

# 3D Long-Term Memories (LtRAMs) for Model-Resident Inference Architectures

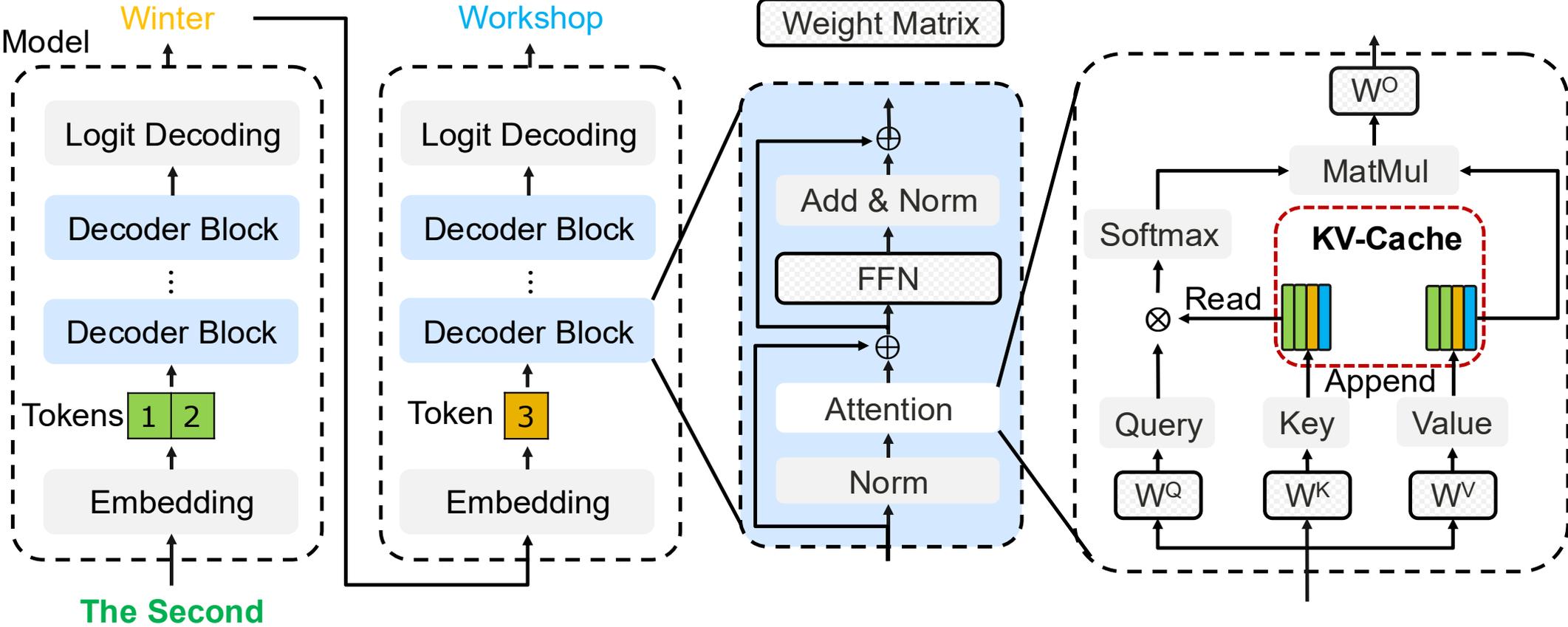
## A Framework for Quantitative Analysis

Samuel A. Dayo<sup>1</sup>, Colin H. Drewes<sup>1</sup>, Alex Aranburu\*, Victor J.B. Jung<sup>2</sup>, Shridhar Mukund\*  
Prof. Robert M. Radway<sup>3</sup>, Prof. Christos Kozyrakis<sup>1,4</sup>, Prof. Luca Benini<sup>2</sup>, Prof. H.-S Philip Wong<sup>1</sup>,  
Prof. Subhasish Mitra<sup>1</sup>

<sup>1</sup>Stanford University <sup>2</sup>ETH Zurich <sup>3</sup>University of Pennsylvania <sup>4</sup>NVIDIA Research \*EMD Electronics

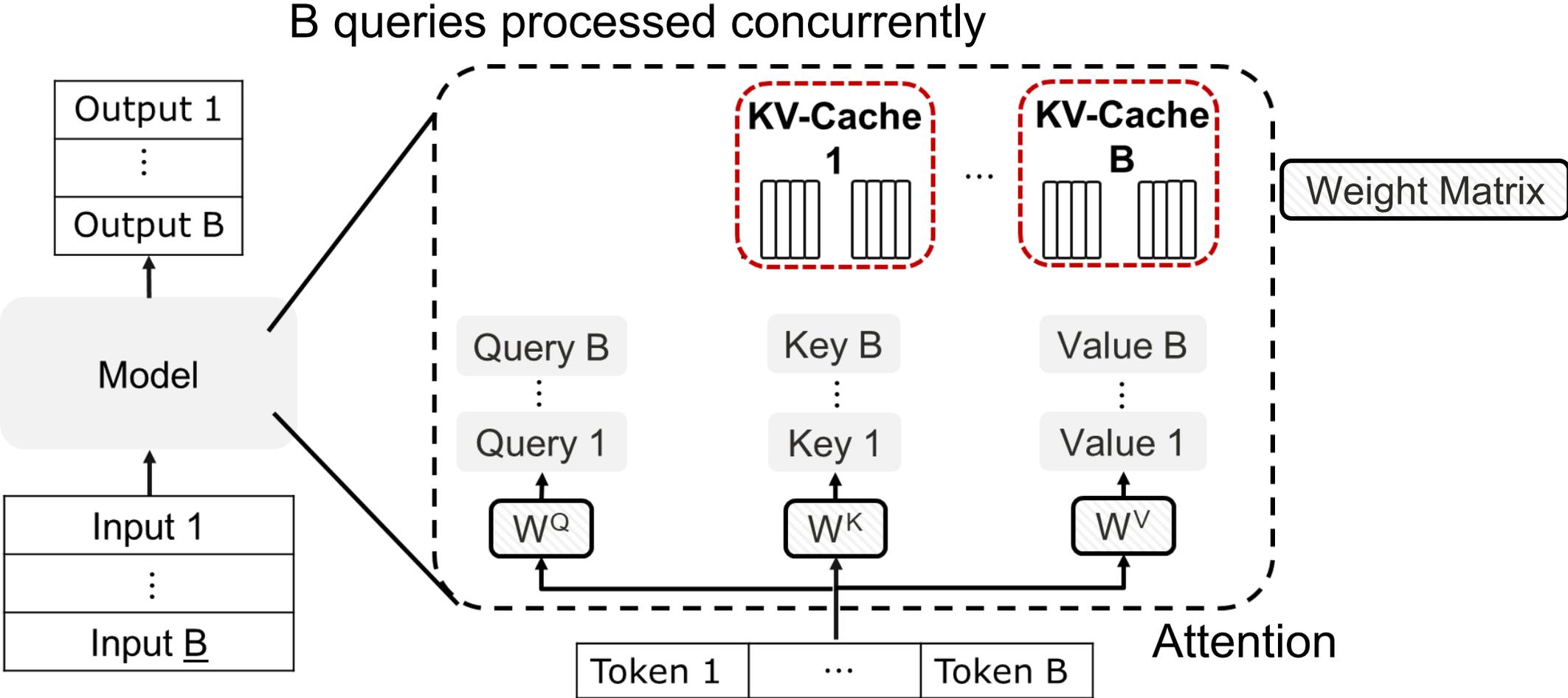
# Background: LLM Inference

Weights and KV-Cache read every pass



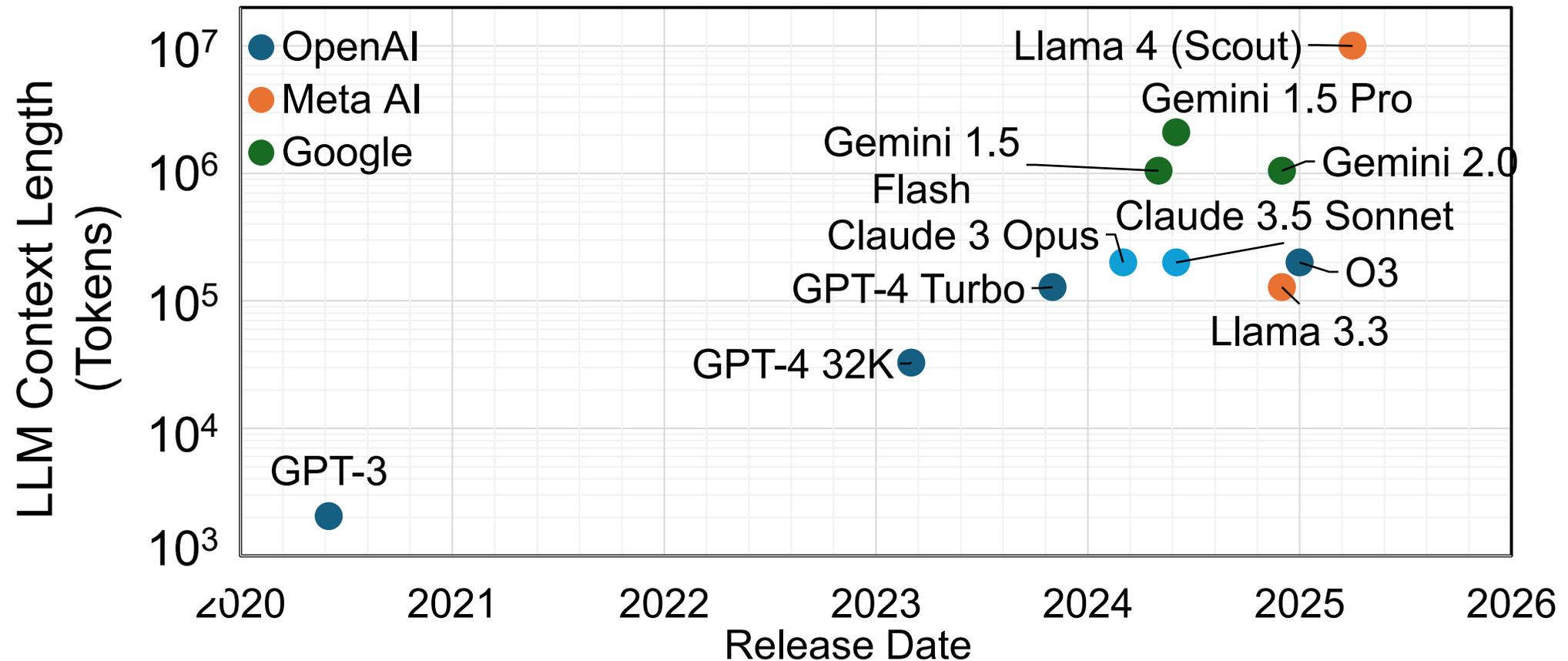
KV-Cache grows linearly with total tokens (context length)

# Background: Batch Inference



Weight access but not KV\$ access amortized

# Background: Explosive Growth of Context Lengths



Larger context  $\rightarrow$  more weight + KV-Cache reads  
Memory needs: high-bandwidth, low-energy access

# Model-Resident Architectures

---

## Model-Streaming

Weight/KV-cache materialized during execution

Continuous re-allocation of weight/KV-cache

Memory is remote from compute node

## Model-Resident

Weight/KV-cache allocated during model load

Allocation persist through model inference

Memory co-located with compute node

Minimize: Weight/KV-Cache data movement

# Long-Term RAM (LtRAM)

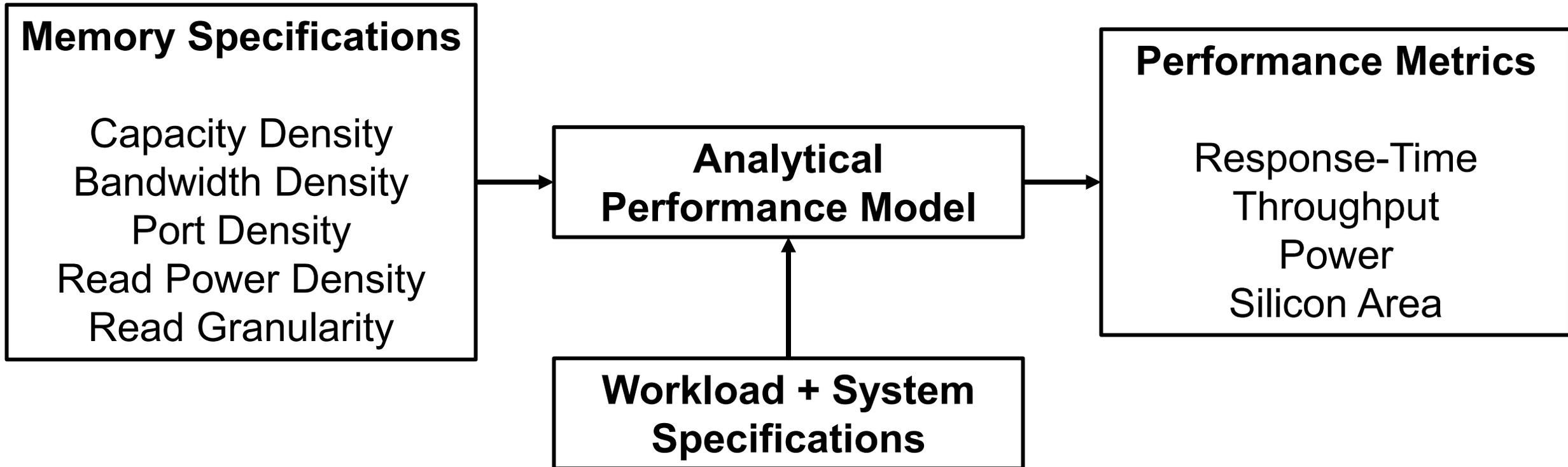
---

	<b>Long-Term RAM</b>	<b>Short-Term RAM</b>
Desired Properties	High density Low read energy	High density Low read/write energy
Acceptable Compromise	Higher write energy Higher write latency Limited write endurance	Retention Refresh overhead
Use Case	Read-mostly weights Read-mostly KV-Cache	Read/write caches Read/write scratchpads

Read-optimized memory weights and KV-Cache

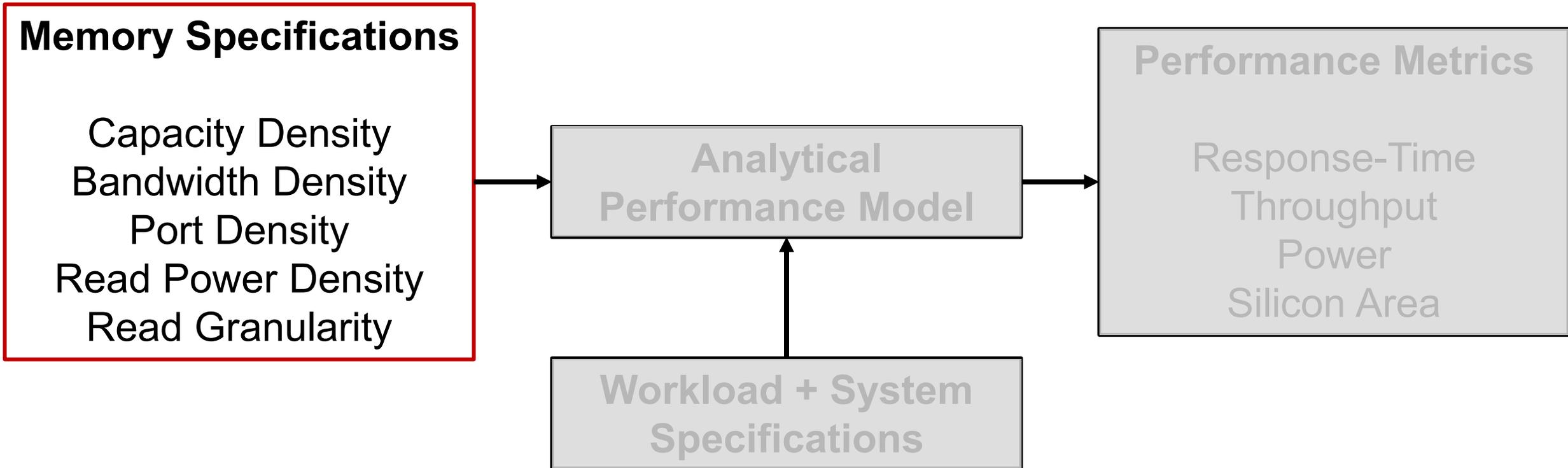
# Quantitative Performance Evaluation Framework

---



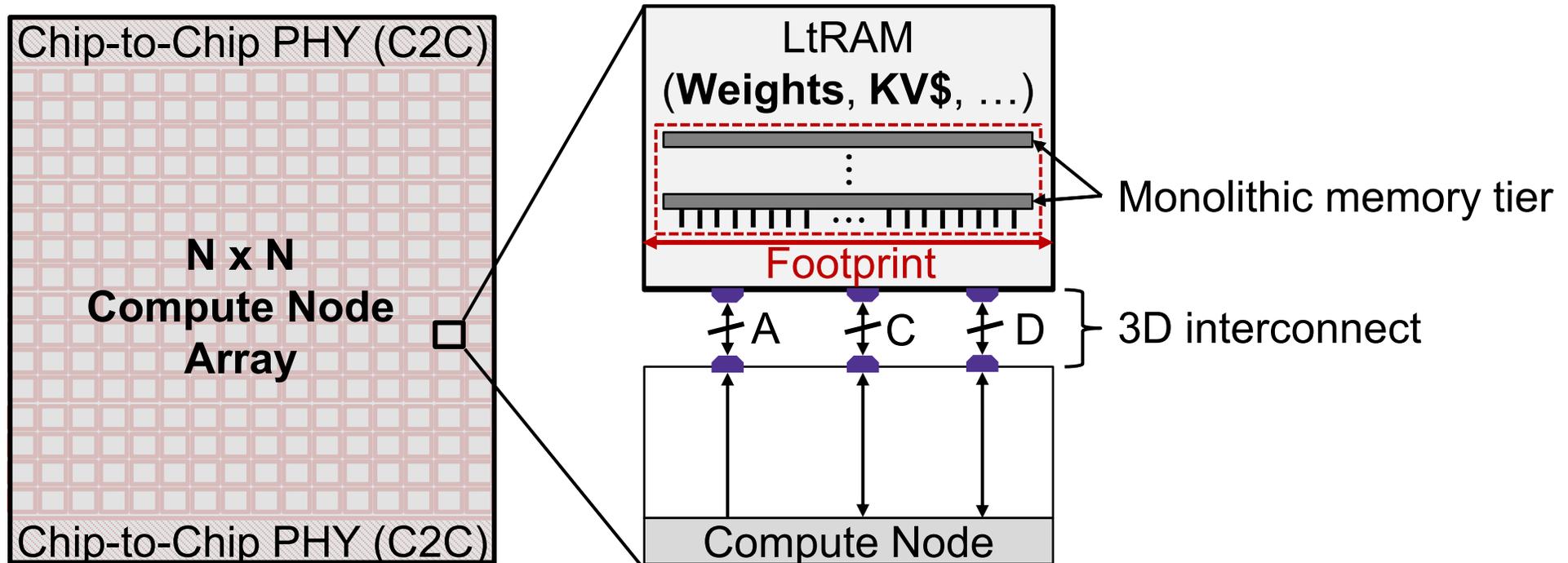
# Quantitative Performance Evaluation Framework

---



# Memory Capacity Density

Memory capacity per unit area (GB/mm<sup>2</sup>)

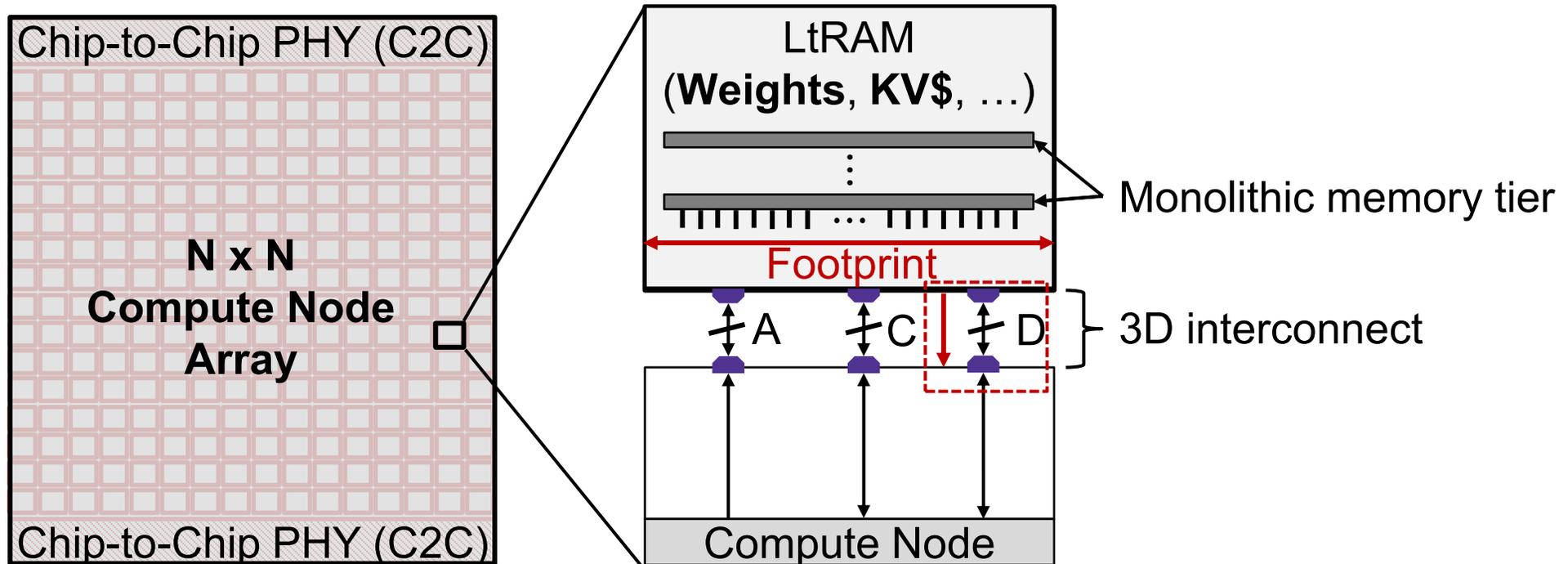


A: Address width; C: Control width; D: Data width

$$(\text{Per-layer capacity} \times \text{Monolithic Layers}) / \text{Footprint}$$

# Memory Bandwidth Density

Peak data transfer throughput per unit area (TB/s/mm<sup>2</sup>)

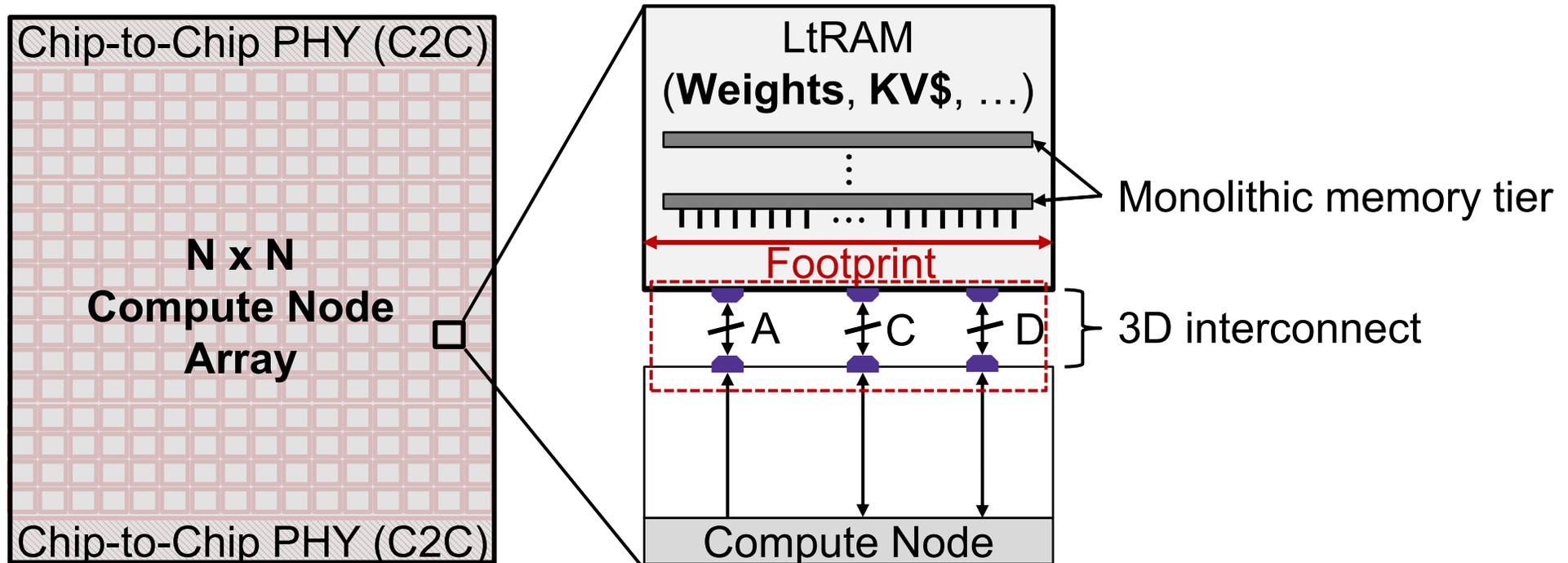


A: Address width; C: Control width; D: Data width

$$(\text{Data width} \times \text{clock frequency}) / \text{Footprint}$$

# Memory Port Density

Independent memory access ports per unit area (ports/mm<sup>2</sup>)

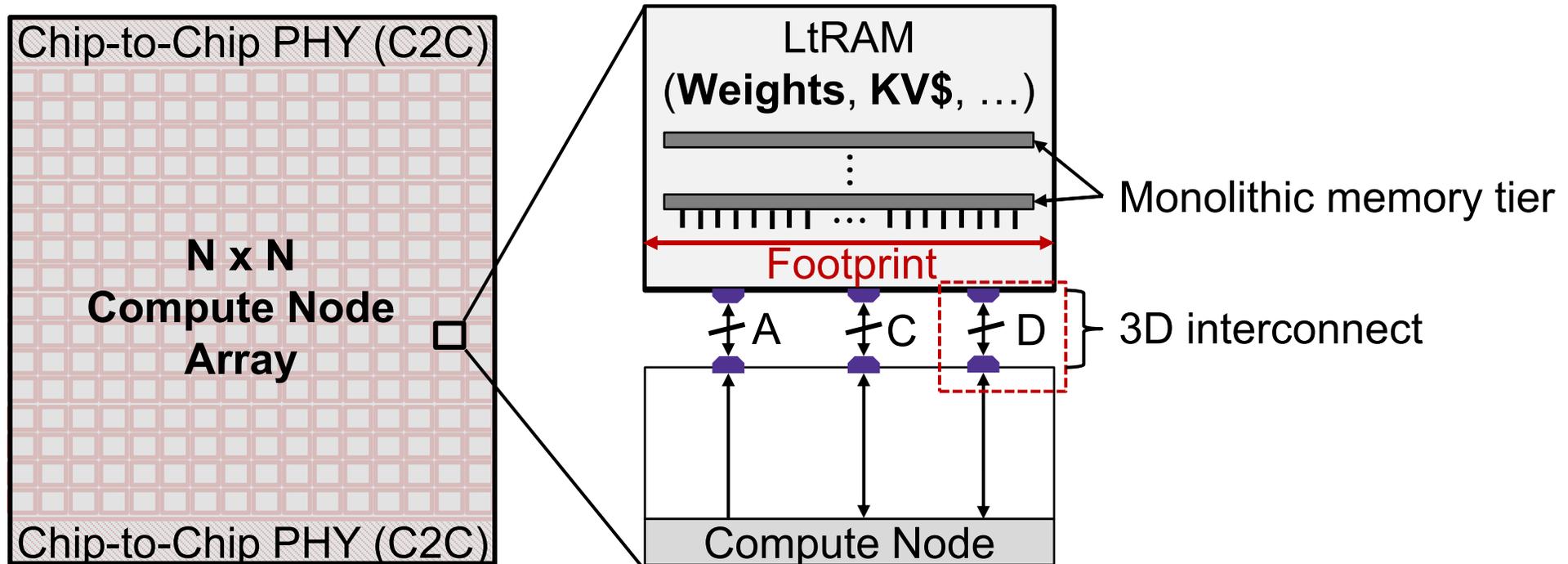


A: Address width; C: Control width; D: Data width

1/ Footprint

# Memory Read Power Density

Memory read power per unit area ( $\text{W}/\text{mm}^2$ )

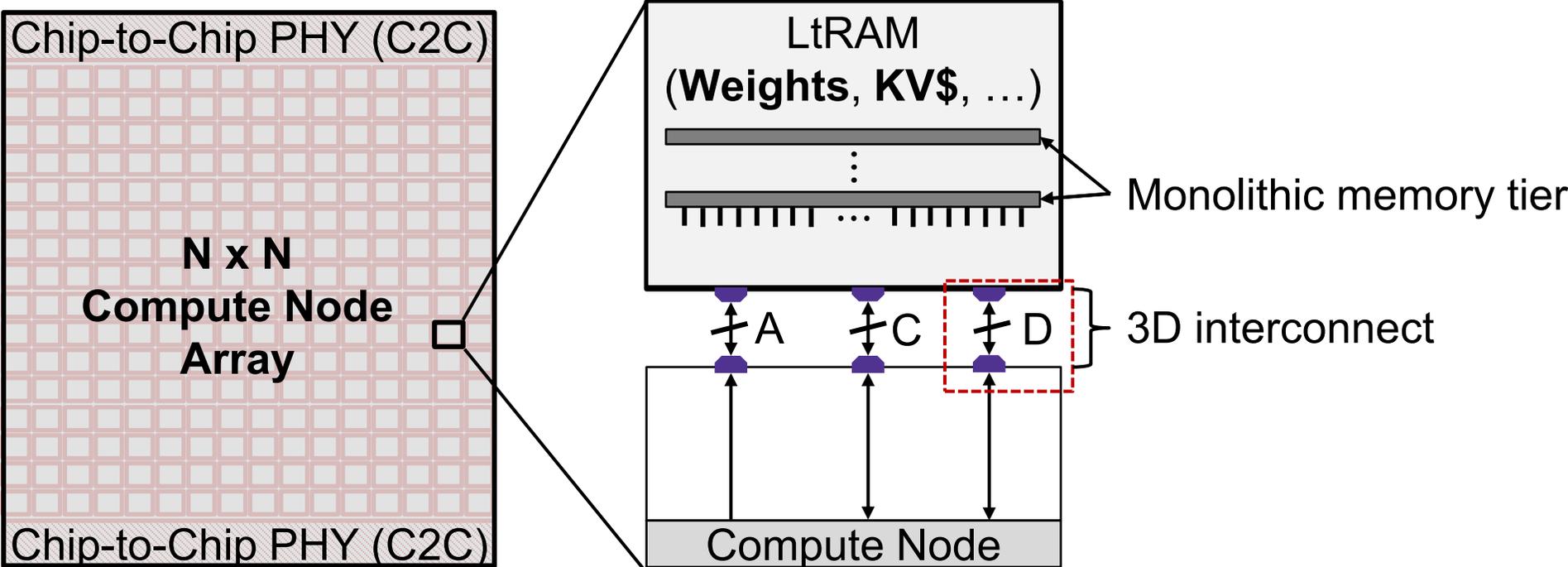


A: Address width; C: Control width; D: Data width

$$(\text{Data width} \times \text{clock frequency} \times \text{energy-per-bit}) / \text{Footprint}$$

# Memory Read Granularity

Random access size (bytes/access)

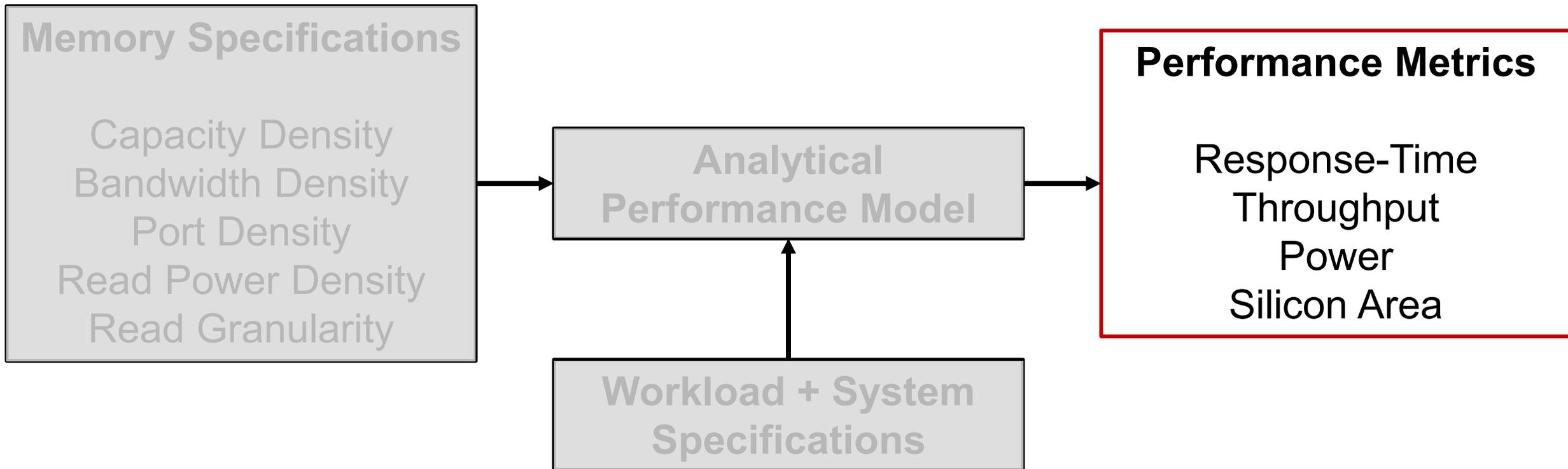


A: Address width; C: Control width; D: Data width

Data width

# Quantitative Performance Evaluation Framework

---



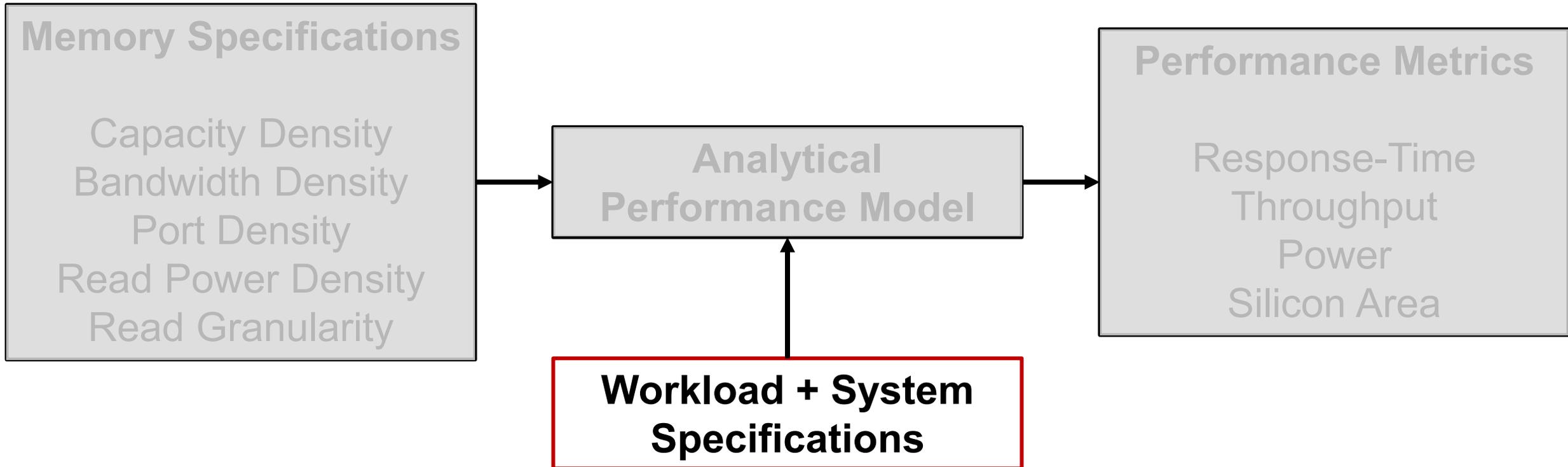
# Review: System Performance Metrics

---

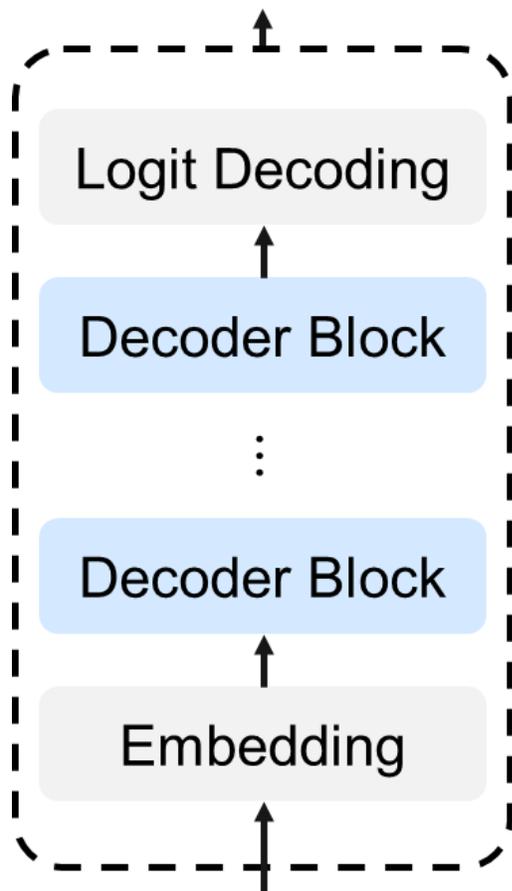
System Performance Metric	Definition	Better
User Query Response Time (s)	End-to-end latency per user query	↓
System Token Response Time (s)	Average time per generated token	↓
User Tokens Throughput (token/s)	Token generation rate per user	↑
System Tokens Throughput(token/s)	Aggregate token generation rate across all users	↑
System Power-Efficiency (token/s/watt)	Represents energy cost to generate token (OPEX)	↑
Silicon Area-Efficiency (token/s/mm <sup>2</sup> )	Silicon area utilization efficiency (CAPEX)	↑

# Quantitative Performance Evaluation Framework

---



# Workload Parameterization: LLM Decode



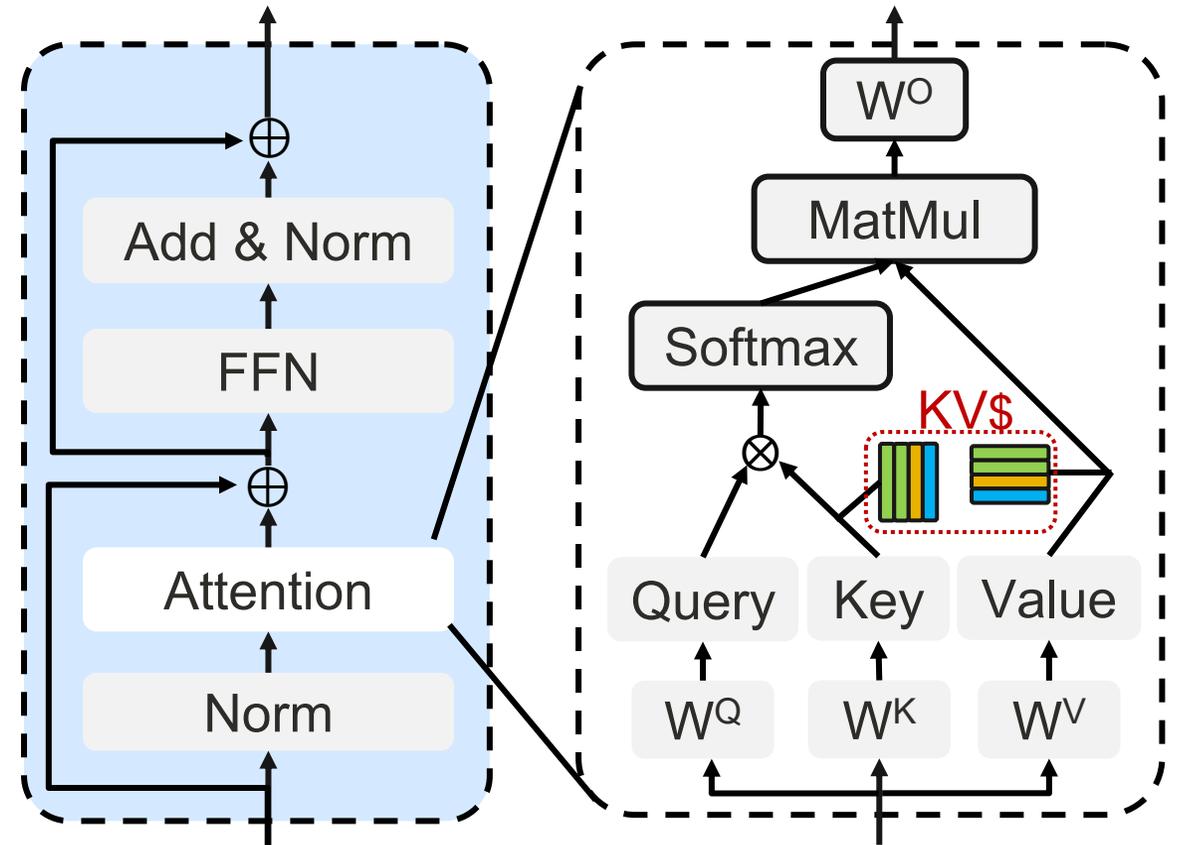
Component	Compute FLOPS(%)	Memory Capacity (%)
Embedding	0.48	1.28
Decoder Blocks	<b>99.47</b>	<b>97.43</b>
Output Logits	0.04	1.28

Dominates workload compute and memory capacity

# Workload Parameterization

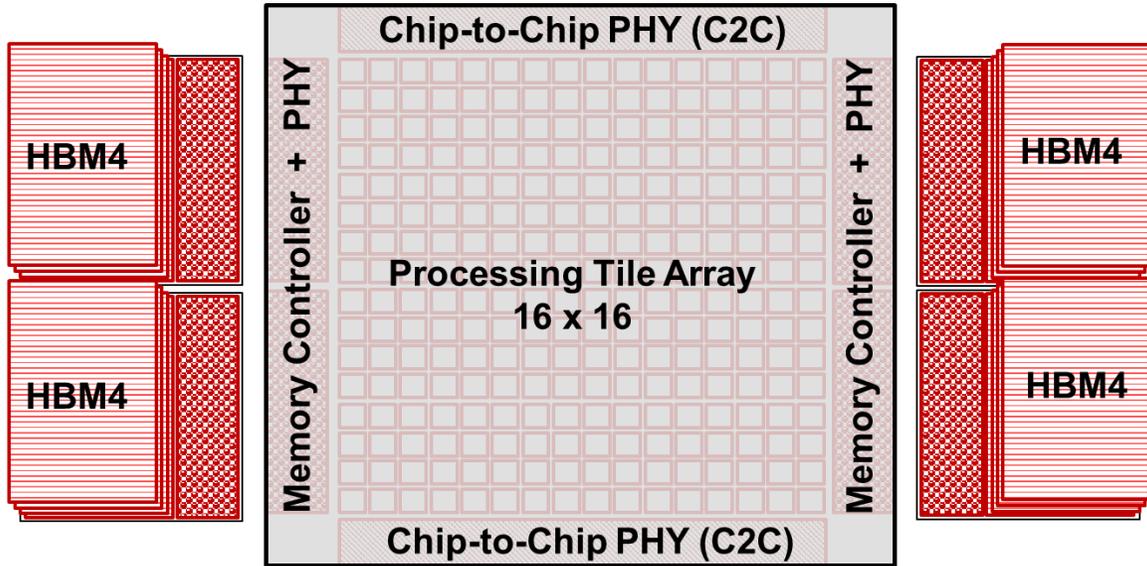
Focus only on performance significant components

Variables	Variable
Emb. Dimension	$d_{model}$
Head Dimension	$d_h$
# Query Heads	$H_q$
# KV Heads	$H_{kv}$
# Decoder Layers	$L$
# Experts	$E$
Top-K (Active Experts)	$K$
Expert Dimension	$d_{ffn}$
Batched Users	$B$
Context Length	$C$

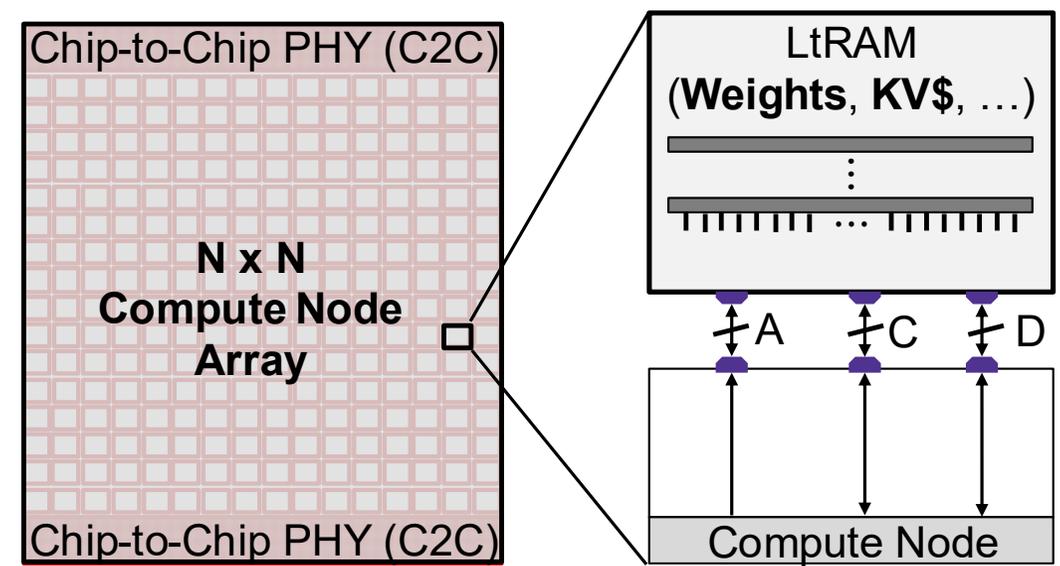


# System Parameterization Assumptions

## Baseline



## Model-Resident



Tile-based spatiotemporal computation

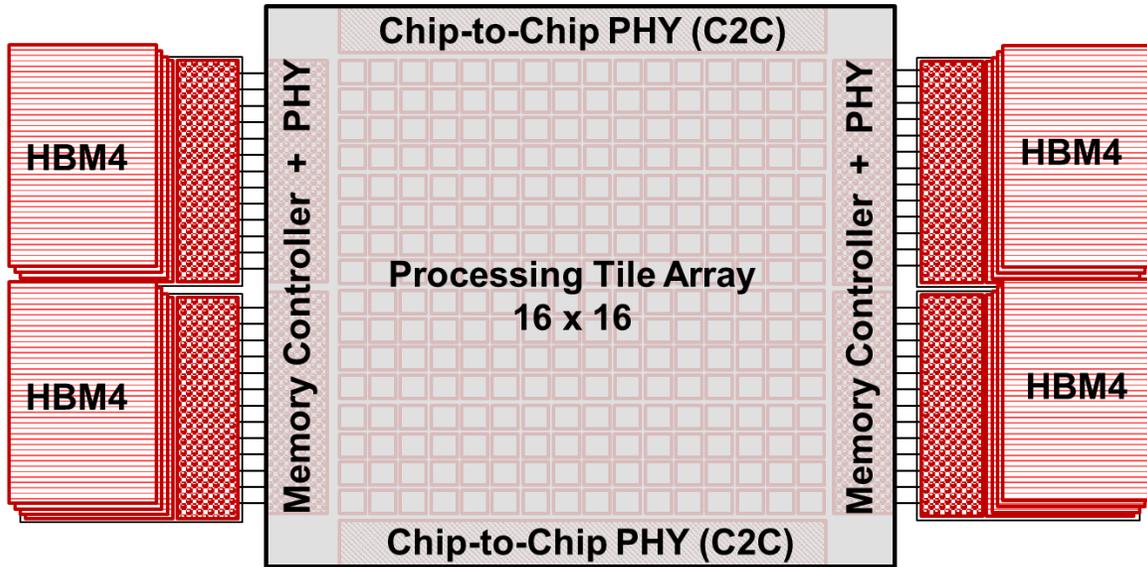
2D Mesh NoC with dimension-ordered routing

Inter-chip pipeline parallelism

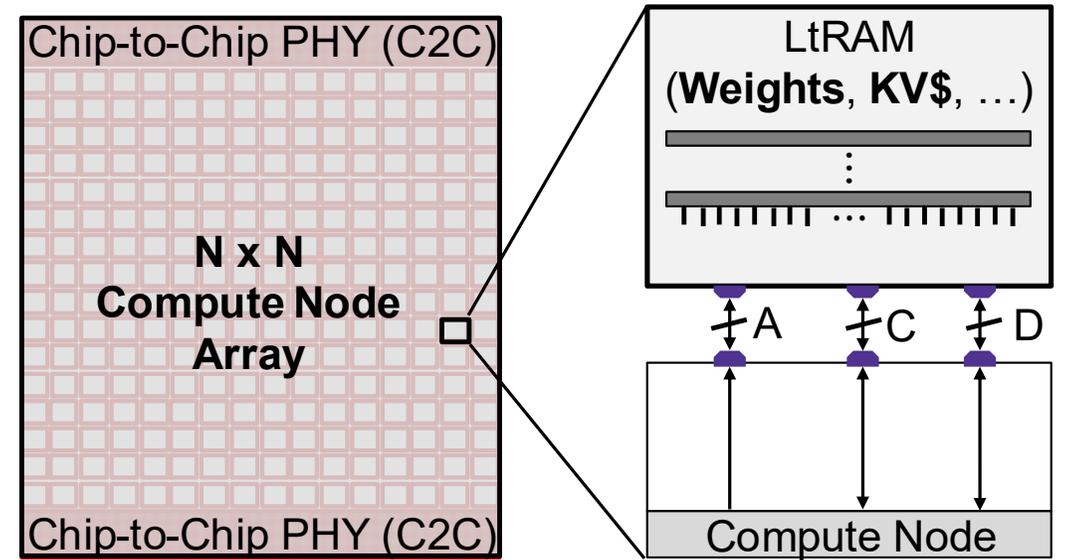
Intra-chip tensor parallelism

# System Parameters

## Baseline



## Model-Resident



### Compute

Parameter	Unit
Area	mm <sup>2</sup>
Energy	pJ/Op
Throughput	TFLOPS

### NoC

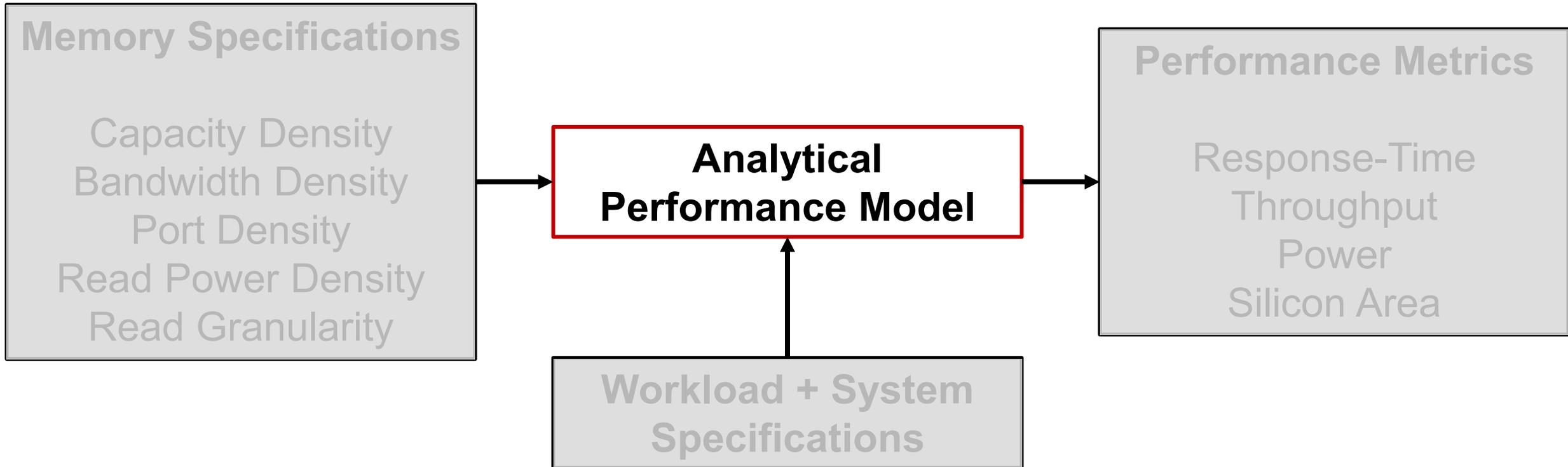
Parameter	Unit
Hop Latency	ps
Hop Energy	pJ/bit
Bandwidth	GB/s

### Chip-to-Chip

Parameter	Unit
Hop Latency	ps
Hop Energy	pJ/bit
Bandwidth	GB/s

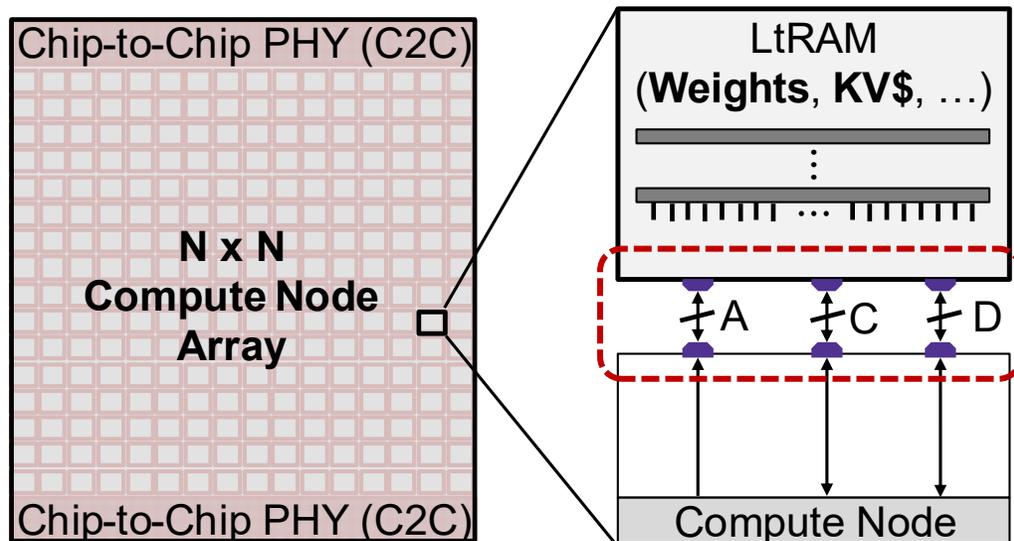
# Quantitative Performance Evaluation Framework

---

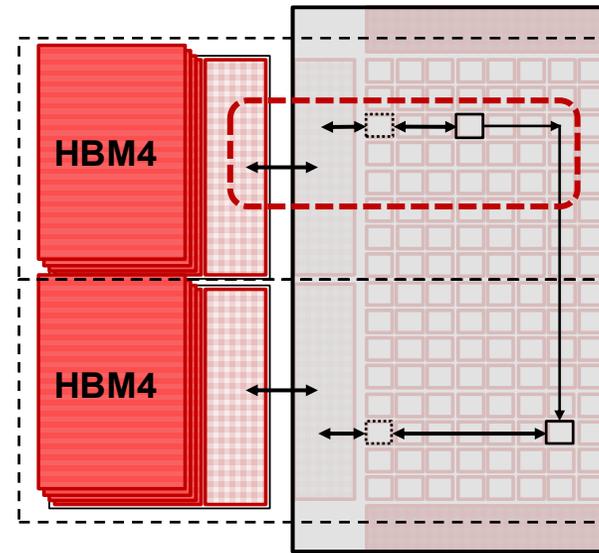


# System Modeling: Memory Access Assumptions

## Model-Resident



## Baseline



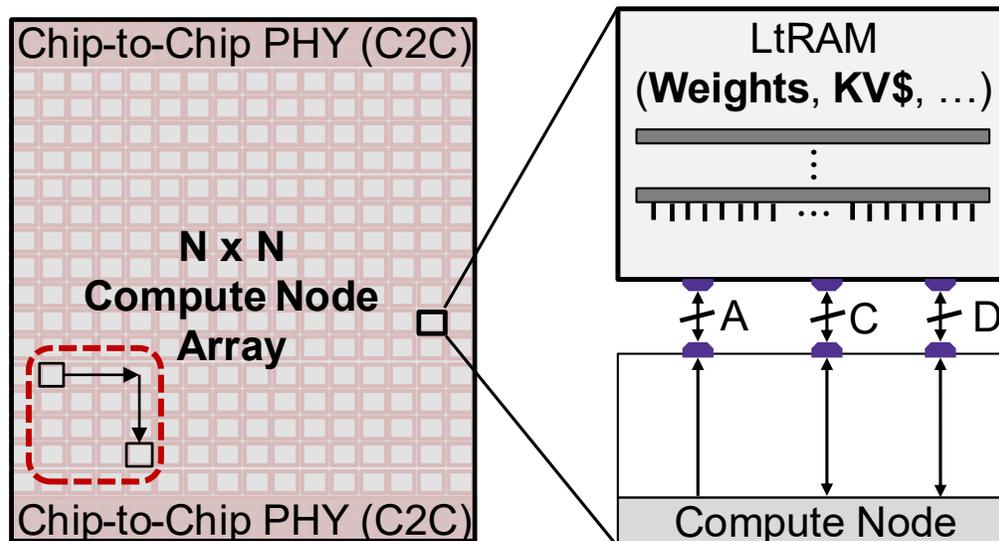
Alpha-beta access model: fixed + streaming

Fixed latency & energy cost: memory access + transport

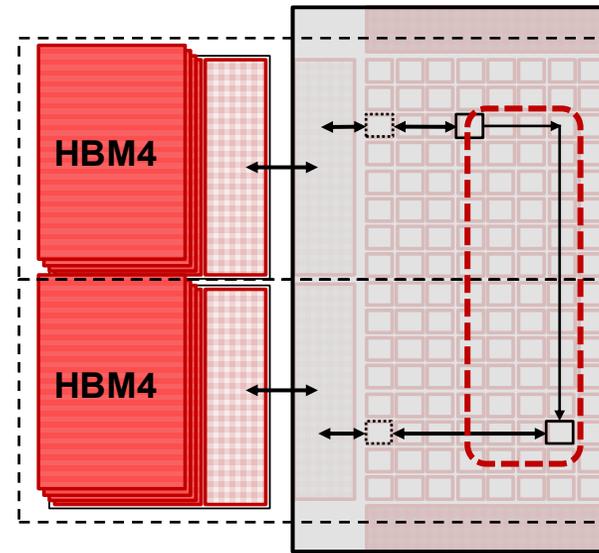
Streaming latency & energy cost: memory + network access

# System Modeling: Network-on-Chip Assumptions

## Model-Resident



## Baseline



Alpha-beta model of NoC transfers: fixed + streaming

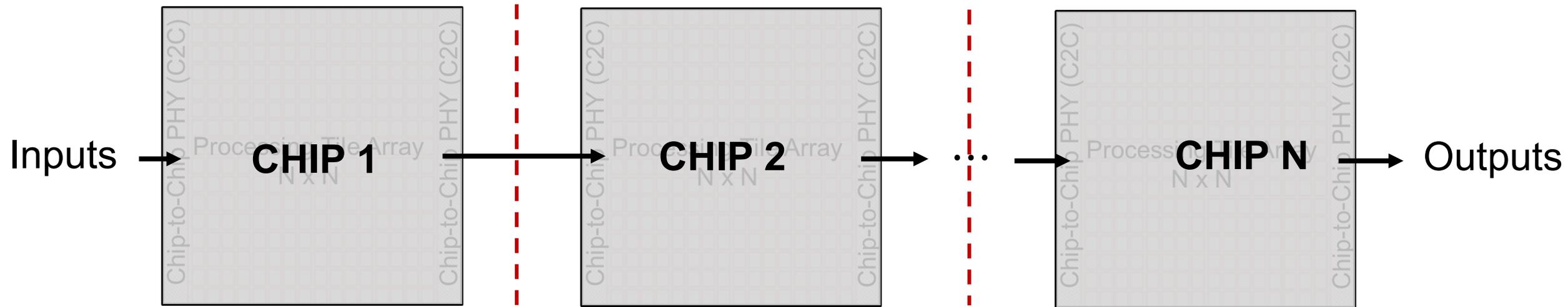
Fixed latency & energy cost: router + hop

Streaming latency & energy cost: NoC bisection

Average-hop approximation

# System Modeling: Multi-Chip Parallelism

---



Inter-chip pipeline parallelism

Workload split by per-chip memory capacity

Integer number of layers per chip

Intra-chip tensor parallelism

# Analytical Performance Model

---

- Layers processed as set of sequential stages
  - Input projection → Self-Attention → MoE
- Overlapped compute/ memory + exposed network synchronization
- Total latency & energy aggregated across stages & layers
  - $T_{\text{Total}} = \text{Layers} \times \left( \sum_{\text{stage}} \max \left( T_{\text{stage}}^{\text{comp}}, T_{\text{stage}}^{\text{mem}} \right) + T_{\text{layer}}^{\text{net}} \right) + T_{C2C}$
  - $E_{\text{Total}} = \text{Layers} \times \left( \sum_{\text{stage}} E_{\text{stage}}^{\text{comp}} + E_{\text{stage}}^{\text{mem}} + E_{\text{layer}}^{\text{net}} \right) + E_{C2C}$
- Validate vs. *LIMINAL*<sup>1</sup>

# Framework Simplifications

---

- ❑ Single face-to-face bonded memory die
- ❑ Uniform per-layer execution
- ❑ Mean hop-count NoC approximations
- ❑ Deterministic MoE routing
- ❑ Simplified roofline compute model

Capture dominant **first-order** architectural trade-offs

# Case Study: Memory Assumptions

---

	<b>Capacity Density</b> (GB/mm <sup>2</sup> )	<b>Bandwidth Density</b> (TB/s/mm <sup>2</sup> )	<b>Port Density</b> (ports/mm <sup>2</sup> )	<b>Read Granularity</b> (bytes/access)	<b>Read Power Density</b> (W/mm <sup>2</sup> )
<u>HBM4</u> <sup>1</sup>	2.56	0.01	0.16	8	0.017
6T Static <sup>2</sup> (SRAM)	0.05	0.35	12.3	256	0.035
OS 1T1C <sup>3</sup> (TSMC)	0.07	4.10	15.9	256	2.050
FeNOR <sup>4</sup> (EMD)	0.32	2.32	1250.0	4	0.023

<sup>1</sup> HBM4 JESD270-4A; <sup>2</sup> DESTINY simulation (512 KB RAM macro)

<sup>3</sup> TSMC. [Chiang, Symp. VLSI Circuits 25]; <sup>4</sup> EMD Electronics Projections

# Case Study: System Parameter Assumptions

System Parameter		Value
<b>NoC<sup>1</sup></b>	Bandwidth	128 GB/s
	Latency	500 ps
	Energy	0.1 pJ/bit
<b>C2C<sup>2</sup></b>	Bandwidth	320 GB/s
	Latency	4 ns
	Energy	1.17 pJ/bit
<b>Compute</b>	Area	800 mm <sup>2</sup>
	Energy <sup>3</sup>	0.160 pJ/op
	Throughput <sup>4</sup>	2250 TFLOPS/s

<sup>1</sup> Orenes-Vera *et. al*, “Massive Data-Centric Parallelism in the Chiplet Era.”

<sup>2</sup> Open Compute Group, “Bunch of wires PHY specification.”

<sup>5</sup> B. Dally *et. al*, “On the model of computation.”

<sup>6</sup> M. Davies *et.al*, *LIMINAL*. [NVIDIA Research].

# Case Study: Representative Workload Assumptions

---

<b>Workload Variable</b>	<b>Llama 4 Scout</b>	<b>GPT-OSS-120B</b>	<b>Llama 3 405B</b>
$d_{model}$	4096	4096	16384
$d_n$	128	128	128
$H_q$	32	32	128
$H_{kv}$	8	8	8
$L$	32	36	126
$E$	16	128	1
$K$	2	4	—
$d_{ffn}$	14336	14336	53248

# Review: System Performance Metrics

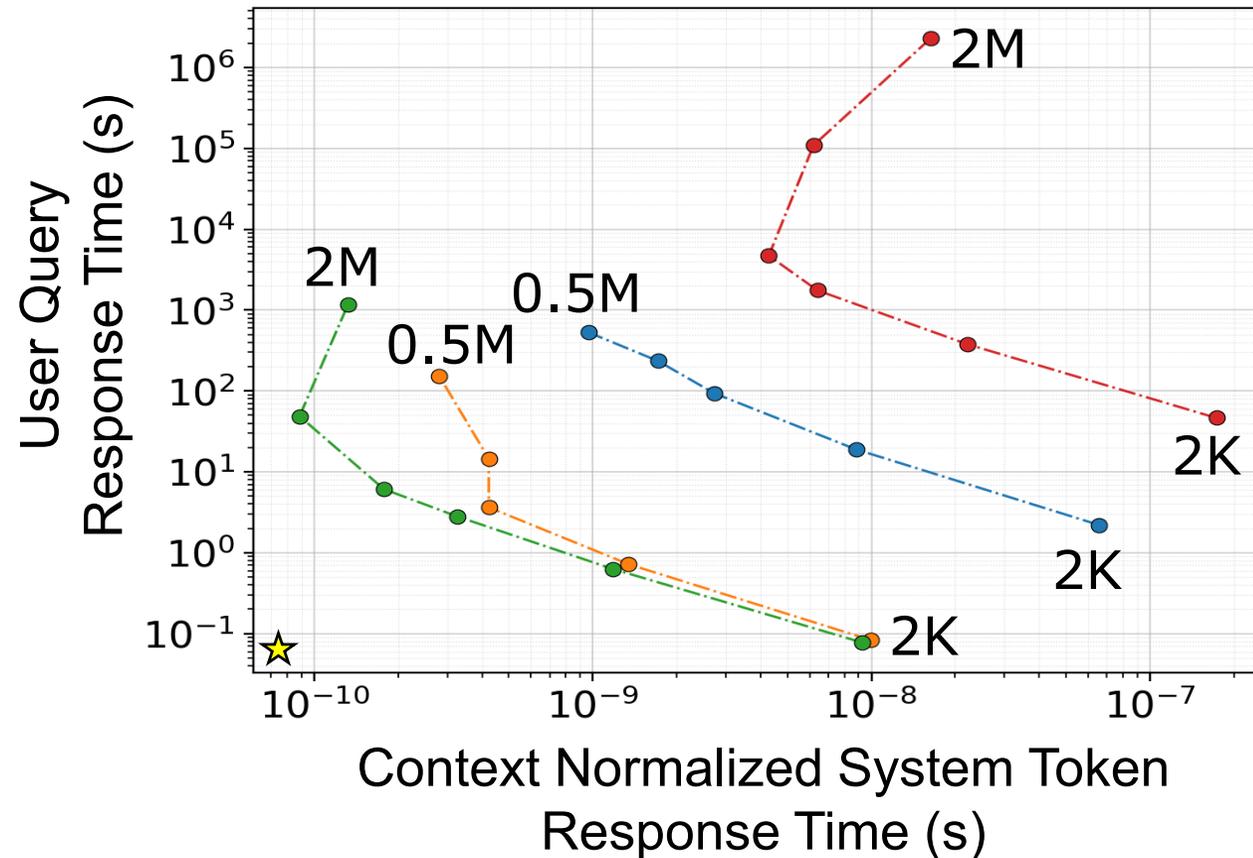
---

System Performance Metric	Definition	Better
User Query Response Time (s)	End-to-end latency per user query	↓
System Token Response Time (s)	Average time per generated token	↓
User Tokens Throughput (token/s)	Token generation rate per user	↑
System Tokens Throughput(token/s)	Aggregate token generation rate across all users	↑
System Power-Efficiency (token/s/watt)	Represents energy cost to generate token (OPEX)	↑
Silicon Area-Efficiency (token/s/mm <sup>2</sup> )	Silicon area utilization efficiency (CAPEX)	↑

# Case-Study: Response-Time Performance

Model: Llama 4 Scout

Context Sweep: [2K, 16K, 64K, 128K, 0.5M, 2M]



Batch Sweep: B = 1 to 128

- Static (optimal B = 8)
- OS 1T1C (optimal B = 2)
- FeNOR (optimal B = 2)
- HBM4 (optimal B = 64)

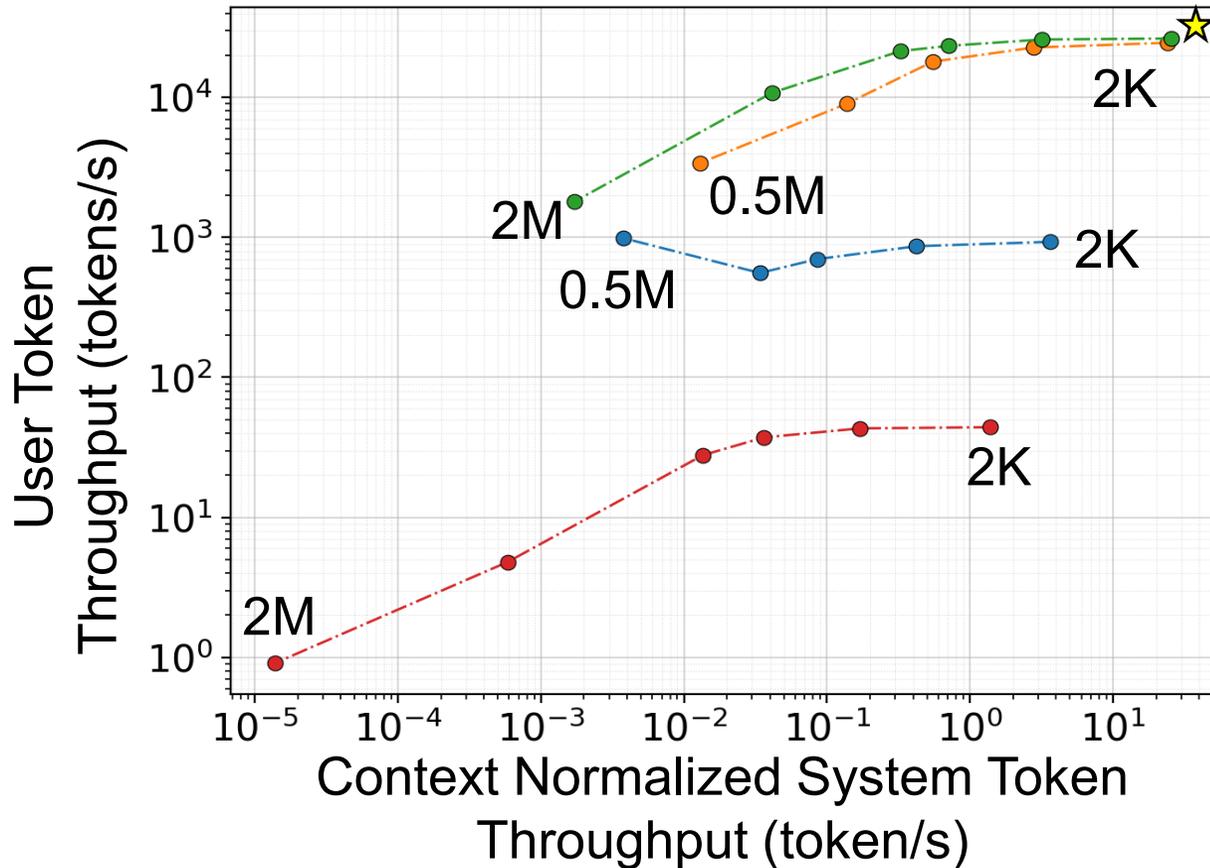
★ Preferred Corner

Model-residency shifts latency limit to longer context

# Case-Study: Throughput Performance

Model: Llama 4 Scout

Context Sweep: [2K, 16K, 64K, 128K, 0.5M, 2M]



Batch Sweep: B = 1 to 128

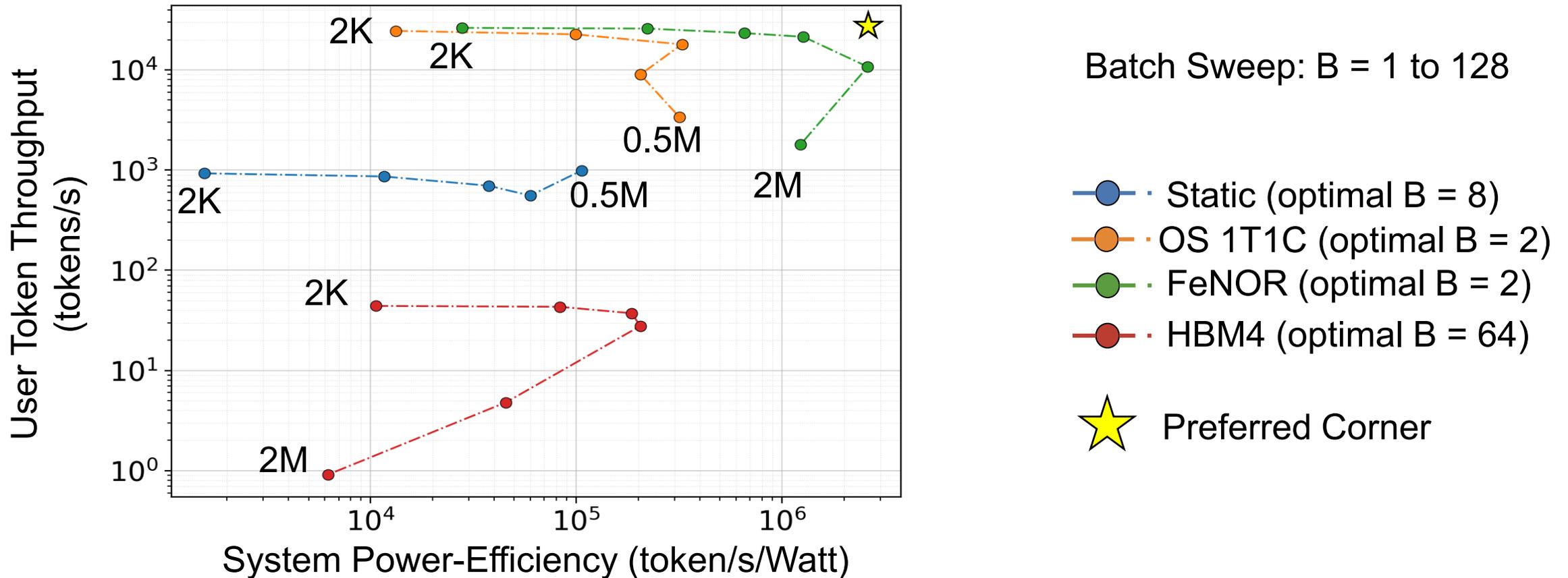
- Static (optimal B = 8)
- OS 1T1C (optimal B = 2)
- FeNOR (optimal B = 2)
- HBM4 (optimal B = 64)
- ★ Preferred Corner

Bandwidth and capacity density determine usable throughput

# Case-Study: Power Performance

Model: Llama 4 Scout

Context Length Sweep: [2K, 16K, 64K, 128K, 0.5M, 2M]

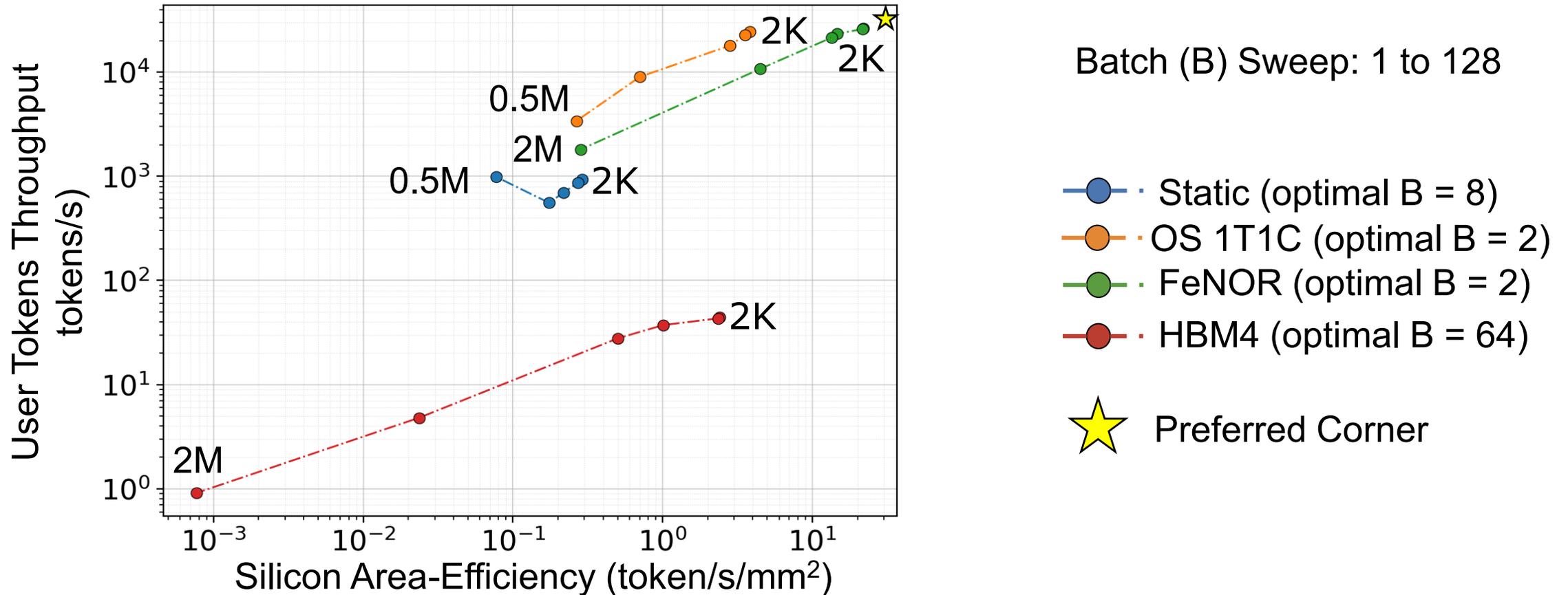


Must co-optimize energy-efficiency with bandwidth density

# Case-Study: Silicon-Area Performance

Model: Llama 4 Scout

Context Length Sweep: [2K, 16K, 64K, 128K, 0.5M, 2M]



Area efficiency requires bandwidth–capacity co-scaling

# Conclusion

---

## Take Away:

- ❑ Model-resident architecture key for AI inference
- ❑ Requires high-density and high read bandwidth

## Future Work:

- ❑ Multi-die-stacking
- ❑ Workload modeling (parallelism, routing skew)
- ❑ Compute & NoC micro-architecture detail
- ❑ Publish open-source

## Call for Action:

- Memory makers to provide high-density, high read bandwidth LtRAMs