

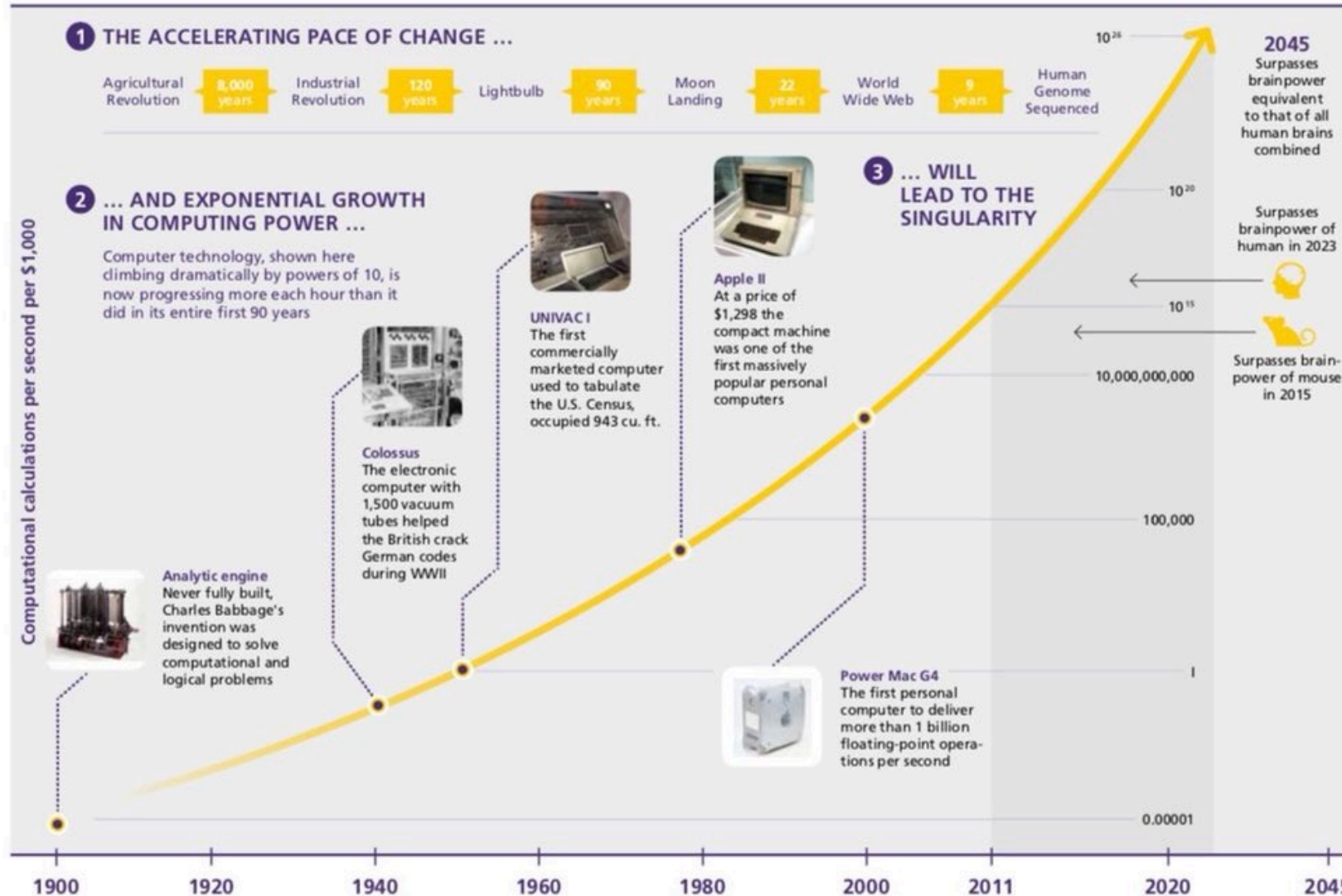
Superoptimizing Machine Learning Systems

Zhihao Jia

Computer Science Department
Carnegie Mellon University

What are the **fundamental driving forces**
behind the success of ML?

Compute Per Second Per Dollar



Surpass human brainpower in 2023

Scaling Law in ML

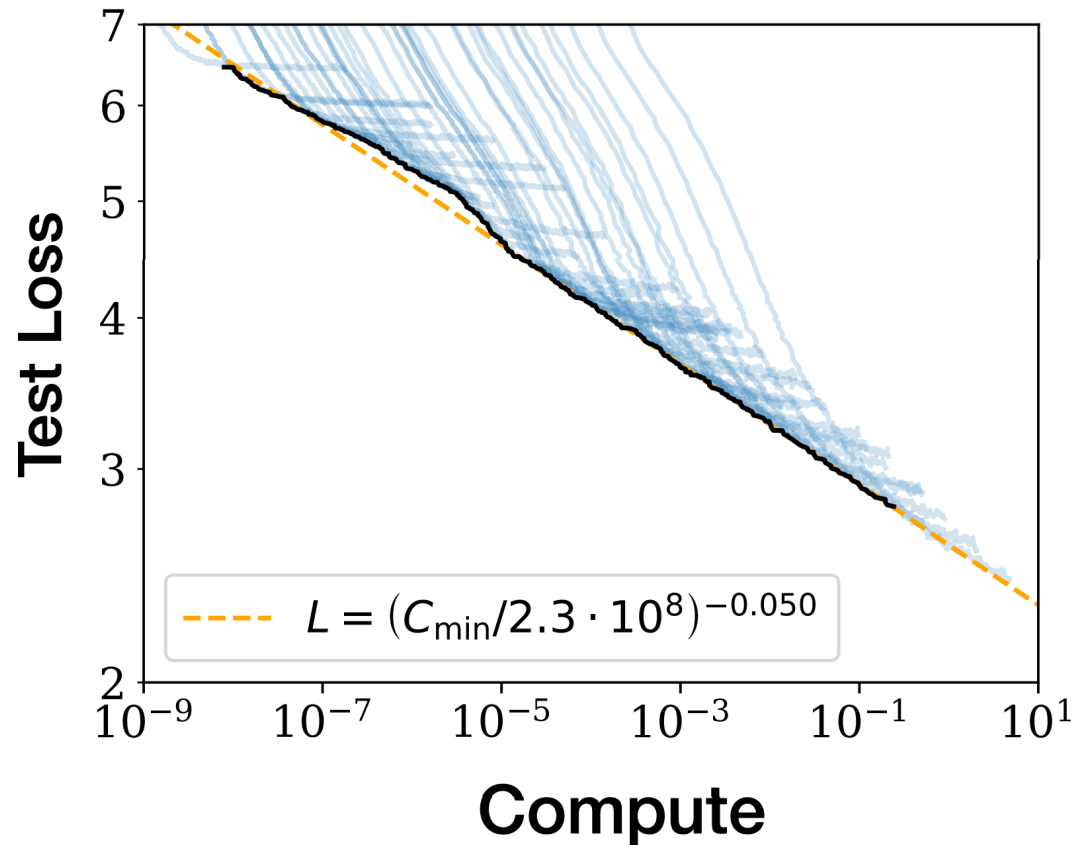
Improving model accuracy & capability



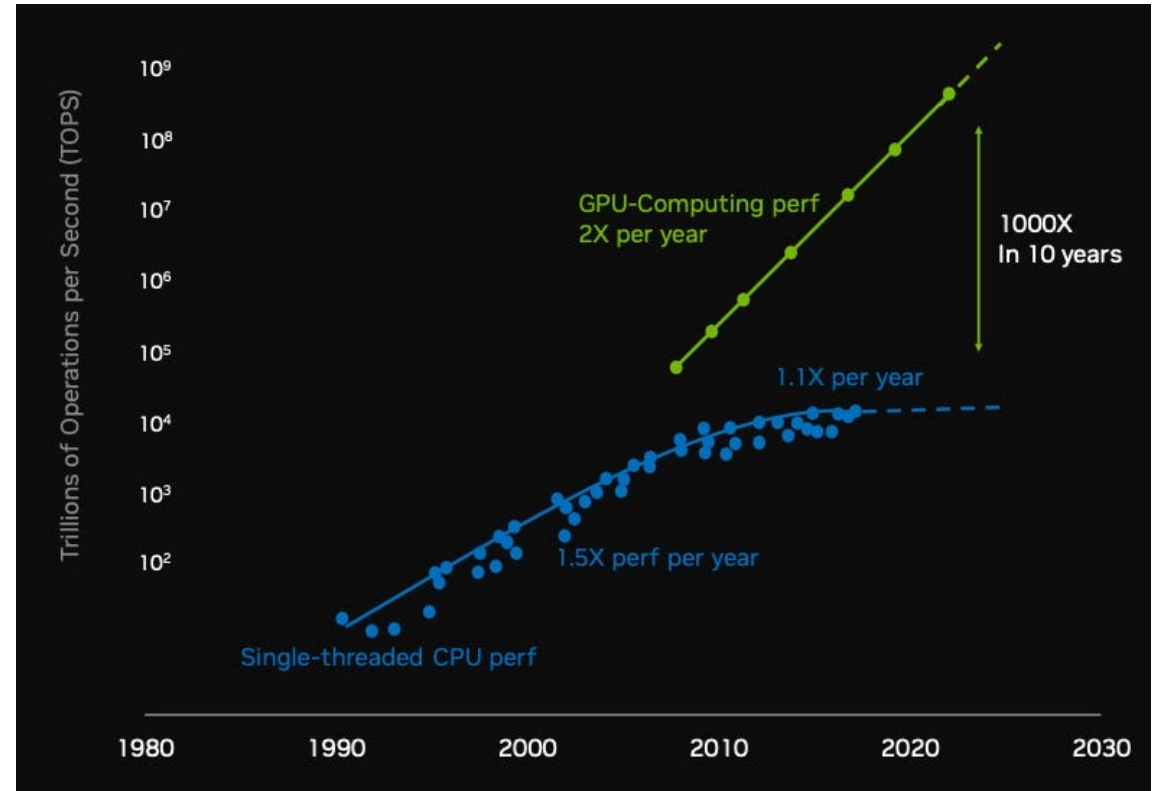
Scaling training & inference compute



Hardware parallelization and specialization

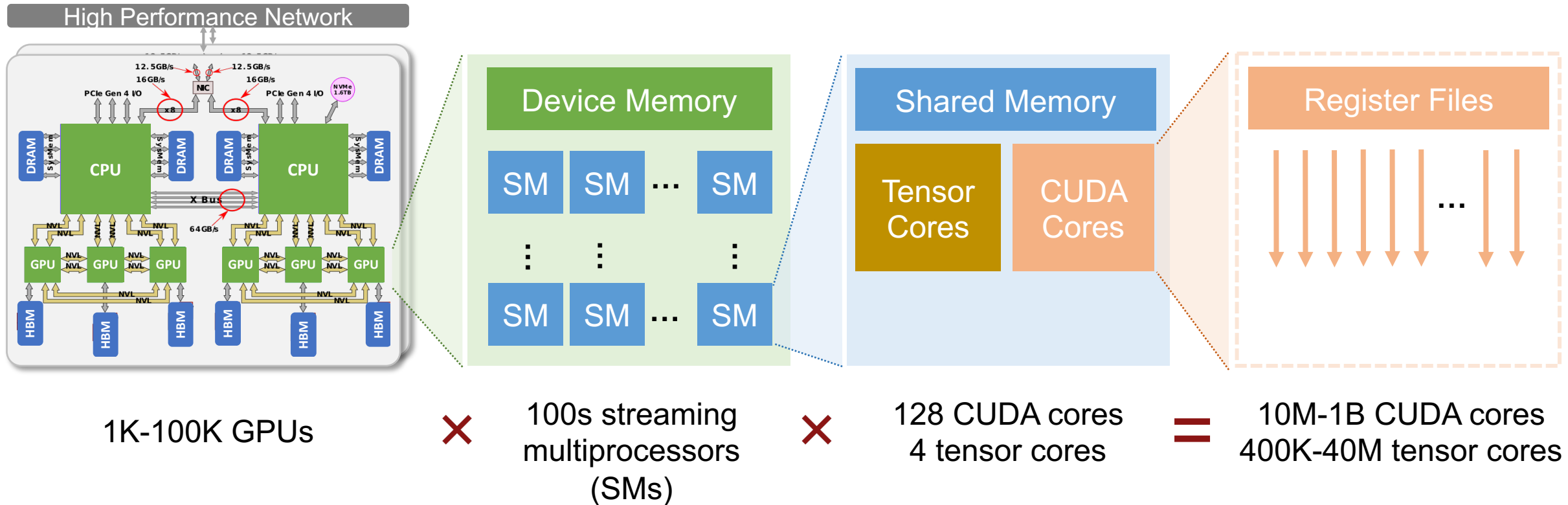


Source: OpenAI



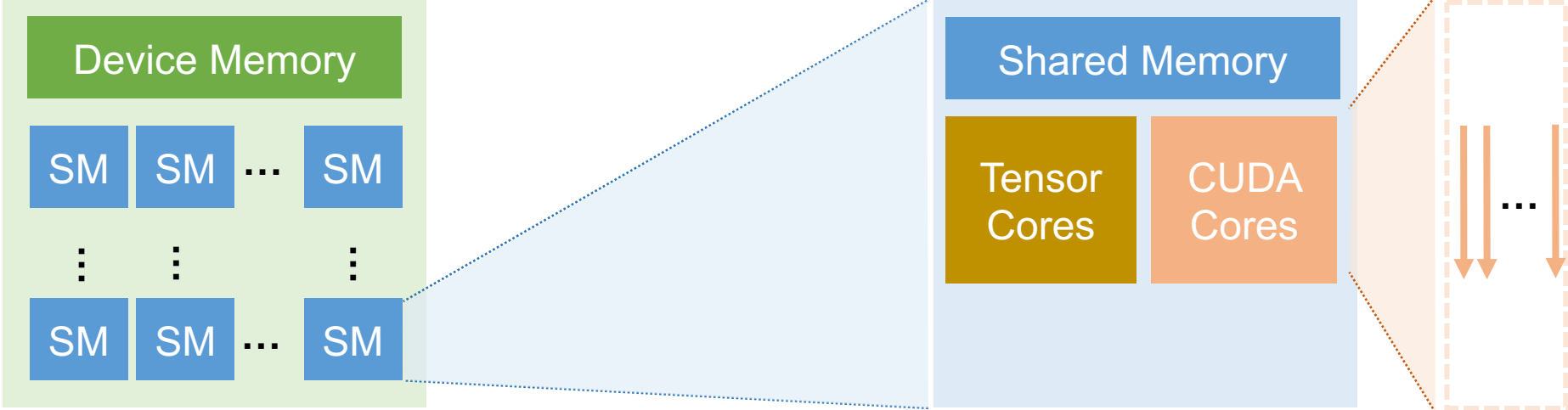
Source: NVIDIA

ML Hardware is Massively Parallel, Highly Heterogeneous

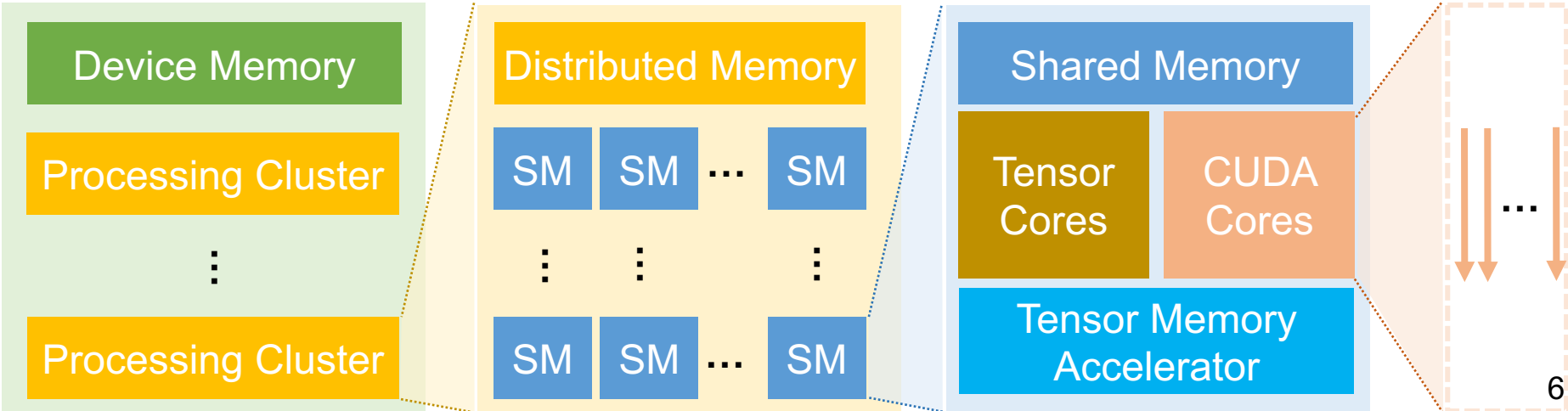


ML Hardware is Quickly Evolving

NVIDIA A100 GPU (2020)

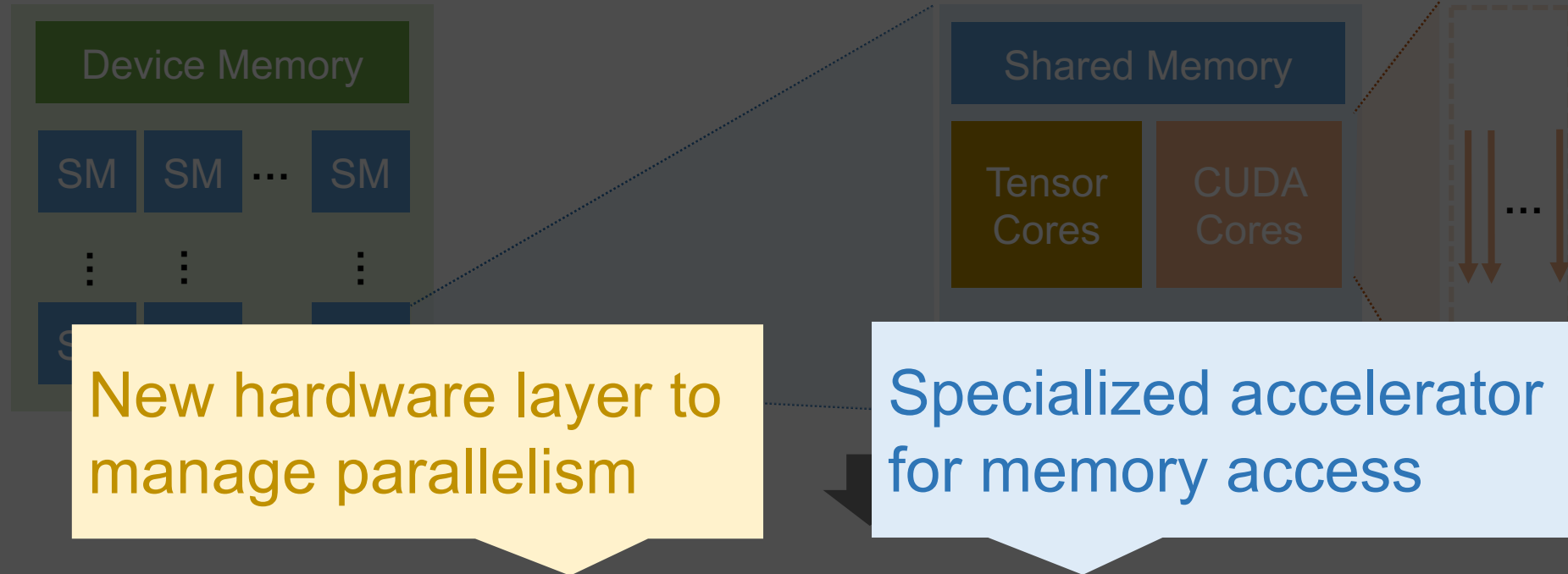


NVIDIA H100 GPU (2022)



ML Hardware is Quickly Evolving

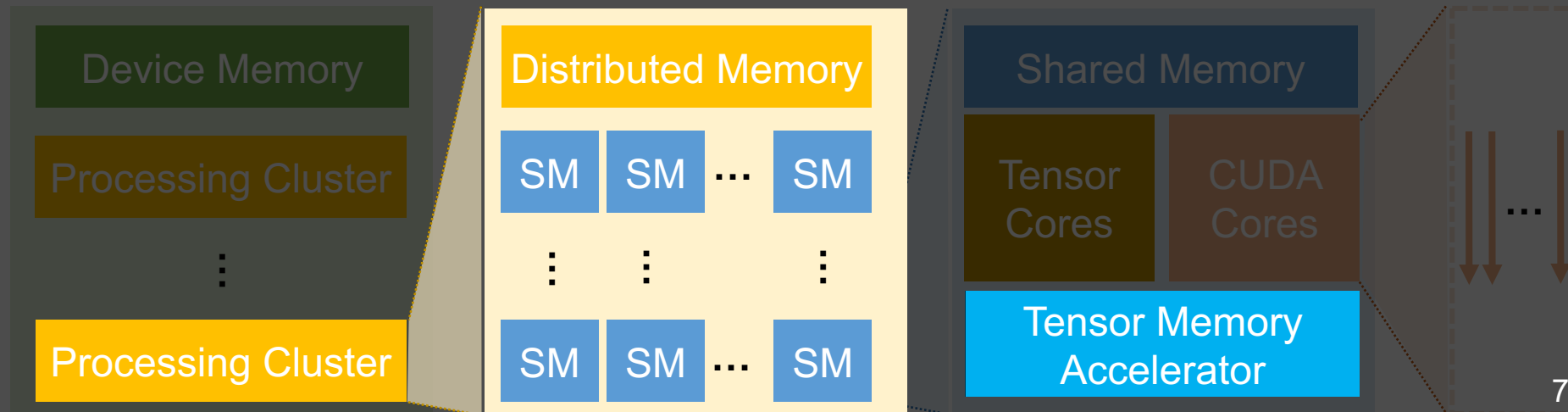
NVIDIA A100 GPU (2020)



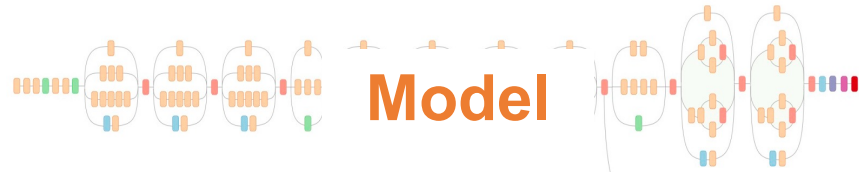
New hardware layer to manage parallelism

Specialized accelerator for memory access

NVIDIA H100 GPU (2022)



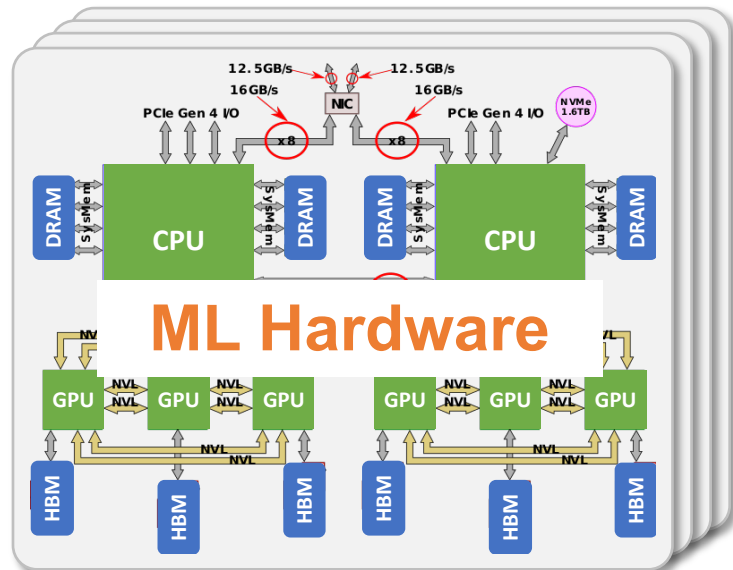
Our Research: ML Systems



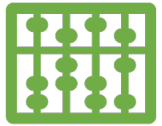
ML Systems

Goal:

Efficiently deploying ML applications on massively parallel, increasingly heterogeneous, rapidly evolving hardware platforms



Key Challenges for Developing ML Systems



Massively Parallel

Billions of compute units on modern ML hardware

How can we find the best way to parallelize ML computation?



Increasingly Heterogeneous

CUDA cores, tensor cores, tensor memory accelerators, processing clusters, and more

How can we handle different accelerator types and complex memory hierarchy?

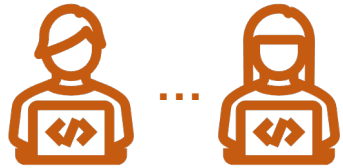
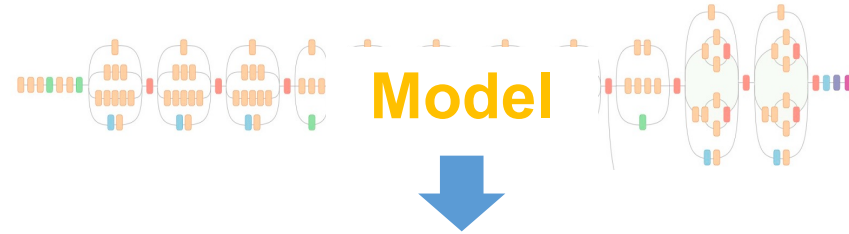


Rapidly Evolving

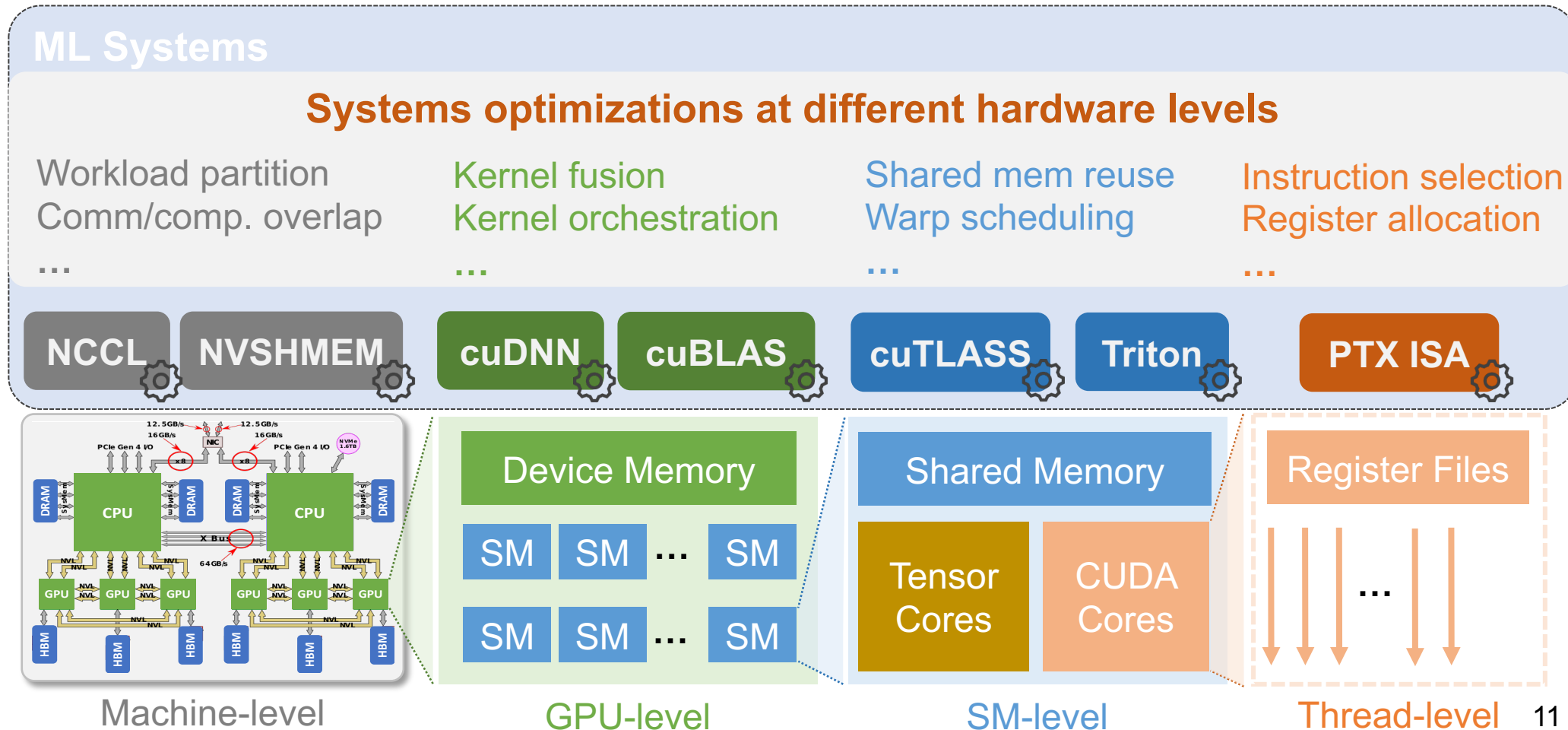
New generation every 2-3 years; but building high-quality systems & compilers takes much longer

How can we deal with the rapid evolution of ML hardware?

Current Practice: Rely on Engineers to Handle HW Complexity

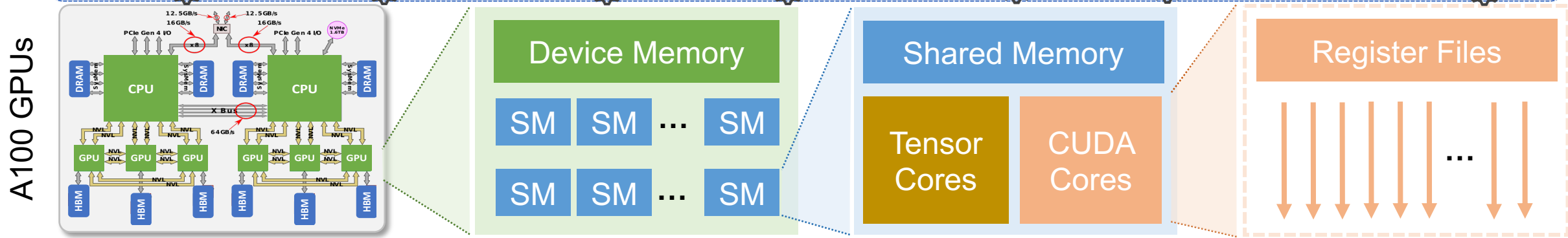
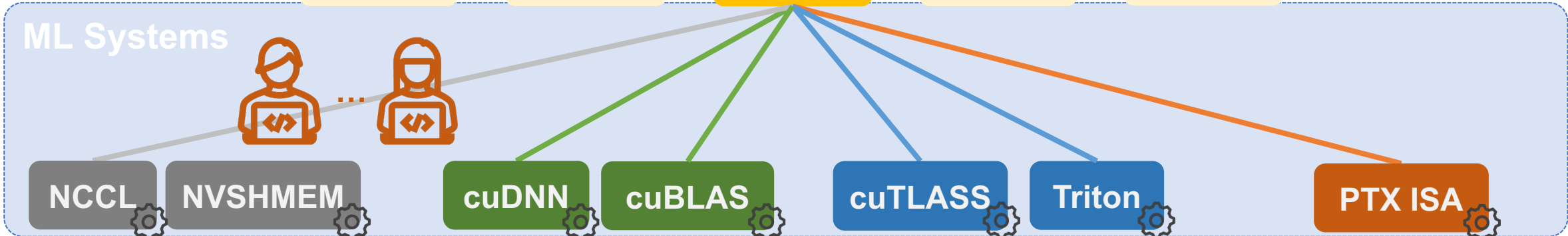
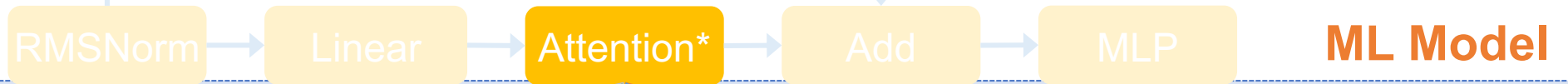


Manually design & implement using HW libraries, compilers, ISAs



Issue 1: Time Intensive to Manually Design and Implement

- Complex interactions between HW levels
- Optimizing attention takes many months; 1000+ calls to HW libraries; 60K LOC



5 NCCL APIs

103 cuBLAS APIs
64 cuDNN APIs

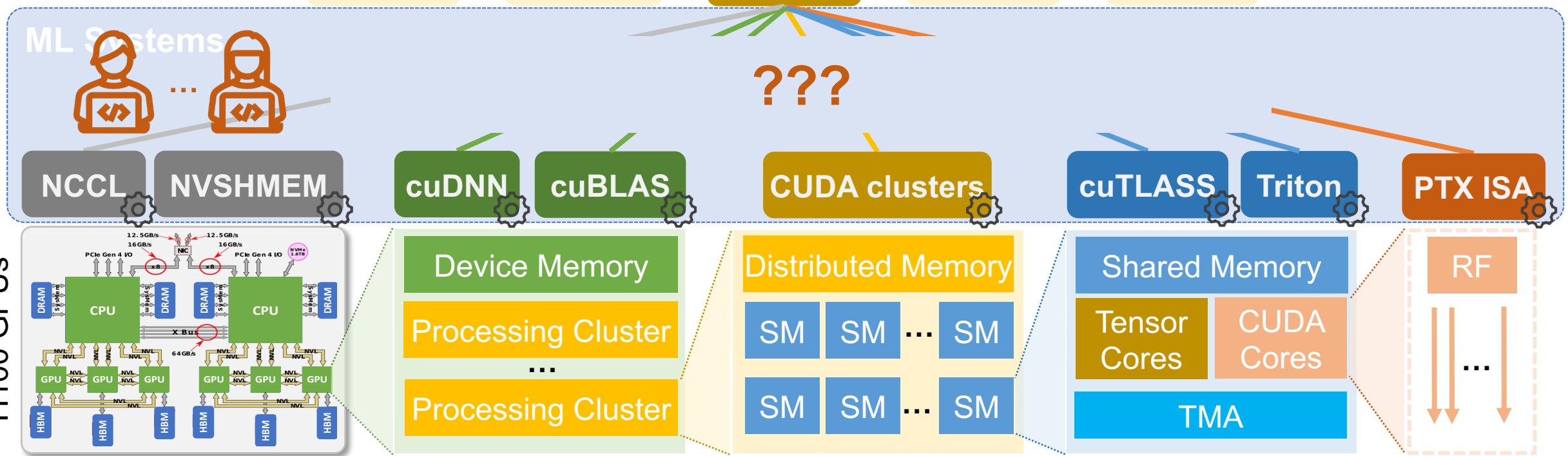
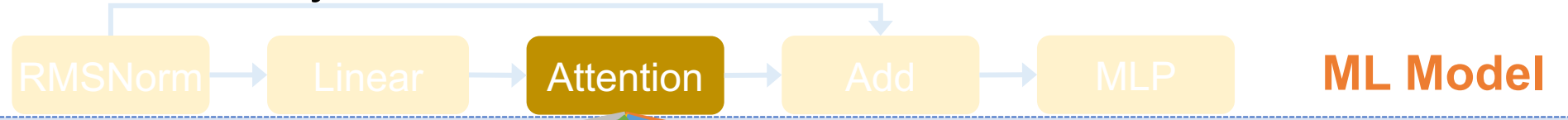
719 cuTLASS APIs
191 Triton APIs

26 PTX ISA APIs

* FlashAttention: Fast and Memory-Efficient Exact Attention

Issue 2: New HW Requires New Design and Implementation

- New HW features affect optimization landscape of others
- Reimplement kernels to optimize attention for H100 (15K LOC)
- Not available until two years after H100's release



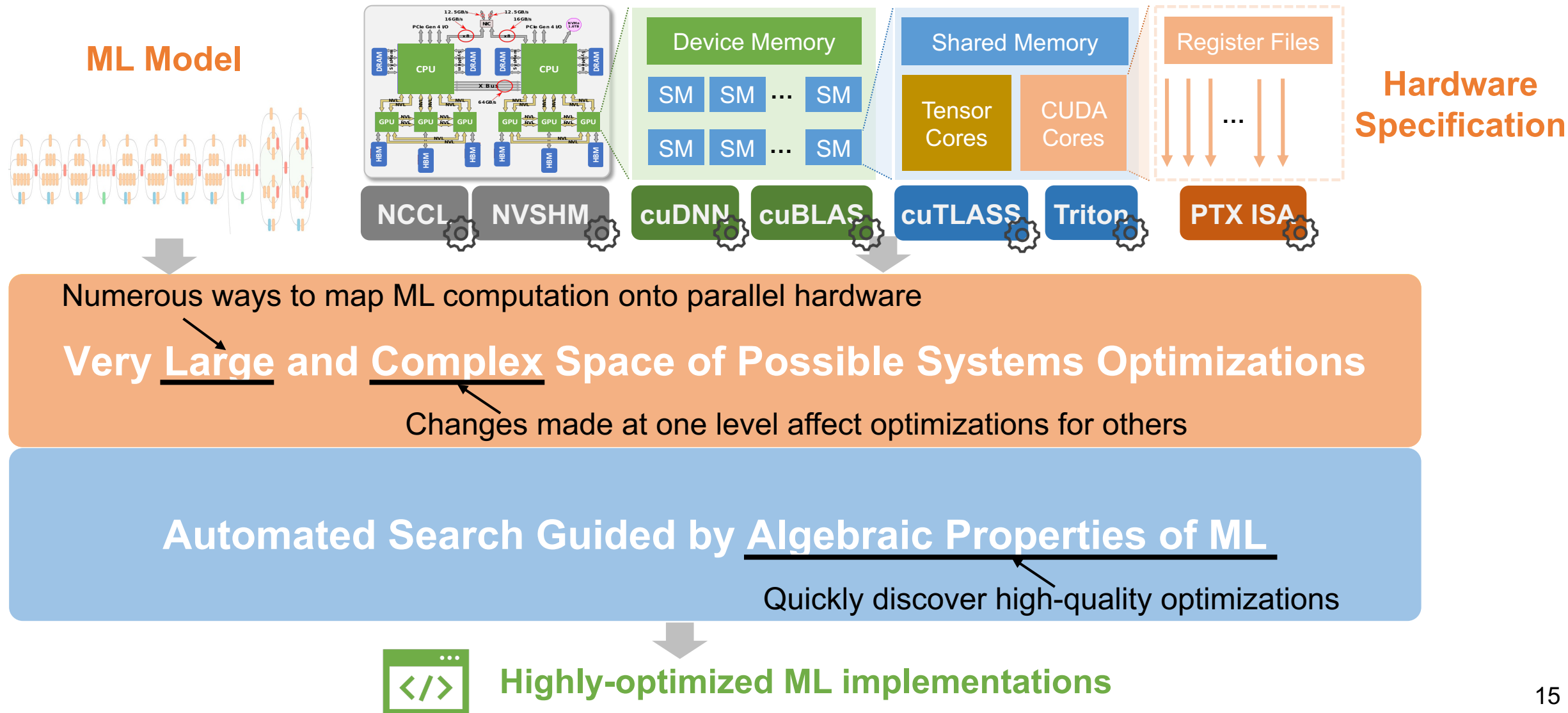
Can ML systems discover and deploy these optimizations **automatically**?

~~Many engineering months, 60K LOC
to optimize attention on A100~~

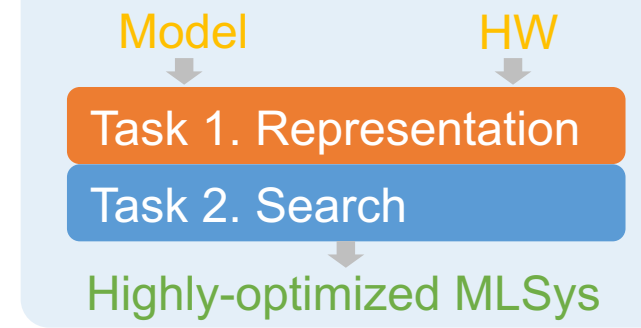
~~Reimplementing kernels, modifying
15K LOC for H100~~

ML systems automate these optimizations
with minimal or even no manual effort

Our Vision: ML Systems *Automatically* Discover and Deploy Optimizations at All Levels

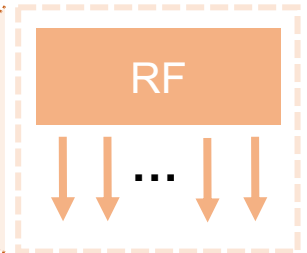
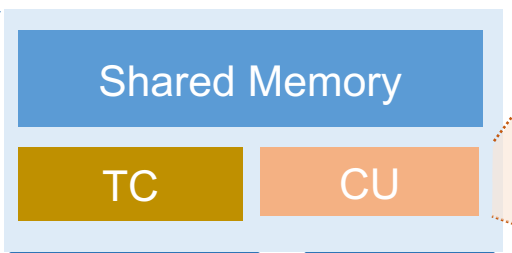
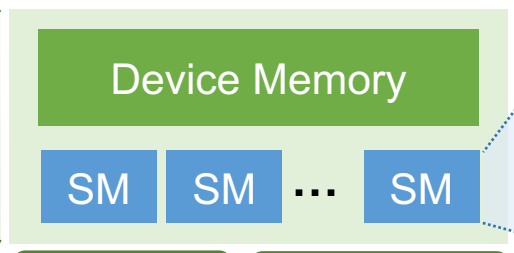
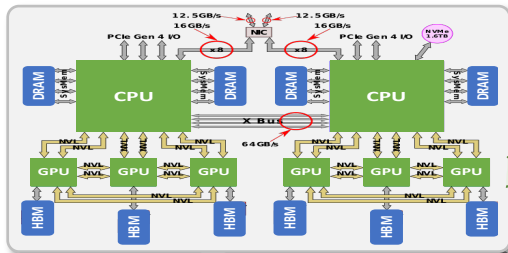


Task 1: How to Represent Optimizations?

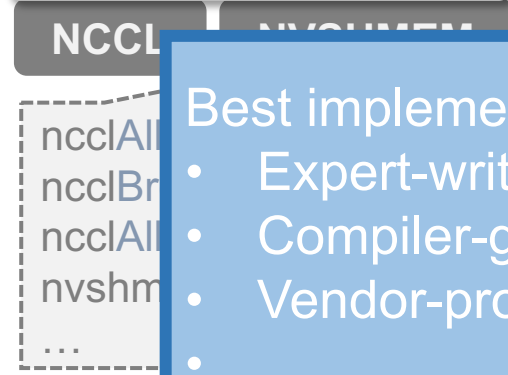


Desiderata: a unified representation across HW layers and architectures

A100 GPUs



Current library APIs



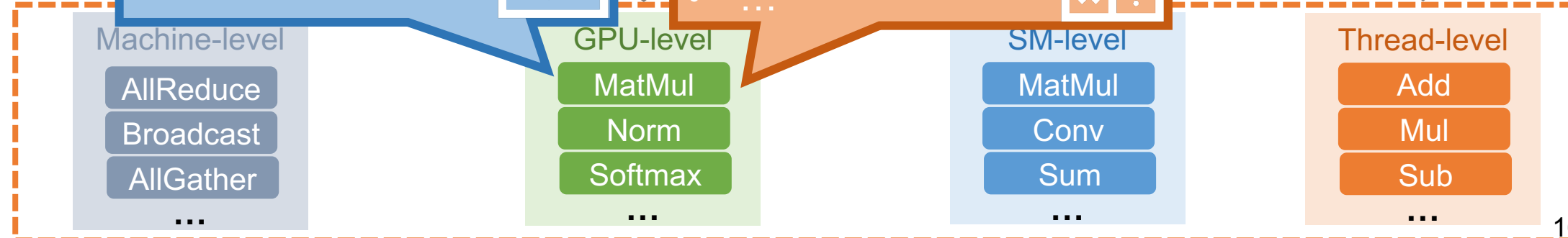
Best implementations:

- Expert-written
- Compiler-generated
- Vendor-provided
- ...

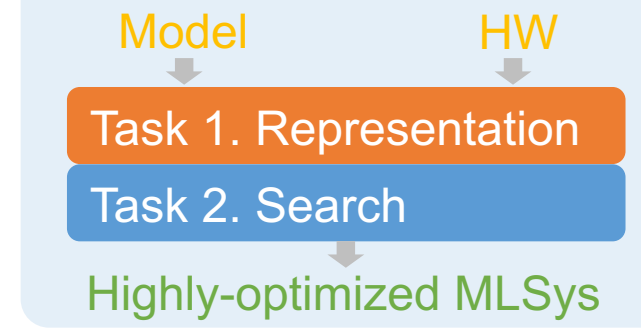
Algebraic properties:

- Linearity
- Commutativity
- Associativity
- Distributivity
- ...

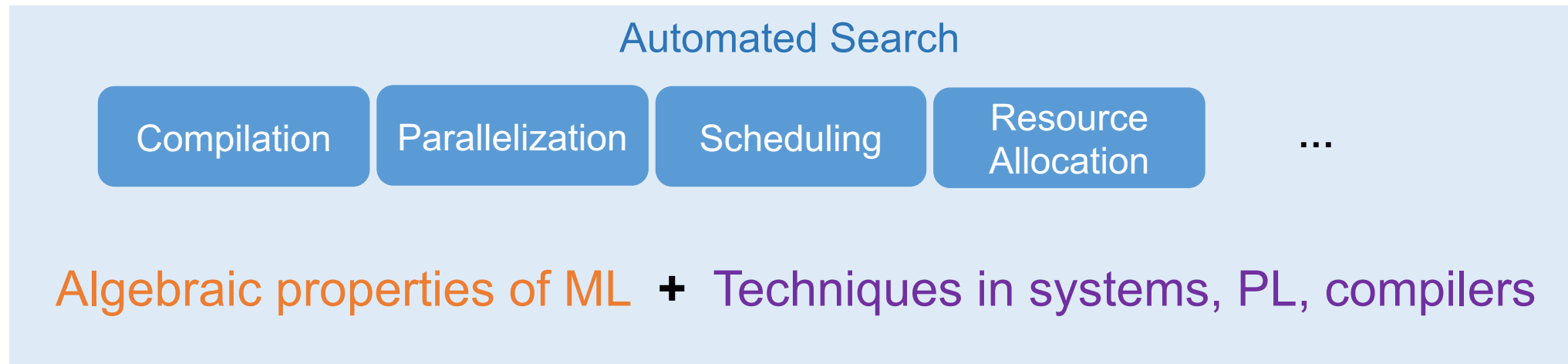
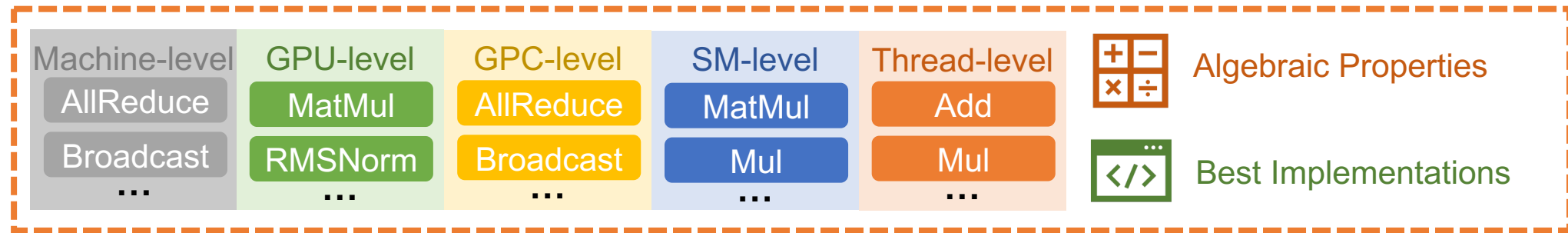
Multi-Level Abstraction



Task 2: How to Find Performant Systems?

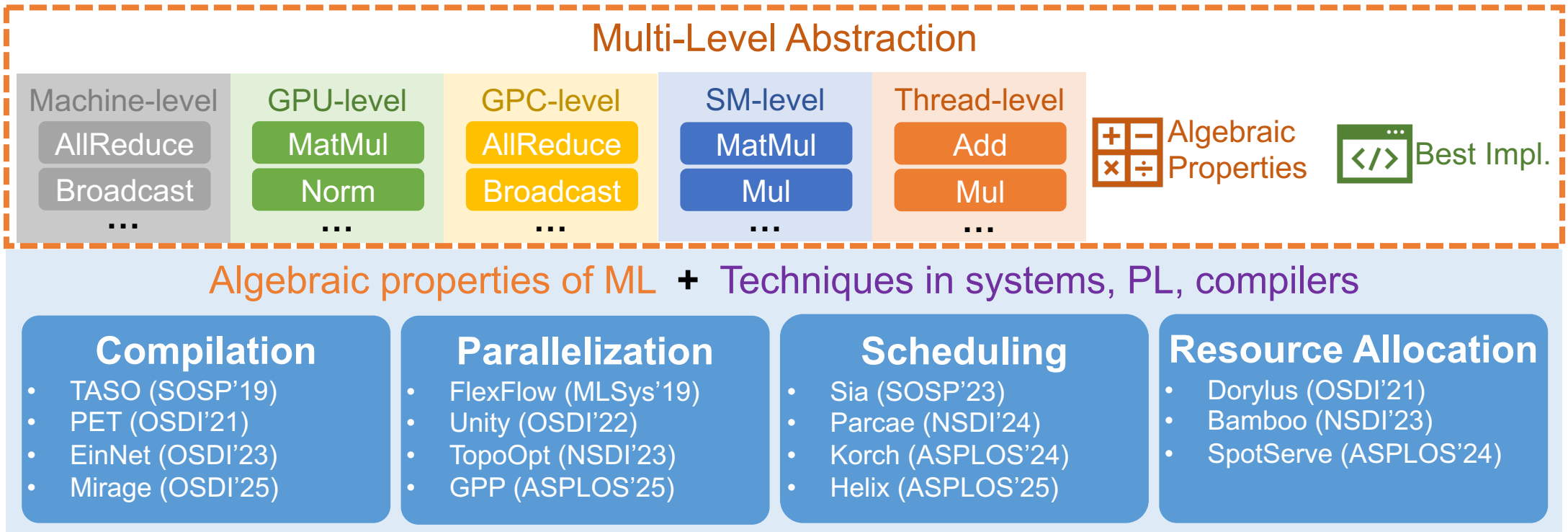


- Need to simultaneously consider many tasks



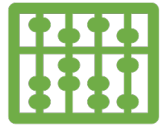
Highly optimized ML systems

Our Research: Automated End-to-end ML Systems



How to Address Three Challenges?

Capture model- and HW-specific optimizations across all levels



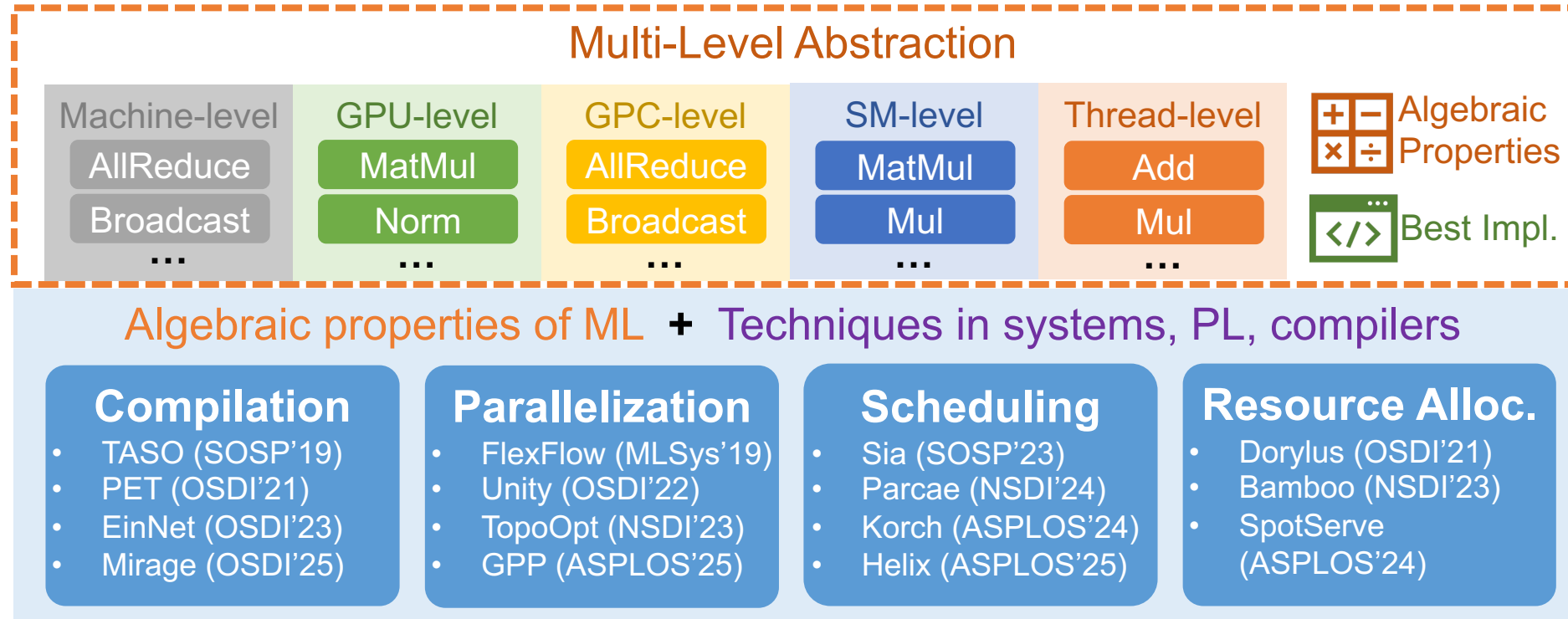
Massively Parallel



Increasingly Heterogeneous



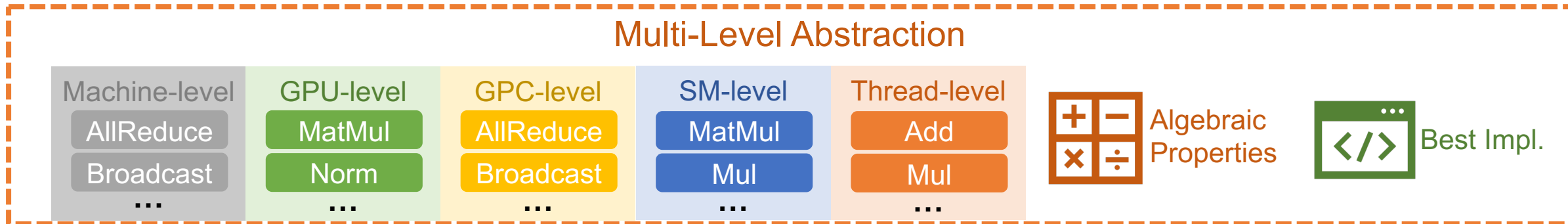
Rapidly Evolving



Automatically discover efficient ways to parallelize ML

Much less engineering efforts; faster adaptation to new HW

Our Techniques Generalize Beyond ML Systems



Algebraic properties of ML + Techniques in systems, PL, compilers

Compilation

- TASO (SOSP'19)
- PET (OSDI'21)
- EinNet (OSDI'23)
- Mirage (OSDI'25)

Parallelization

- FlexFlow (MLSys'19)
- Unity (OSDI'22)
- TopoOpt (NSDI'23)
- GPP (ASPLOS'25)

Scheduling

- Sia (SOSP'23)
- Parcae (NSDI'24)
- Korch (ASPLOS'24)
- Helix (ASPLOS'25)

Resource Allocation

- Dorylus (OSDI'21)
- Bamboo (NSDI'23)
- SpotServe (ASPLOS'24)

Quantum Circuit Optimizer

- Quartz (PLDI'22)
- Quarl (OOPSLA'24)



Quantum Algebra

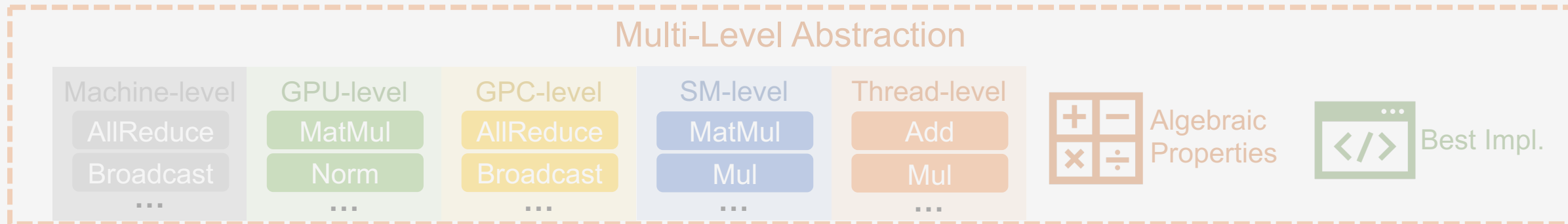
Database Query Optimizer

- TOD (VLDB'22)
- SDPipe (VLDB'23)



Relational Algebra

Our Techniques Generalize Beyond ML Systems



Algebraic properties of ML + Techniques in systems, PL, compilers

Compilation

- TASO (SOSP'19)
- PET (OSDI'21)
- EinNet (OSDI'23)
- Mirage (OSDI'25)

Parallelization

- FlexFlow (MLSys'19)
- Unity (OSDI'22)
- TopoOpt (NSDI'23)
- GPP (ASPLOS'25)

Scheduling

- Sia (SOSP'23)
- Parcae (NSDI'24)
- Korch (ASPLOS'24)
- Helix (ASPLOS'25)

Resource Allocation

- Dorylus (OSDI'21)
- Bamboo (NSDI'23)
- SpotServe (ASPLOS'24)

Quantum Circuit Optimizer

- Quartz (PLDI'22)
- Quarl (OOPSLA'24)



Quantum Algebra

Database Query Optimizer

- TOD (VLDB'22)
- SDPipe (VLDB'23)



Relational Algebra

Mirage: A Multi-Level SuperOptimizer for ML

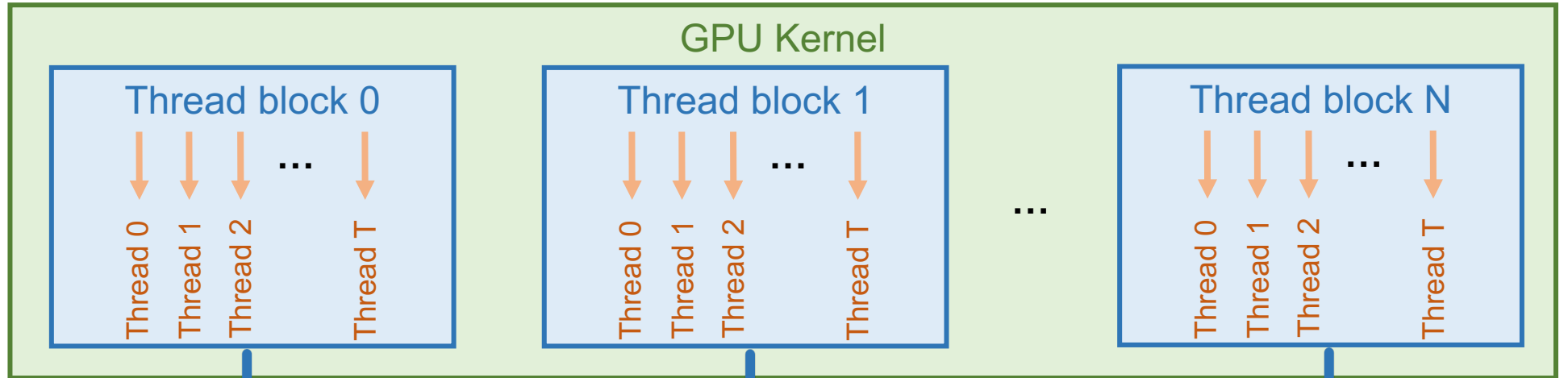
Key idea: automatically translates pure math definition of ML models into highly optimized GPU code



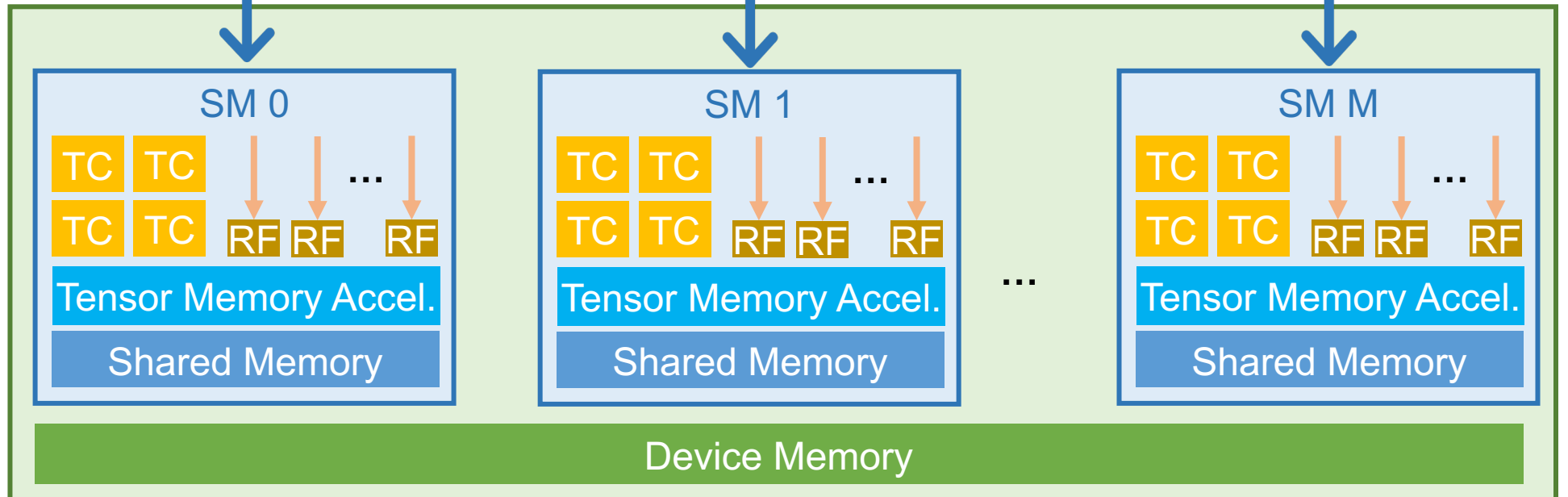
- **Less engineering effort:** thousands of lines of CUDA code → a few lines of code in Mirage
- **Better performance:** outperform existing systems by 1.1-2.9x
- **Faster adaptation:** day-0 support for new models; no manual effort

GPU Programming 101

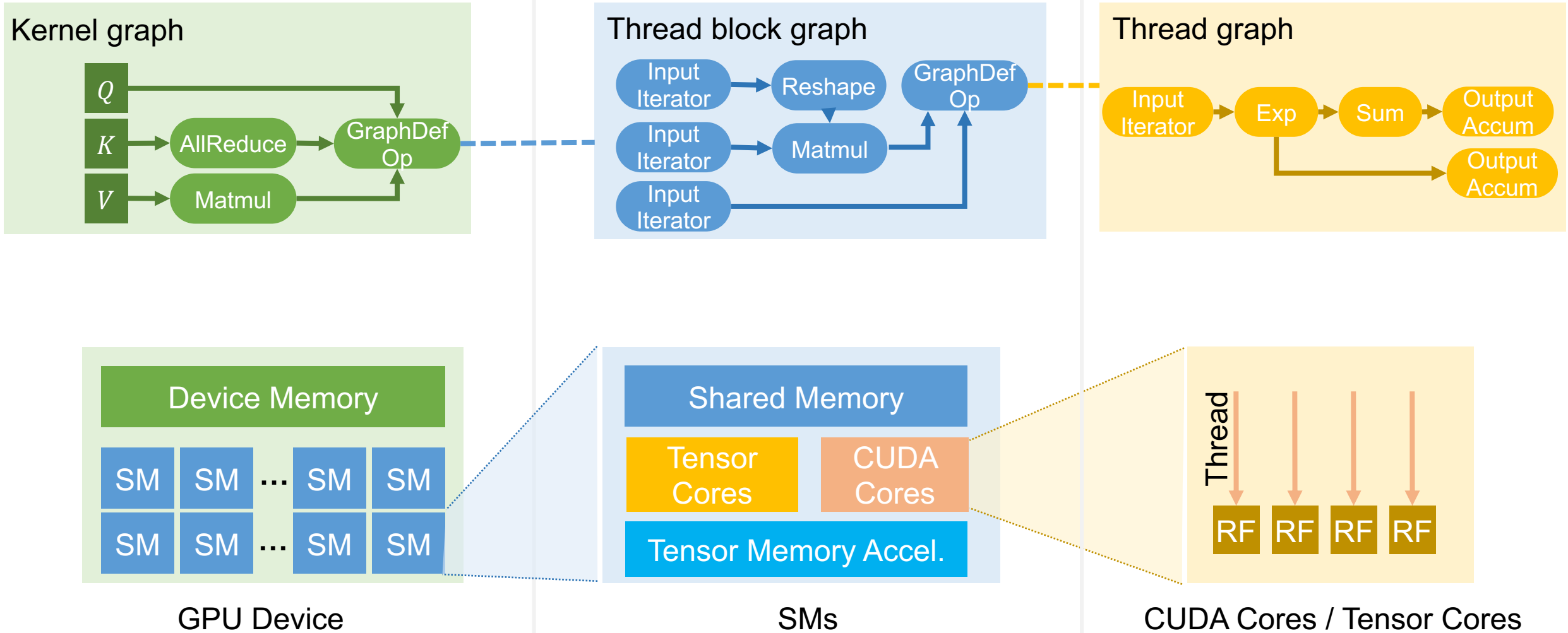
Programming
Abstraction



Hardware
Architecture

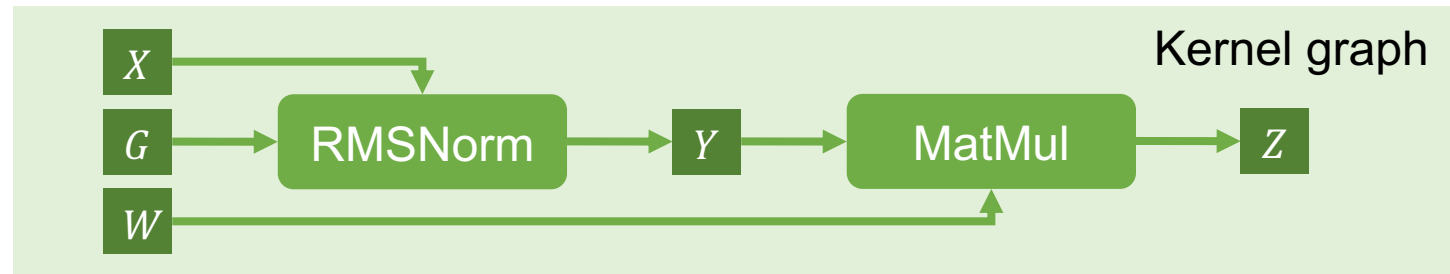


μ Graphs: Hierarchical Graph Representation



Motivating Example: RMSNorm & MatMul in LLMs

Existing systems launch two kernels since Y does not fit in shared memory



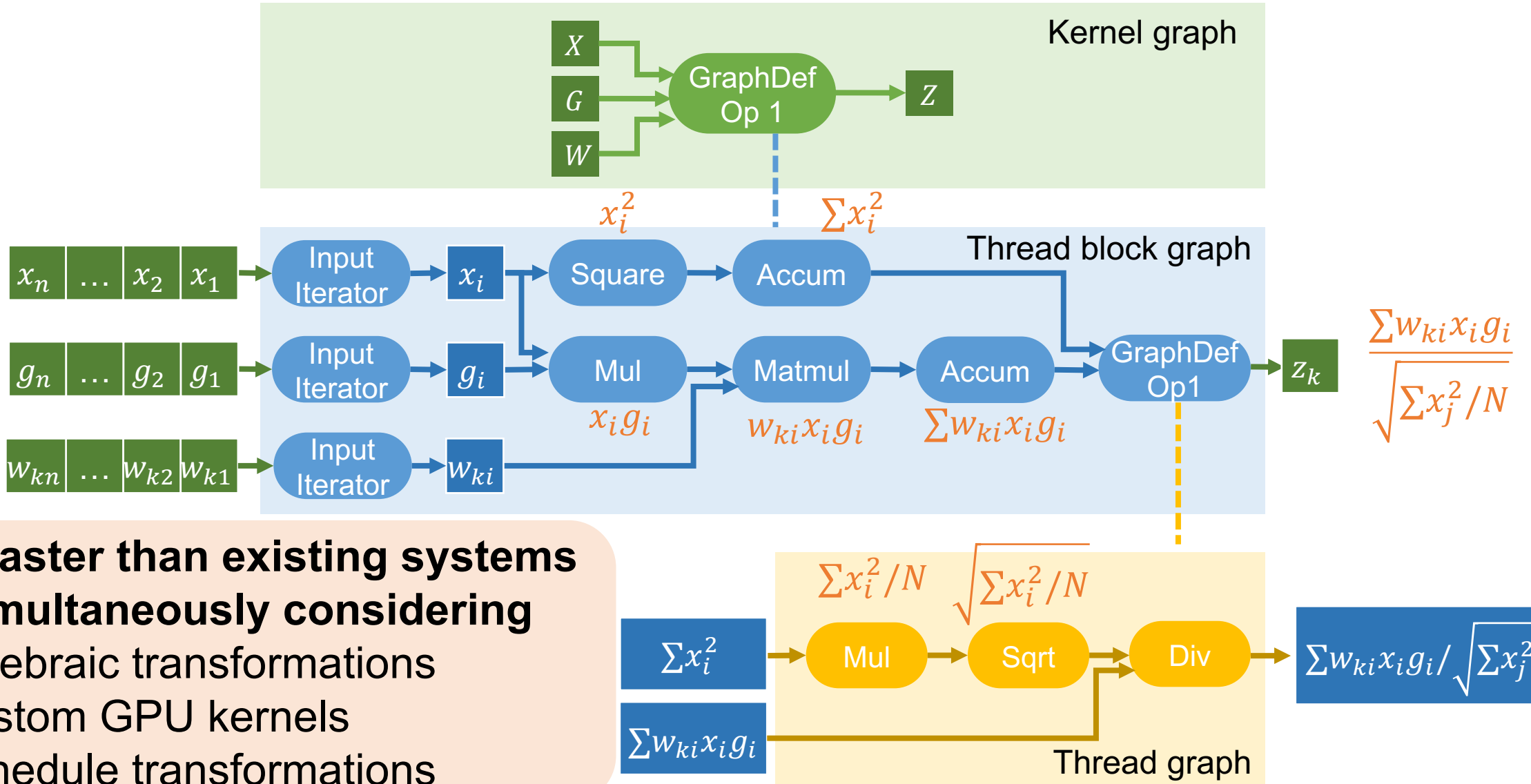
$$y_i = \frac{x_i g_i}{\sqrt{\frac{1}{N} \sum_j x_j^2}}$$

$$z_i = \sum_k w_{ik} y_k$$

Performance issues:

1. No shared memory reuse
2. Kernel launch overhead

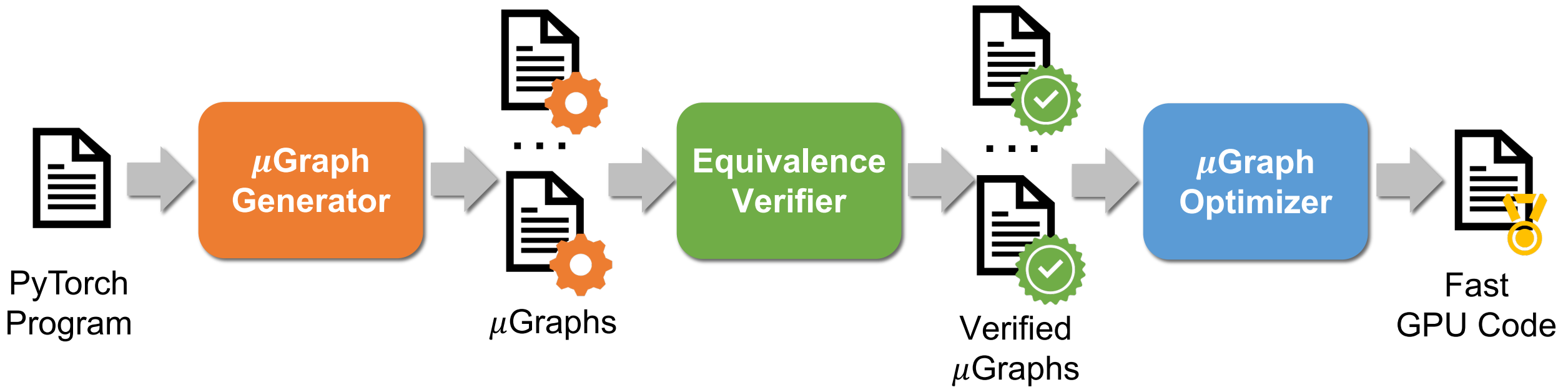
Best Discovered μ Graph for RMSNorm & MatMul



2.2x faster than existing systems by simultaneously considering

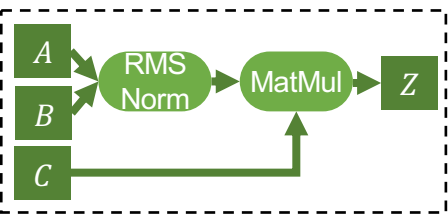
- Algebraic transformations
- Custom GPU kernels
- Schedule transformations

Mirage Overview

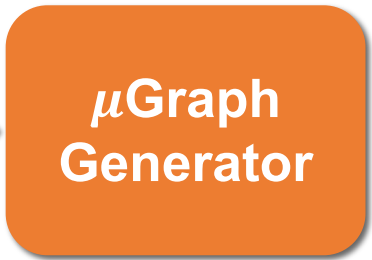


μ Graph Generator

Consider all possible μ Graphs using available operators

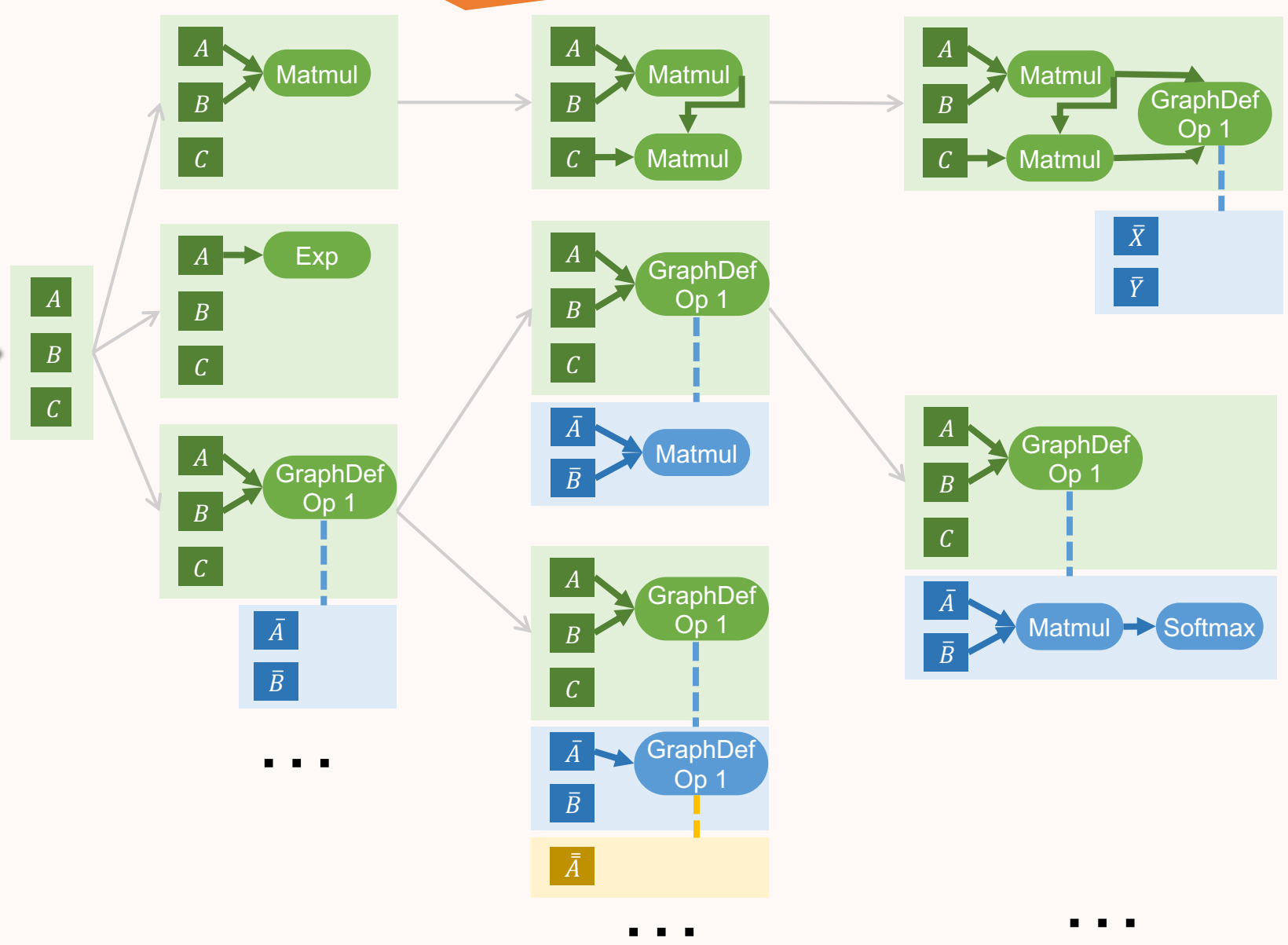


PyTorch Program



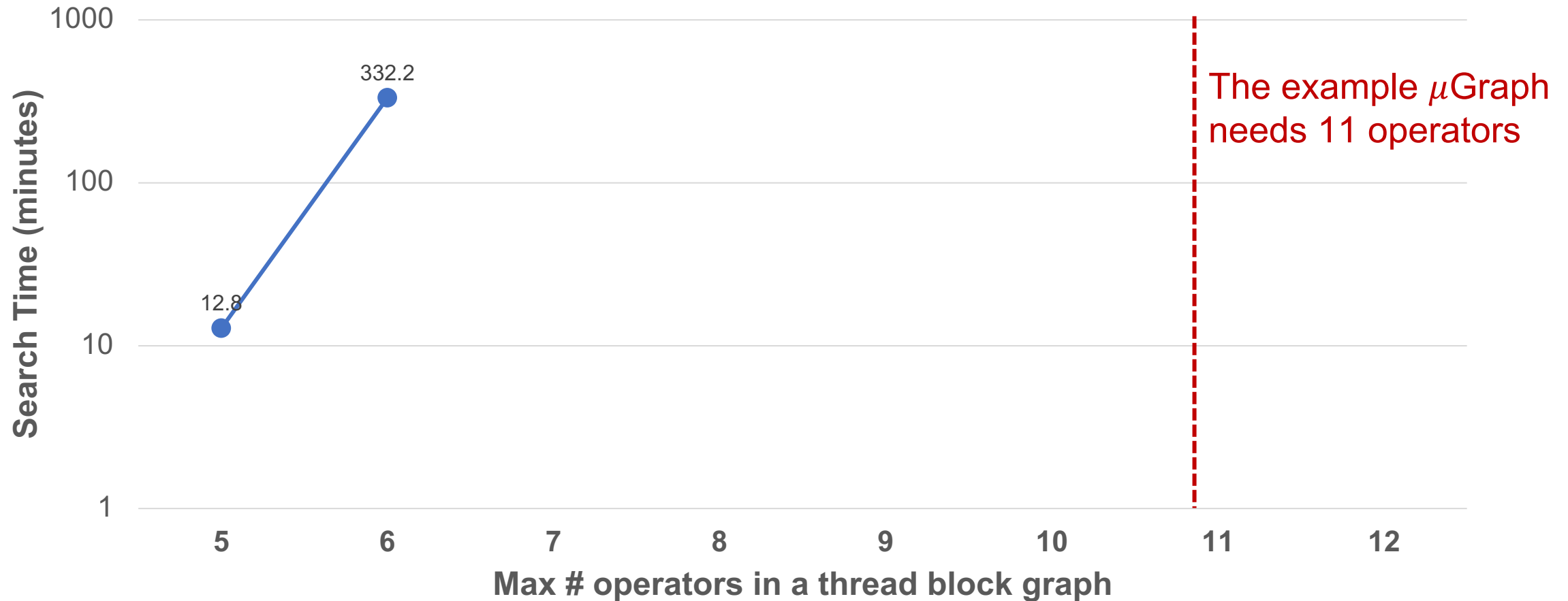
GPU-level	SM-level	Thread-level
Matmul	Matmul	Add
Norm	Mul	Mul
Exp	Sum	Div
...

Operators at the kernel, thread block, and thread levels





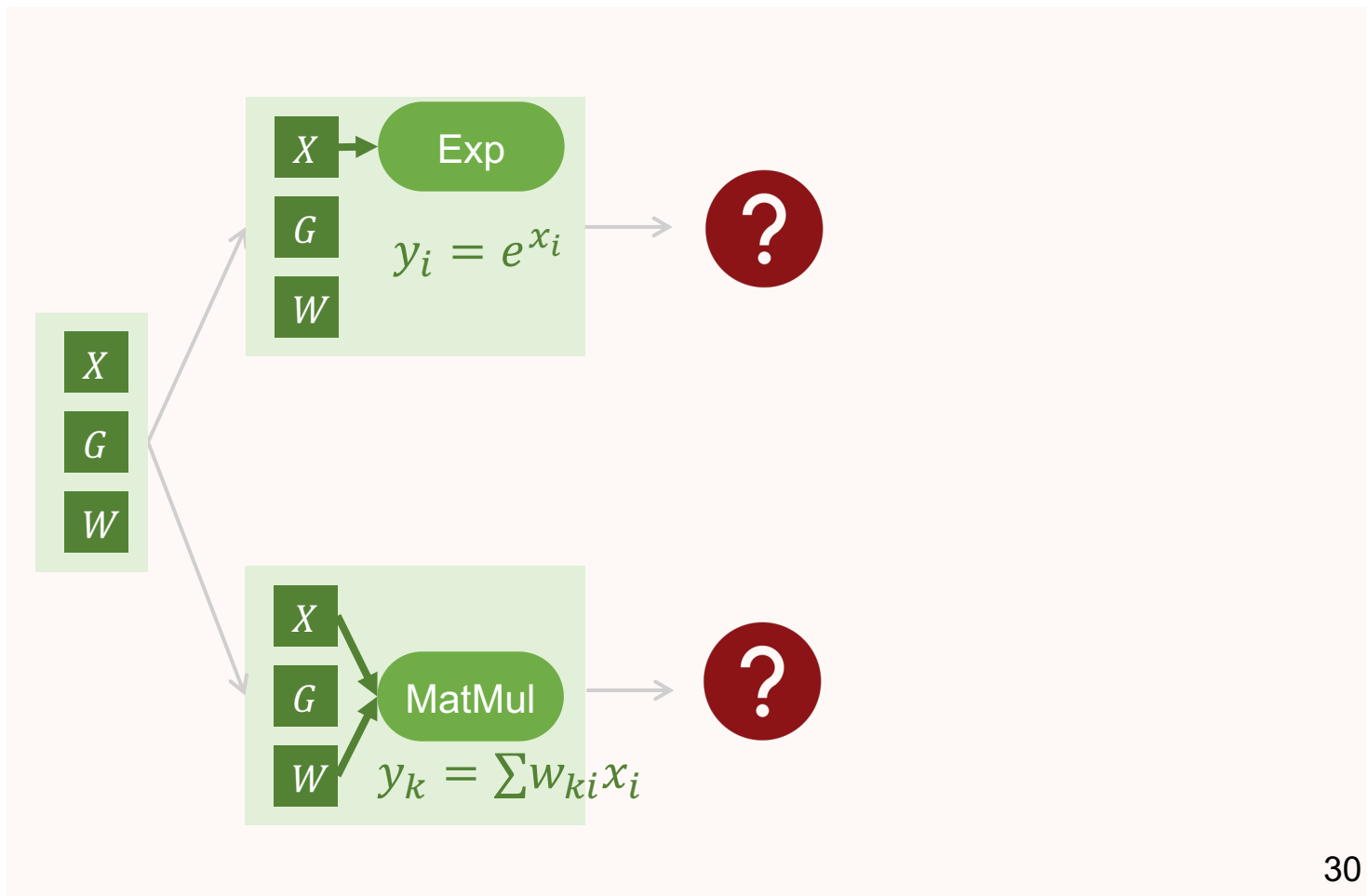
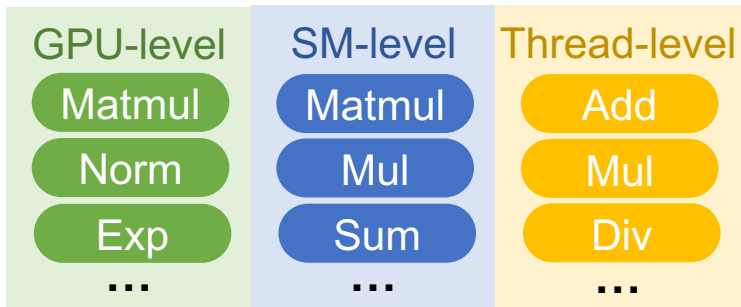
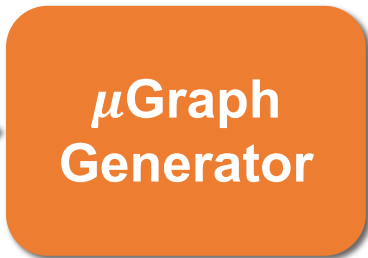
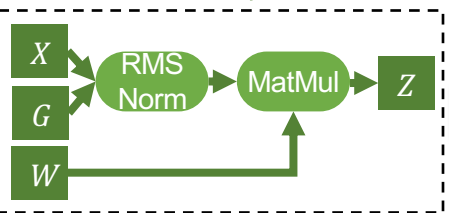
Challenge: Extremely Large Search Space





Can Algebraic Properties Guide Search?

$$z_k = \frac{\sum w_{ki} x_i g_i}{\sqrt{\sum x_j^2 / N}}$$

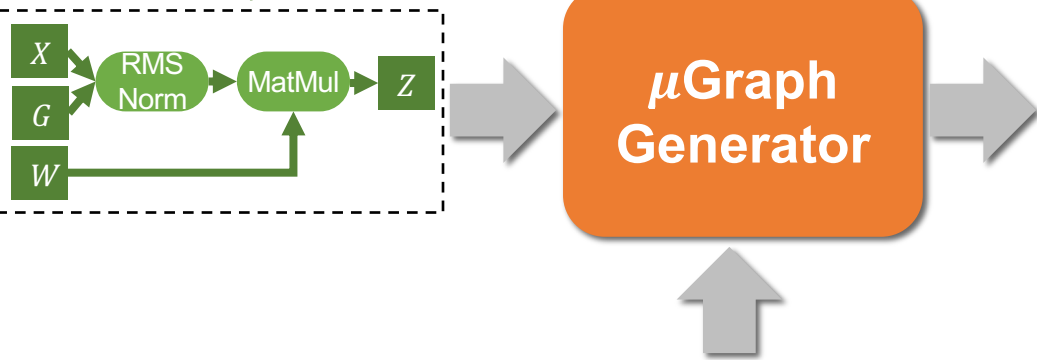




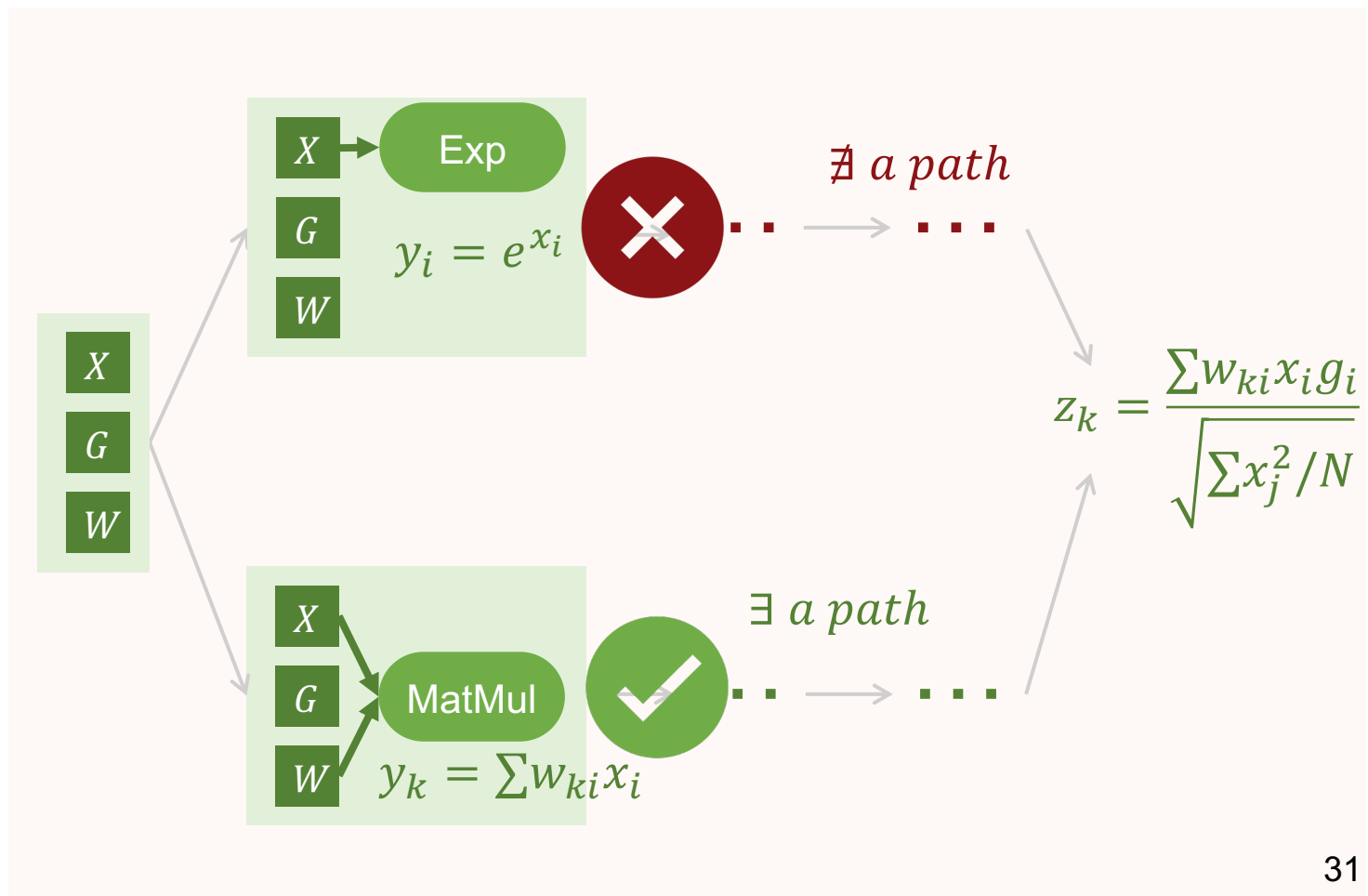
Can Algebraic Properties Guide Search?

Key idea: prune candidates that do not have a path to desired computation using given operators

$$z_k = \frac{\sum w_{ki} x_i g_i}{\sqrt{\sum x_j^2 / N}}$$



GPU-level	SM-level	Thread-level
Matmul	Matmul	Add
Norm	Mul	Mul
Exp	Sum	Div
...

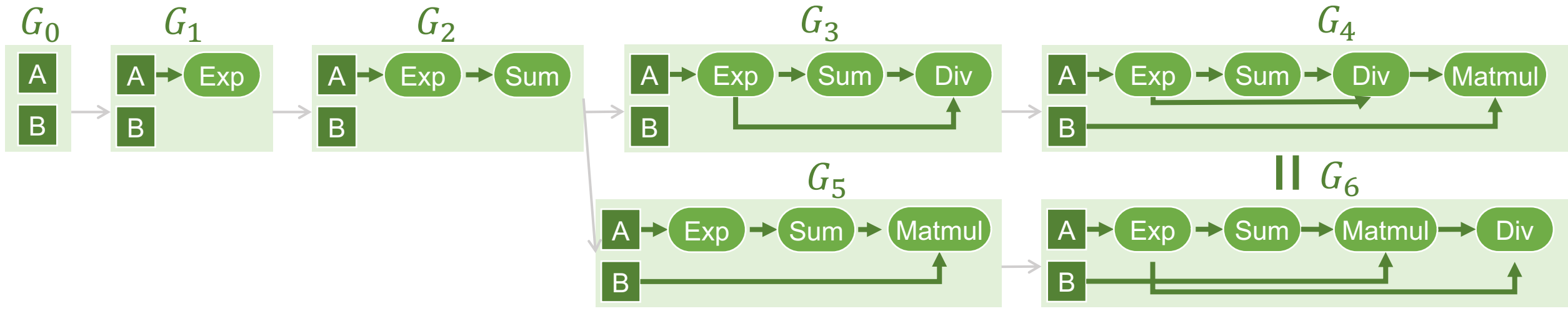




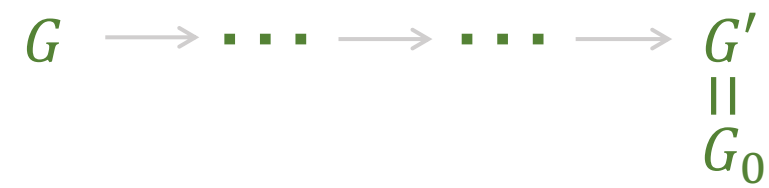
Abstract Expression

Goal #1. capture subgraph relations: $E(G_0) \preceq E(G_1) \preceq E(G_2) \preceq E(G_3) \preceq E(G_4)$

Goal #2. capture graph equivalence: $E(G_4) = E(G_6)$



\exists a path from G to a μ Graph G' equivalent to the input G_0 ?



$\Rightarrow E(G) \preceq E(G_0) ?$

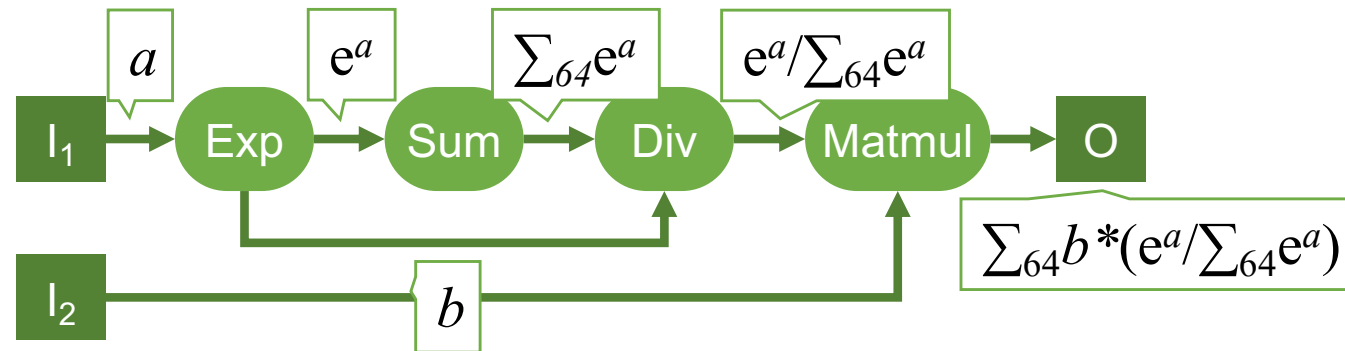


Abstract Expression: An Implementation

Represent a tensor's computation by abstracting away index details

- E.g., $C = A \times B$: $c_{ij} = \sum_{k=1}^{64} a_{ik} b_{kj} \Rightarrow$ abstract expression is $c = \sum_{64} ab$

Recursively compute abstract expressions





Abstract Expression: An Implementation

Mirage uses first-order logic to reason about two relations

Goal 1: subexpression

Subexpression Axioms A_{sub}

$\forall x, y. \text{subexpr}(x, \text{add}(x, y))$	
$\forall x, y. \text{subexpr}(x, \text{mul}(x, y))$	
$\forall x, y. \text{subexpr}(x, \text{div}(x, y))$	
$\forall x, y. \text{subexpr}(y, \text{div}(x, y))$	
$\forall x. \text{subexpr}(x, \text{exp}(x))$	
$\forall x, i. \text{subexpr}(x, \text{sum}(i, x))$	
$\forall x. \text{subexpr}(x, x)$	reflexivity
$\forall x, y, z. \text{subexpr}(x, y) \wedge \text{subexpr}(y, z) \rightarrow \text{subexpr}(x, z)$	transitivity

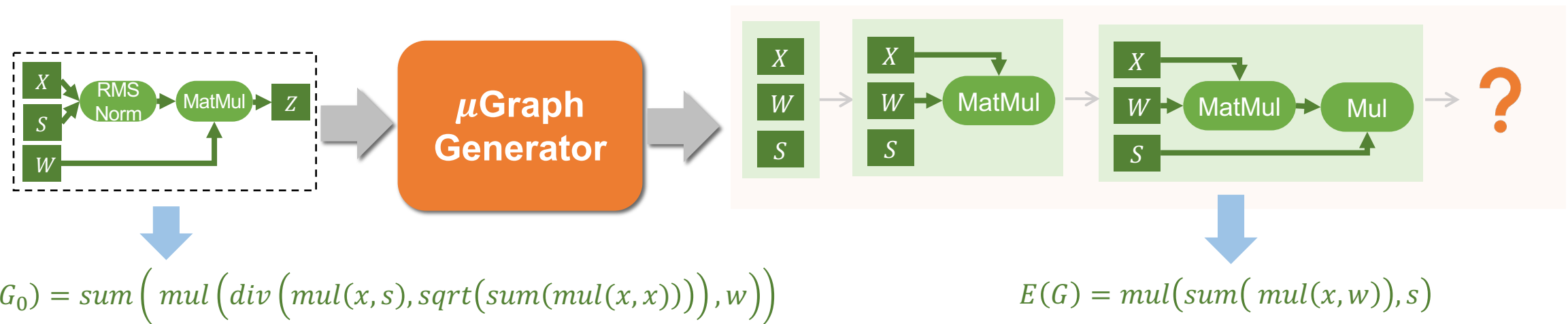
Goal 2: equivalence

Equivalence Axioms A_{eq}

$\forall x, y. \text{add}(x, y) = \text{add}(y, x)$	commutativity
$\forall x, y. \text{mul}(x, y) = \text{mul}(y, x)$	commutativity
$\forall x, y, z. \text{add}(x, \text{add}(y, z)) = \text{add}(\text{add}(x, y), z)$	associativity
$\forall x, y, z. \text{mul}(x, \text{mul}(y, z)) = \text{mul}(\text{mul}(x, y), z)$	associativity
$\forall x, y, z. \text{add}(\text{mul}(x, z), \text{mul}(y, z)) = \text{mul}(\text{add}(x, y), z)$	distributivity
$\forall x, y, z. \text{add}(\text{div}(x, z), \text{div}(y, z)) = \text{div}(\text{add}(x, y), z)$	associativity
$\forall x, y, z. \text{mul}(x, \text{div}(y, z)) = \text{div}(\text{mul}(x, y), z)$	associativity
$\forall x, y, z. \text{div}(\text{div}(x, y), z) = \text{div}(x, \text{mul}(y, z))$	associativity
$\forall x. x = \text{sum}(1, x)$	identity reduction
$\forall x, i, j. \text{sum}(i, \text{sum}(j, x)) = \text{sum}(i * j, x)$	associativity
$\forall x, y, i. \text{sum}(i, \text{add}(x, y)) = \text{add}(\text{sum}(i, x), \text{sum}(i, y))$	associativity
$\forall x, y, i. \text{sum}(i, \text{mul}(x, y)) = \text{mul}(\text{sum}(i, x), y)$	distributivity
$\forall x, y, i. \text{sum}(i, \text{div}(x, y)) = \text{div}(\text{sum}(i, x), y)$	distributivity



Abstract Expression-Guided Search

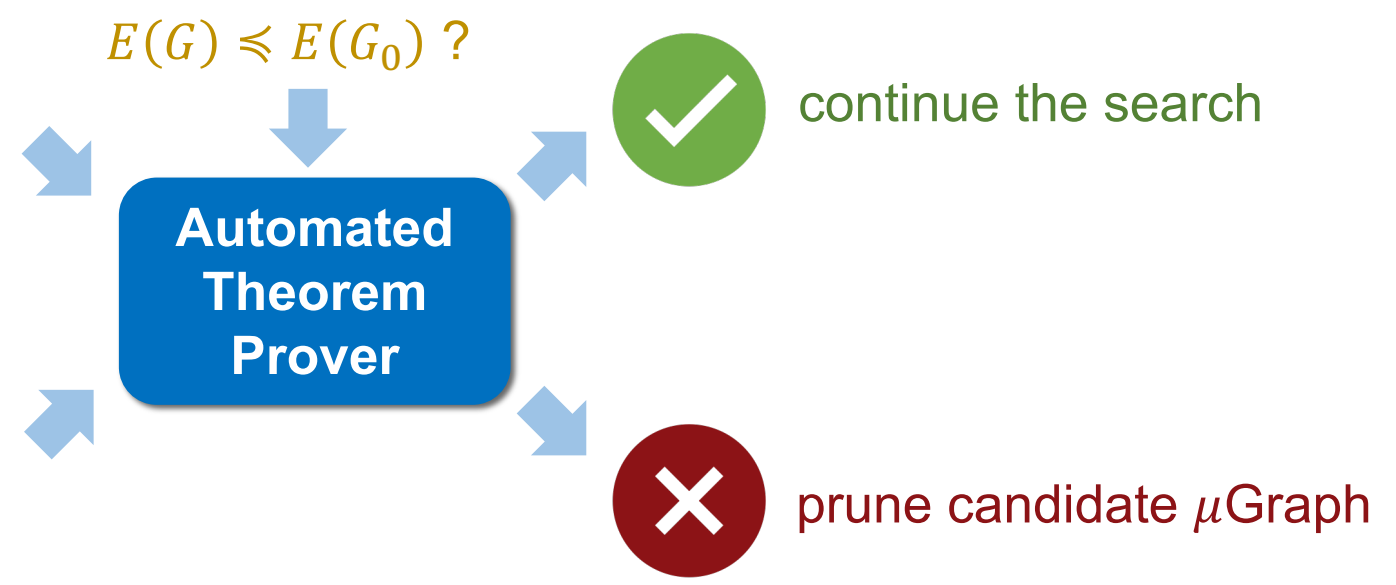


A1. $\forall x, y. x \leqslant mul(x, y)$
 A2. $\forall x, i. x \leqslant sum(i, x)$
 A3. ...

Subexpression axioms

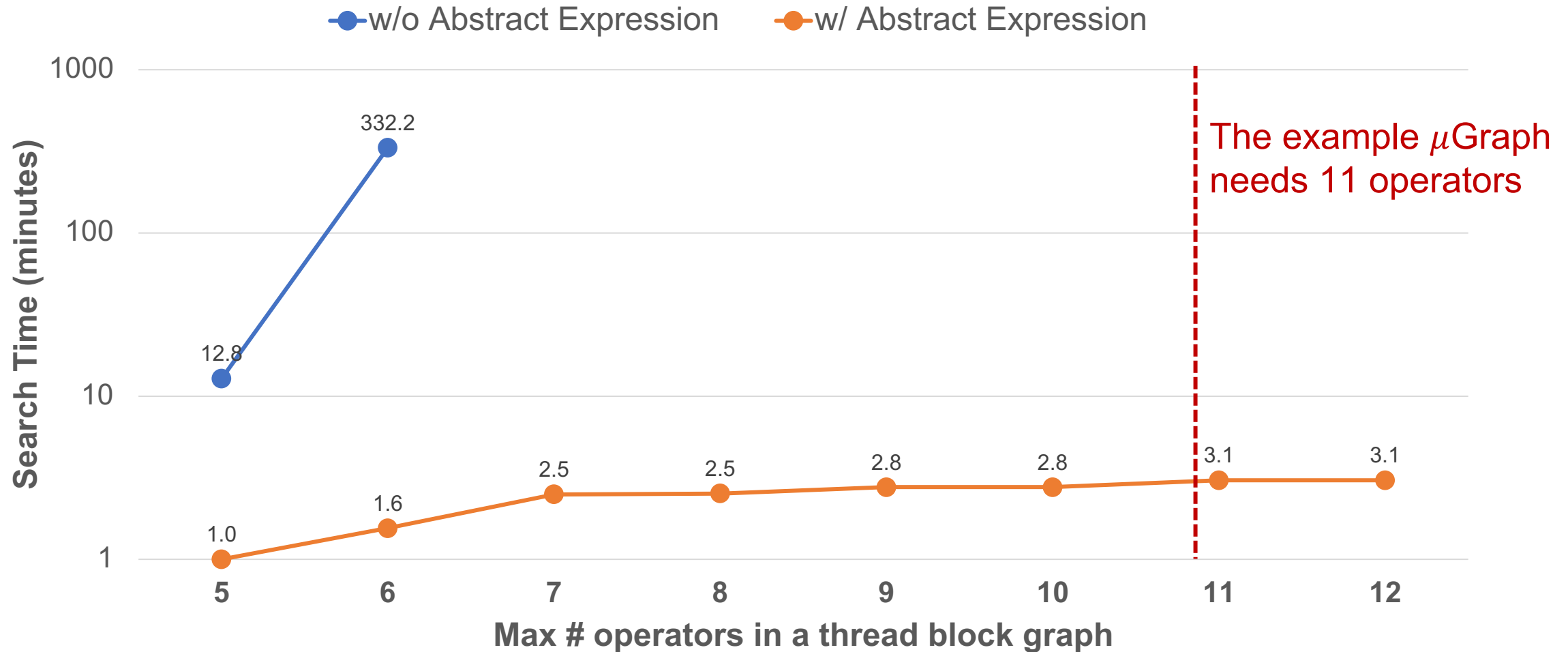
E1. $\forall x, y, z. mul(x, div(y, z)) = div(mul(x, y), z)$
 E2. ...

Equivalence axioms



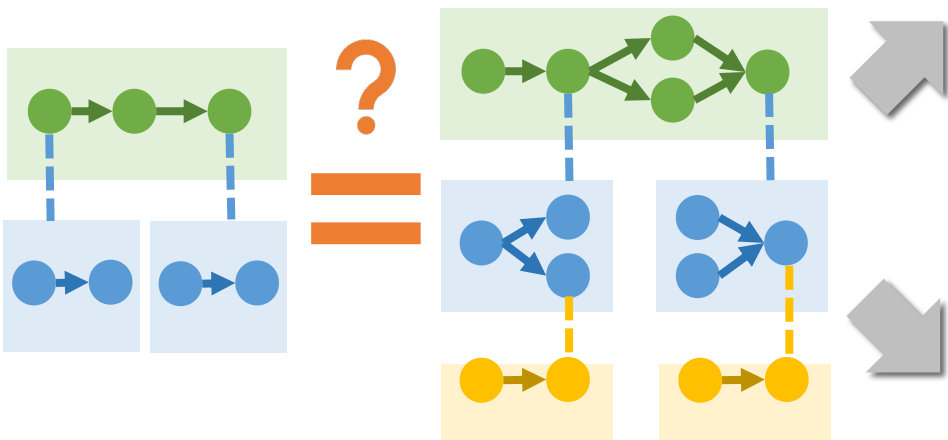


Abstract Expression Significantly Improves Scalability





μGraph Verifiers



Probabilistic
Equivalence
Verifier

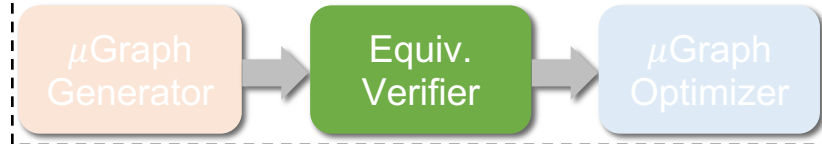
Random tests over finite fields

- Pro: directly support new operators
- Con: limited to a subset of models (i.e., those w/o ReLU)

Solver-based
Equivalence
Verifier

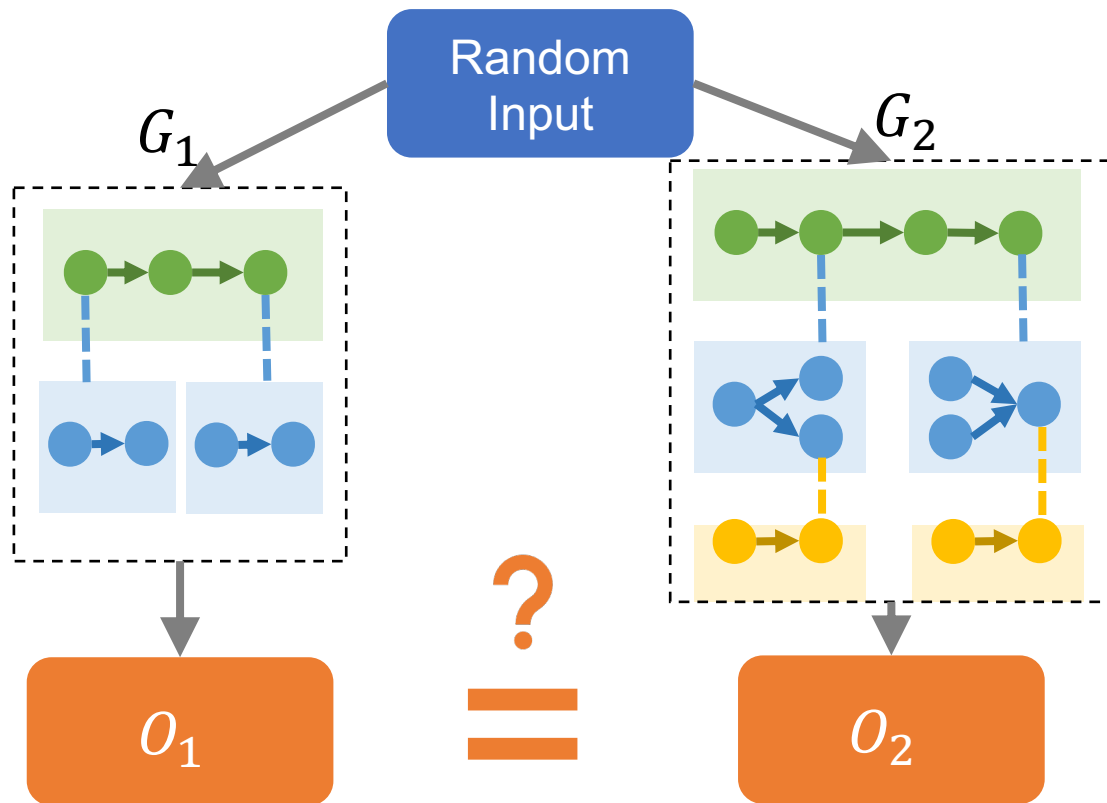
Formal verification

- Pro: support arbitrary ML operators
- Con: manual effort to support new operators



Probabilistic Equivalence Verifier

Idea: use random inputs in *finite fields* to examine μ Graph equivalence

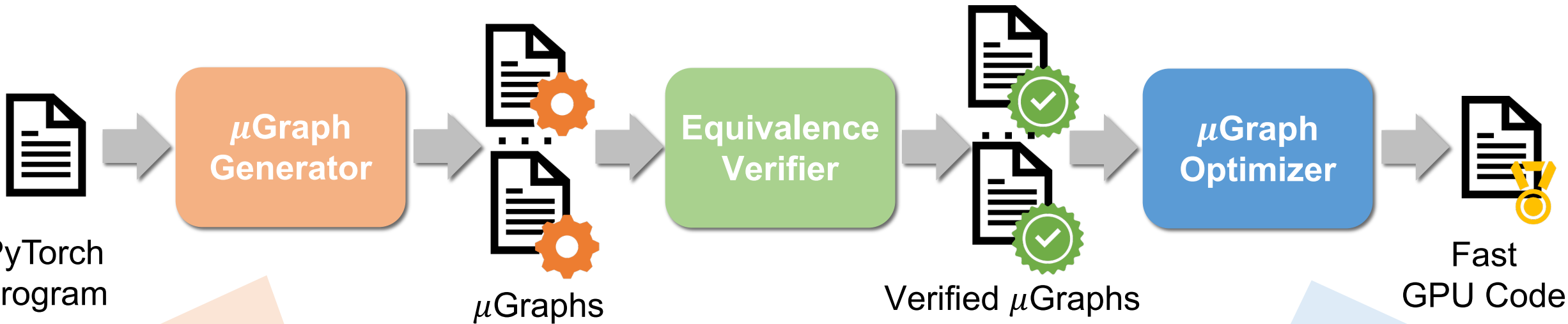


Theorem 1: if G_1 is equivalent to G_2 , then $O_1 = O_2$

Theorem 2: if G_1 is not equivalent to G_2 , then $O_1 \neq O_2$ with a certain probability p^*

* $p \geq \frac{1}{T}$, where T is the size of intermediate tensors

μ Graph Optimizer



Only consider output-alternating optimizations:

- Algebraic transformations
- Kernel instantiation
- Compute organization

Reduce generator's search space

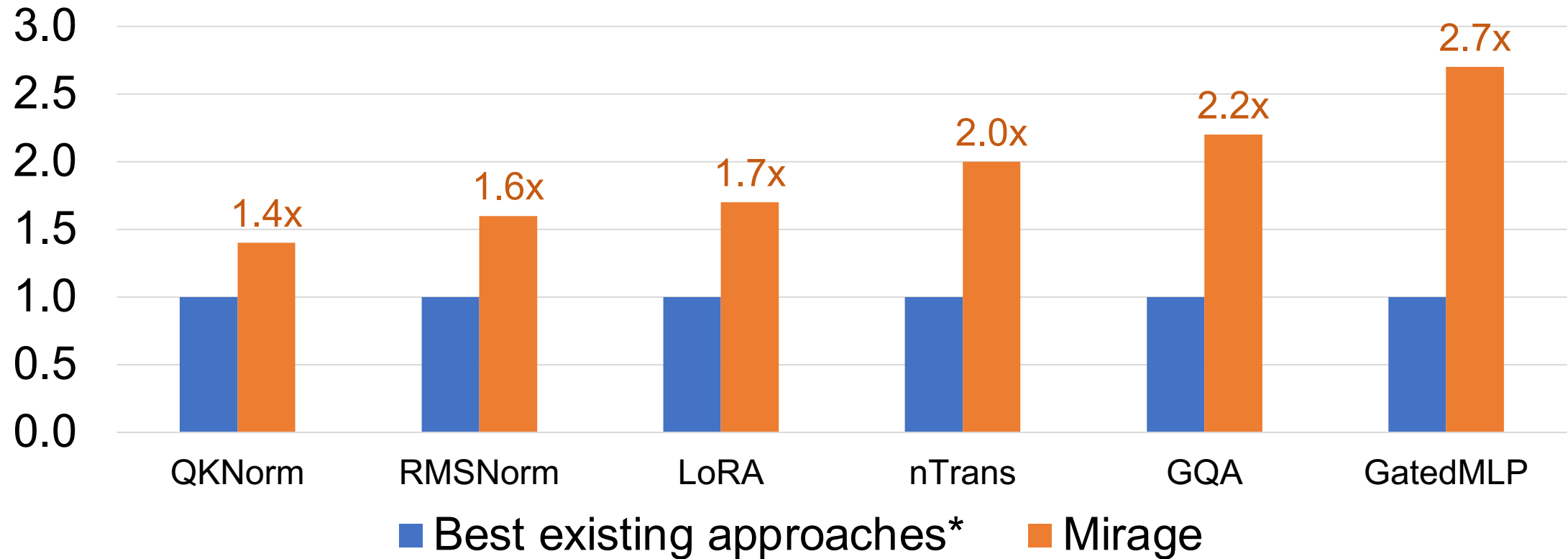
Other optimizations are deferred to μ Graph Optimizer:

- Tensor layouts
- Memory planning
- Operator scheduling

Solve these tasks optimally

Mirage Outperforms Existing Approaches

Relative performance on H100 (higher is better)

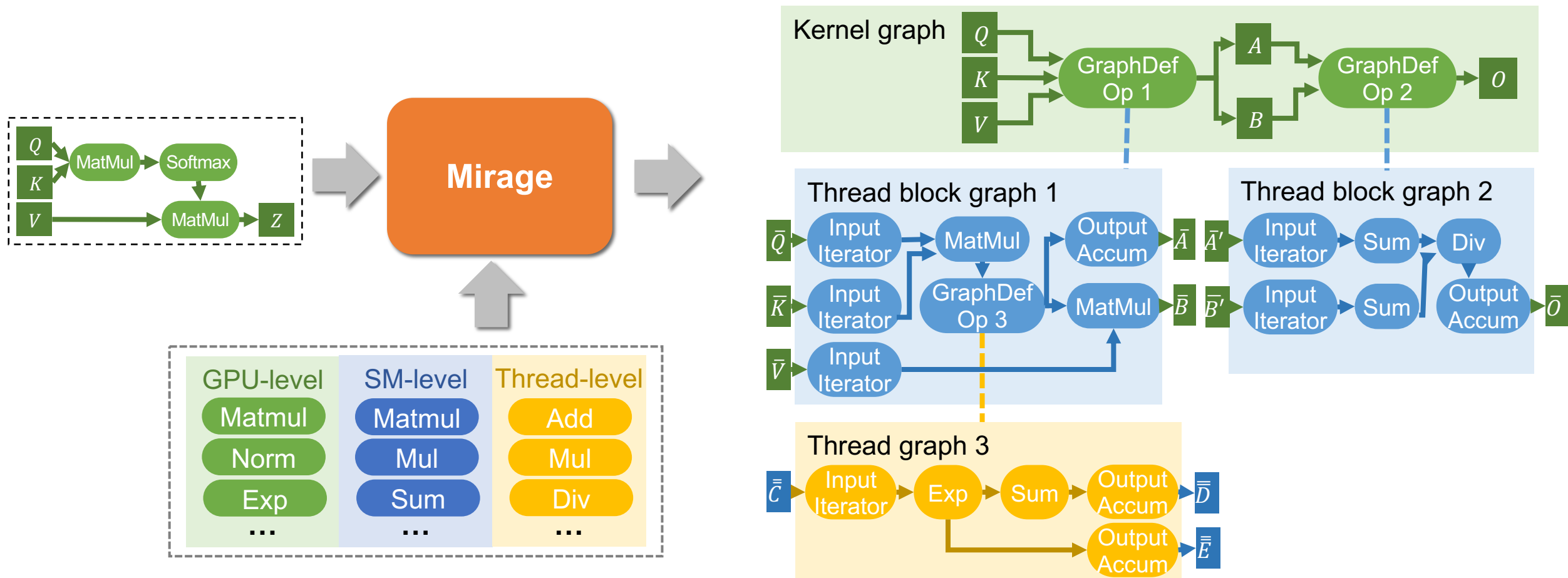


Vender-provided Expert-written Compiler-generated

*Best existing approaches are the best of cuDNN/cuBLAS, FlashAttention, PyTorch, TensorRT, Triton, and TVM.

Mirage Discovers Hardware-Customized μ Graphs

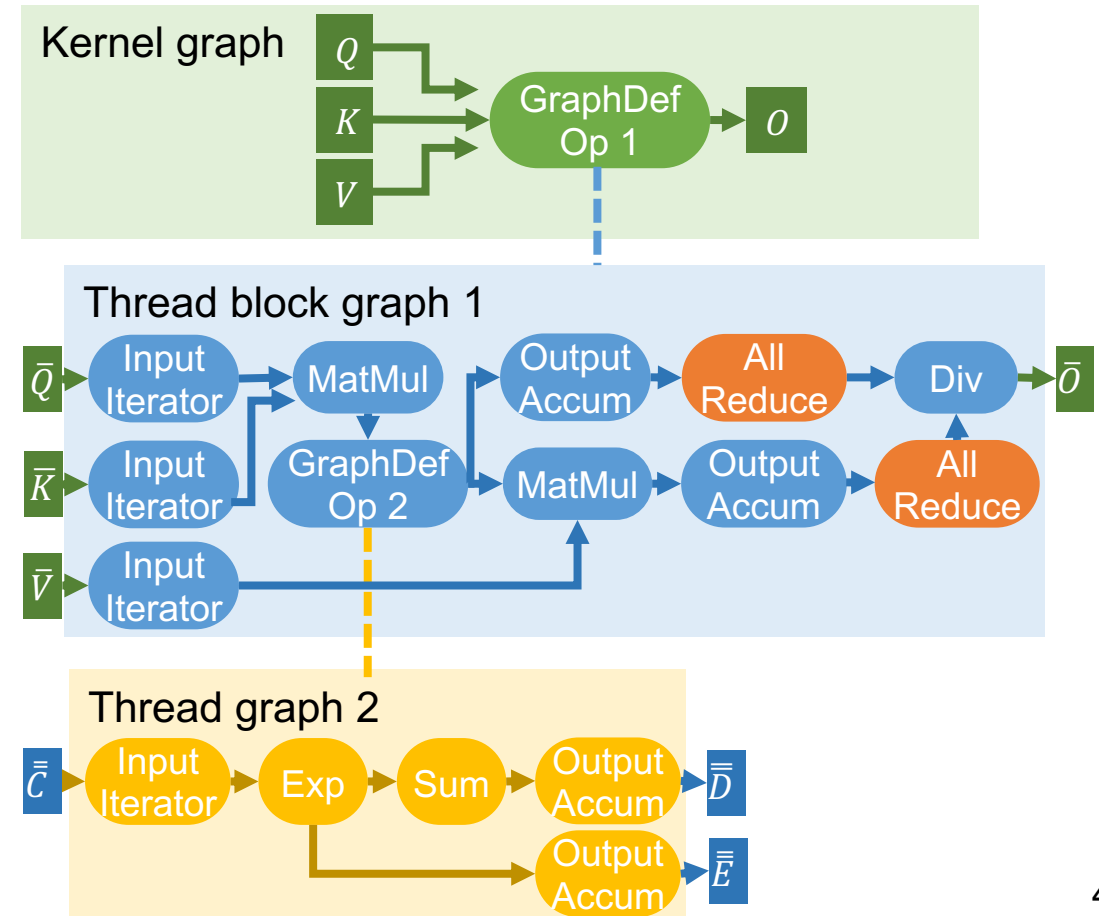
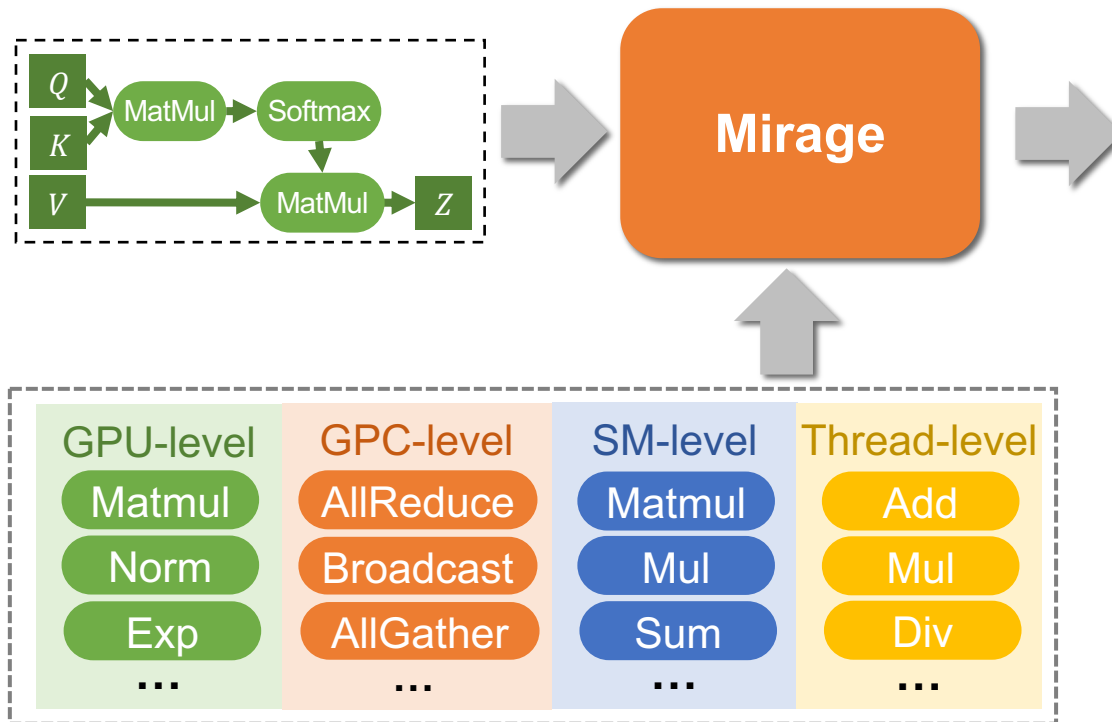
Find μ Graphs similar to expert-written implementations for attention on A100



Mirage Discovers Hardware-Customized μ Graphs

Leverage **GPC-level AllReduce** to accelerate attention on H100

- **2.2x faster** than best existing kernels



Our Research: Superoptimizing ML Systems

