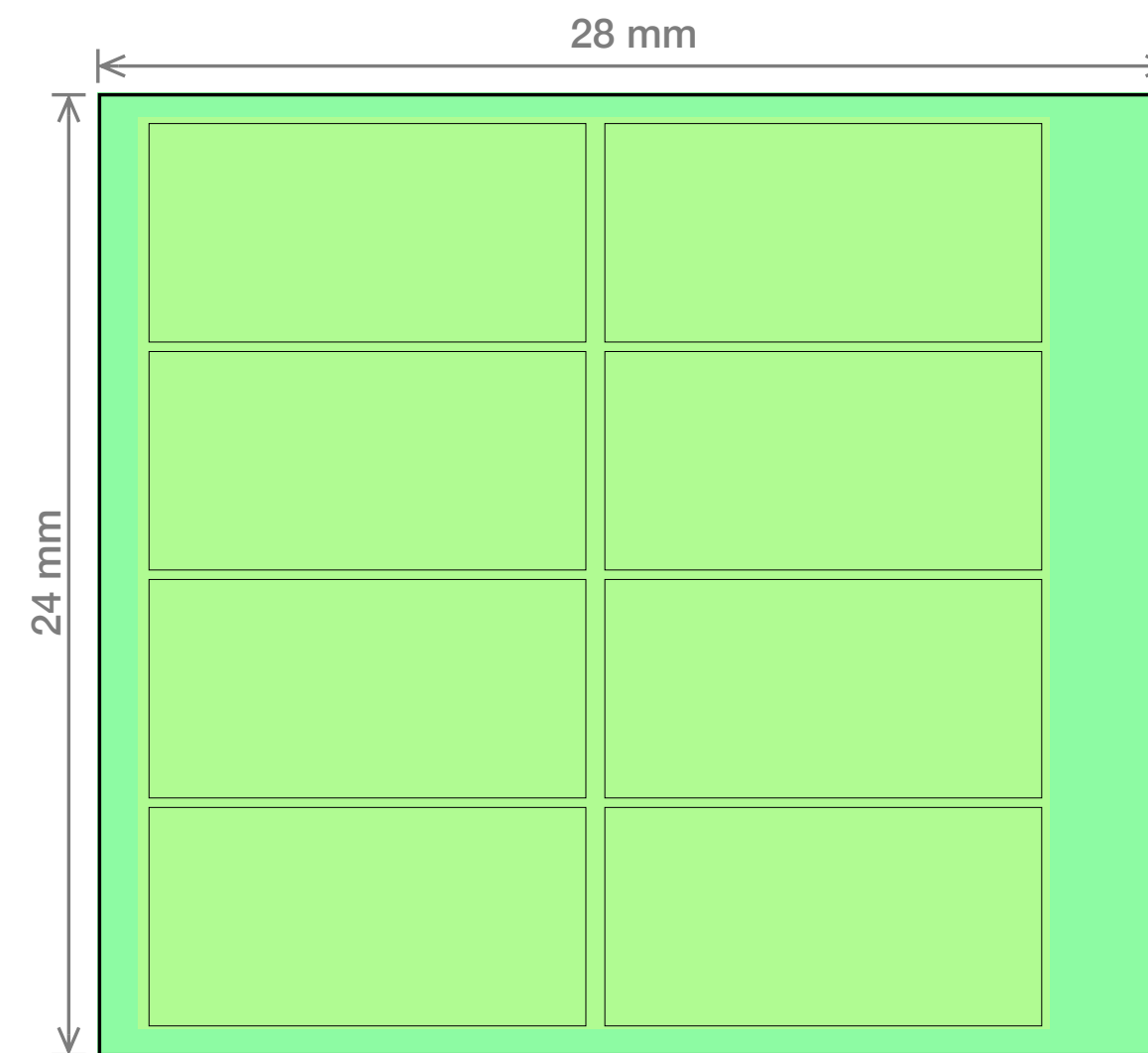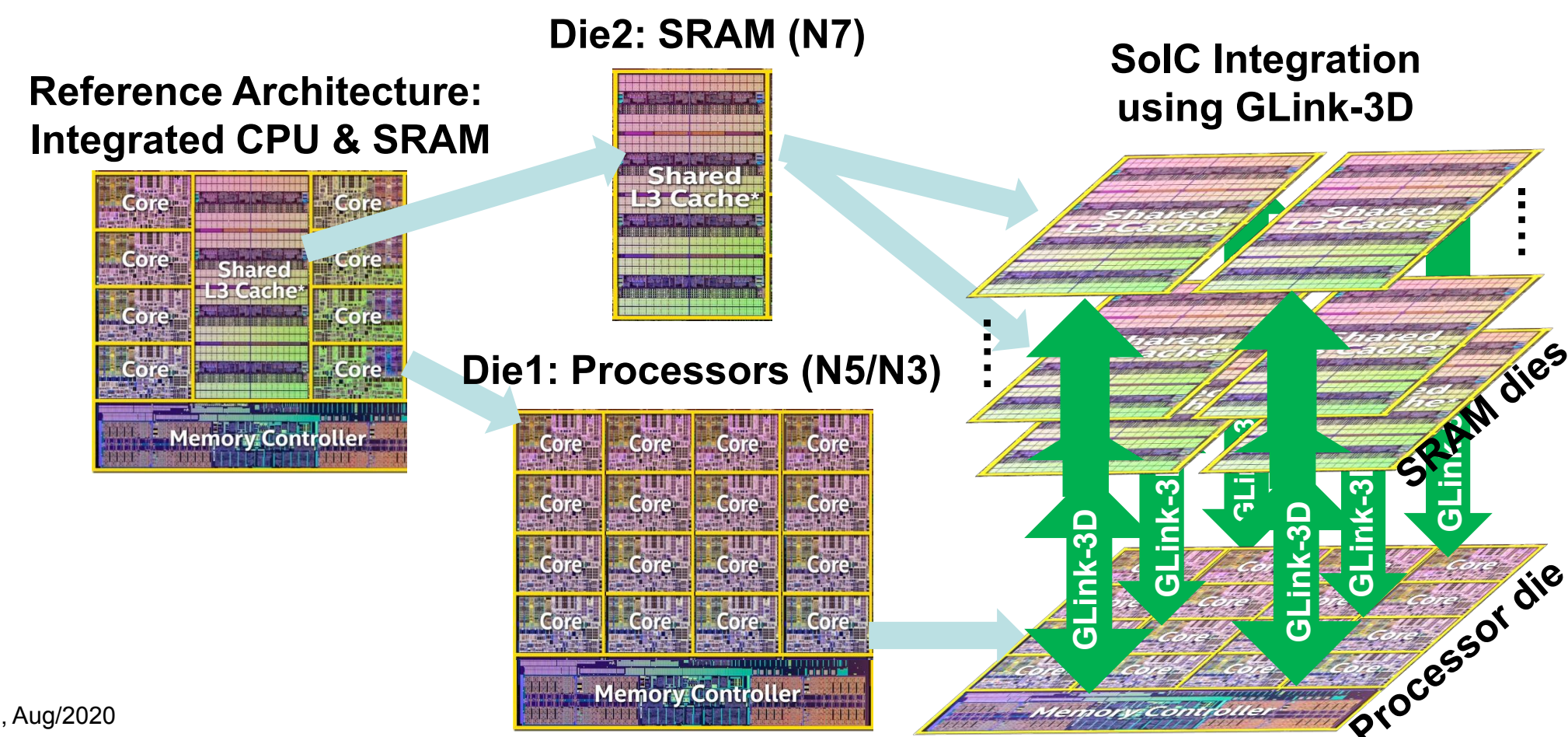- GLink-3D interface allows low area/power/latency connection

Reference Architecture: Integrated CPU & SRAM

Die2: SRAM (N7)

Shared L3 Cache

Die1: Processors (N5/N3)

SoIC Integration using GLink-3D

SRAM dies

Processor die

GLink-3D

Core / Shared L3 Cache / Memory Controller

Igor Elkanovich, Aug/2020

Copyright © 2020 GUC

Uncompromising Performance    P·12

**AMD Ryzen 7 5800X3D shipped out of factory, first CPU with 3D V-Cache**

AMD's first consumer CPU with its next-gen 3D V-Cache technology is now shipping, Ryzen 7 5800X3D will be in-hands this month.

**Anthony Garreffa**
**@anthony256**

PUBLISHED WED, MAR 2 2022 10:12 PM CST

# Agenda

1. Economics & technology

2. Accelerator architecture

3. Examples

4. Programming model & evaluation

5. Memory technologies

6. Conclusion

We illustrate with examples chosen for simplicity of explanation

Real commercial designs could do better

# HPC Example (SpMV++)
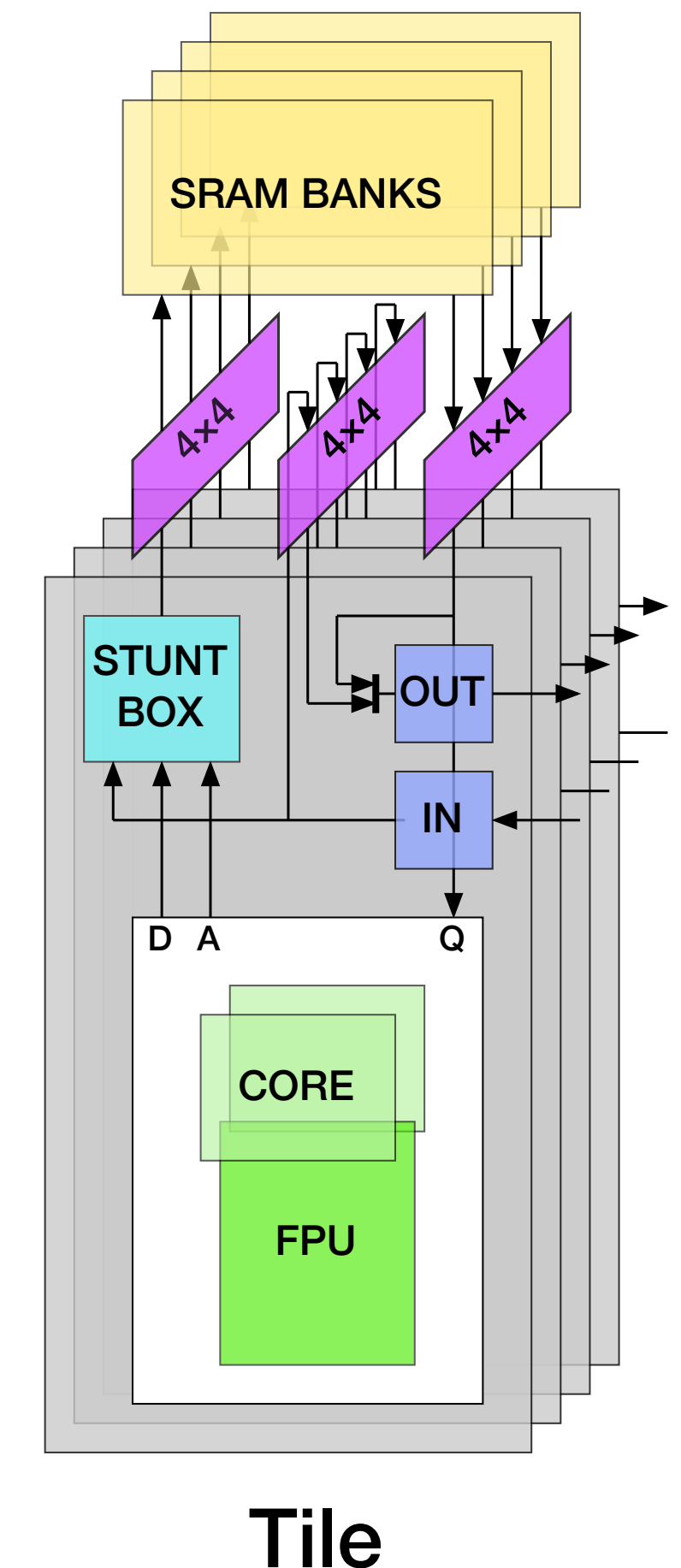
Specialized sparse matrix accelerator

- 8 GB, 4 TFLOPS (DP), 250 W PCIe card

- 8 bytes/FLOP memory bandwidth

- Logic 147.5 W (144 pJ per tile, 22nm)

- SRAM 25.6 W (25.6 pJ per tile, 7nm)
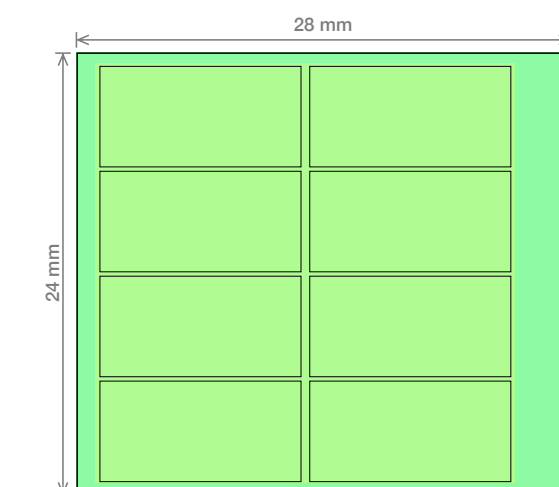
Evaluate using HPCG benchmark

Tile

# SpMV++ vs. GPU

NVIDIA A100 is today's premier HPC accelerator

- 7 nm

- 826 mm$^2$

- HBM2E

- 250 W

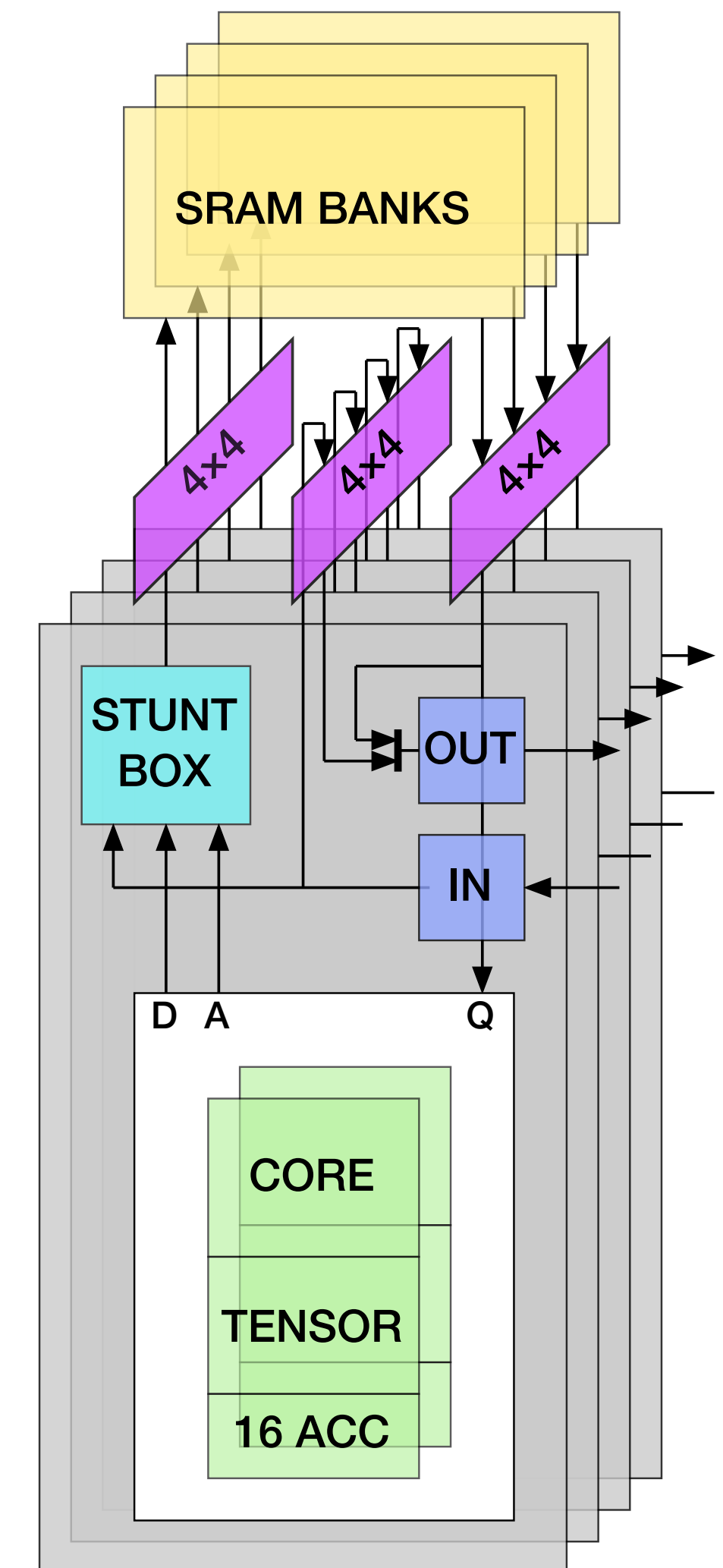| Technology | A100 | SpMV++ | Improve |
|---|---|---|---|
| Capacity (GB) | 80 | 8 | 0.1x |
| Peak TFLOPS (DP) | 19.5 | 4.1 | 0.21x |
| Bandwidth (TB/s) | 2 | 32.8 | 16x |
| Bytes/FLOP | 0.10 | 8.0 | 78x |
| HPCG TFLOPS | 0.227 | 2.9 | 13x |
| FPU Utilization* | 1.16% | 70% | 60x |
| Power (W) | 250 | 250 | 1x |
| GFLOPS/W | 0.91 | 11.5 | 13x |
| | *estimate for SpMV++ | | |



**NVIDIA A100 GPU**



**SpMV++**

# ML Example

Specialized machine learning accelerator

- ## 4-bit precision (logarithmic numbers)
  "Ultra-Low Precision 4-bit Training of Deep Neural Network," IBM Research

- ## 8 GB, 262 TOPS
  4-bit multiply, higher precision accumulate

- ## 500 MHz, low VDD
  0.45V vs. 0.85V for 1 GHz operation

- ## 3.9 TOPS/W energy efficiency



SRAM BANKS

4x4  4x4  4x4

STUNT BOX

OUT

IN

D A          Q

CORE

TENSOR

16 ACC

# Mainstream Programming Model
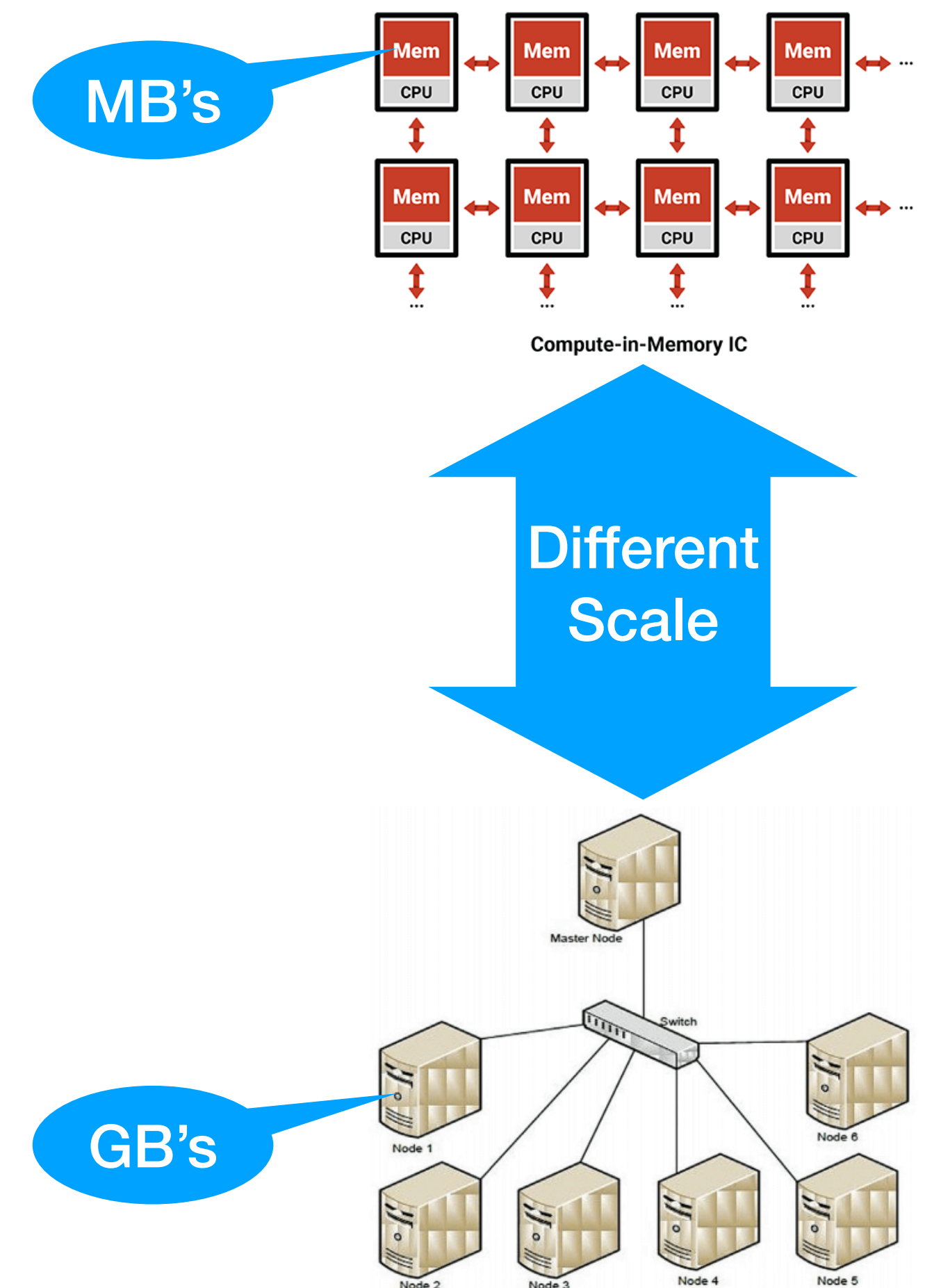
Cluster of computers ➡ array of tiles

- Multicores with shared coherent memory (SMP)

Network ➡ network on chip (NoC)

- Application specific interconnection topology

- Front-end computers ➡ host processors

Same application programming interface (API)

- Linux threads, sockets, RDMA…

MB's

Mem / CPU | Mem / CPU | Mem / CPU | Mem / CPU ...

Mem / CPU | Mem / CPU | Mem / CPU | Mem / CPU ...

Compute-in-Memory IC

Different Scale

Master Node

Switch

Node 1

Node 6

Node 2  Node 3  Node 4  Node 5

GB's

# Codesign Methodology

1. Develop algorithm on standard cluster of SMP servers

   • Standard Linux API, but respecting target node memory size

2. Simulate near-memory RISC-V SMP cluster

   • Parallel "thread per core, process per node" simulator (next slide)

3. Develop co-processor with custom instructions

   • Refine codesign and validate performance improvement

# Cavatools

***Caveat*** — RISC-V user-mode Linux virtual machine

- Thread-per-core execution-driven simulator, ≈100 MIPS

- Shared memory (eg. OpenMP), multiple nodes (eg. MPI)

Custom instruction definition

- Spec ➞ compiler intrinsic, asm, sim

```
fd += fs3 * load(rs1+imm)
```

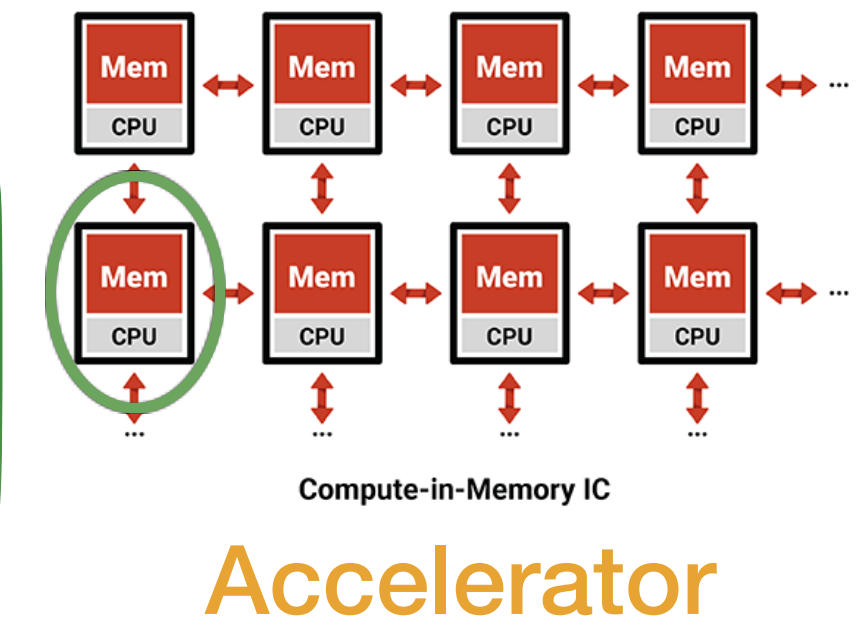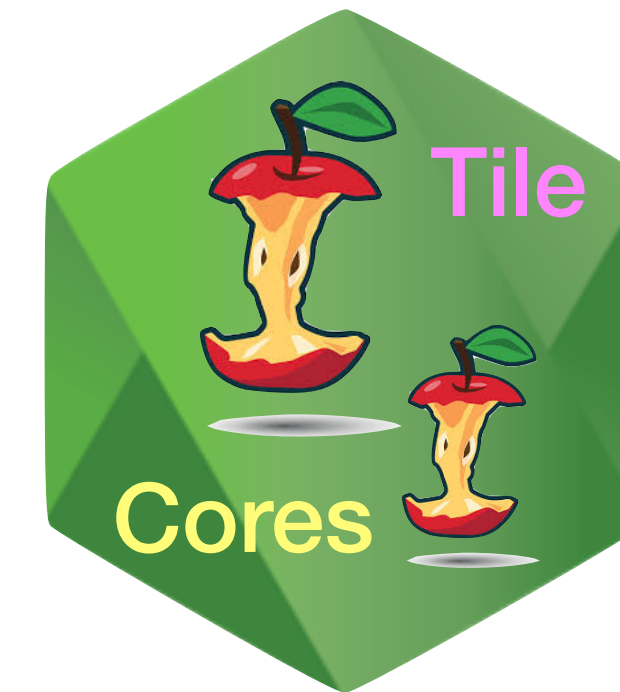| fs3 | imm[7:2] | op2 | rs1 | op3 | fd | opcode |
|-----|----------|-----|-----|-----|----|--------|

Open source (this work was partially supported by BSC)

> More details in ICS Conference presentation "Cavatools:  Parallel Architecture Simulator for RISC-V" (PDF)

# *Caveat* Simulation Paradigm

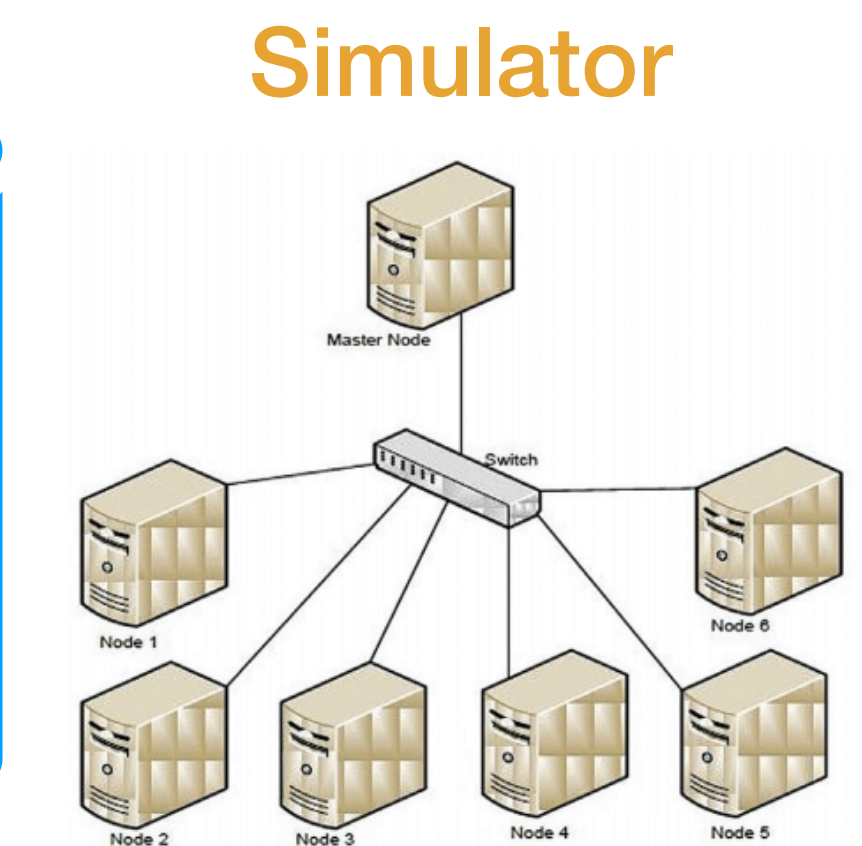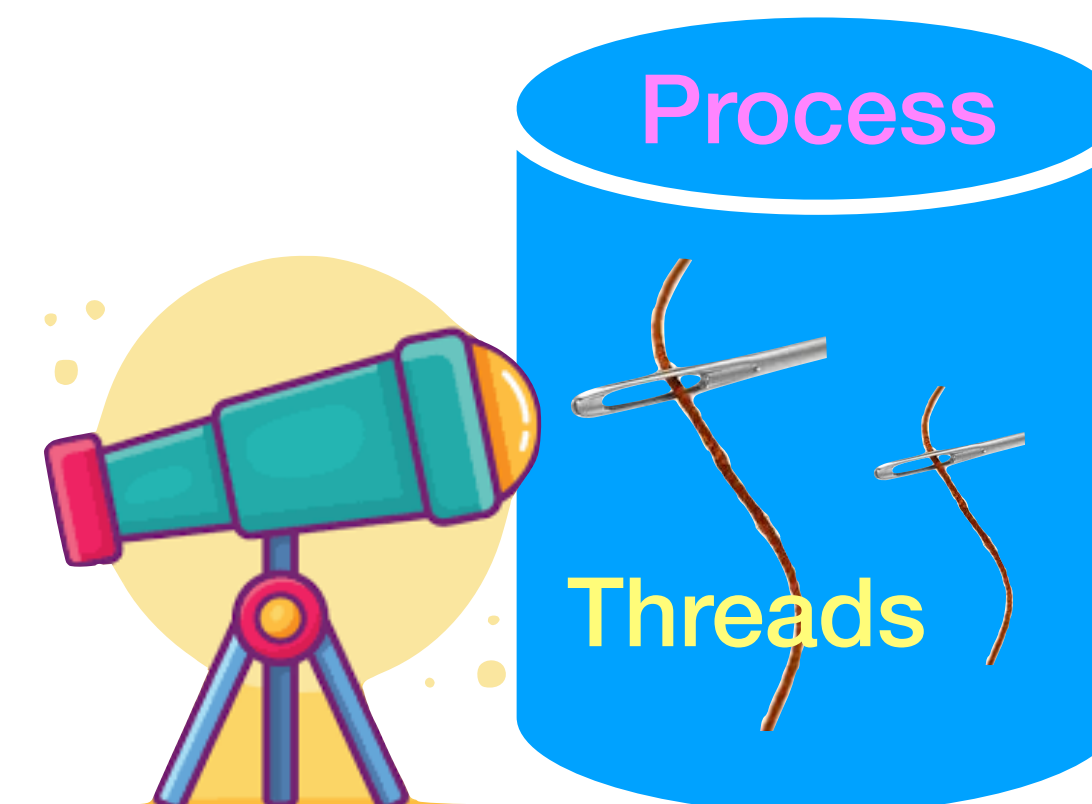Array of tiles ➜ Linux processes

- SMP cores ➜ Linux threads

- RISC-V AMO ➜ x86 CMPXCHG

Network on Chip ➜ Linux sockets

- Messages ➜ Linux read(), write()

*Erised* — realtime visualization



Tile

Cores

Compute-in-Memory IC

Accelerator



Simulator

Process

Threads

# *Erised* Performance Visualizer

Within a tile

- Pipeline stalls

- Instruction Buffer Misses

Across chip

- System calls (IPC)

- Message queues

# Agenda

1. Economics & technology

2. Accelerator architecture

3. Examples

4. Programming model & evaluation

5. Memory technologies

6. Conclusion

> Optimized SRAM process is less expensive than ASIC, but 7nm SRAM main memory is still quite expensive
>
> We need a path to more affordable memory technology

# Magnetoresistive Memory

Leverage mature process for low cost

- MTJ extra module in standard CMOS

- Similar to RF, Analog process modules

Attractive characteristics

- Read time, energy ≈ SRAM (but write ↑)

- Wear-out is no longer a problem



**SOT-MRAM To Challenge SRAM**

567 Shares    f 54    🐦 37    in 444    <

*Spin-orbit torque memory adds endurance and faster write speeds, but displacing existing memories is still not easy.*

JANUARY 13TH, 2022 - BY: **BRYON MOYER**

Conferences > 2018 IEEE International Solid... ❓

**A 1Mb 28nm STT-MRAM with 2.8ns read access time at 1.2V VDD using single-cap offset-cancelled sense amplifier and in-situ self-write-termination**

# Reducing Memory Cost

SRAM 7nm ≈ 2 MB/mm² (TSMC)

- Scaling less than logic gates

MRAM 28nm > 1 MB/mm² today

- Nonvolatile → mobile devices

- 3D (like XPoint™) in future

**Samsung Demonstrates the World's First MRAM Based In-Memory Computing**

Audio 🔇 Sh

| Technology | DRAM | MRAM | SRAM |
|---|---|---|---|
| Cell Size (F²) | 6 | 16 | 120 |
| Relative Cost | $1.00 | $2.67 | $20.00 |



The Worldwide MRAM Industry is Expected to Reach $4.9 Billion by 2026

| Everspin 4th Gen. MRAM (DDR4) |
|---|
| EMD4E001G DDR4 ST-MRAM |
| 105.1 mm² (12.29 mm x 8.55 mm) |
| 28 nm |
| 1,024 Mb (1 Gb) |
| 9.75 Mb/mm² |
| 0.0396 µm² |
| 110 nm / 180 nm |
| pMTJ |
| Between M3 and M4 |
| 7 |

# Summary

We provide RISC-V platform architecture
High bandwidth, low latency SRAM memory
Standard Linux programming environment
Simulation tools to validate performance

SRAM BANKS

4x4  4x4  4x4

STUNT BOX

OUT

IN

CUSTOMIZED CORES

You design arithmetic, interconnect, software for your application

# Conclusion

- RISC-V enables hardware innovation but development cost limits creativity

- Proposed paradigm for near-memory accelerator using state-of-the-art memory with mature logic technology

- Enable commercially competitive accelerators based on novel ideas by smaller organizations with limited resources

# Thank You

**Abstract:** The RISC-V architecture has opened new opportunities for many people to innovate in computer design. However to design a chip that can compete in the marketplace against veteran industry computer designers with their vast resources is still a formidable challenge. We propose a solution for specialized accelerators with near-memory processing architectures. We observe the critical technology is the embedded memory because it consumes most of the silicon area and determines the power/bandwidth of the chip. If instead memory is stacked on top of the logic chip, then a less dense, lower cost mature technology can be used for the logic. Communication wire power will be lower because the through-silicon via (TSV) interconnect traverses a much smaller distance, offsetting the lower power efficiency of mature logic technology. Design cost of the re-useable hi-tech memory chip is amortized across multiple accelerators. We believe this approach can help smaller organizations with limited resources design commercially competitive novel accelerators.

**Bio:** Peter Hsu received his Ph.D. from University of Illinois Urbana-Champaige. He started work at IBM T. J. Watson Research Center on the 801 Project. He joined SGI in 1990 as architect of MIPS R8000 TFP microprocessor; the chip powered fifty TOP500 supercomputers in 1994. Peter co-founded ArtX in 1997 to develop the Nintendo GameCube video game console. He joined Oracle Labs in 2011 as Architect and built a fifty thousand core parallel SQL database accelerator. Peter moved to Europe in 2018 and was visiting professor at EPFL University in Switzerland, then senior researcher at the Barcelona Supercomputing Center. In 2022 he started a consulting company in Spain, Peter Hsu & Associates, S.L.U.