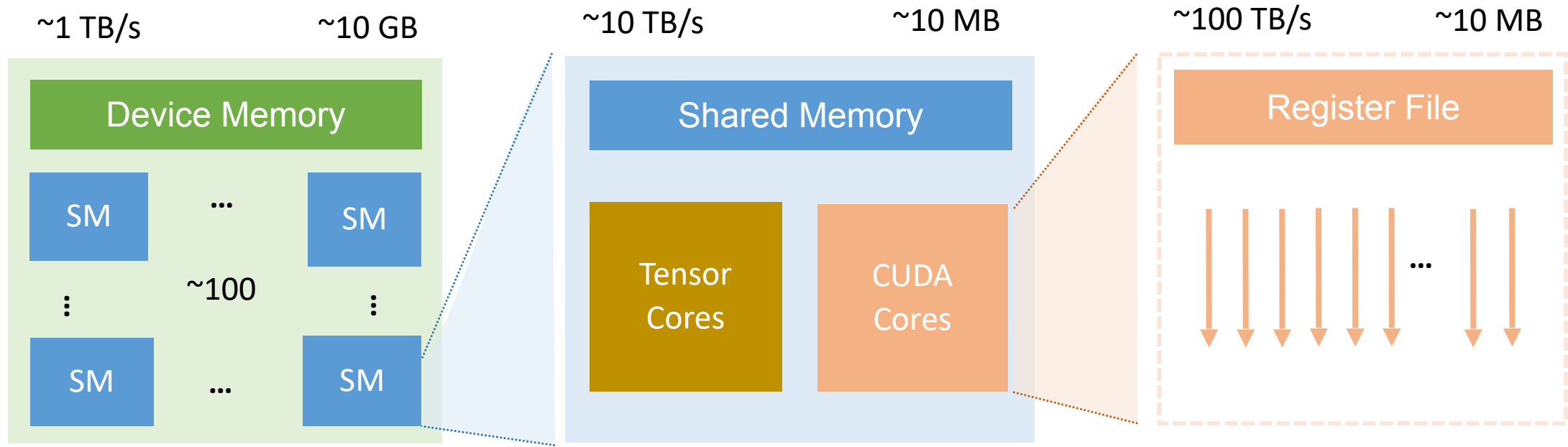


Mirage: A Multi-Level Superoptimizer for Tensor Programs

Mengdi Wu Xinhao Cheng Shengyu Liu Chunan Shi Jianan Ji
Man Kit Ao Praveen Velliengiri Xupeng Miao Oded Padon Zhihao Jia

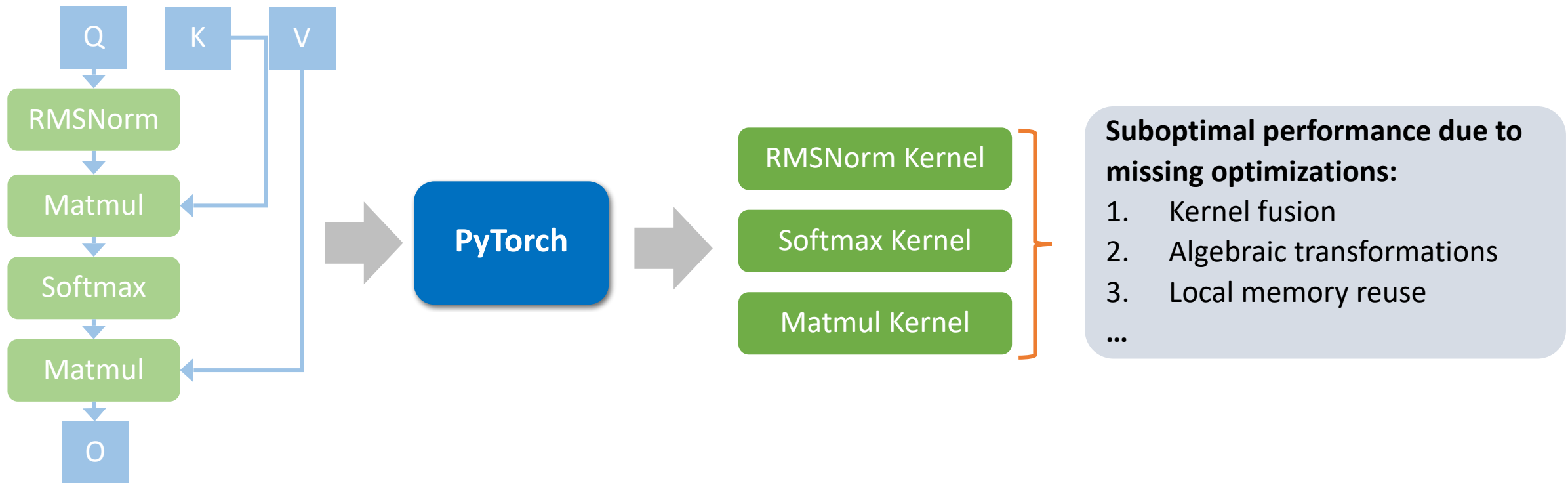


GPUs: Complex Computation and Memory Hierarchy

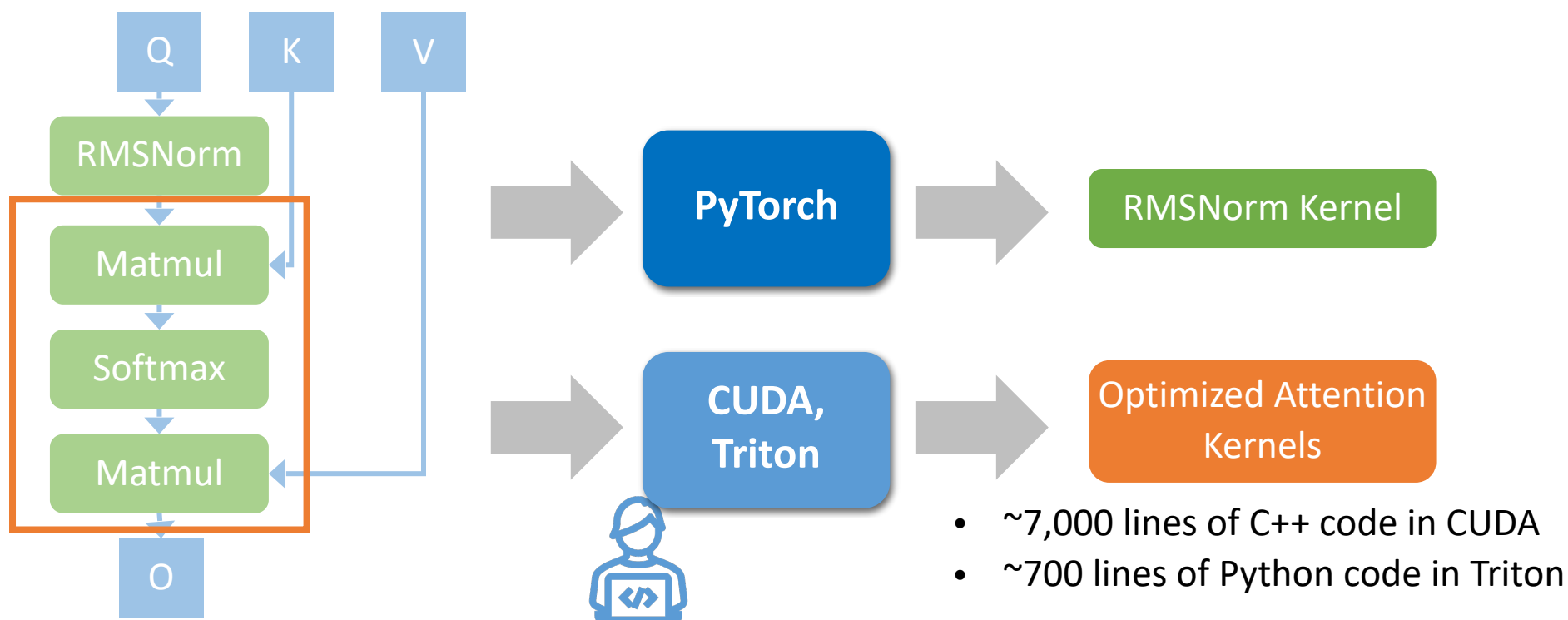


Hard to design high-performance GPU kernels

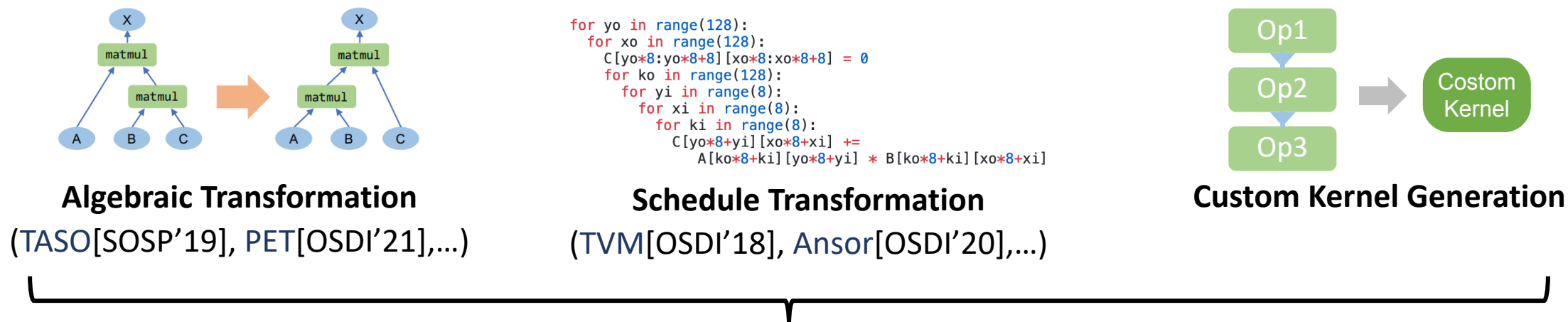
Existing Systems Launch Kernels for Individual Operators



Manually Implement Optimized Kernels for Certain Computation Patterns

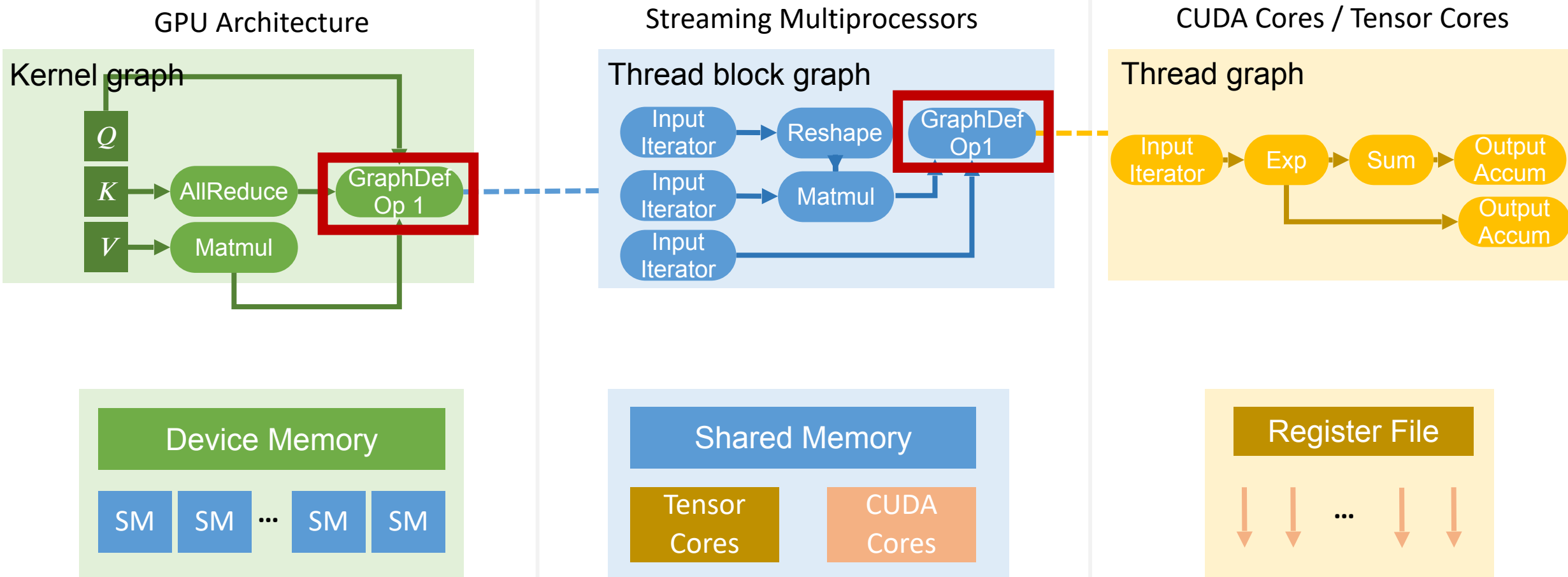


Mirage: A *Multi-Level* Superoptimizer



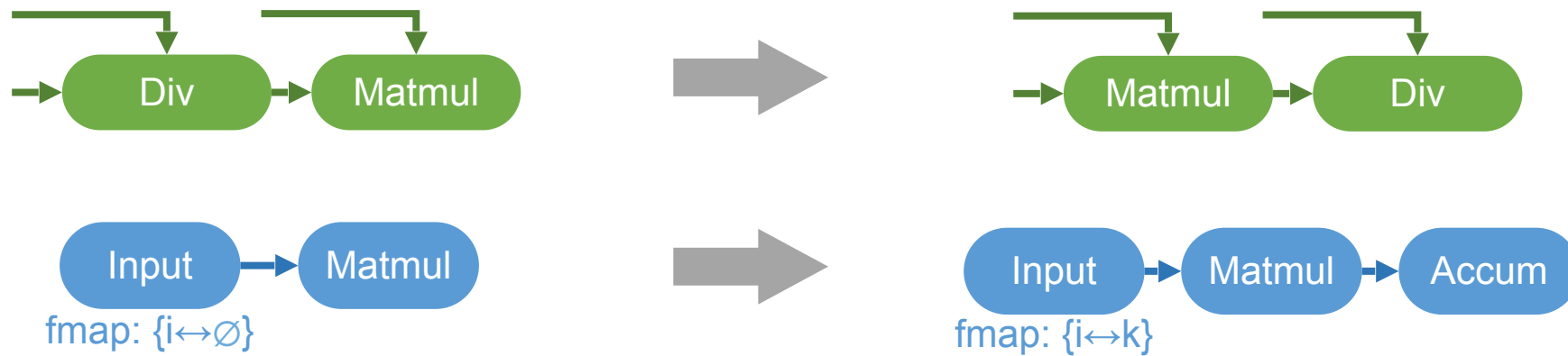
- **Less engineering effort:** 7,000 lines of CUDA code / 700 lines of Triton code → a few lines of Python code in Mirage
- **Better performance:** outperform existing systems by up to 3.3x
- **Easy adaptation:** do not rely on manual implementation

μ Graph: Multi-Level Graph Representation



Challenge: Search Space Exploration

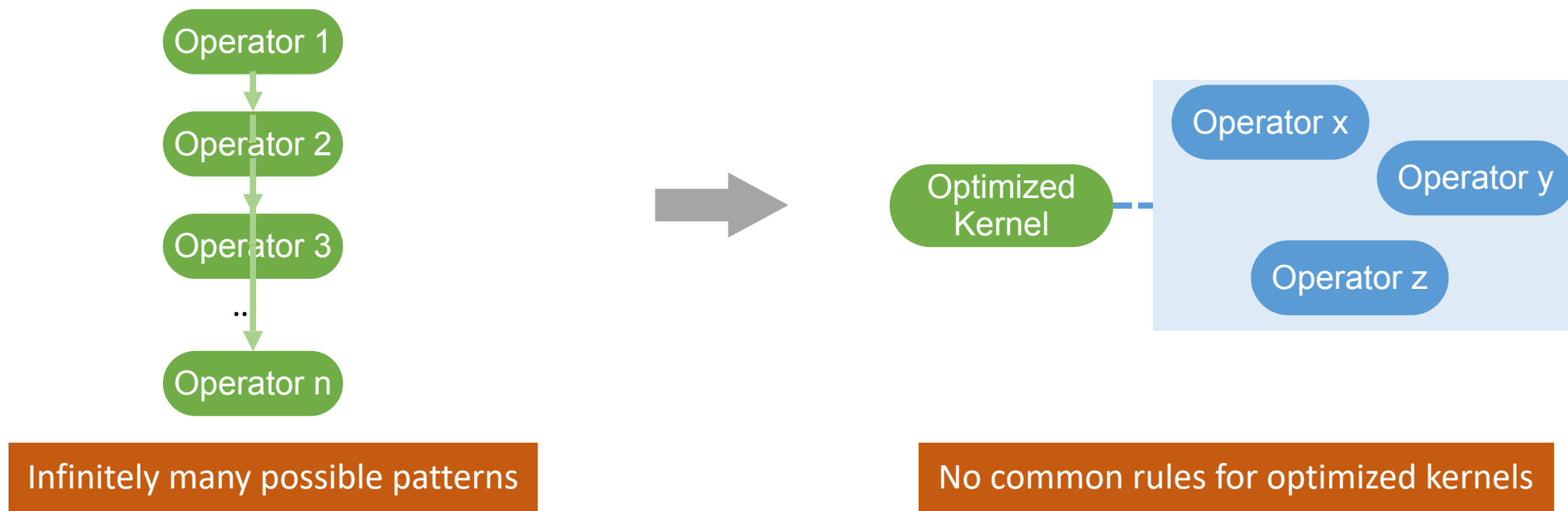
Transformation-based approach: define *transformation rules* $\mathcal{T} : G_{pattern} \mapsto G_{opt}$



Easy for optimizations within one level,
(Transformations can be summarized as *limited* rules)
but ...

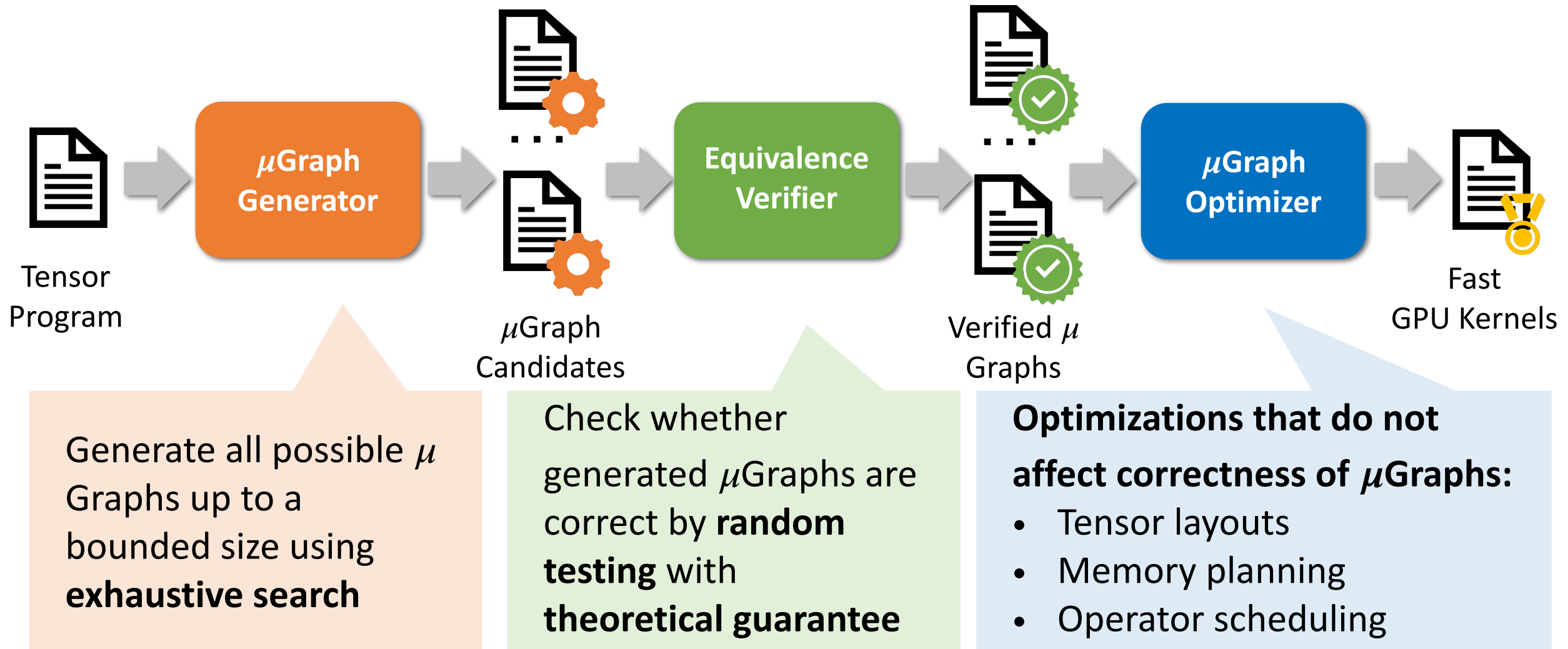
Challenge: Search Space Exploration

Transformation-based approach: define *transformation rules* $\mathcal{T} : G_{pattern} \mapsto G_{opt}$



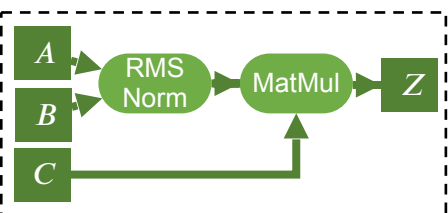
Hard for multi-level optimizations
(Transformations are *dense* and *irregular* functions)

Mirage Overview



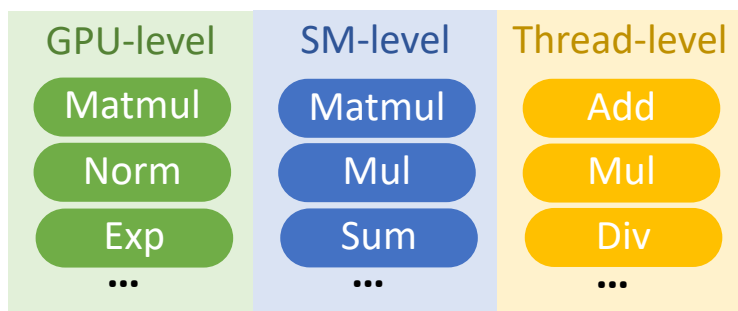
μ Graph Generator

Exhaustively search for all possible μ Graphs using available operators up to a size

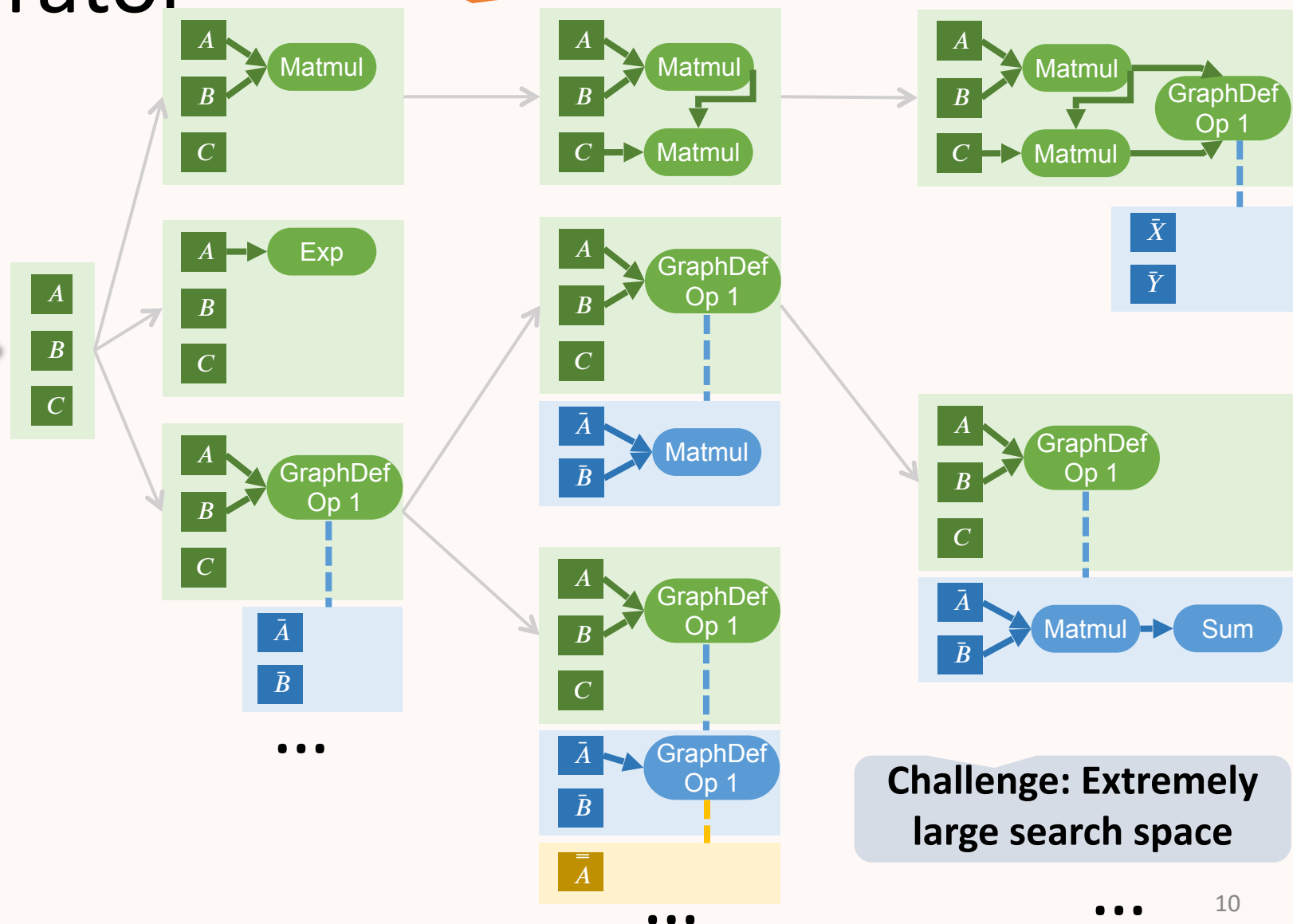


Tensor Program

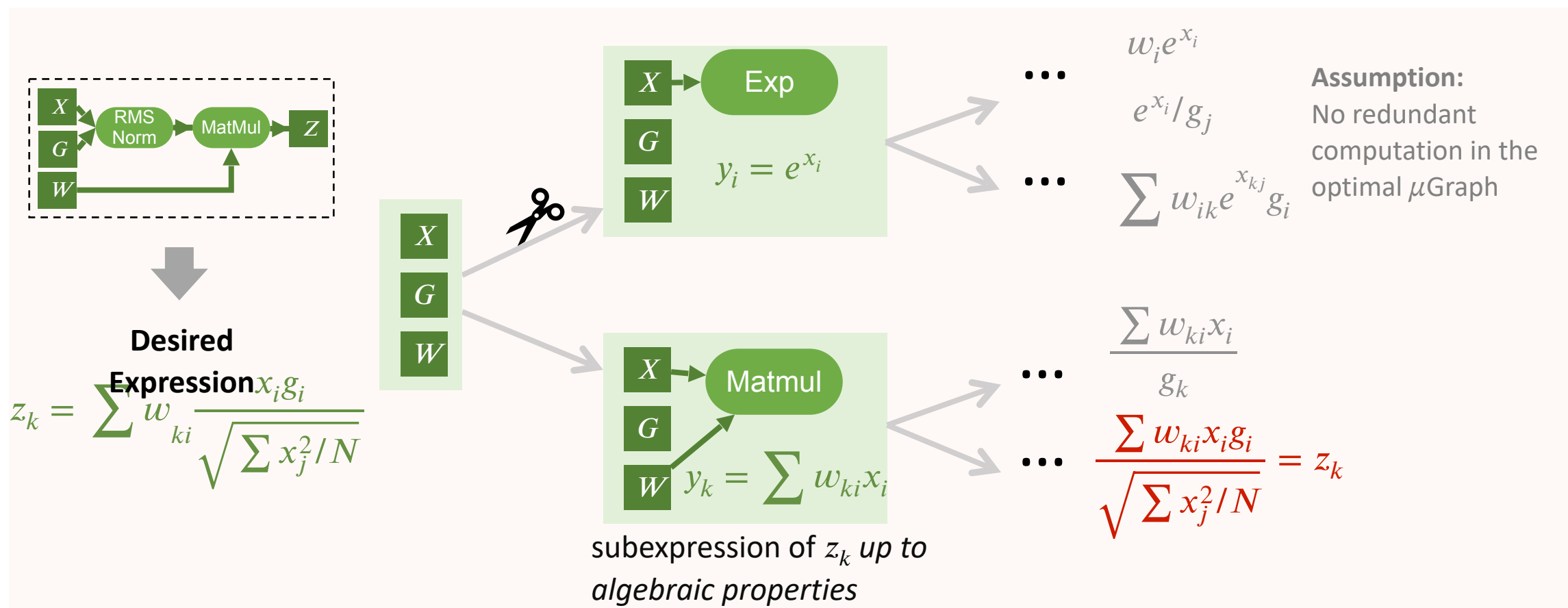
μ Graph Generator



Operators at the kernel, thread block, and thread levels



Expression-Guided Pruning



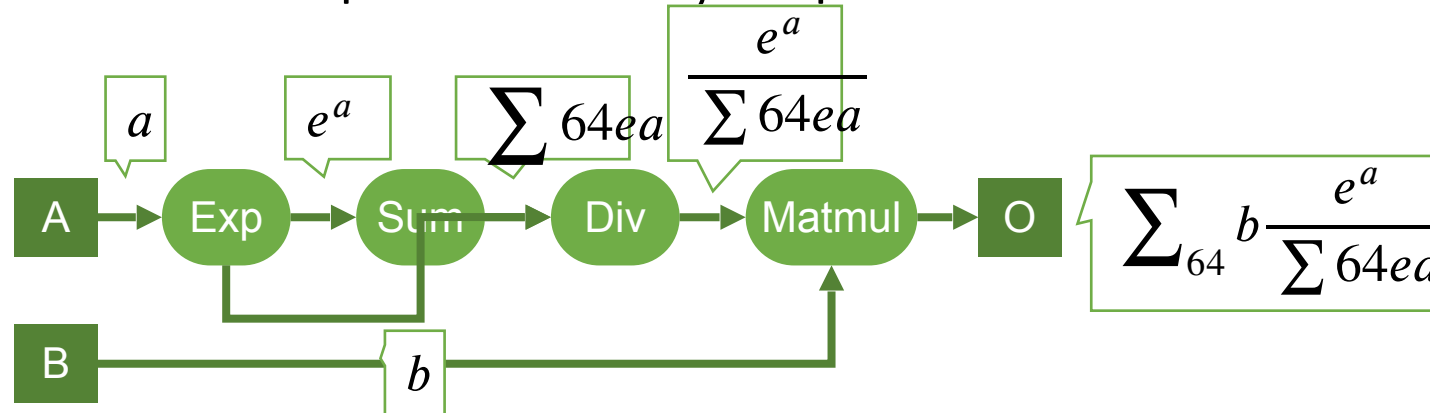
Full-information expressions are too complicated to reason about

Abstract Expression

Key idea: abstract away the index details

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

Recursively compute abstract expressions for a μ Graph



- Capture most semantic information
- Easy to reason about

Axioms for Abstract Expression

Mirage uses first-order logic to reason about abstract expression relations

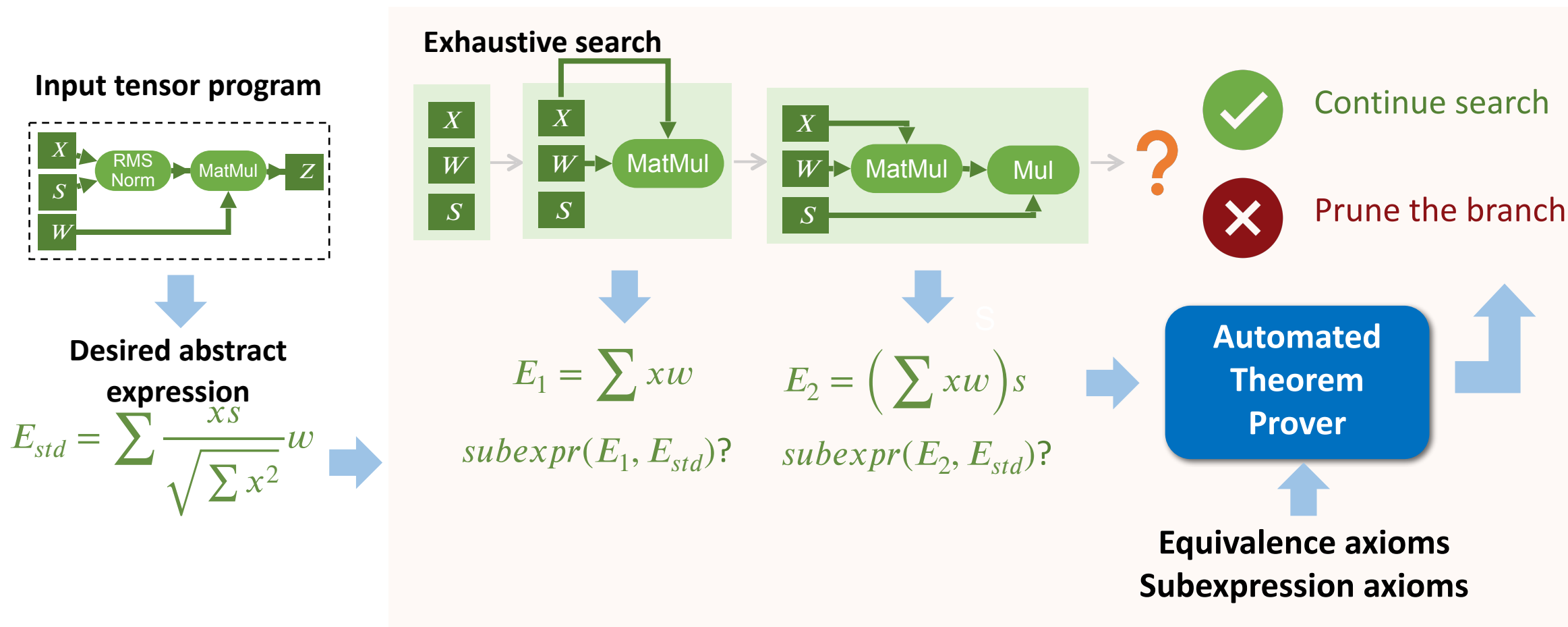
Equivalence axioms

$\forall x, y. \text{add}(x, y) = \text{add}(y, x)$	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

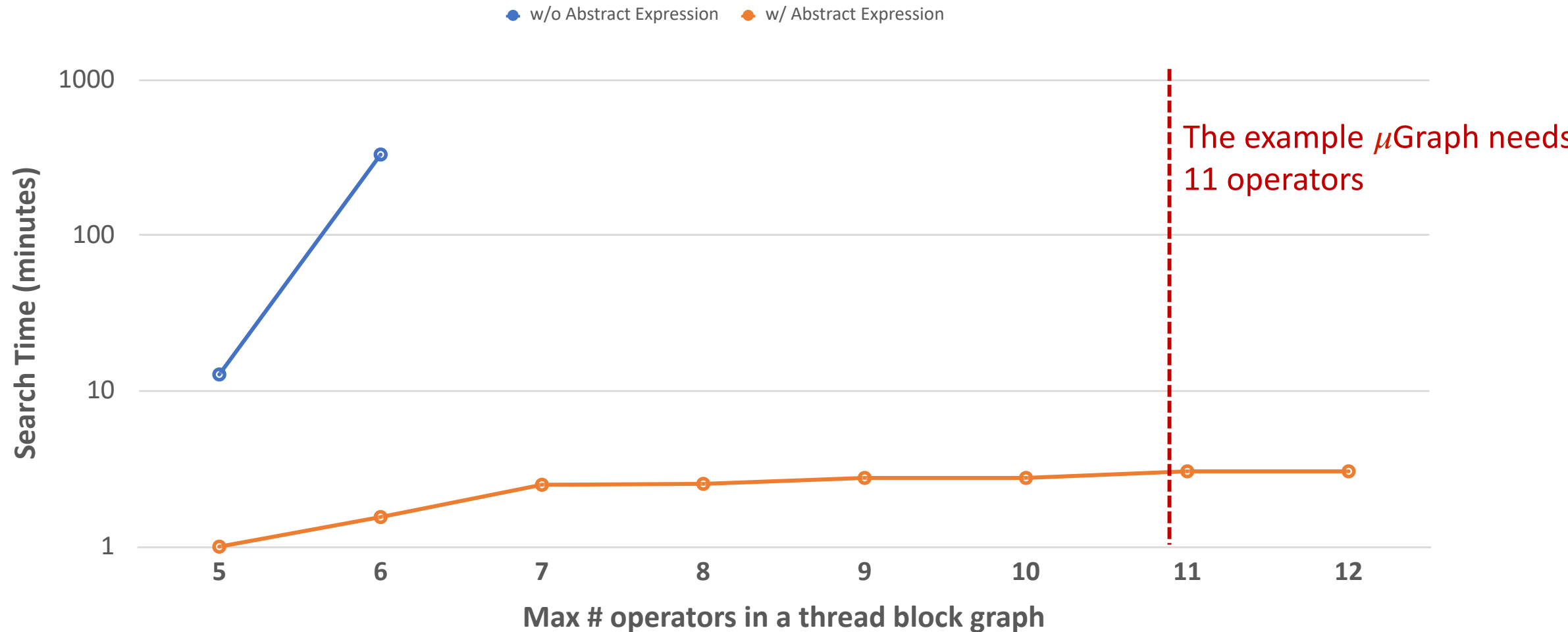
Subexpression axioms

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

Abstract Expression-Guided Search

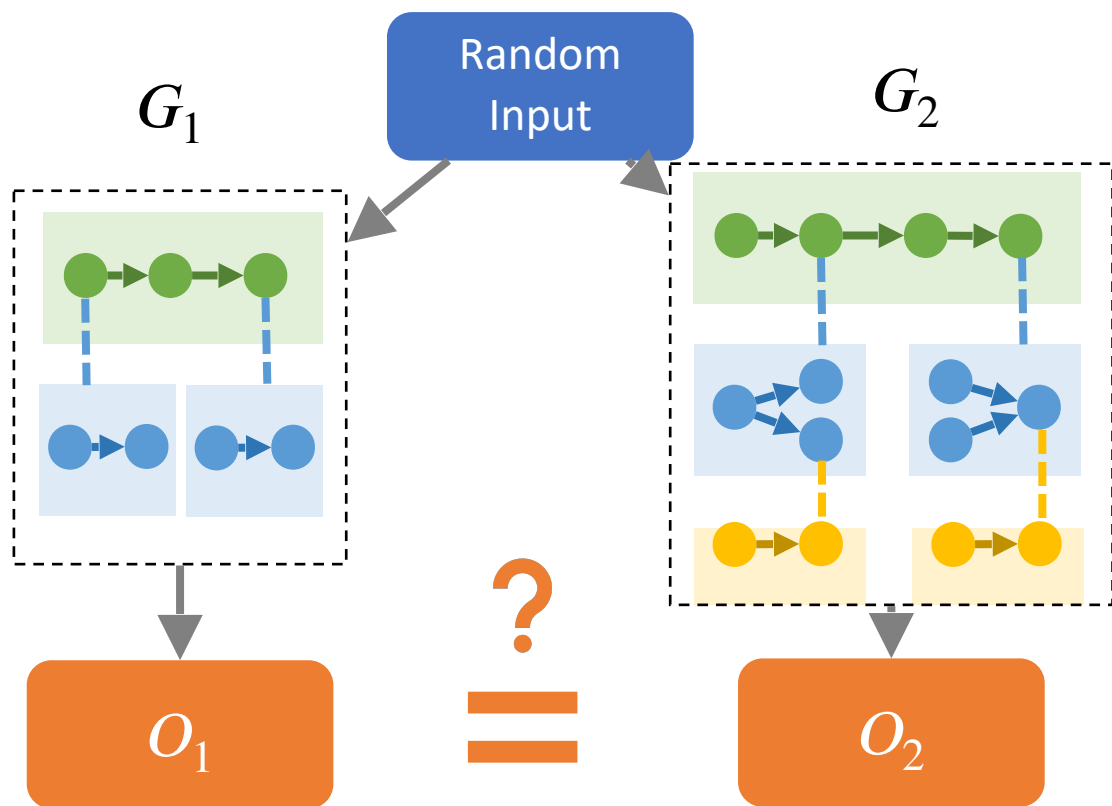


Abstract Expression Significantly Improves Scalability



Probabilistic Equivalence Verifier

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

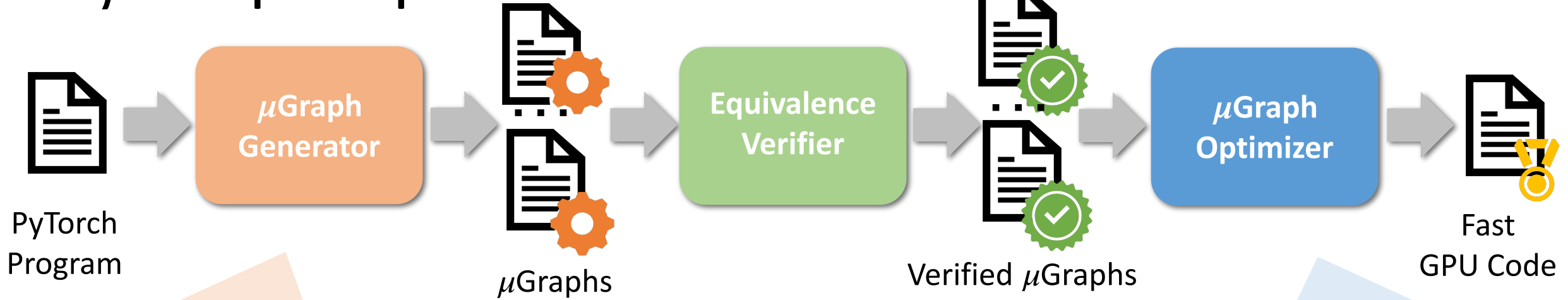


Theorem 1: if G_1 is equivalent to G_2 ,
then $O_1 = O_2$

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

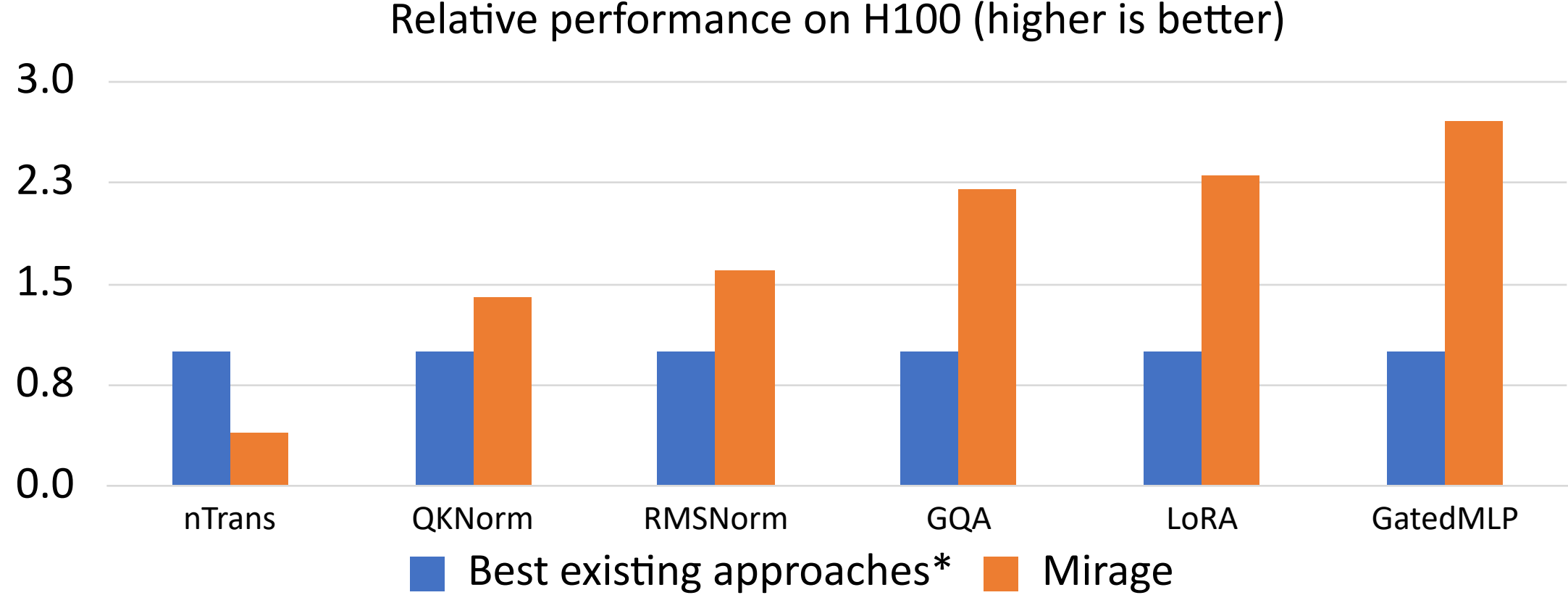
Reduce generator's search space

Other optimizations are deferred to μ Graph Optimizer:

- Tensor layouts
- Memory planning
- Operator scheduling

Solve these tasks optimally

Mirage Outperforms Existing Approaches



Vender-provided

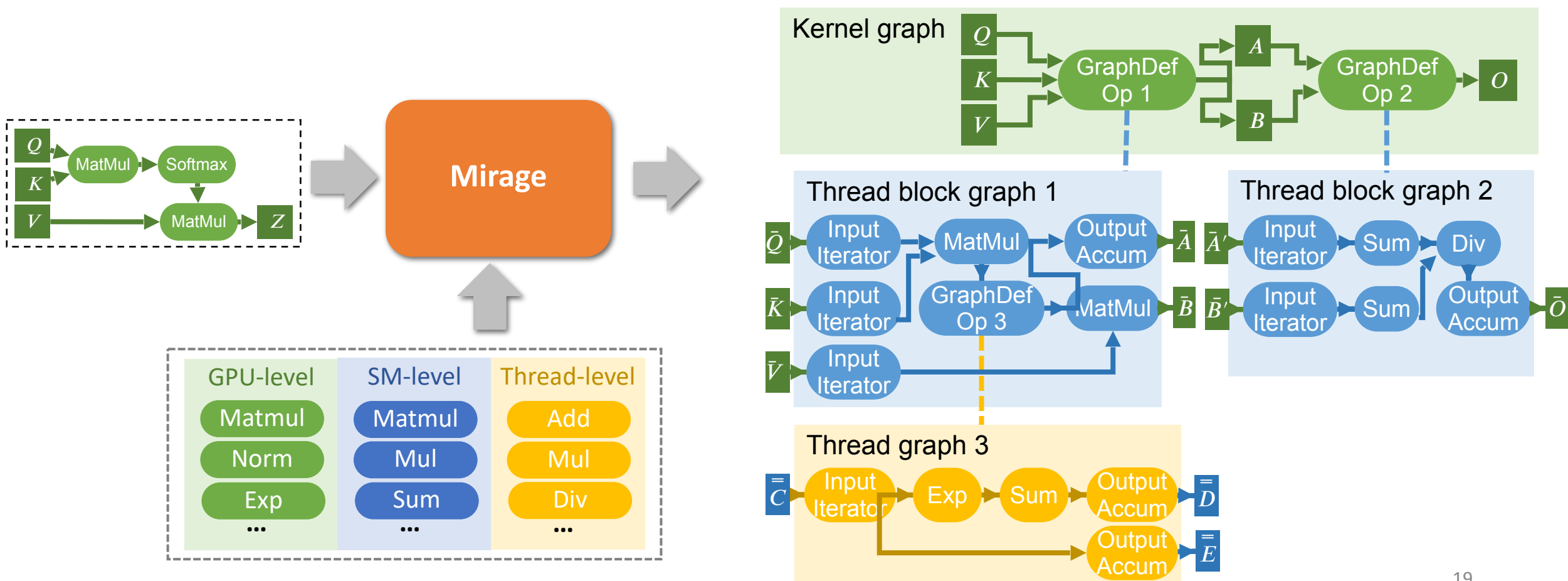
Expert-written

Compiler-generated

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

Mirage Discovers Hardware-Customized μ Graphs

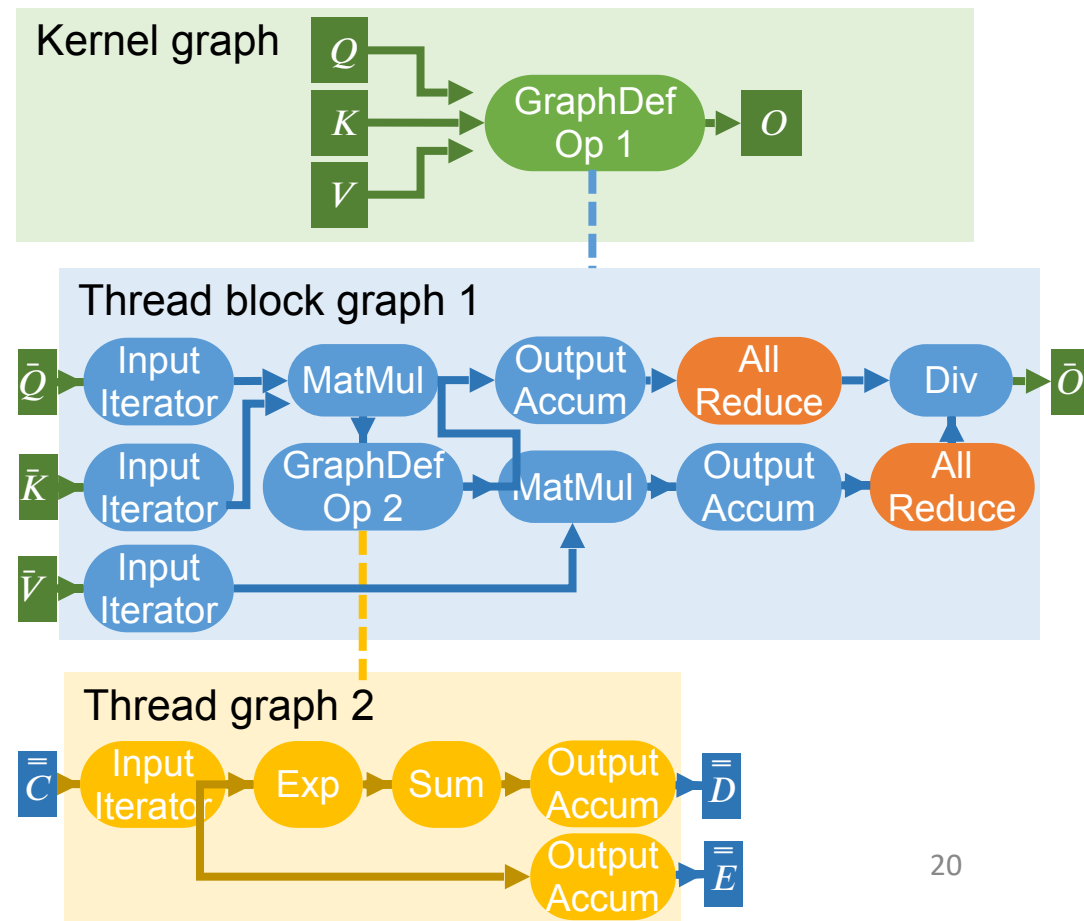
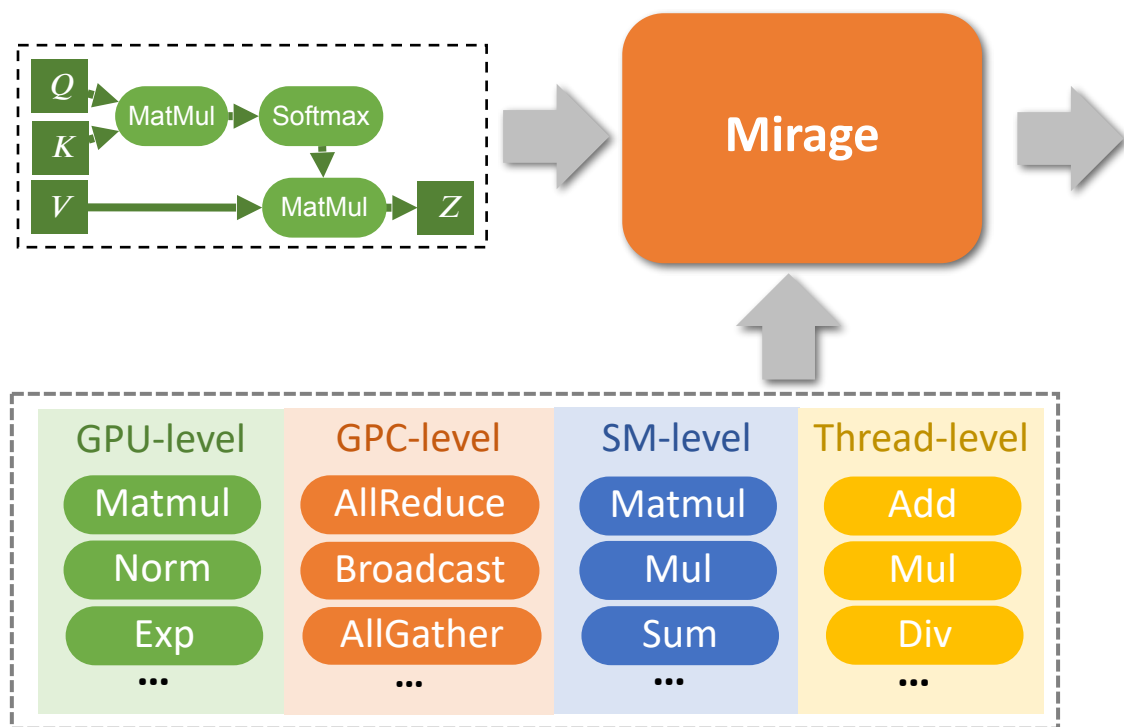
Find μ Graphs similar to expert-written implementations* for attention on A100



Mirage Discovers Hardware-Customized μ Graphs

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

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



Mirage: A Multi-Level Superoptimizer

- *Algebraic transformation + Schedule transformation + New kernel generation*
- **Minimal engineering effort:** A few lines of Python codes from users
- **High performance:** Outperform existing systems by up to 3.3x

Our Paper @ OSDI 2025



<https://www.usenix.org/conference/osdi25/presentation/wu-mengdi>

 Our Repo



<https://github.com/mirage-project/mirage>

