



EINNET: Optimizing Tensor Programs with Derivation-Based Transformations

OSDI23 @ Boston, July 12, 2023

Liyan Zheng, Haojie Wang, Jidong Zhai,
Muyan Hu, Zixuan Ma, Tuwei Wang, Shuhong Huang,
Xupeng Miao, Shizhi Tang, Kezhao Huang, Zhihao Jia

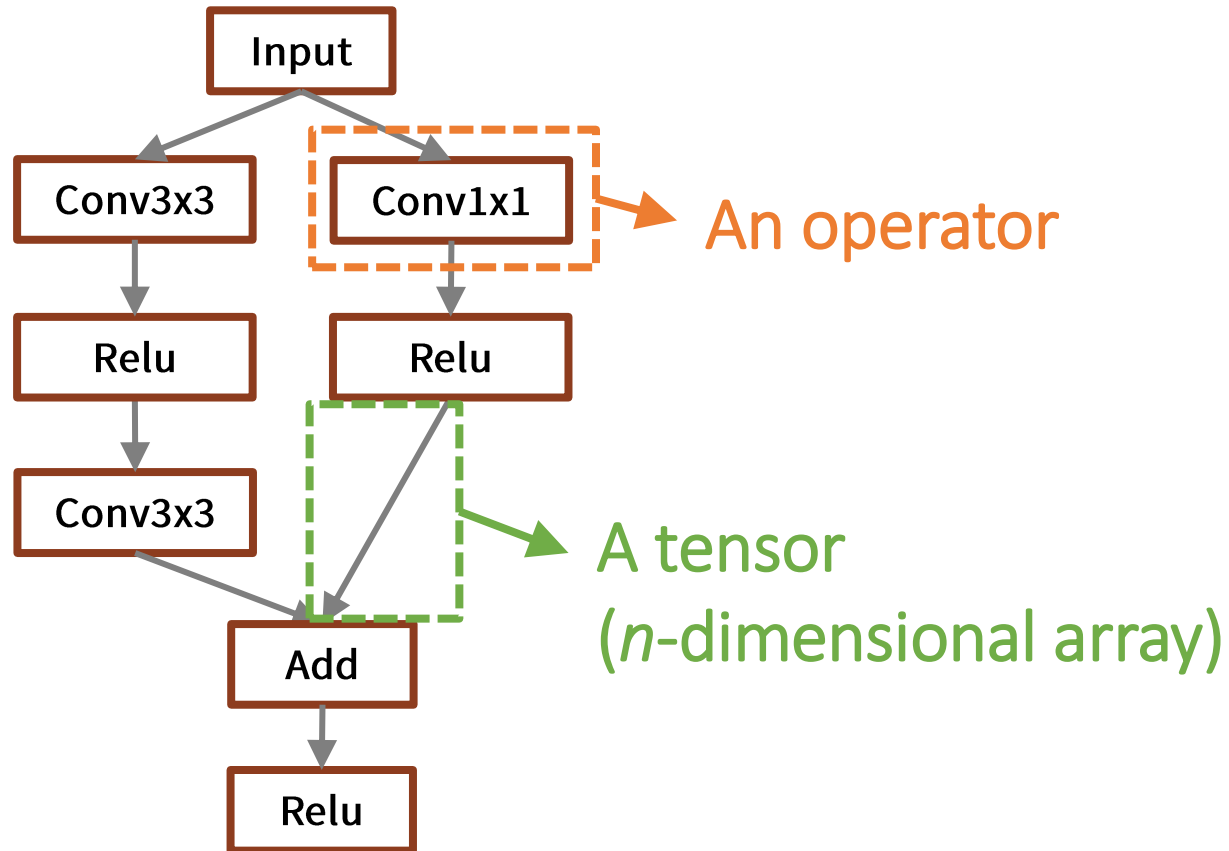
Tsinghua University

Carnegie Mellon University



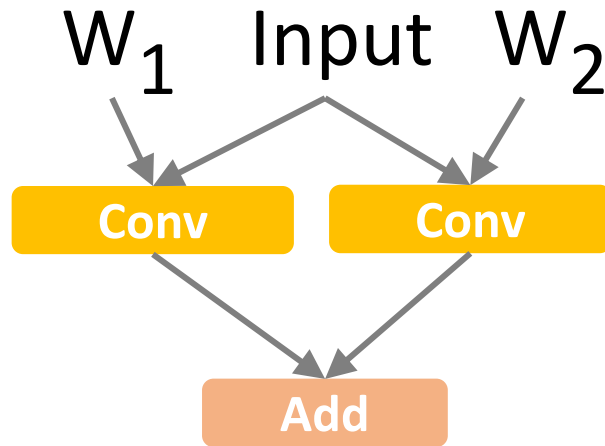
Tensor Programs

- Widely used in deep learning
- Represented as computation graphs

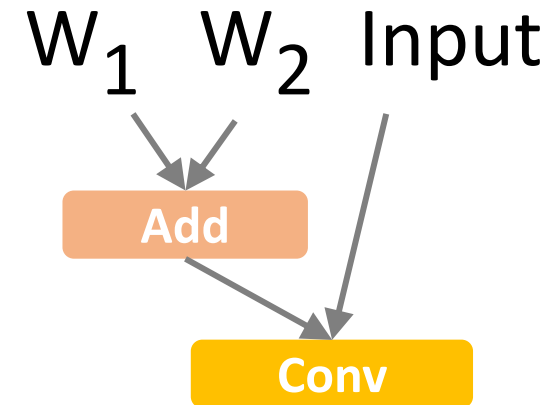
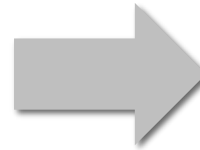


Tensor Program Transformations

- **Goal:** optimizing program performance



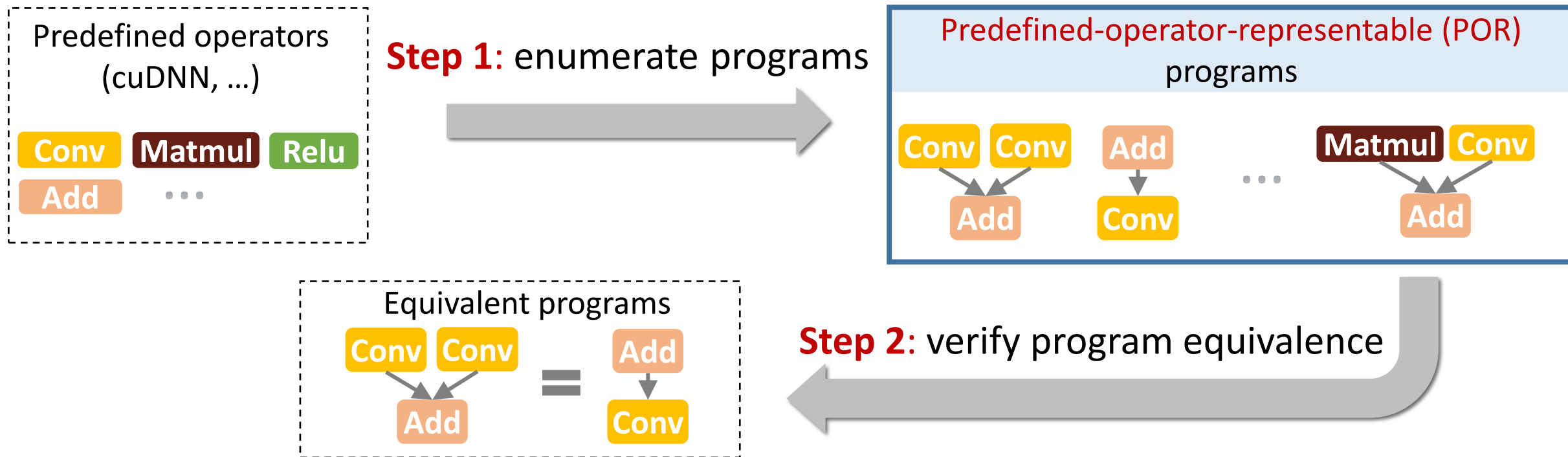
Original Program



Optimized Program

Automatic Transformations

- Superoptimization-based approaches (TASO¹ and PET²)



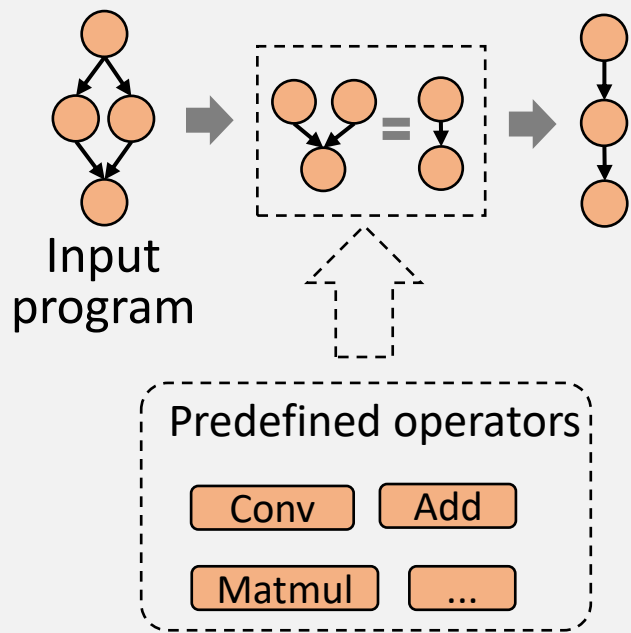
😊 Automatic exploration 😭 Only POR transformations 😭 Time-consuming verifications

1. TASO: Optimizing Deep Learning Computation with Automated Generation of Graph Substitutions. SOSP'19
2. PET: Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections. OSDI'21
EinNET: Optimizing Tensor Programs with Derivation-Based Transformations

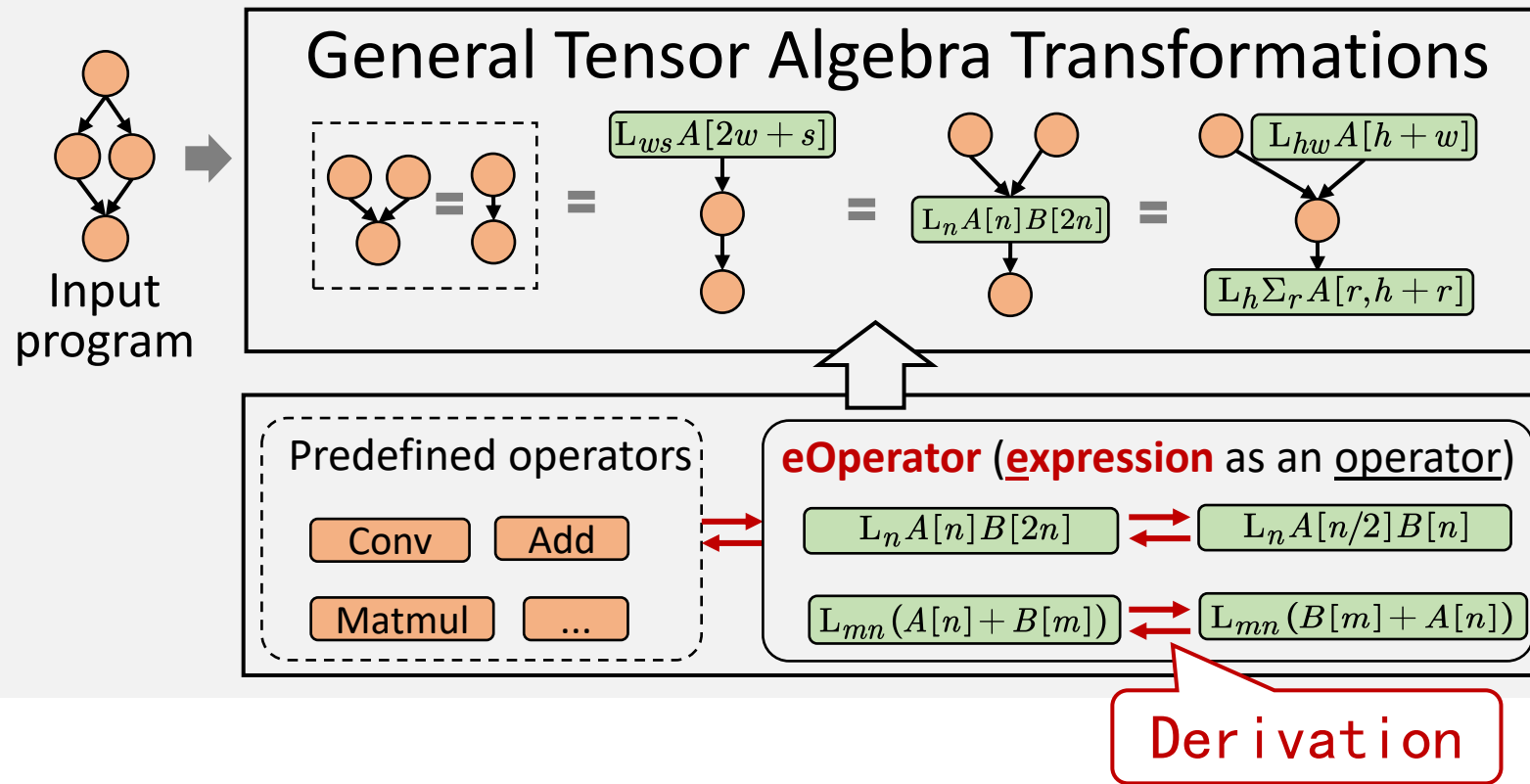
EINNET

- A tensor program optimizer based on **tensor expression derivation**

Prior work

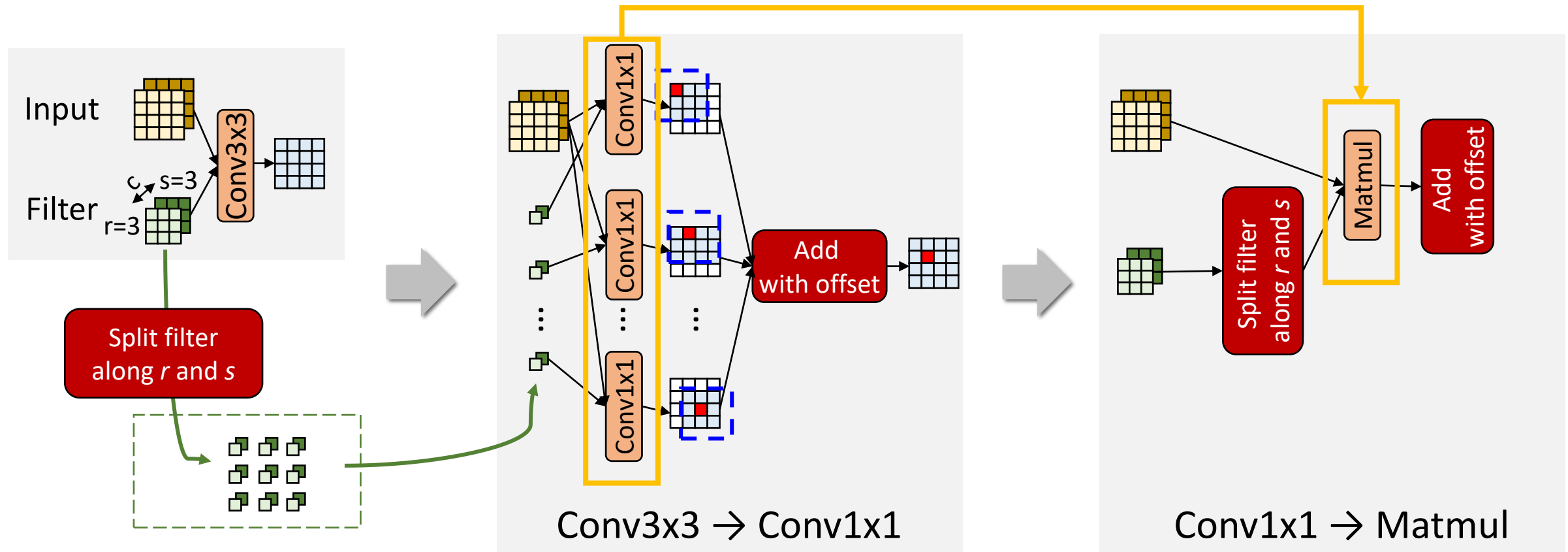


EINNET



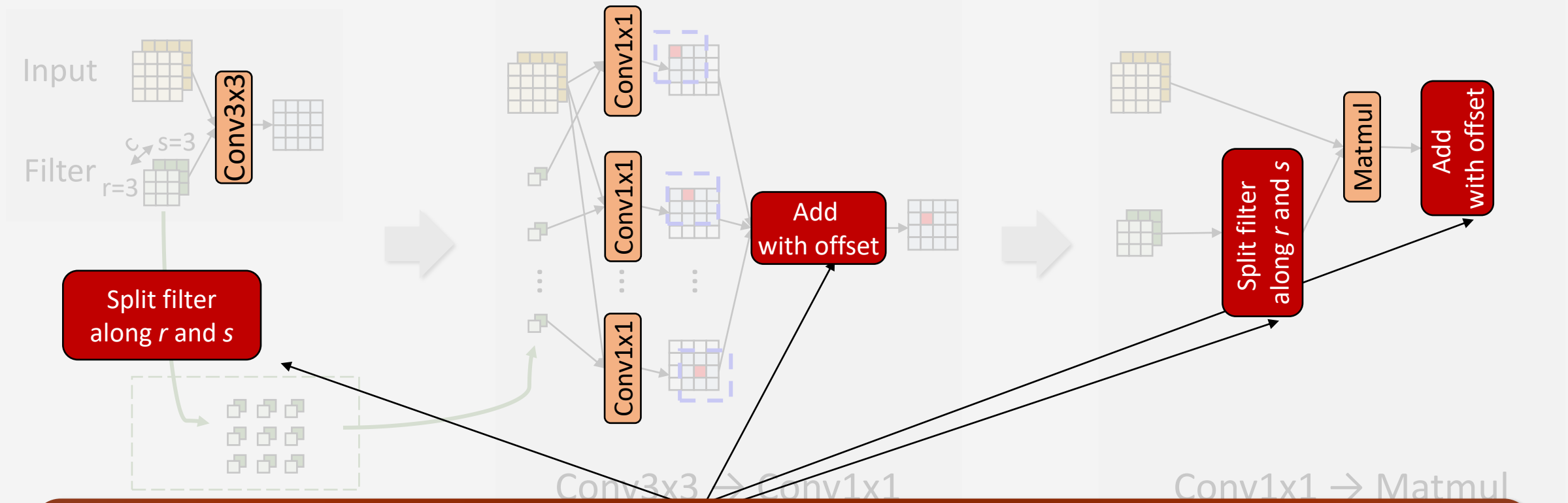
Larger optimization space and up to 2.7x speedup

Motivating Example: Convolution to Matmul



2x speedup on Nvidia A100 GPUs

Motivating Example: Convolution to Matmul

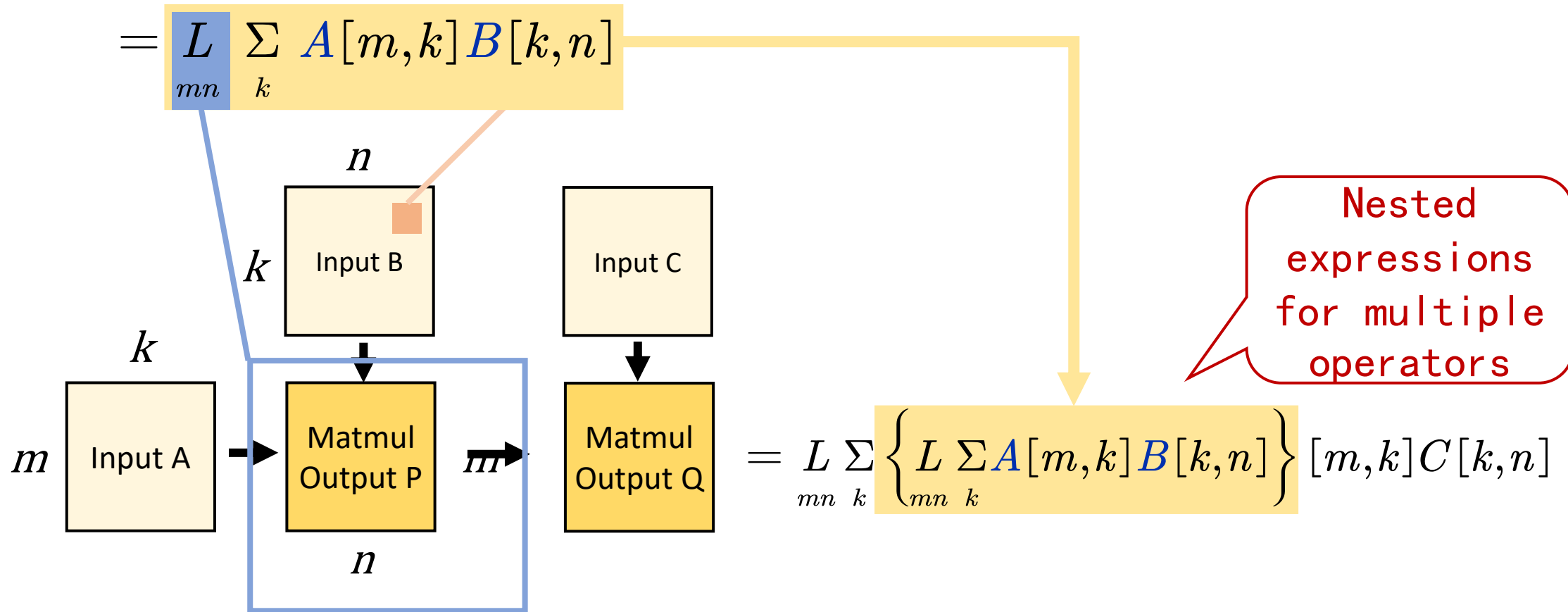


Non-predefined operators enable more optimizations

Tensor Expressions

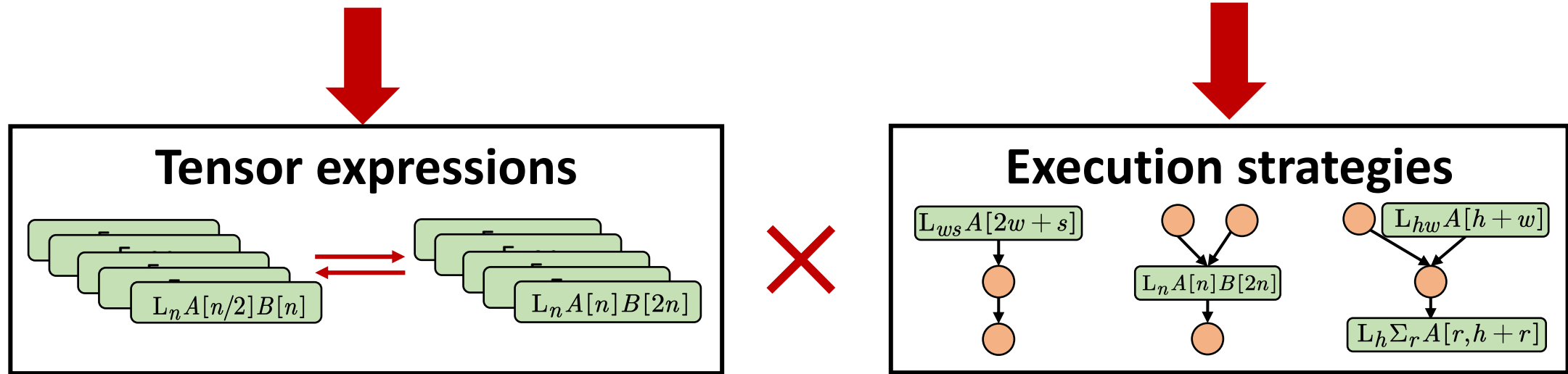
- Specify the computation semantics

$$P[m, n] = \text{Matmul}(A, B)$$



Key Challenges

1. discover transformations
2. execute expressions on HW



3. explore the large search space

Challenge#1: Expression Transformation

A strawman approach: **superoptimization for expressions**

Input expression

$$2A[i] + B[j]$$

Enumerate



Candidates

$$2A[i] \times B[j]$$

$$A[i] + B[j]$$

$$B[j] + 2A[i]$$

...

Verified

- **Limitation**

- Infinitely many expressions
- Hard to verify equivalence

EinNet discovers equivalent expressions by **derivations**

Derivation Rules

- Mathematically equivalent rewrites

Intra-expression derivations

Summation splitting
$$L_{\vec{x}, \vec{s}_1, \vec{s}_2} \Sigma f(\vec{T}, \vec{x}, \vec{s}_1, \vec{s}_2) \implies L_{\vec{x}, \vec{s}_1} \Sigma \left\{ L_{\vec{x}, \vec{s}_1, \vec{s}_2} \Sigma f(\vec{T}, \vec{x}, \vec{s}_1, \vec{s}_2) \right\} [\vec{x}, \vec{s}_1]$$

Variable substitution
$$L_{\vec{x}} f(\vec{T}, \vec{x}) \implies L_{\vec{x}} \left\{ L_{\vec{y}} f(\vec{T}, \Phi^{-1}(\vec{y})) \right\} [\Phi(\vec{x})]$$

Traversal merging
$$L_{\vec{x}} \left\{ L_{\vec{y}} f(\vec{T}, \vec{y}) \right\} [\Phi^{-1}(\vec{x})] \implies L_{\vec{x}} f(\vec{T}, \vec{x})$$

Boundary relaxing
$$L_{\vec{x}} f(\vec{T}, \vec{x}) \implies L_{\vec{x}}^{\mathbb{X} \cup \mathbb{Y}} f(\vec{T}, \vec{x})$$

Boundary tightening
$$L_{\vec{x}}^{\mathbb{X} \cup \mathbb{Y}} f(\vec{T}, \vec{x}) \implies L_{\vec{x}} f(\vec{T}, \vec{x})$$

Inter-expression derivations

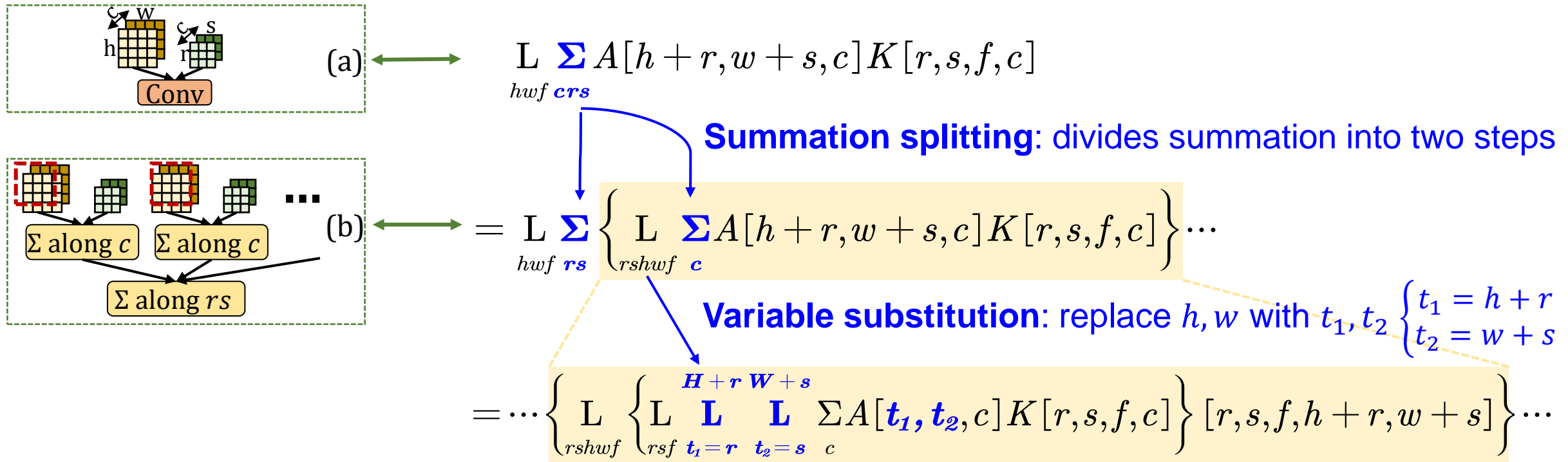
Expression splitting
$$L_{\vec{x}}^{\mathbb{X}} f(\vec{T}, \vec{x}) \implies L_{\vec{x}}^{X_0} f(\vec{T}, \vec{x}) \sim L_{\vec{x}}^{X_1} f(\vec{T}, \vec{x})$$

Expression merging
$$L_{\vec{x}}^{X_0} f(\vec{T}, \vec{x}) \sim L_{\vec{x}}^{X_1} f(\vec{T}, \vec{x}) \implies L_{\vec{x}}^{\mathbb{X}} f(\vec{T}, \vec{x})$$

Expression fusion
$$L_{\vec{x}}^{\mathbb{X}} f(\vec{T}, \vec{x}) \rightarrow L_{\vec{x}}^{X'} g(\vec{T}', \vec{x}) \implies L_{\vec{x}}^{\mathbb{X}} g(f(\vec{T}, \vec{x}))$$

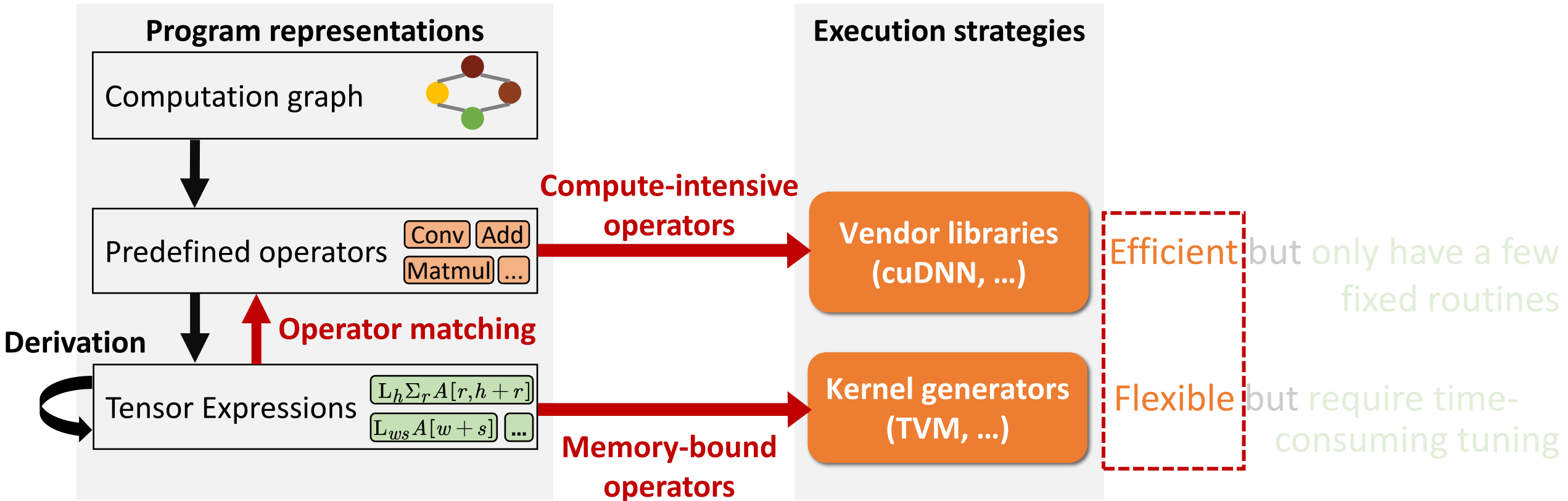
- Support custom derivation rules

Derivations in the Motivating Example



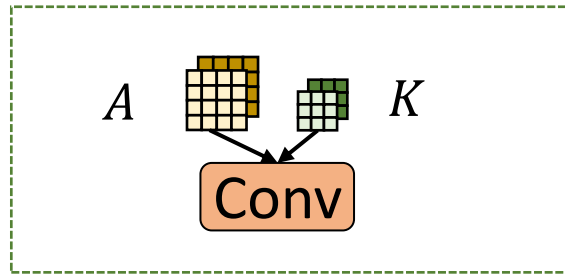
Derivation creates new equivalent expressions

Challenge#2: Executing Expressions



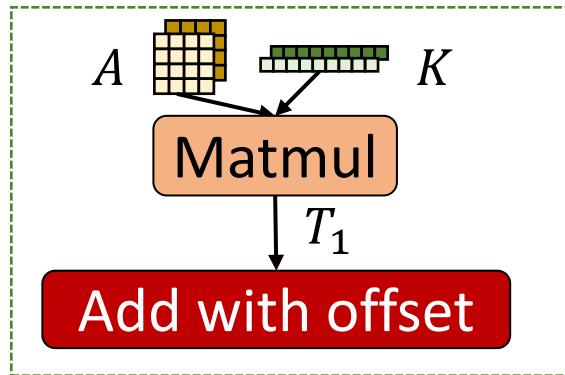
Combine the benefits of vendor libraries and kernel generators

Operator Matching & Kernel Generation



$$\mathbb{L} \sum_{hwf\ crs} A[h+r, w+s, c] K[r, s, f, c]$$

Several steps of derivations



$$= \mathbb{L} \sum_{hwf\ rs} \left\{ \mathbb{L} \sum_{rsf\ t_1=0\ t_2=0\ c} \mathbb{L} \sum_{t_1=0\ t_2=0\ c} A[t_1, t_2, c] K[r, s, f, c] \right\} [r, s, f, h+r, w+s]$$

$$= \mathbb{L} \sum_{hwf\ rs} T_1[r, s, f, h+r, w+s]$$

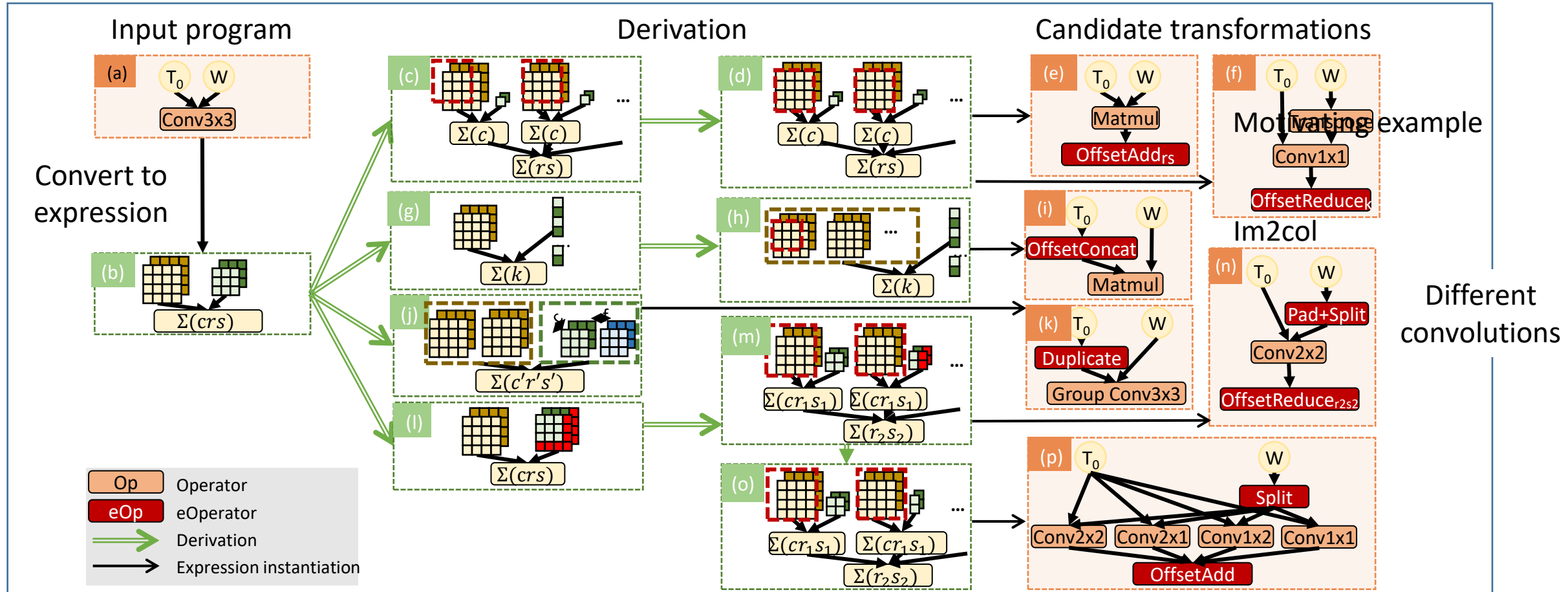
$$= T_2$$

Operator matching: replace compute-intensive expressions with *Matmul(K, A)*

Kernel generation: tune kernels for memory-bound expressions

* Details available in the paper

Challenge#3: Large search space



More optimization opportunities

Challenge#3: Large search space

- A computation graph of a single convolution
 - ~10 steps of derivation
 - $\sim 10^8$ candidates
 - ~10 hours

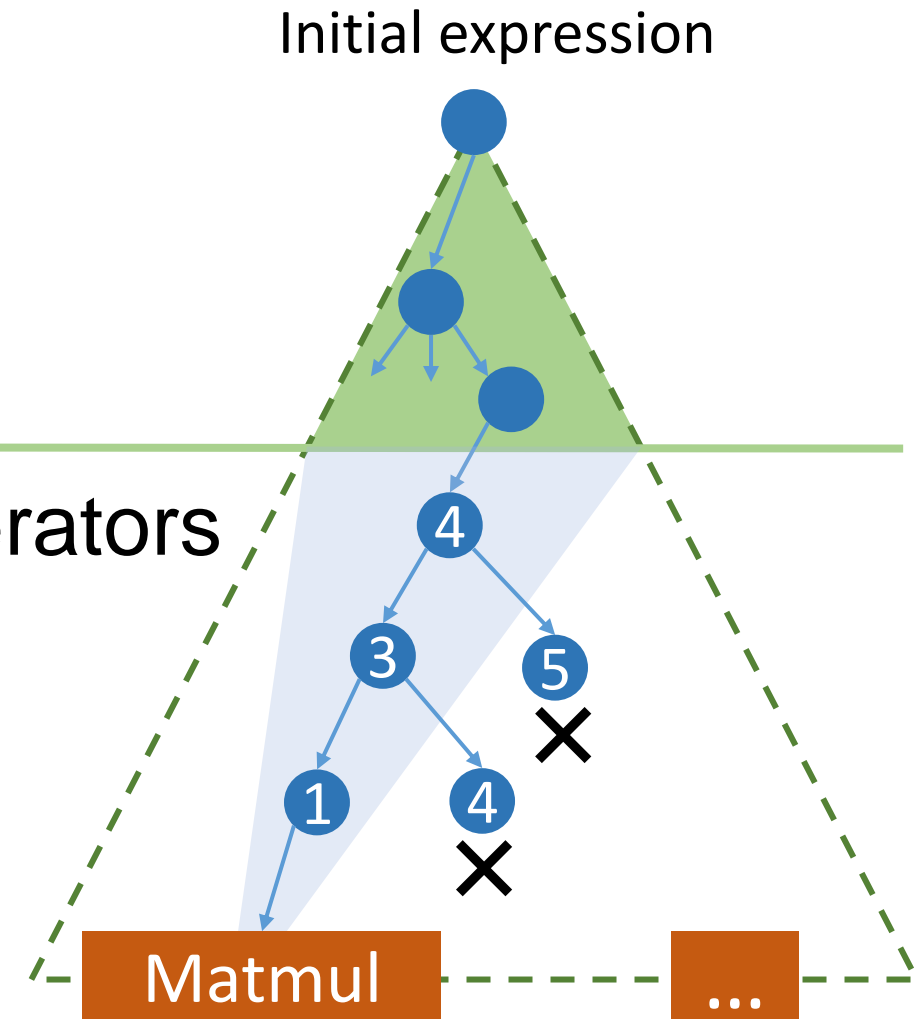
Solution: **expression distance** to guide search

- Measure similarity between expressions

Expression-Distance-Guided Search

- Two search stages
- **Stage I**: explore search space
 - Apply all possible derivations

- **Stage II**: converge to performant operators
 - Guided by expression distance



Evaluation

Platforms:

Nvidia A100 & V100 GPUs

Backends:

cuBLAS + cuDNN, AutoTVM, Ansor

Models:

Language model: Longformer

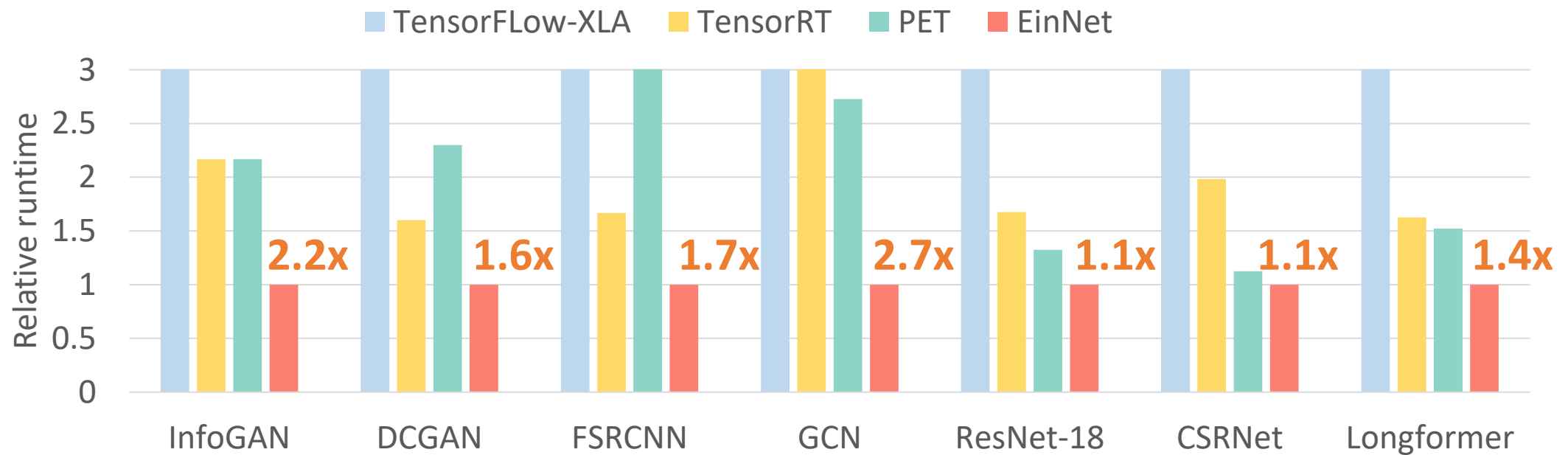
Image generation: InfoGAN, DCGAN, FSRCNN

Image understanding: GCN, ResNet-18, CSRNet

Baseline:

TensorRT, PET, Tensorflow, Tensorflow-XLA, Nimble, TVM

End-to-End Inference Evaluation (Nvidia A100 GPU)



EinNet outperforms existing optimizers by up to 2.7x

Conclusion



EINNET is a **derivation-based** tensor program optimizer



Proposed technique: expression derivation



Larger search space: general tensor algebra transformations



Better performance: up to 2.7x speedup

Available at: <https://github.com/InfiniTensor/InfiniTensor>

Q&A
Thank you!