# **EINNET: Optimizing Tensor Programs** with Derivation-Based Transformations

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

#### MOTIVATION

#### Tensor program transformations

- Optimize program performance
- Preserve program outputs

#### <u>Automatic program optimizers</u>

Superoptimization-based approaches

- Step I: enumerate candidate programs
  by predefined operators → limited space
- Step II: verify candidate programs with the original program → time-consuming

### **Derivation-based optimizer (our work)**

Proposed technique: tensor expression derivation
 Larger search space: tensor algebra transformations
 Better performance: up to <u>2.7x</u> speedup

#### 





#### Approach

#### Tensor algebra expressions

Specify computation semantics mathematically Nested expressions for multiple operators

#### **Expression derivation**

Mathematically equivalent rewrite  $\Delta \mathcal{F}$ 

#### **Expression execution**

Different execution strategies

- Math libraries: efficient but only fixed routines
- Kernel generators: flexible but require timeconsuming tuning

Get an ideal combination via operator matching

- Compute-intensive operations  $\rightarrow$  libraries
- Memory-bound operations  $\rightarrow$  generators

## Search expression transformation space





Stage I: enlarge search space

- Apply all derivation rules under a search depth limit
- Stage II: opportunistically leverage math libraries
- Reduce expression distance to target



**Effective on different backends** 

#### EVALUATION

#### End-to-end inference (up to 2.7x speedup)

