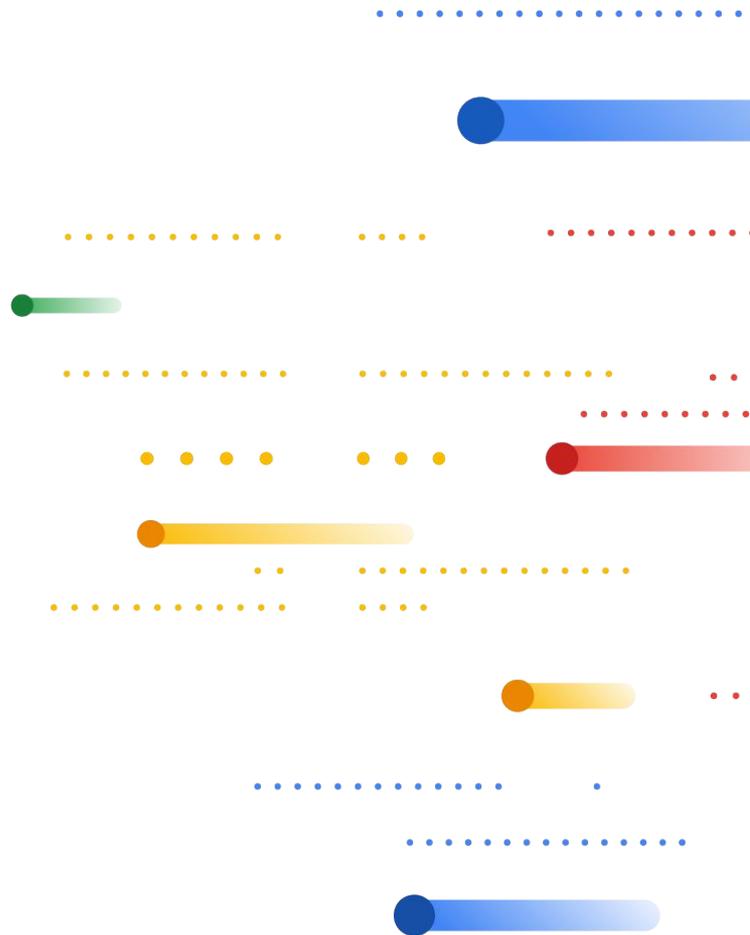


# ML for Autotuning Production ML Compilers

Phitchaya Mangpo Phothilimthana  
mangpo@google.com



# Search-Based ML Compilers

optimization scope	<i>graph</i>	TASO PET	DeepCuts
	<i>subgraph</i>		TVM Halide TensorComp... FlexTensor Anso AdaTune Chameleon

# Search-Based ML Compilers

optimization scope	<i>graph</i>	TASO PET	DeepCuts
optimization scope	<i>subgraph</i>		TVM Halide TensorComp... FlexTensor Anso AdaTune Chameleon

# Search at Subgraph Level is Suboptimal

A common strategy **partitions** a graph into subgraphs **according to the neural net layers**, ignoring cross-layer optimization opportunities.

Empirical result: a **regression** of **up to 2.6x** and **32% on average** across 150 ML models by limiting fusions in XLA to be within layers.

# Search-Based ML Compilers

optimization scope	<i>graph</i>	TASO PET	DeepCuts
	<i>subgraph</i>		TVM Halide TensorComp... FlexTensor Anso AdaTune Chameleon

# Search Approaches: Long Compile Time

optimization scope	<i>graph</i>	XLA	TASO PET	DeepCuts
	<i>subgraph</i>			TVM Halide TensorComp... FlexTensor Anso AdaTune Chameleon
		<b><i>seconds</i></b>	<b><i>minutes</i></b>	<b><i>hours</i></b>
	<b>compile time (for ResNet like inference)</b>			

# Production Compilers: Multi-Pass

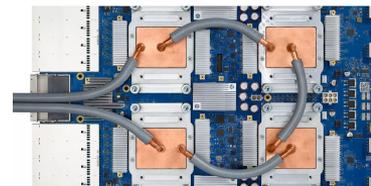
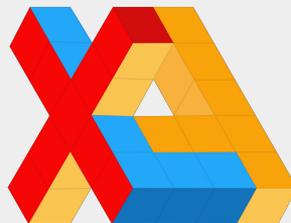
- Models evaluated by research compilers: up to 1,000 nodes
- Industrial-scale models: up to **500,000 nodes!**
- That's why **production ML compilers** still decompose the compilation into **multiple passes**.
- **None** of the existing approaches **support** autotuning different optimizations in a **multi-pass compiler**.
  - **Challenge:** search space of a pass is highly dependent on decisions made in prior passes.

# Our Goal

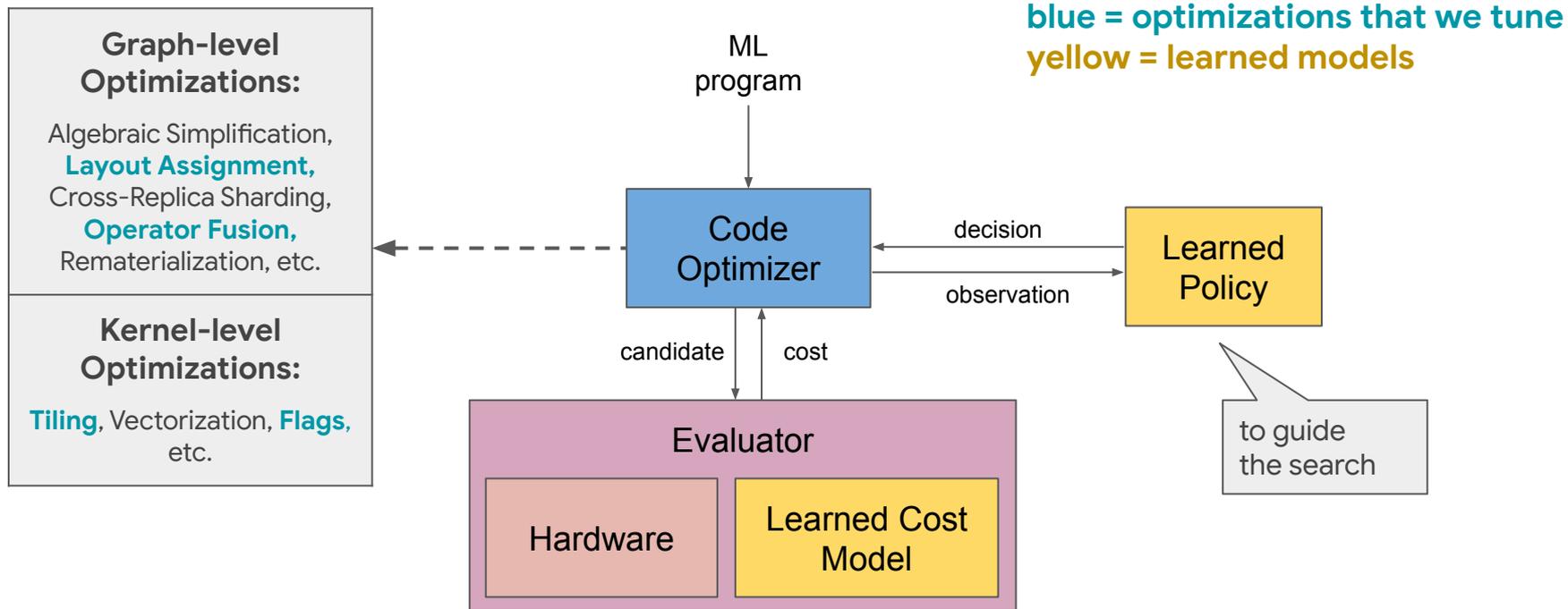
Bring the benefits of **search-based** exploration to **multi-pass compilers**:

- for both graph and subgraph levels
- with flexibility via configurable search to tune subset of optimizations of interest

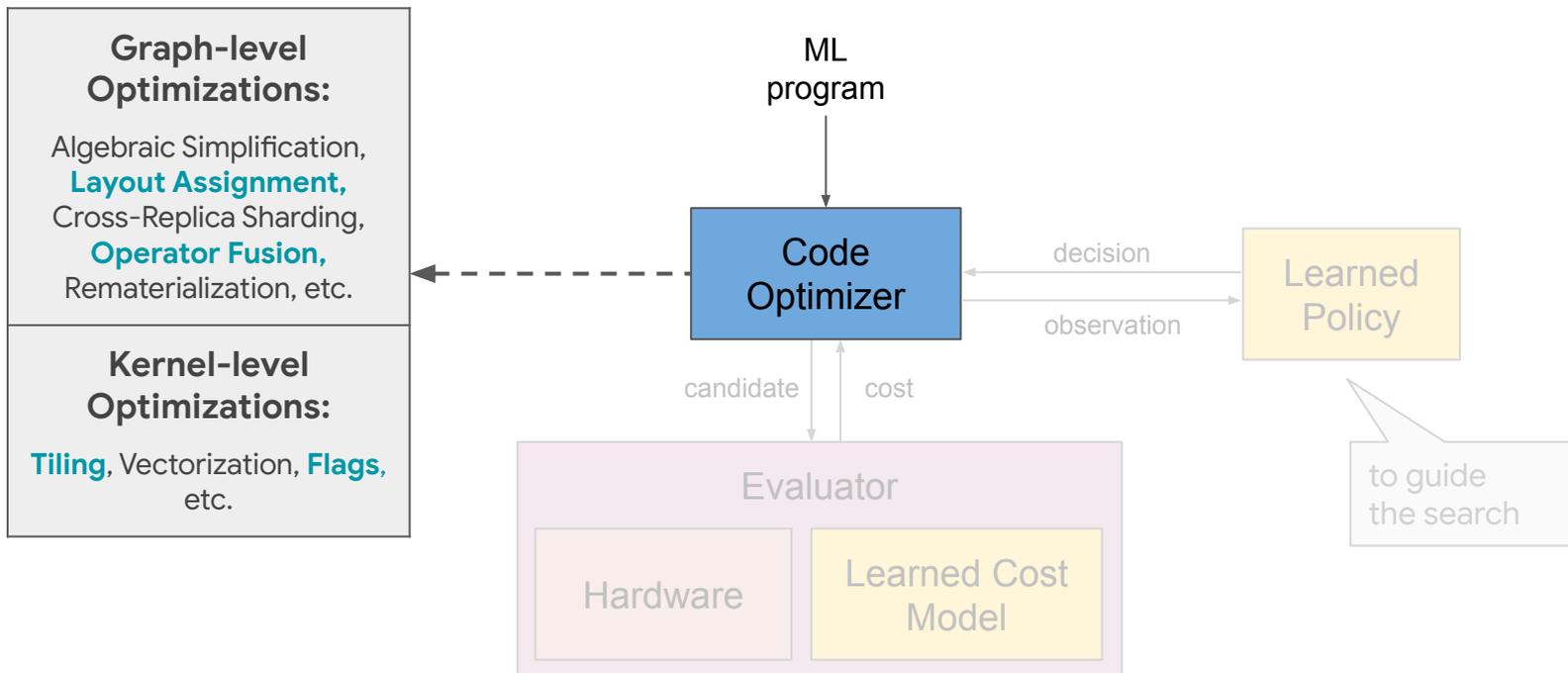
# Production ML Compilation Stack at Google



# XTAT: XLA TPU Autotuner



# XTAT: XLA TPU Autotuner



# Pass Configuration

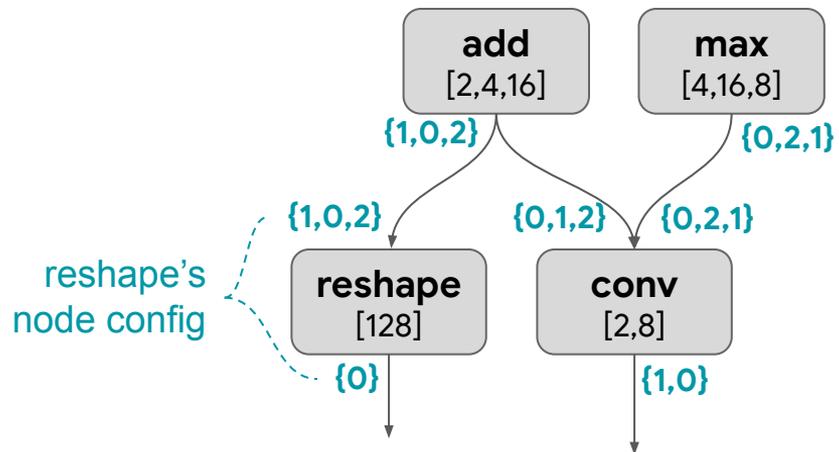
**configuration on a tensor graph  
for an optimization pass**

is

**a collection of per-node configurations** that control  
how the pass transforms each node in the graph

# Layout Assignment

Example:





# Layout Search Space

## Option #1: Naive

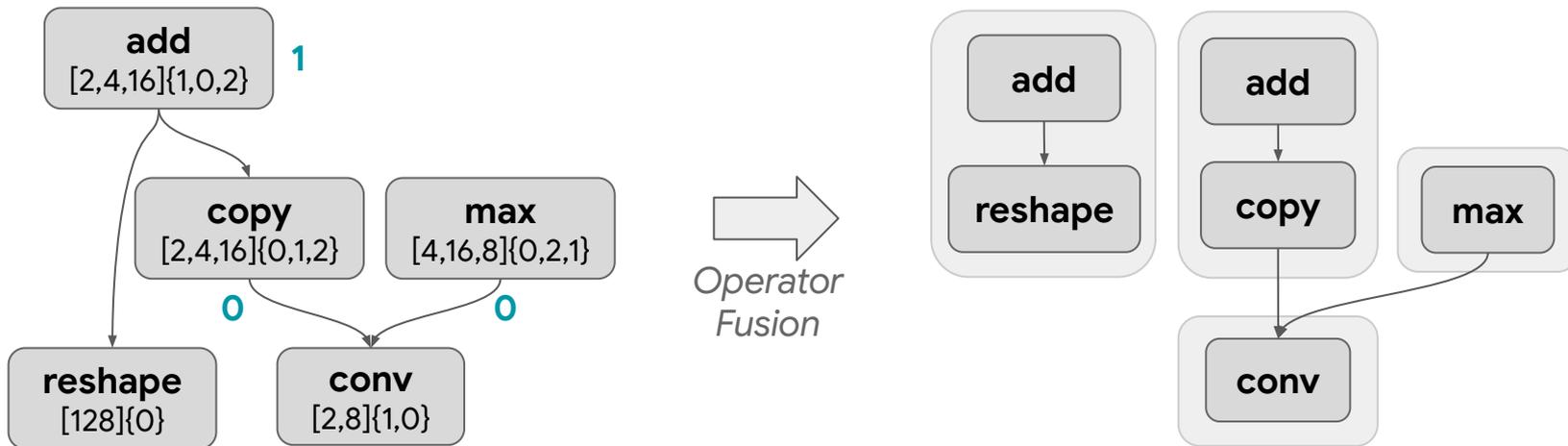
- Layout options for **each input/output** are **permutation** of its dimensions.
- Many **invalid configs** because there are constraints between tensors.

## Option #2: Proposed

- Tune **layout options for important ops** (convolution and reshape).
- For each important op, get valid input-output layouts from XLA.
- Leverage XLA **layout propagation** algorithm.

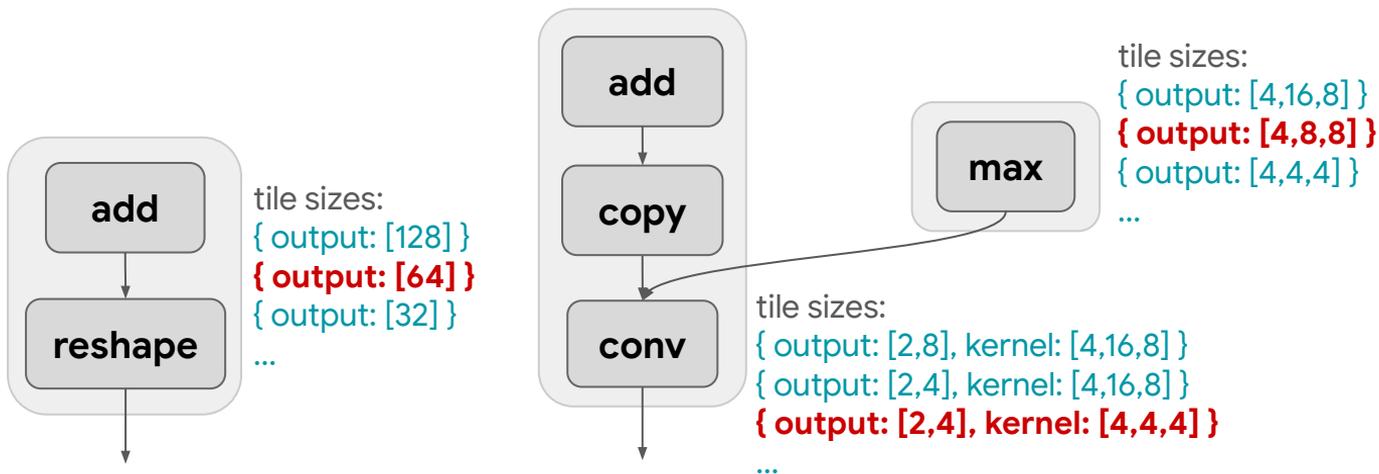
# Operator Fusion

Example:



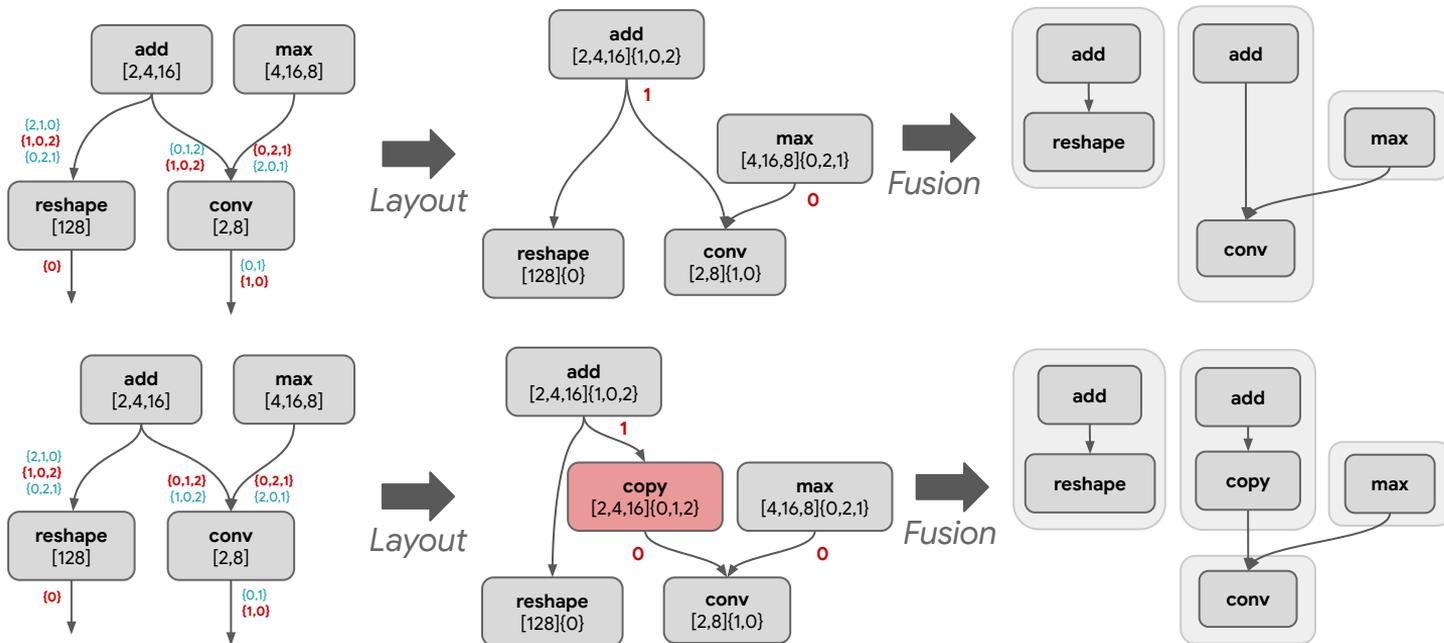
# Tile Size & Code Gen Flags Search Space

Tune config for each fused node (kernel) independently.



# Joint Autotuning: Challenges

$$g_A \text{ --- } A(\text{config}_A) \text{ ---} \rightarrow g_B \text{ --- } B(\text{config}_B) \text{ ---} \rightarrow g_{out}$$



$config_A$  determines the input graph  $g_B$  to pass B and its search space

When we change  $config_A$  to  $config'_A$ ,  $g_B$  is changed, and  $config_B$  is no longer valid.

How to not start the search for B from scratch when  $config_A$  is changed?

# Methodology for Joint Autotuning

```
function SEARCHSTEP(C)  
  optid ← SelectOpt(Opts)  
  C' ← GenerateCandidates(optid, C)  
  for c : C' do  
    UpdateAndApplyCandidate(optid, c)  
    Evaluate(c)  
  end for  
  return SelectCandidates(C, C')  
end function
```

## Returns:

A, B, C, A, B, C, ... (joint tuning)

A, A, ..., B, B, ..., C, C, ... (sequential)

or some combinations of them

# Methodology for Joint Autotuning

```

function SEARCHSTEP(C)
  optid ← SelectOpt(Opts)
  C' ← GenerateCandidates(optid, C)
  for c : C' do
    UpdateAndApplyCandidate(optid, c)
    Evaluate(c)
  end for
  return SelectCandidates(C, C')
end function

```

## Candidate *c*:

*c*.graphs = [*g<sub>A</sub>*, *g<sub>B</sub>*, *g<sub>out</sub>*]  
*c*.configs = [*config<sub>A</sub>*, *config<sub>B</sub>*]

## Change *config<sub>A</sub>*:

*c*.graphs = [*g<sub>A</sub>*, *g<sub>B</sub>*, *g<sub>out</sub>*]  
*c*.configs = [*config<sub>A</sub>'*, *config<sub>B</sub>*]

## Fix *c* to be *well-formed*:

*c*.graphs = [*g<sub>A</sub>*, *g<sub>B</sub>'*, *g<sub>out</sub>'*]  
*c*.configs = [*config<sub>A</sub>'*, *config<sub>B</sub>'*]

# Construct Well-Formed Candidate

## Key ideas:

- Update subsequent graphs
- Update  $\text{config}_B'$  to have configurations for all nodes in  $g_B'$  from:
  - $\text{config}_B$
  - **global configuration store** (maintaining the best config per node)
  - **default value**

### Change $\text{config}_A$ :

```
c.graphs = [gA, gB, gout]  
c.configs = [configA', configB]
```



### Fix $c$ to be *well-formed*:

```
c.graphs = [gA, gB', gout']  
c.configs = [configA', configB']
```

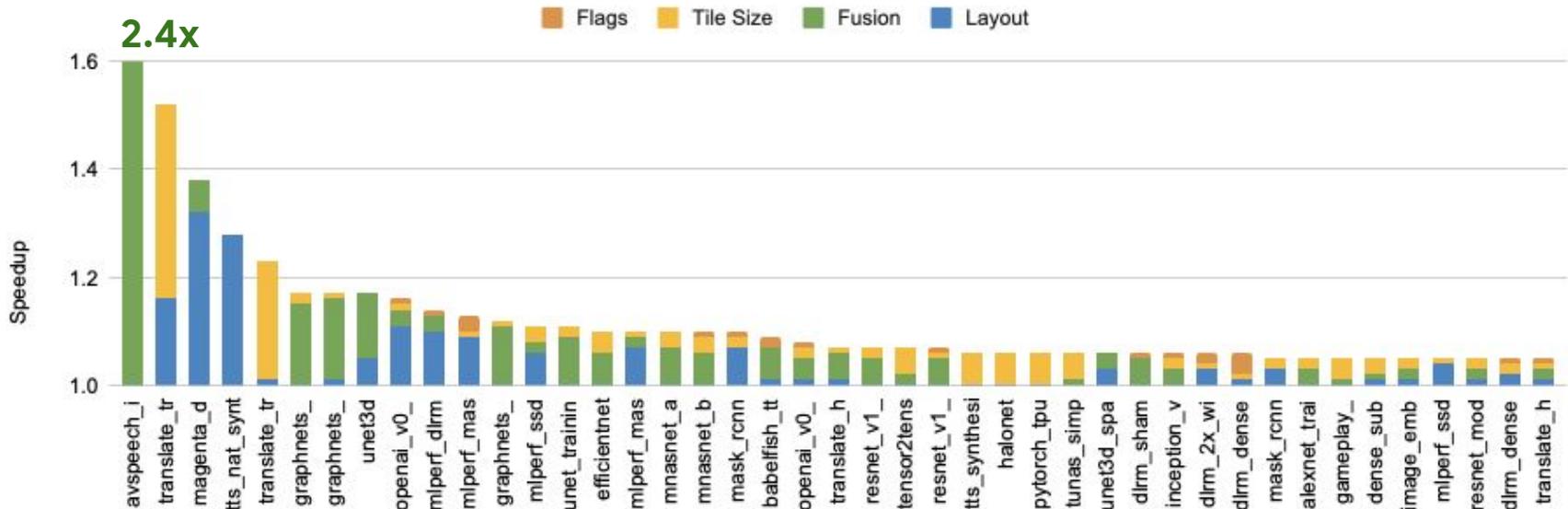
# End-to-End Search Schedule

- Separate tuning graph-level and kernel-level optimizations for scalability
- Tuning **layout + fusion jointly is better** than sequentially
- Tuning **tile size + flag jointly is worse** than sequentially

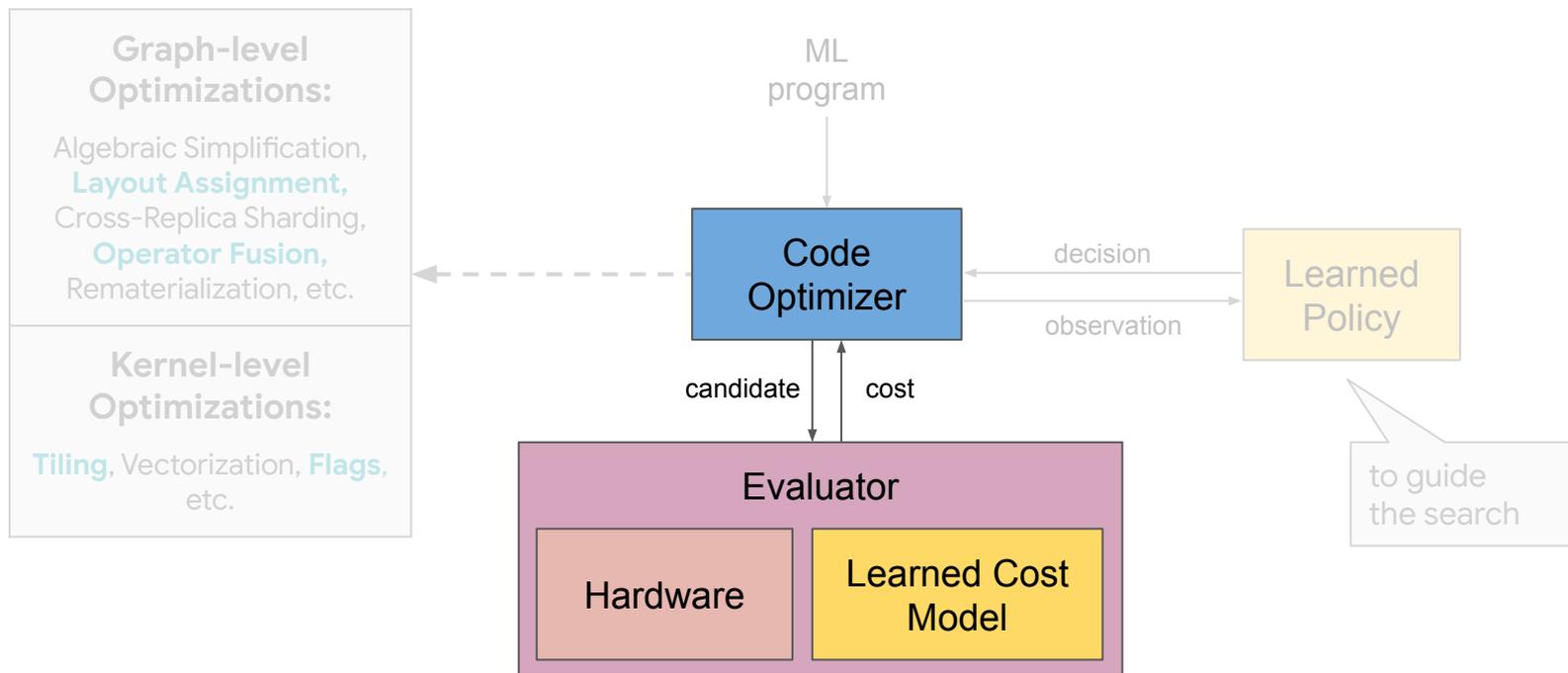
Tune **layout-fusion jointly** (simulated annealing)  
→ then tune **tile size** (exhaustive)  
→ then tune code gen **flags** (exhaustive)

# End-to-End Runtime Speedup

We measured end-to-end model speedups from autotuning **150 ML models**. The figure shows models that achieve 5% or more improvement.

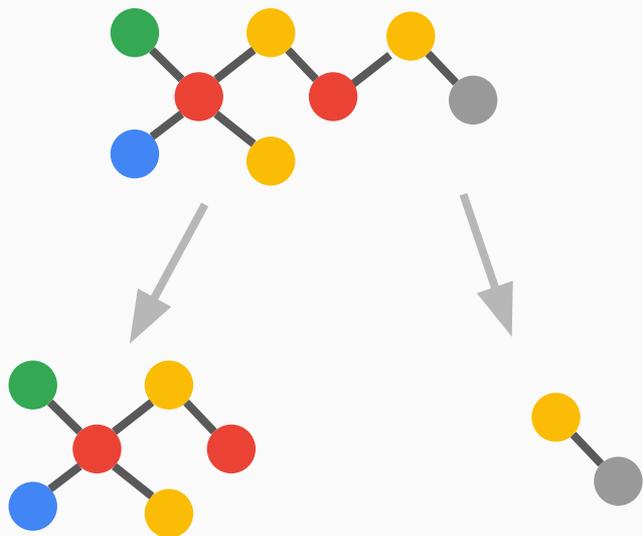


# Learned Cost Model



# Overview of Cost Model

## 1. Decompose Into Kernels

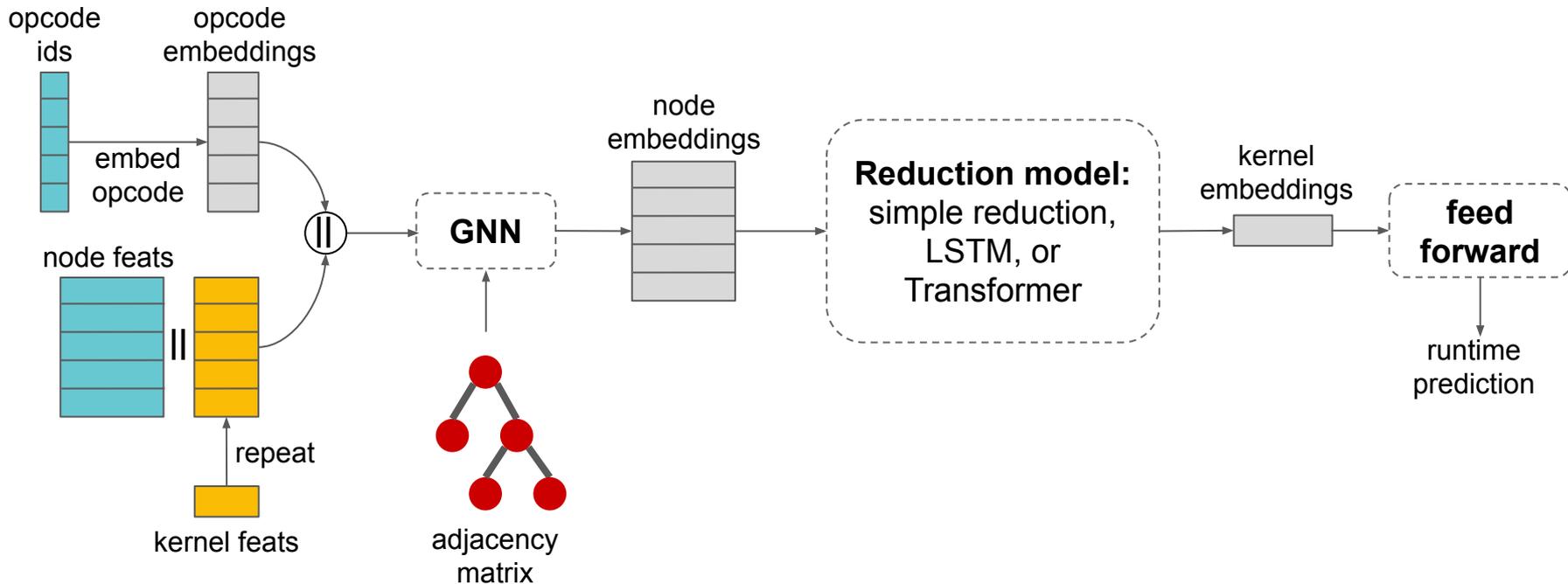


## 2. Regression Per Kernel

$$f(\text{KERNEL}) \approx 5.2s$$

RUNTIME

# Model Architecture



# Losses

## Mean Squared Error

for absolute runtime prediction.  
Targets are log-transformed.

$$L = \sum_{i=1}^n (y'_i - y_i)^2$$

## Pairwise Rank Loss

for relative runtime prediction.

$$L = \sum_{i=1}^n \sum_{j=1}^n \frac{\phi(y'_i - y'_j) \cdot \text{pos}(y_i - y_j)}{n \cdot (n - 1) / 2}$$

$$\phi(z) = \begin{cases} (1 - z)_+ & \text{hinge function} \mathbf{or} \\ \log(1 + e^{-z}) & \text{logistic function} \end{cases}$$

$$\text{pos}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

# Accuracy Evaluation and Baseline

- **Accuracy evaluation tasks**
  - Tile size selection (relative runtimes)
  - Fusion (absolute runtimes)
- **Baseline:** XLA's hand-written, analytical performance model
  - XLA argmins all tile sizes using this performance model
  - Fusion does not use this model. It uses other heuristics.

# Accuracy: Tile Size Selection

Compare **true** runtimes between best predicted and actual best tile size. **APE**:

$$100 \times \frac{\sum_{k \in K} |t_{c'_k}^k - \min_{c \in C_k} t_c^k|}{\sum_{k \in K} \min_{c \in C_k} t_c^k}$$

In **random** split, learned model **~halves** APE.

	Learned	Analytical
<b>ConvDRAW</b>	9.7	3.9
<b>WaveRNN</b>	1.5	2.8
<b>NMT Model</b>	3.1	13.1
<b>SSD</b>	3.9	7.3
<b>RNN</b>	8.0	10.2
<b>ResNet v1</b>	2.8	4.6
<b>ResNet v2</b>	2.7	5.4
<b>Translate</b>	3.4	7.1
<b>Median</b>	3.3	6.2
<b>Mean</b>	3.7	6.1

# Accuracy: Fusion

Compare **Mean Absolute Percentage Error** of kernel runtime predictions.

**Random split:** learned model improves MAPE by **~85%**.

	Learned	Analytical
<b>ConvDRAW</b>	17.5	21.6
<b>WaveRNN</b>	2.9	322.9
<b>NMT Model</b>	9.8	26.3
<b>SSD</b>	11.4	55.9
<b>RNN</b>	1.9	20.5
<b>ResNet v1</b>	3.1	11.5
<b>ResNet v2</b>	2.4	13.3
<b>Translate</b>	2.1	27.2
<b>Median</b>	3.0	24.0
<b>Mean</b>	4.5	31.1

# Ablations: takeaways

- Using a **rank loss** for the tile-size task reduced APE by 10 pts. on average.
- **GraphSAGE** outperformed using a **Graph Attention Networks** or **sequence model** and was less sensitive to hyperparameter selection.
- Replacing the **LSTM/Transformer reduction** with a **non-learned reduction** works almost as well (and improves inference time).

# Training for All Optimization Tasks

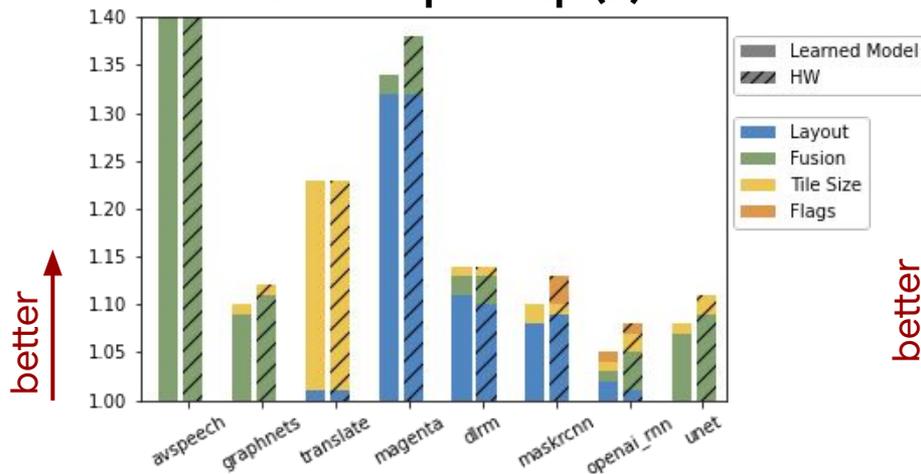
- **Generate training data from 150 ML models** using random layout, fusion, tile size, and flag configurations.
- Train:
  - one model for all **graph-level** optimizations to predict **absolute runtime**
  - one model for **tile-size** to predict **relative runtime**
  - one model for **flags** to predict **relative runtime**
- The graph embedding network is shared between tile-size and flags models.

# Tuning with Learned Cost Model

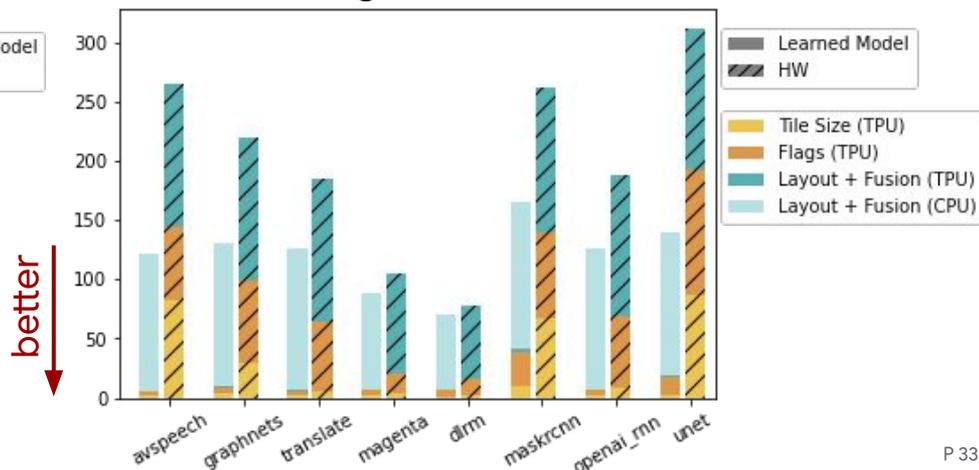
Execute the top k configurations from each worker according to the model on real hardware and pick the best.

- k = 10 for graph-level optimizations
- k = 5 for kernel-level optimizations

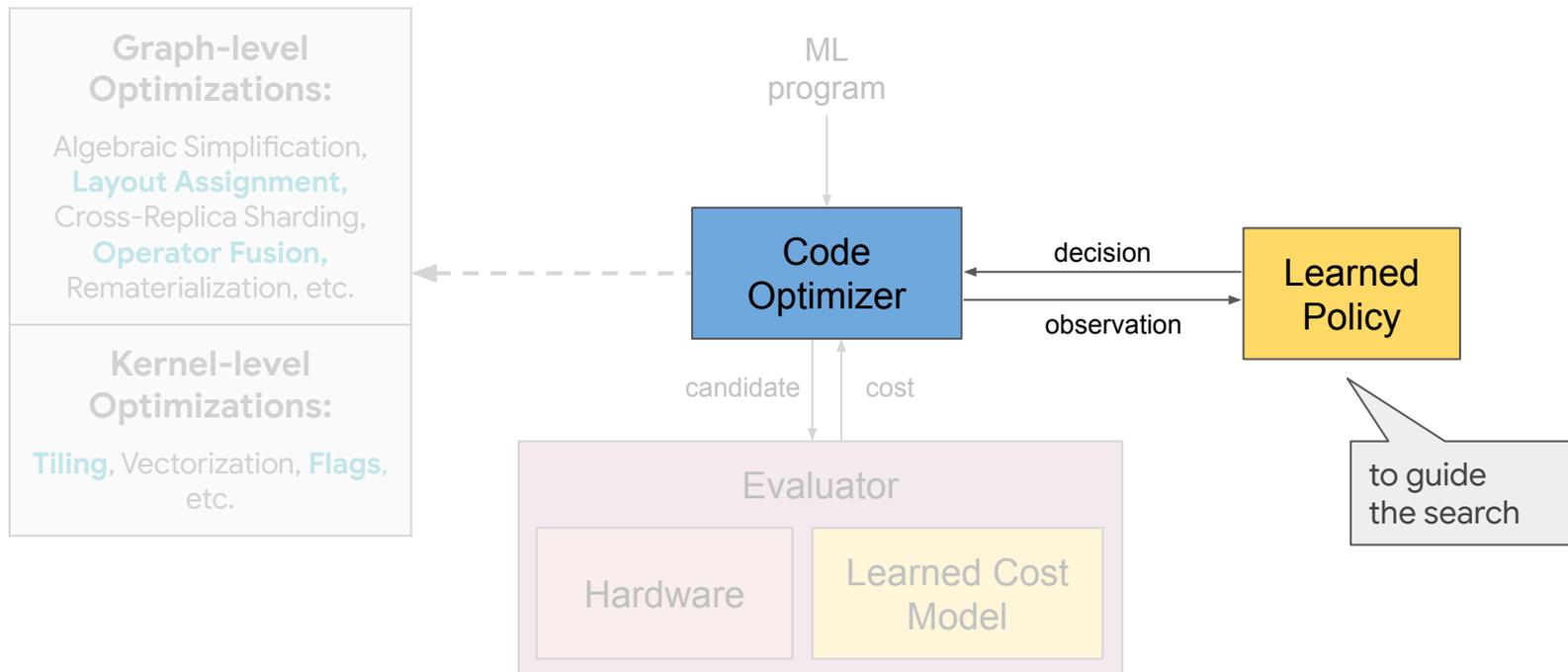
## Runtime Speedup (x)



## Tuning Time (min)



# Search Strategies



# Search Strategies

- Exhaustive
- Simulated annealing (SA)
- Evolutionary (EVO)
- Model-based optimization (MBO)
- Deep reinforcement learning (RL)

# Model-Based Optimization (MBO)

- At each optimization round, a set of candidate **regression models are fit to the acquired data.**
- Good models are **assembled to define an acquisition function.**
- The acquisition function is then **optimized by EVO to generate a new batch of samples.**
- Candidate models: ridge regression, random forests, gradient boosting, and neural networks

*Ref: Angermueller et al., Model-based reinforcement learning for biological sequence design, ICLR 2019*

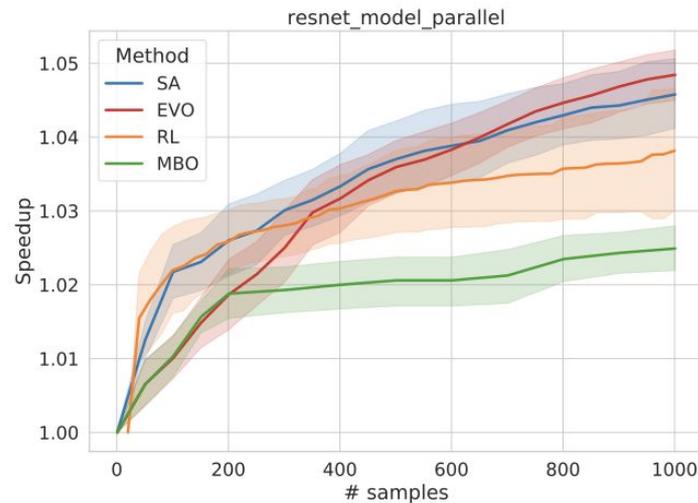
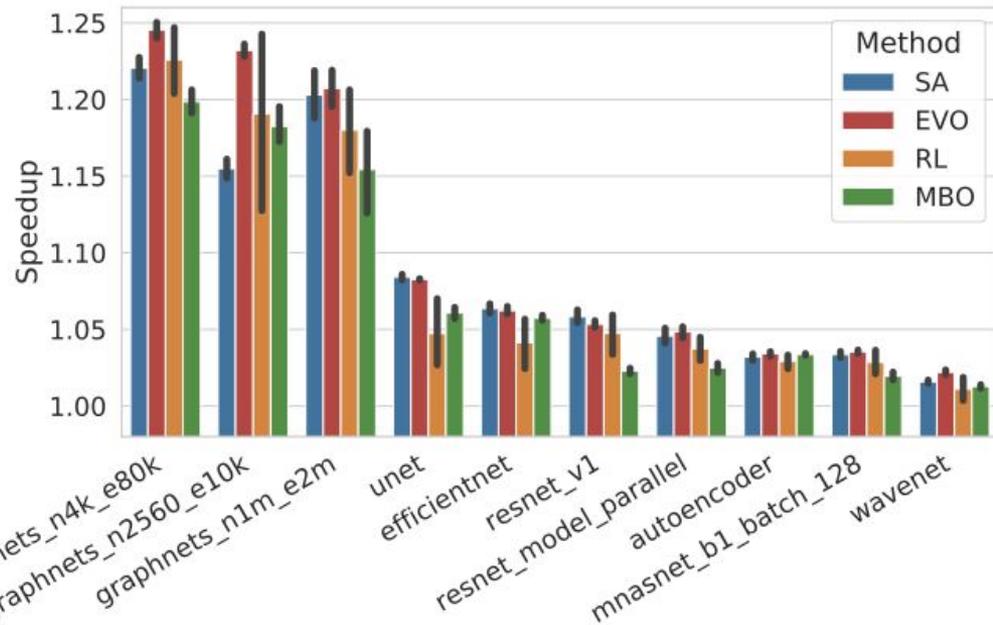
# Deep Reinforcement Learning (RL)

- Designed specifically for ML compiler's graph optimizations
- Uses a **graph neural network** to create node embeddings and segmented **recurrent attention layers** to capture long-range dependencies
- **Non-autoregressive**
  - N node decisions are done in parallel
  - Conventional autoregressive approach is infeasible as N can be as large as 100k

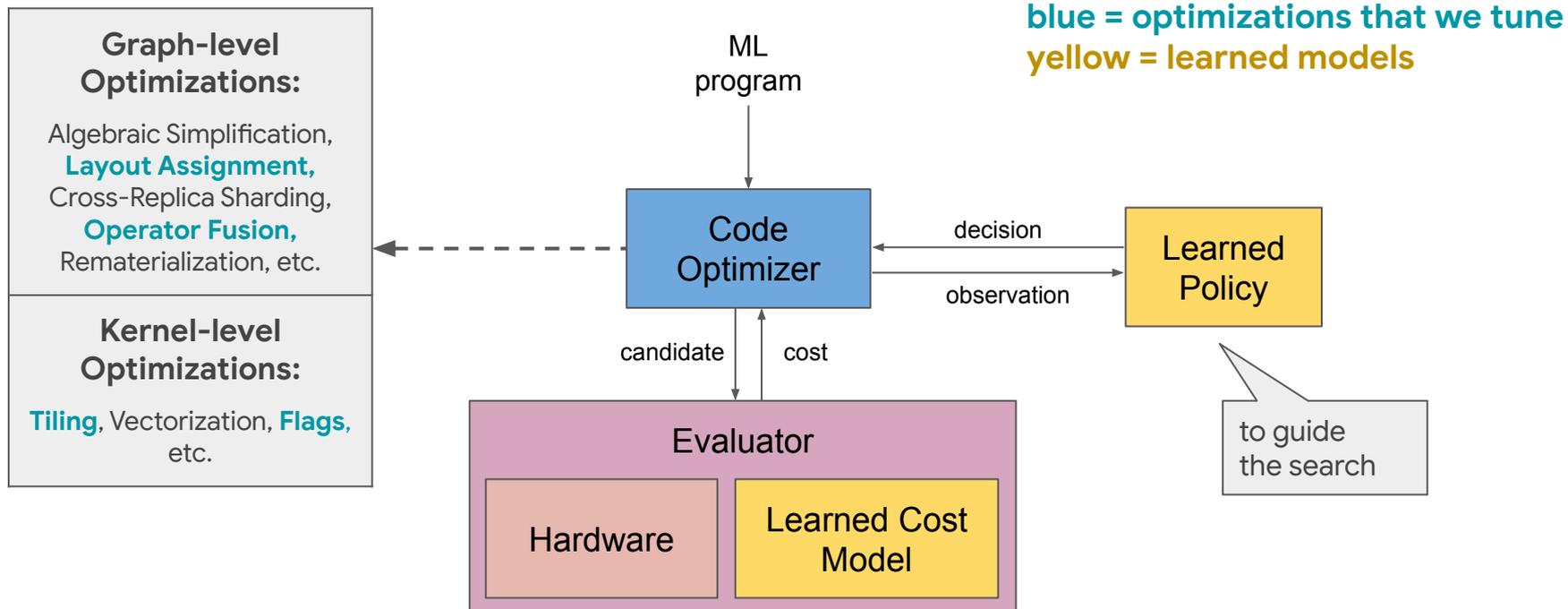
*Ref: Zhou et al., Transferable graph optimizers for ml compilers, NeurIPS 2020*

# Search Strategies: Fusion Autotuning

Average speedup across 10 runs. Each run evaluated 10,000 candidates.

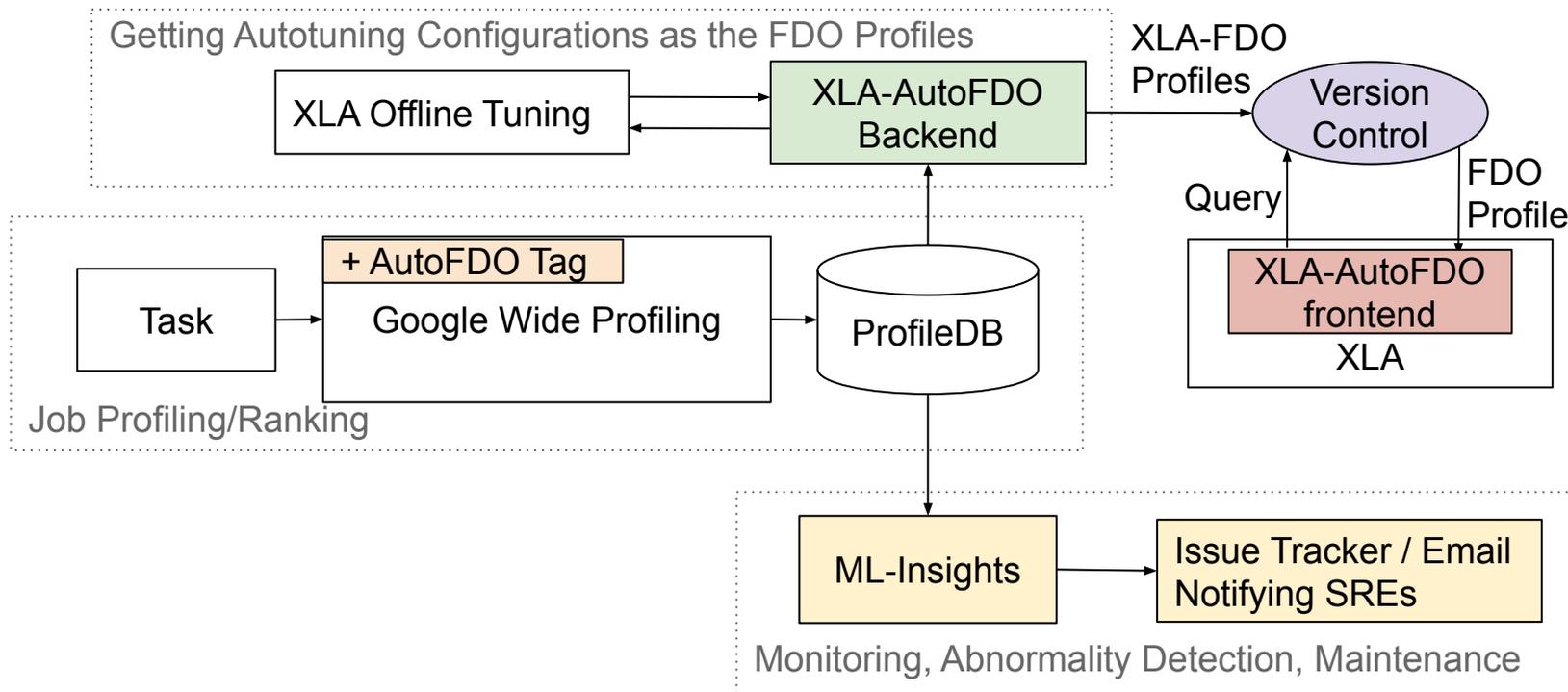


# XTAT: XLA TPU Autotuner



# Data-Center Scale Deployment

# Fleet-Wide TPU Autotuning at Google



# Fleet-Wide TPU Autotuning at Google

- Have deployed the tile size and flags autotuning to optimize top workloads in the TPU fleet daily
- Learned cost model enabled tuning 20x more kernels per day
- Save >1% of total TPU consumption
- Savings / tuning cost: >10x

# References

Phothilimthana et al., **A Flexible Approach to Autotuning Multi-Pass Machine Learning Compilers**, PACT 2021.

Kaufman and Phothilimthana et al., **A Learned Performance Model for Tensor Processing Units**, MLSys 2021.

# Contributors

Mangpo Phothilimthana

Amit Sabne

Karthik Srinivasa Murthy

Nikhil Sarda

Yanqi Zhou

Christof Angermueller

Emma Wang

Sam Kaufman

Charith Mendis

Sudip Roy

Mike Burrows

Berkin Ilbeyi

Bjarke Rouné

Blake Hechtman

Ketan Mandke

Reza Farahani

Shen Wang

Yuanzhong Xu