

LATTE: Layer Algorithm-aware Training Time Estimation for Heterogeneous Federated Learning

Kun Wang¹, Zimu Zhou², Zhenjiang Li¹

¹ Department of Computer Science, City University of Hong Kong

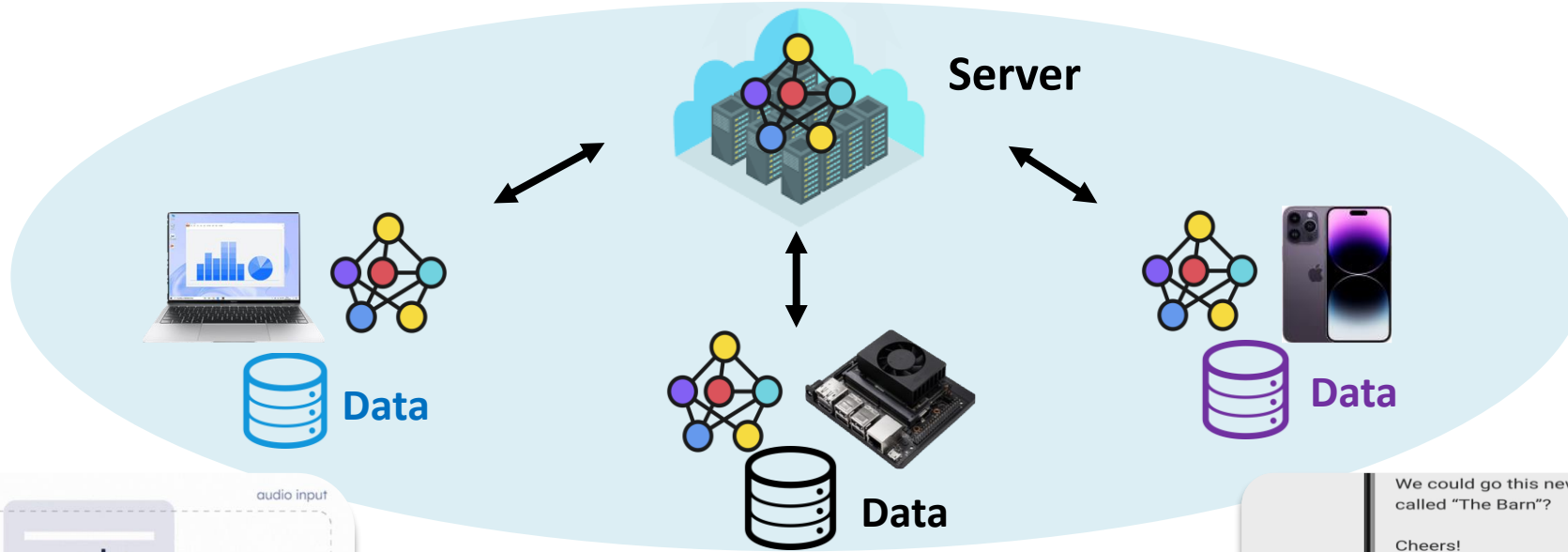
² School of Data Science, City University of Hong Kong



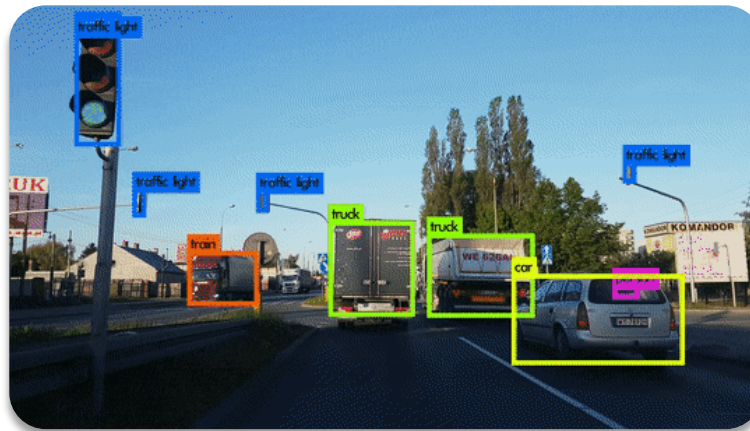
香港城市大學

City University of Hong Kong

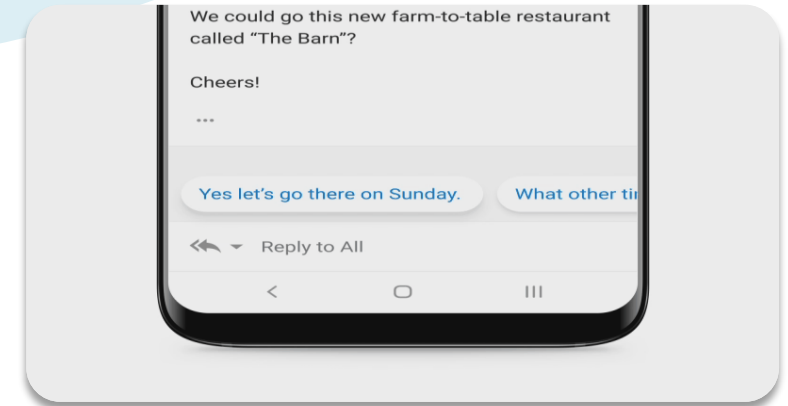
Background - On-Device Federated Training



Speech Recognition

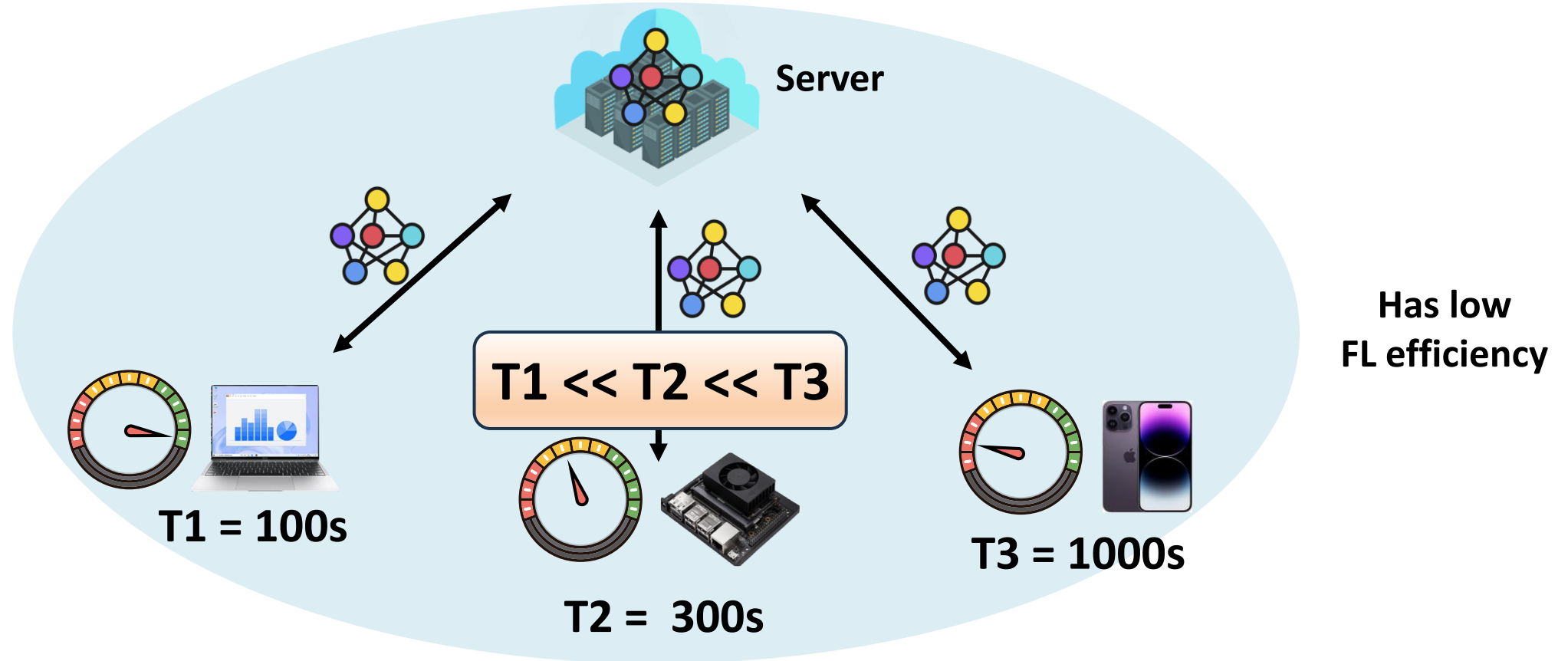


Autonomous Driving



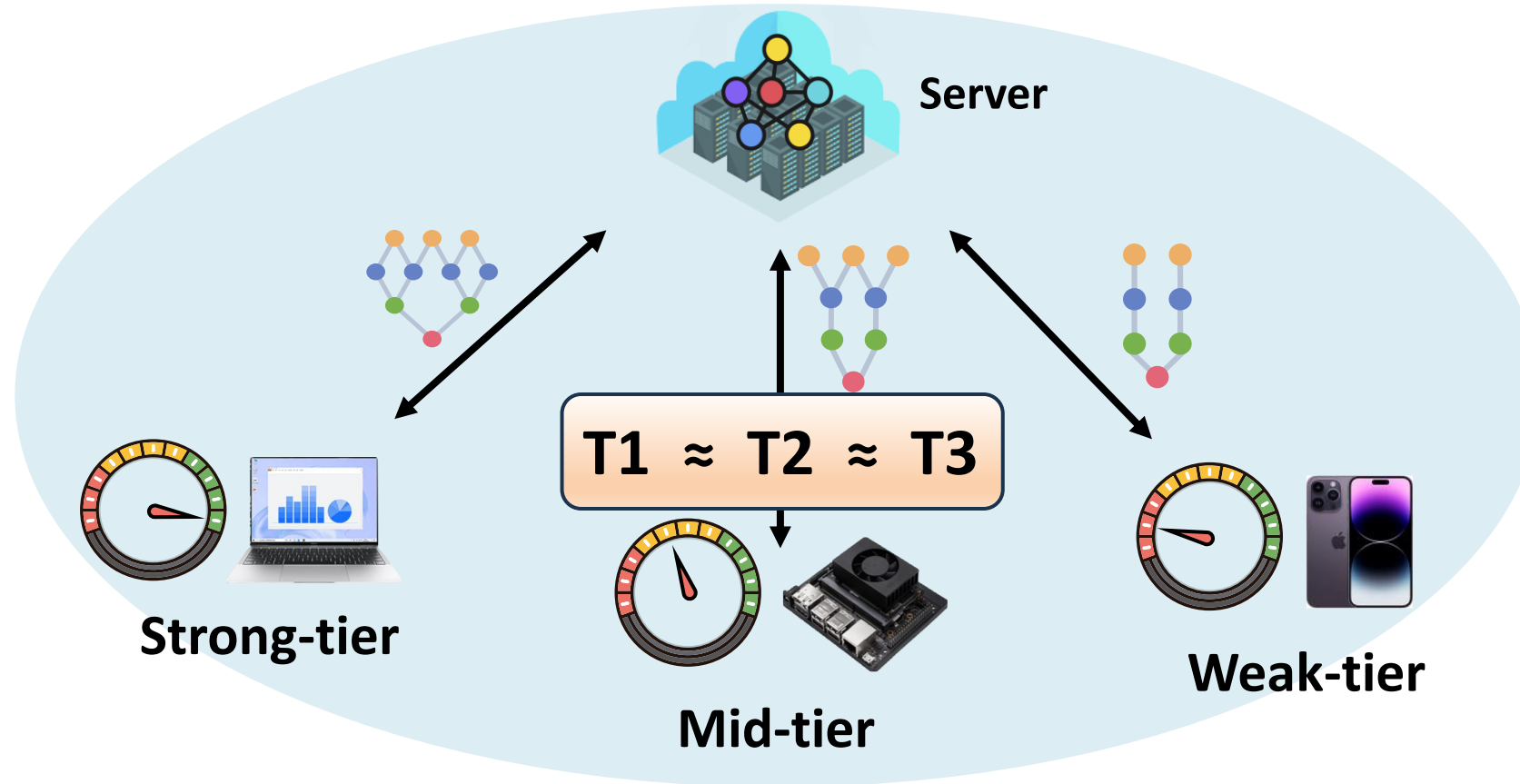
Smart Reply helper

Huge training time gap between devices



Stronger Devices need to **wait** for Weaker Devices

Idea - Heterogeneous Federated Learning

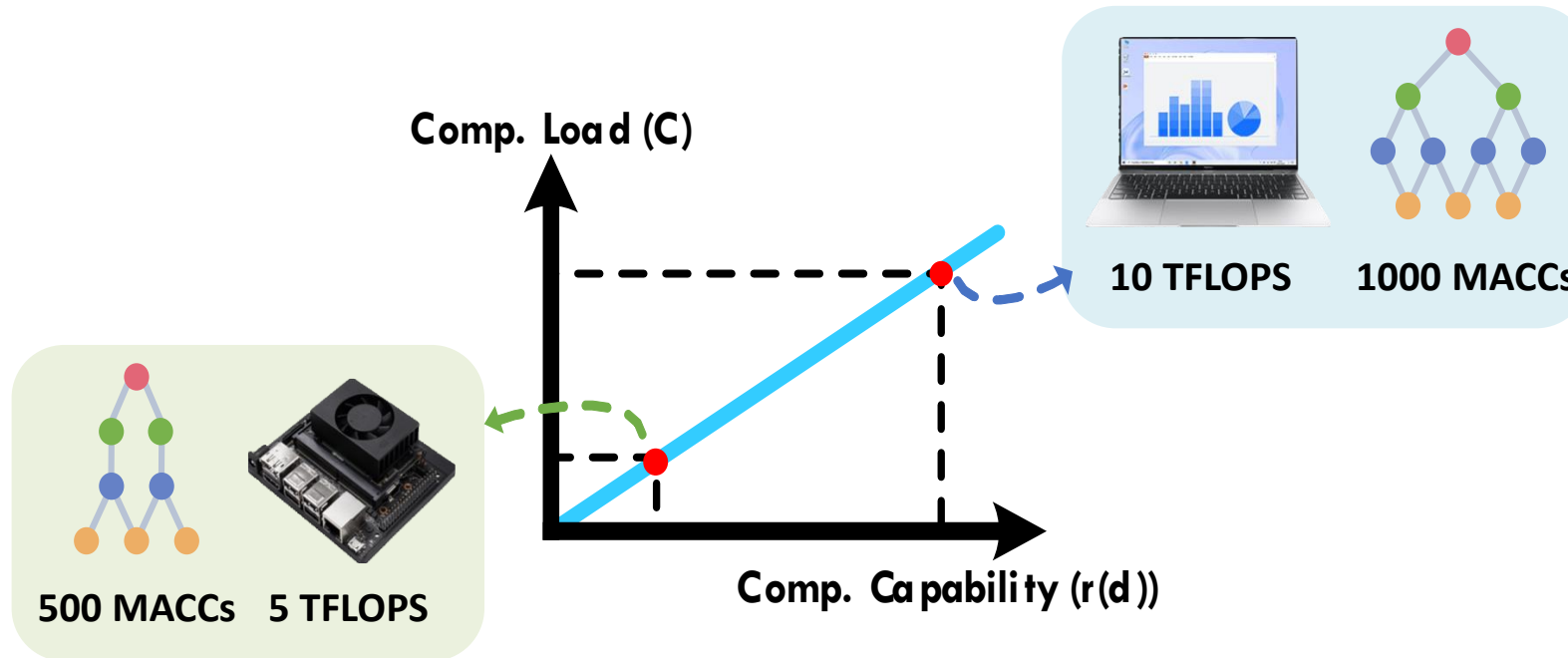


It was proved that
can really improve
FL efficiency !

Different devices but **Similar** Training Time

Recent work - Fine-grained method

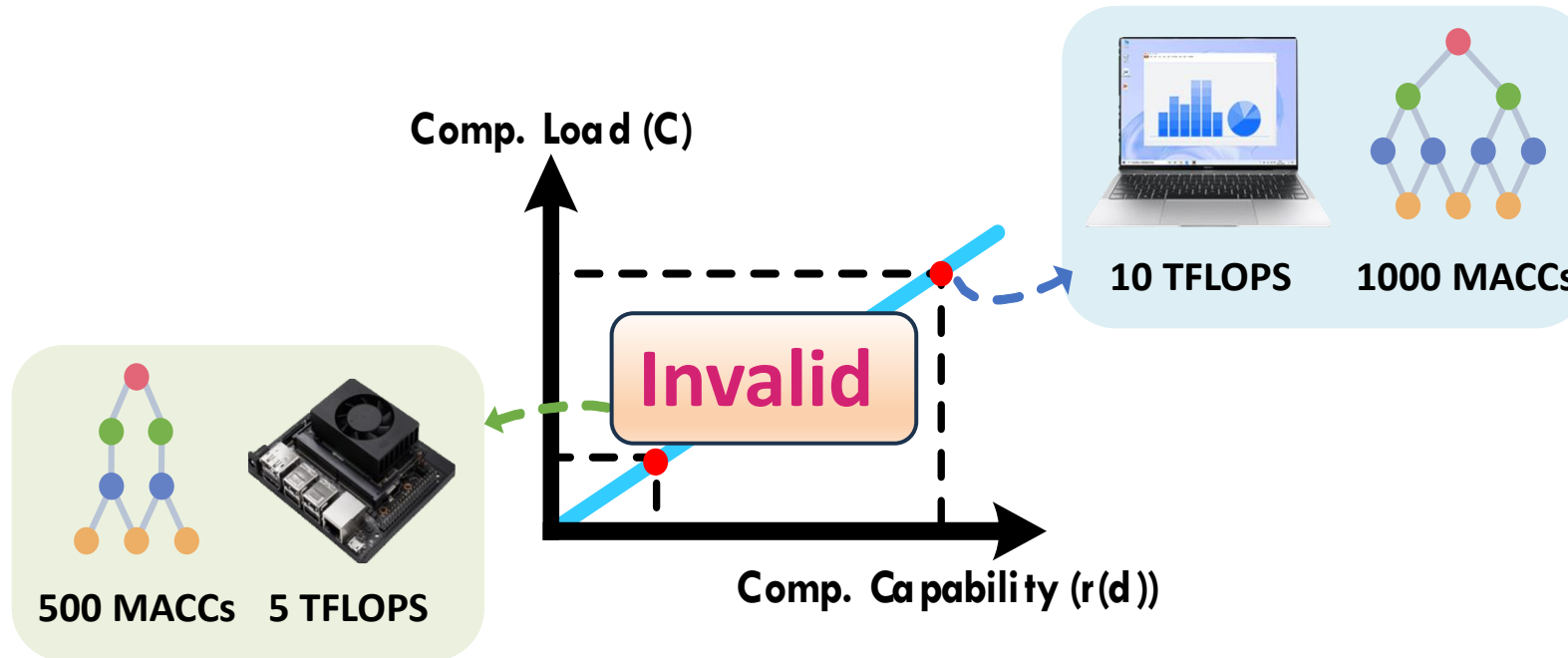
$$T = \frac{C_1}{r(d_1)} = \frac{C_2}{r(d_2)}$$



Allocating sub-models according
device's computing power (i.e., FLOPS)

Recent work - Fine-grained method

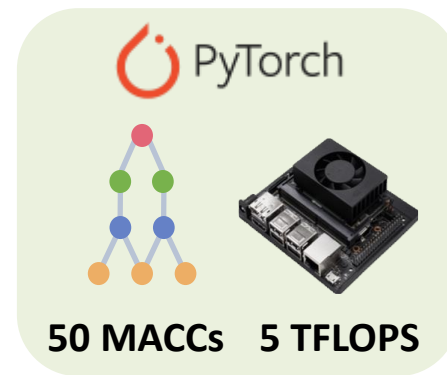
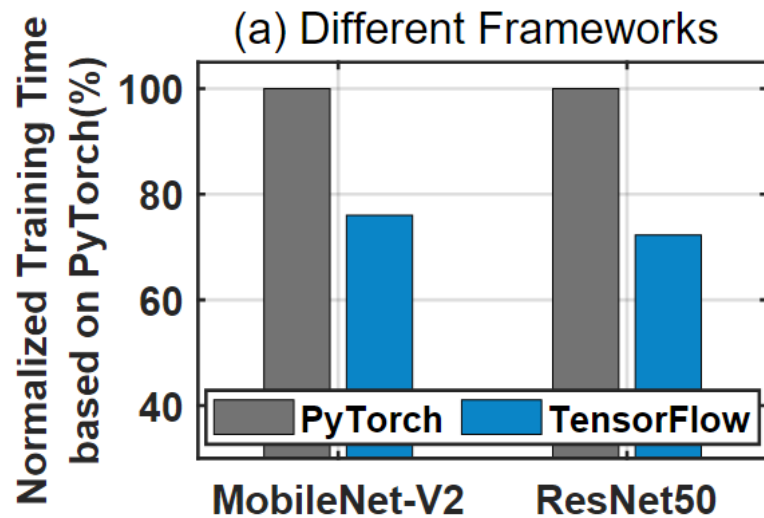
$$T = \frac{C_1}{r(d_1)} = \frac{C_2}{r(d_2)}$$



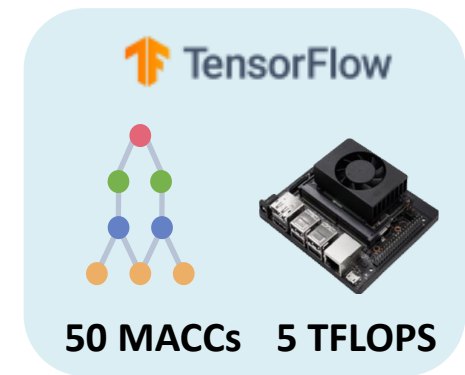
This training time modeling is still **over-simple!**

Problem 1 - Training Time Inconsistency

- Even training same models in the same devices, **training time has huge gap** by using different DL frameworks.



V.S.



Training Time Inconsistency Problem

Key Observation - Layer Algorithm Diversity

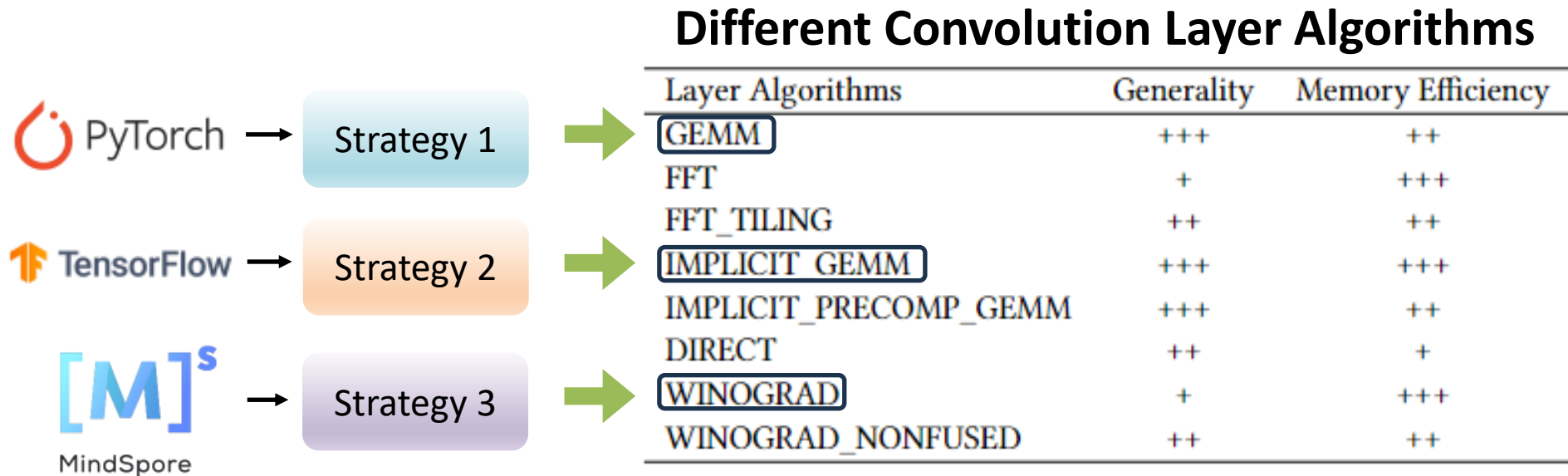
- There are several **candidate layer algorithms implementation** in DL frameworks.

```
1: static const algo_t Fwd_algos[] = {
2:     CUDNN_CONVOLUTION_FWD_ALGO_GEMM,
3:
4:     1: static const algo_t Bwd_algos[] = {
5:         2:     CUDNN_CONVOLUTION_BWD_FILTER_GEMM,
6:         3:     CUDNN_CONVOLUTION_BWD_FILTER_FFT,
7:         4:     CUDNN_CONVOLUTION_BWD_FILTER_FFT_TILING,
8:         5:     CUDNN_CONVOLUTION_BWD_FILTER_IMPLICIT_GEMM,
9:         6:     CUDNN_CONVOLUTION_BWD_FILTER_IMPLICIT_PRECOMP_GEMM,
10:        7:     CUDNN_CONVOLUTION_BWD_FILTER_DIRECT,
11:        8:     CUDNN_CONVOLUTION_BWD_FILTER_WINOGRAD,
12:        9:     CUDNN_CONVOLUTION_BWD_FILTER_WINOGRAD_NONFUSED,
13:    10: };
```

Analyzing source codes of DL framework's training mechanism

Key Observation - Different Algorithm Selecting Strategy

- Different DL frameworks may **select different layer algorithms** as implementation.



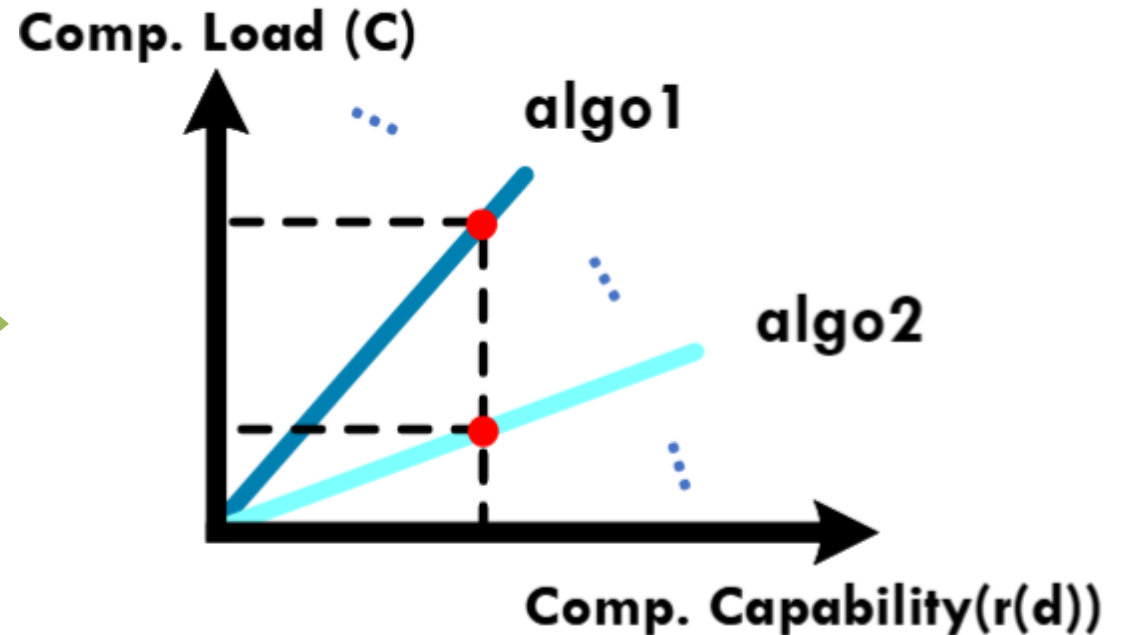
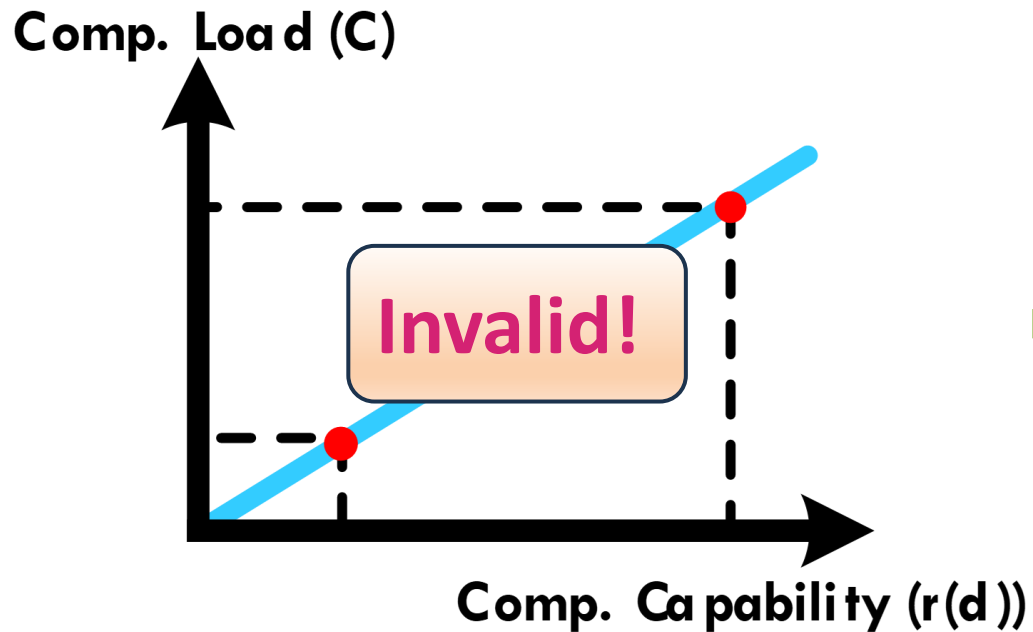
Different layer algorithms have **different computation workloads thus **different training time.****

Training Time Modeling Reformulation

$$T = \frac{C_1}{r(d_1)} = \frac{C_2}{r(d_2)}$$



$$T = \frac{C(\text{algo}_1)}{r(d_1)} = \frac{C(\text{algo}_2)}{r(d_2)}$$



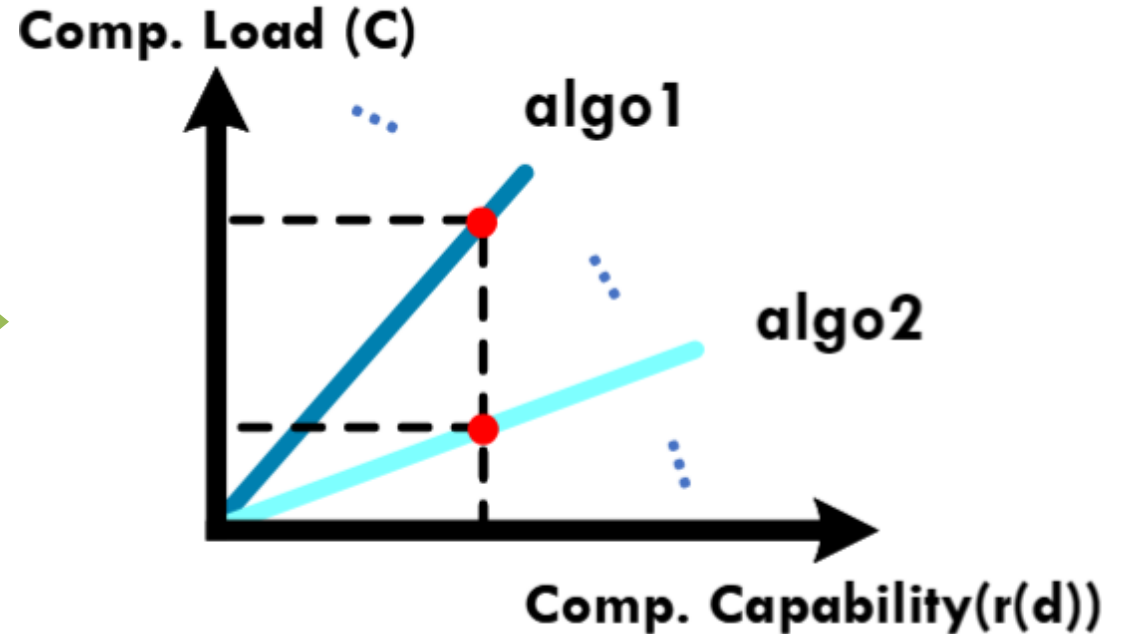
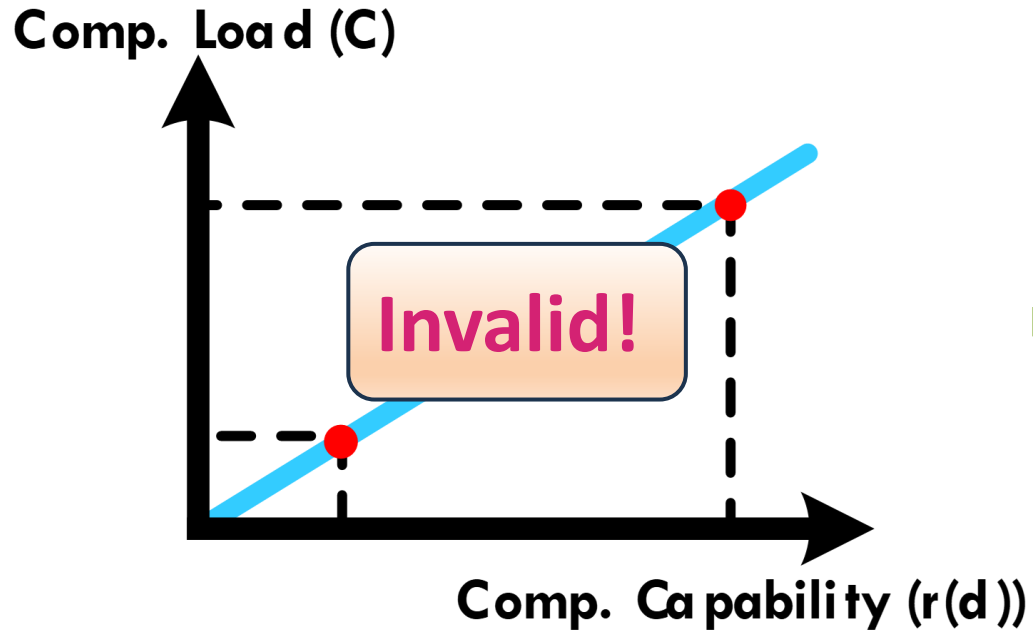
Accurate training time modeling
need to **consider the layer algorithms.**

Training Time Modeling Reformulation

$$T = \frac{C_1}{r(d_1)} = \frac{C_2}{r(d_2)}$$

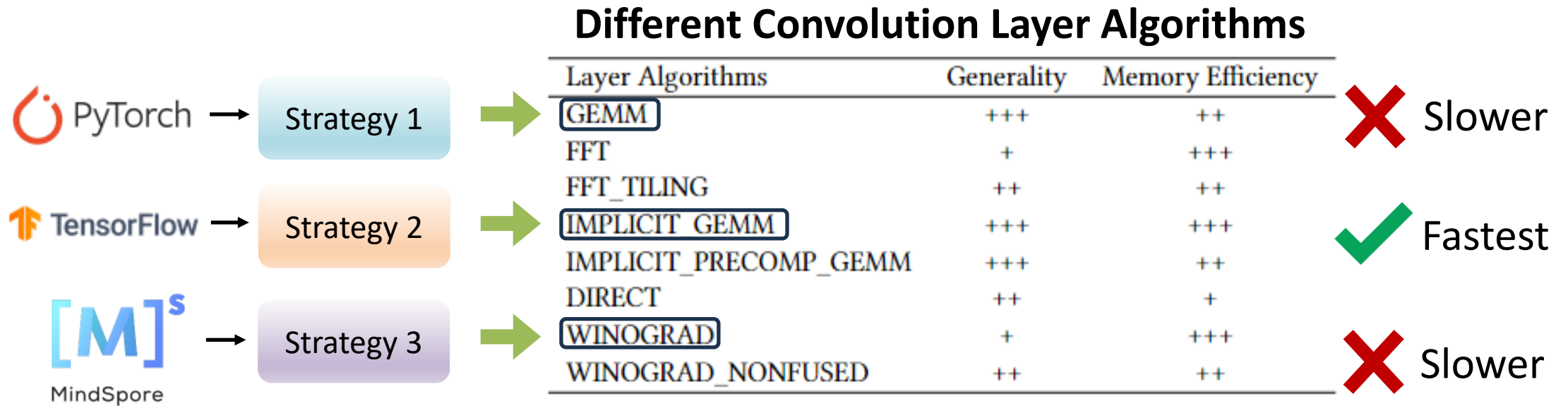


$$T = \frac{C(\text{algo}_1)}{r(d_1)} = \frac{C(\text{algo}_2)}{r(d_2)}$$



Now we indeed can allocate sub-models efficiently, however...

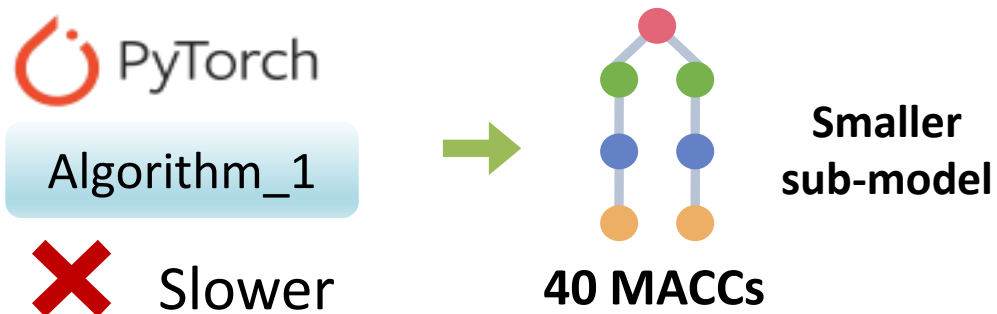
Problem 2 - Most Strategies are not optimal



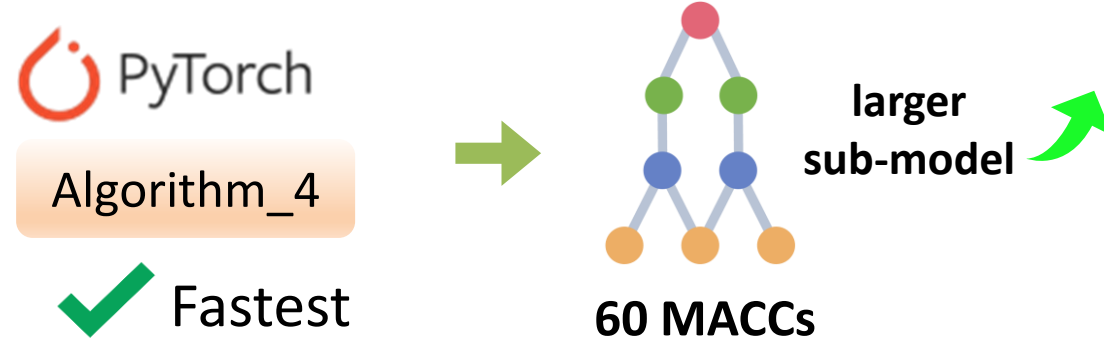
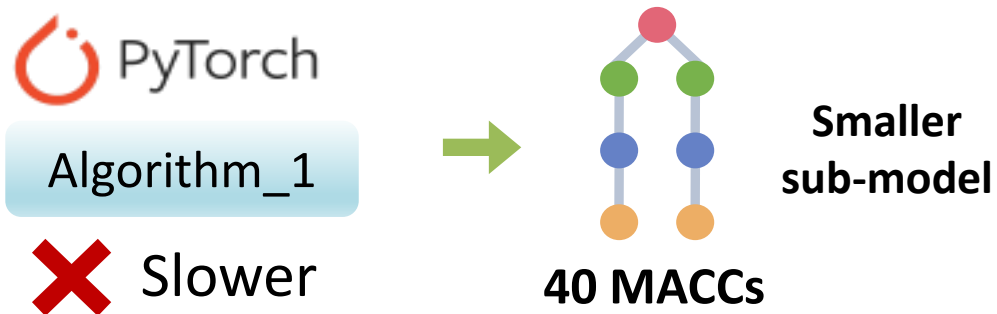
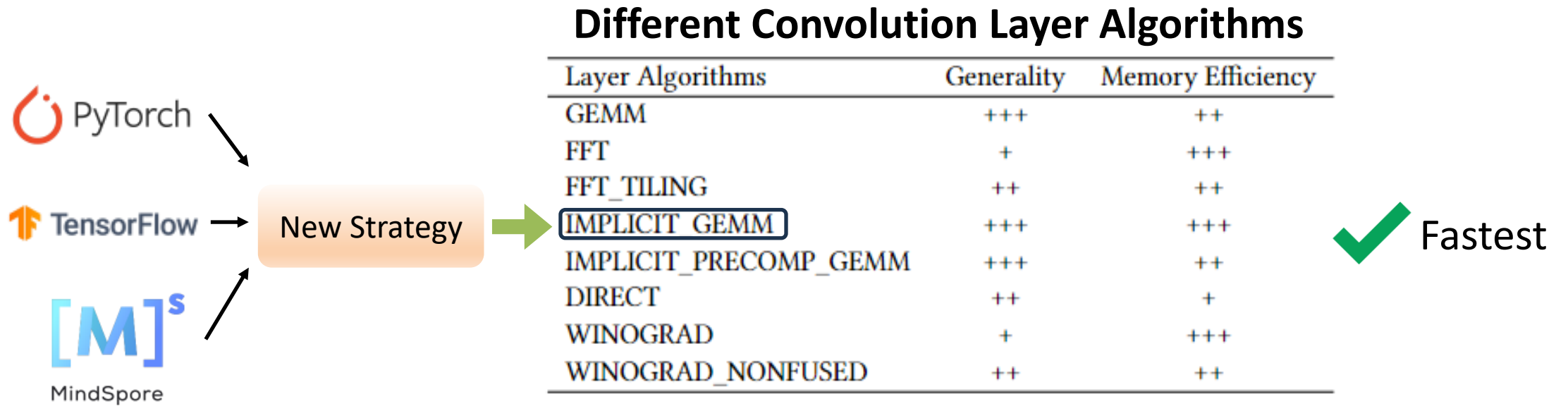
Most strategies **selecting the sub-optimal algorithms** as their layer implementation

Problem 2 - Most Strategies are not optimal

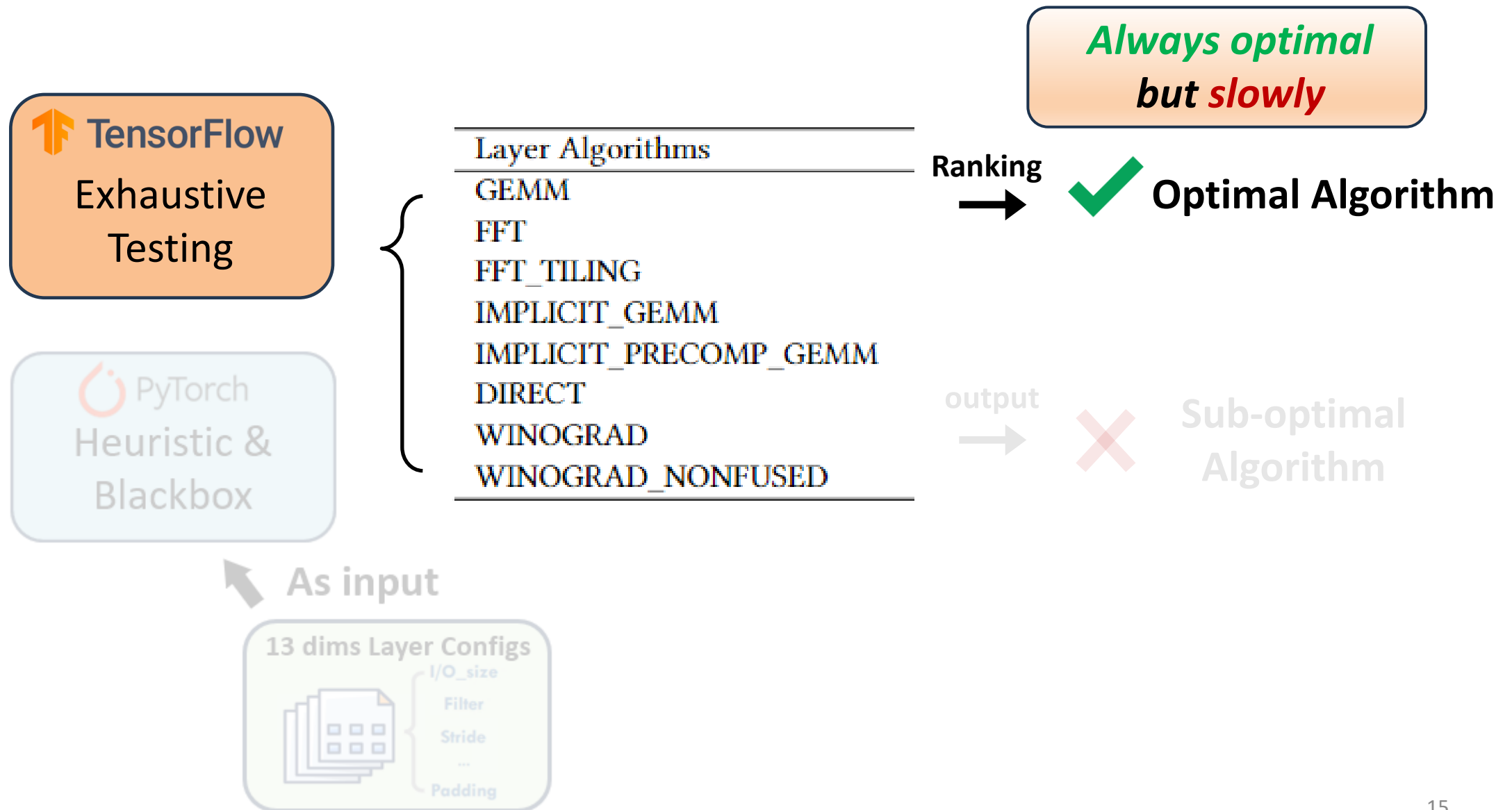
		Different Convolution Layer Algorithms			
		Layer Algorithms	Generality	Memory Efficiency	
PyTorch → Strategy 1	→	GEMM	+++	++	✗ Slower
		FFT	+	+++	
		FFT_TILING	++	++	
TensorFlow → Strategy 2	→	IMPLICIT_GEMM	+++	+++	✓ Fastest
		IMPLICIT_PRECOMP_GEMM	+++	++	
		DIRECT	++	+	
[M] ^s MindSpore → Strategy 3	→	WINOGRAD	+	+++	✗ Slower
		WINOGRAD_NONFUSED	++	++	



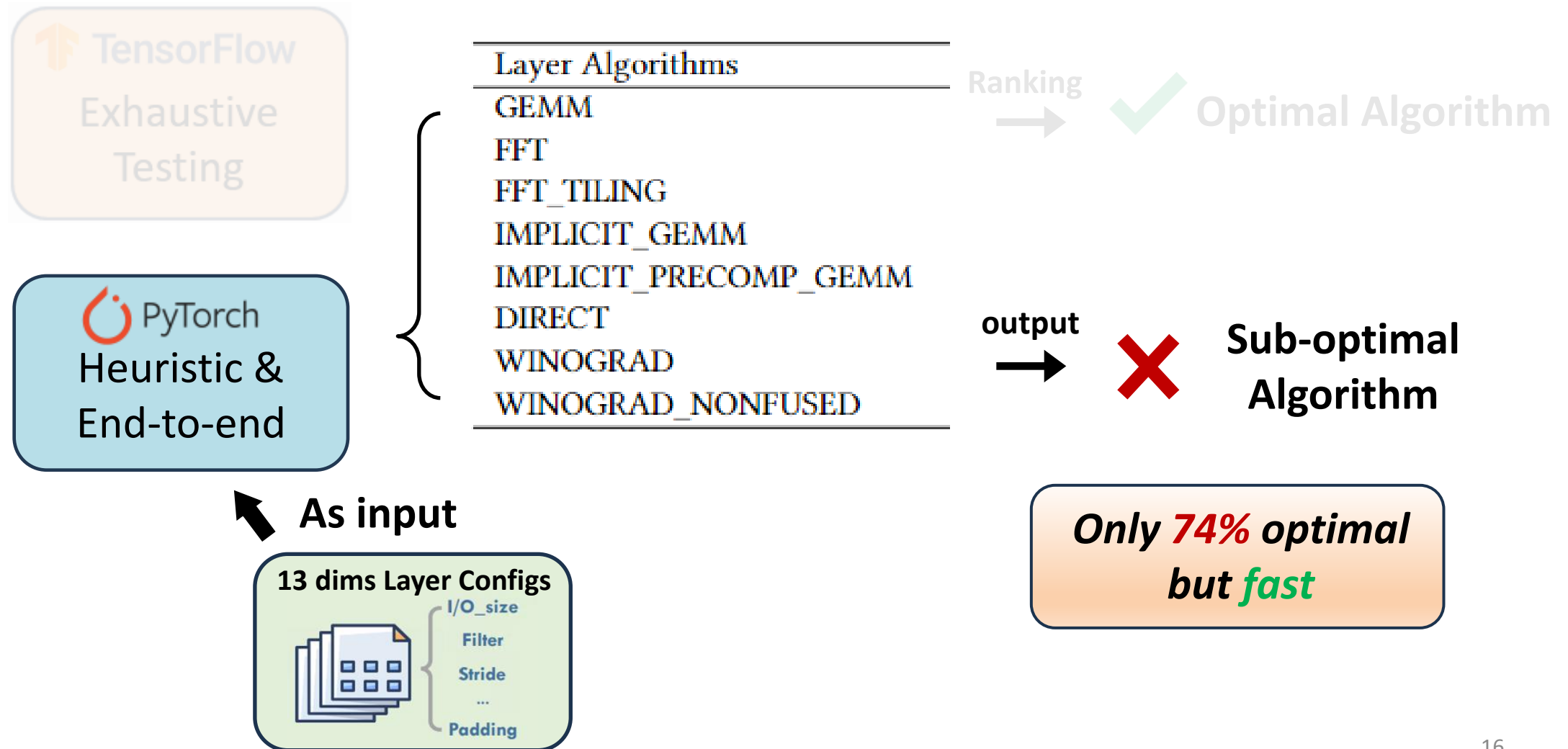
How to design a better selecting strategy?



Observation - TensorFlow's Exhaustive Testing Strategy

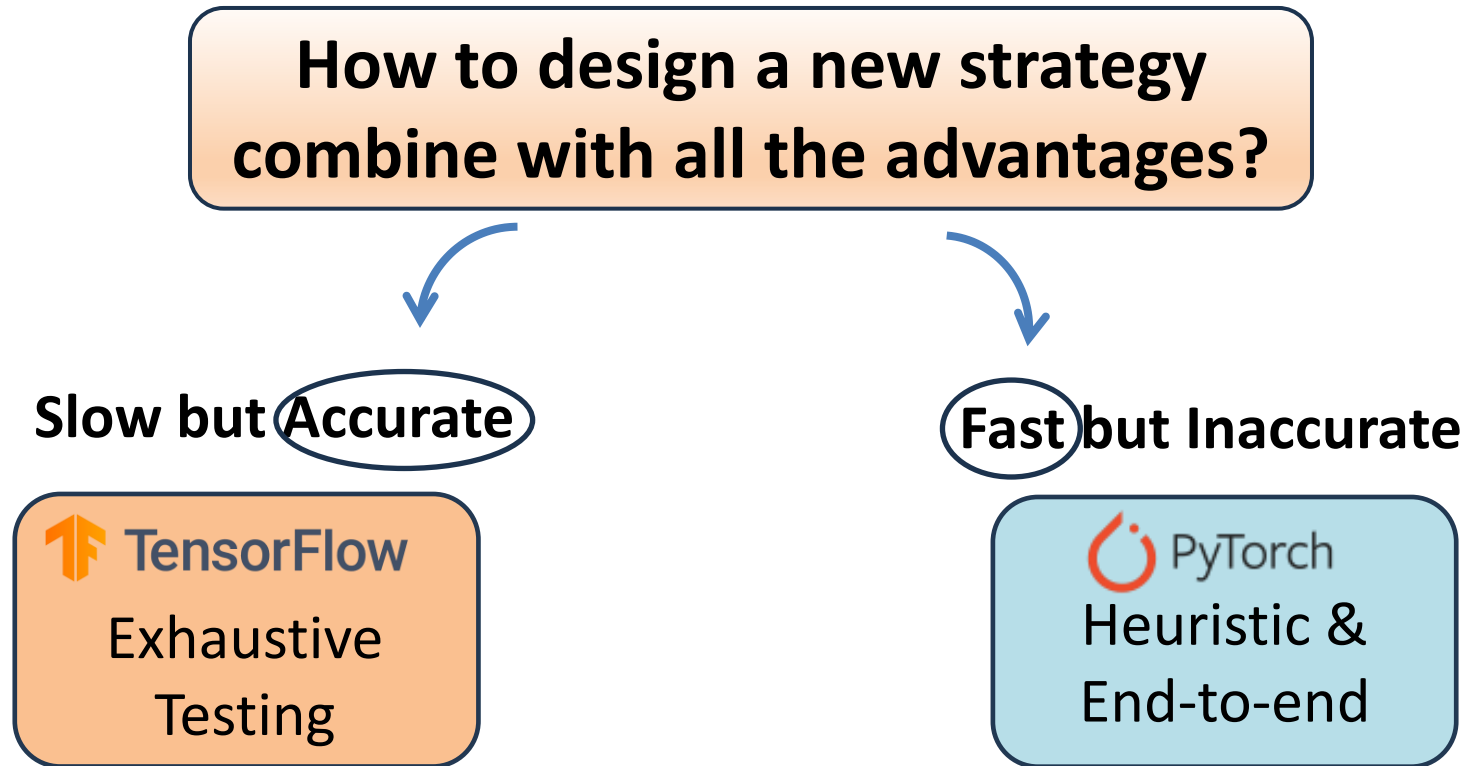


Observation - PyTorch's Heuristic Blackbox Strategy



Thinking

- How to design an **accurate** and **fast** algorithm selecting strategy?



Design 1 - Layer-Algorithm Selector

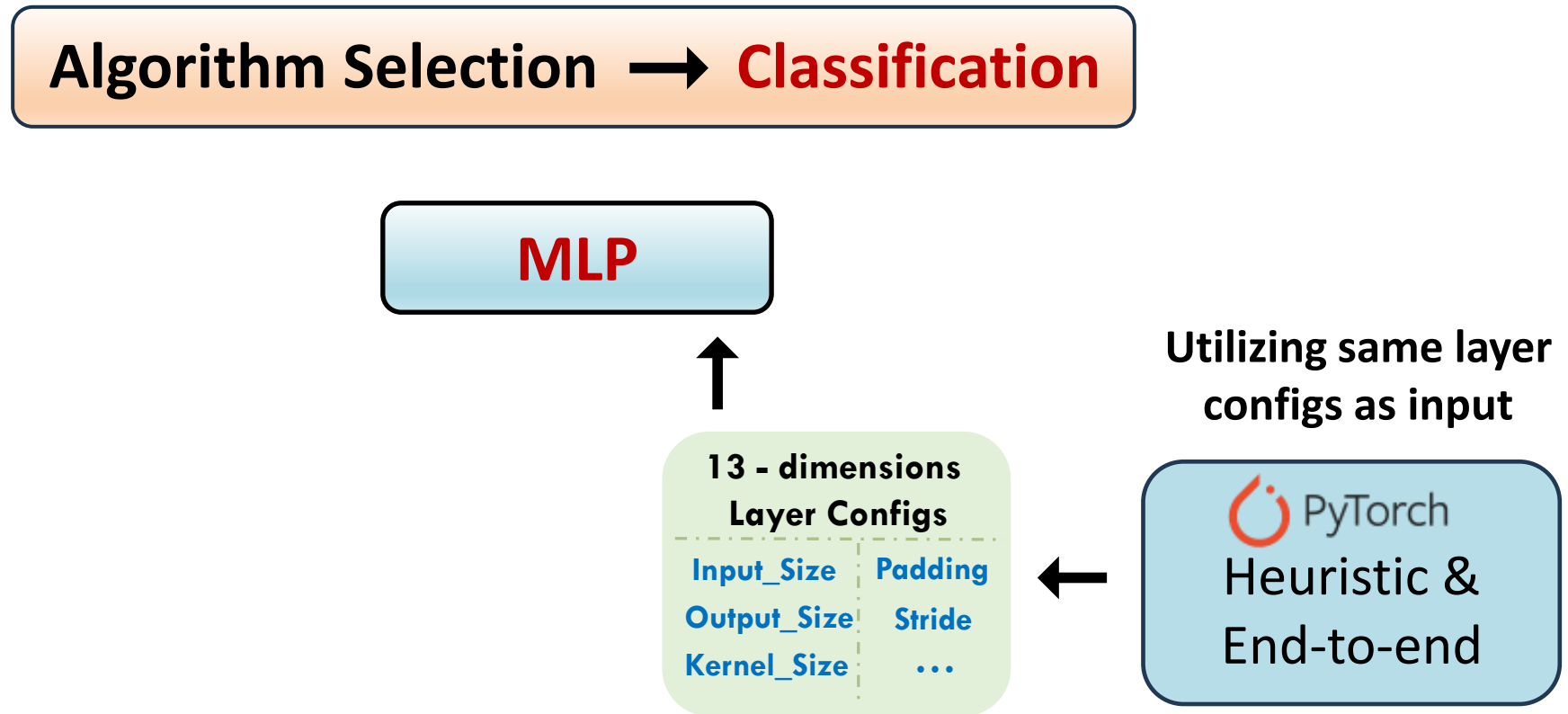
- Data-driven based Layer-Algorithm Selector

Algorithm Selection → Classification

MLP

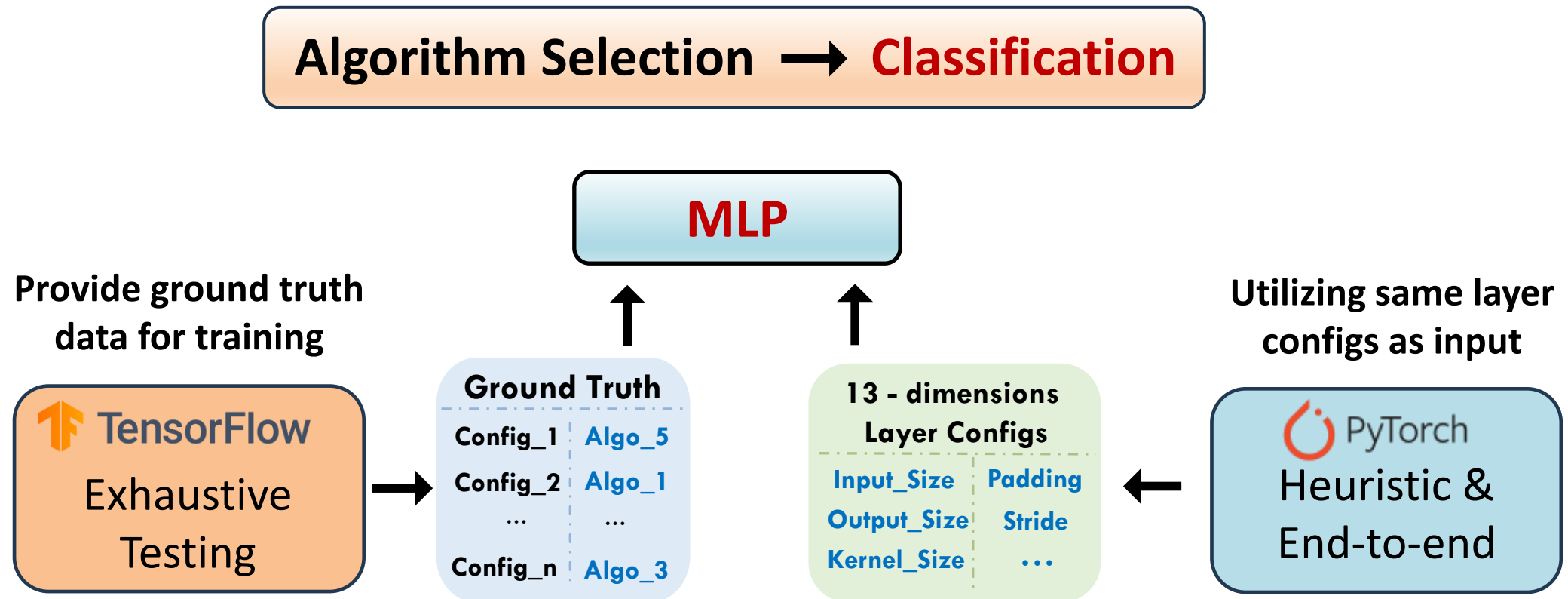
Design 1 - Layer-Algorithm Selector

- Data-driven based Layer-Algorithm Selector



Design 1 - Layer-Algorithm Selector

- Data-driven based Layer-Algorithm Selector



Design 2 - Training Time Estimator

- Training Time Estimation by profiling

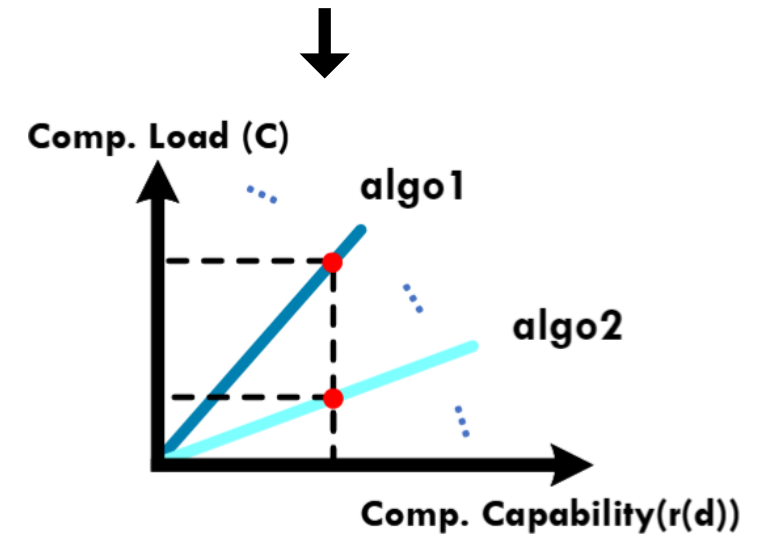


Design 2 - Training Time Estimator

- Training Time Estimation by profiling

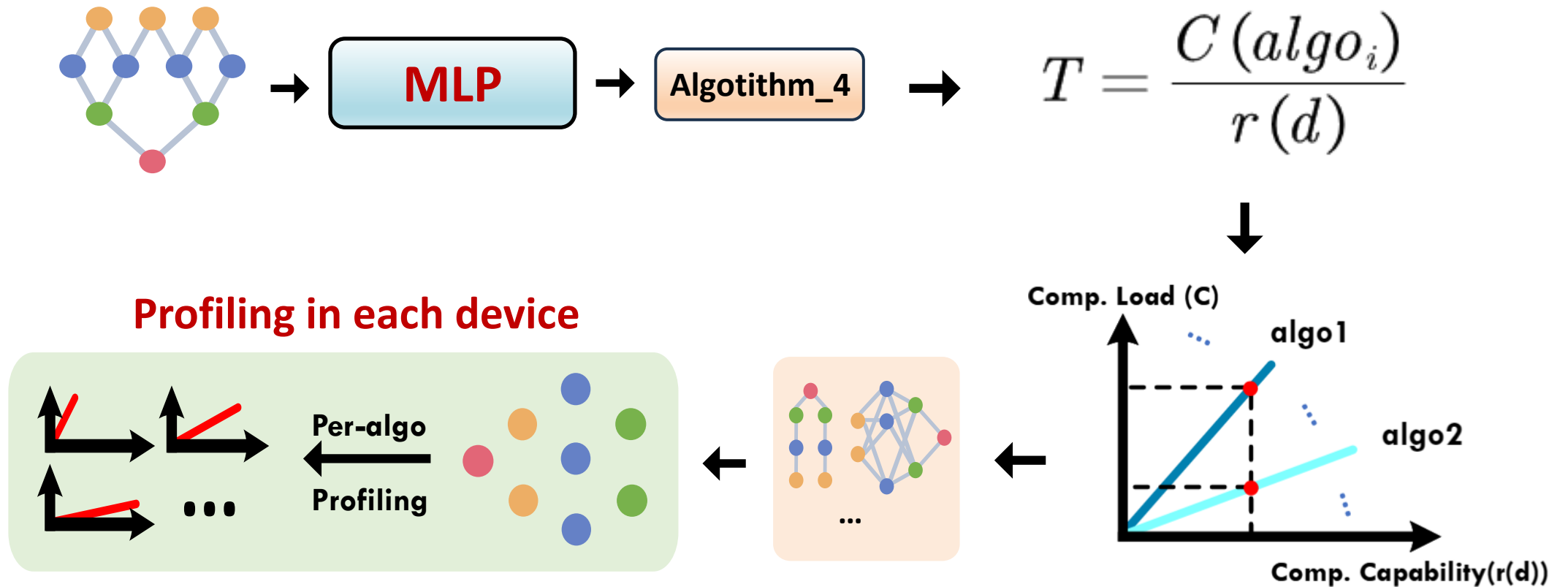


$$T = \frac{C(algo_i)}{r(d)}$$

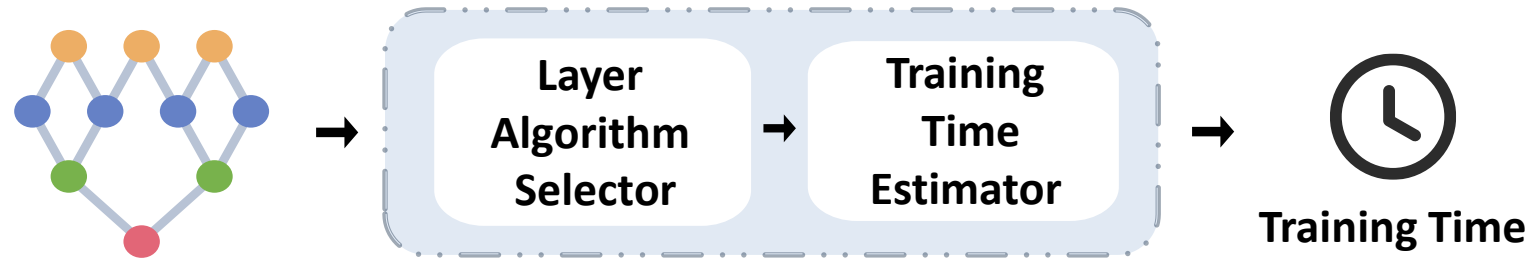


Design 2 - Training Time Estimator

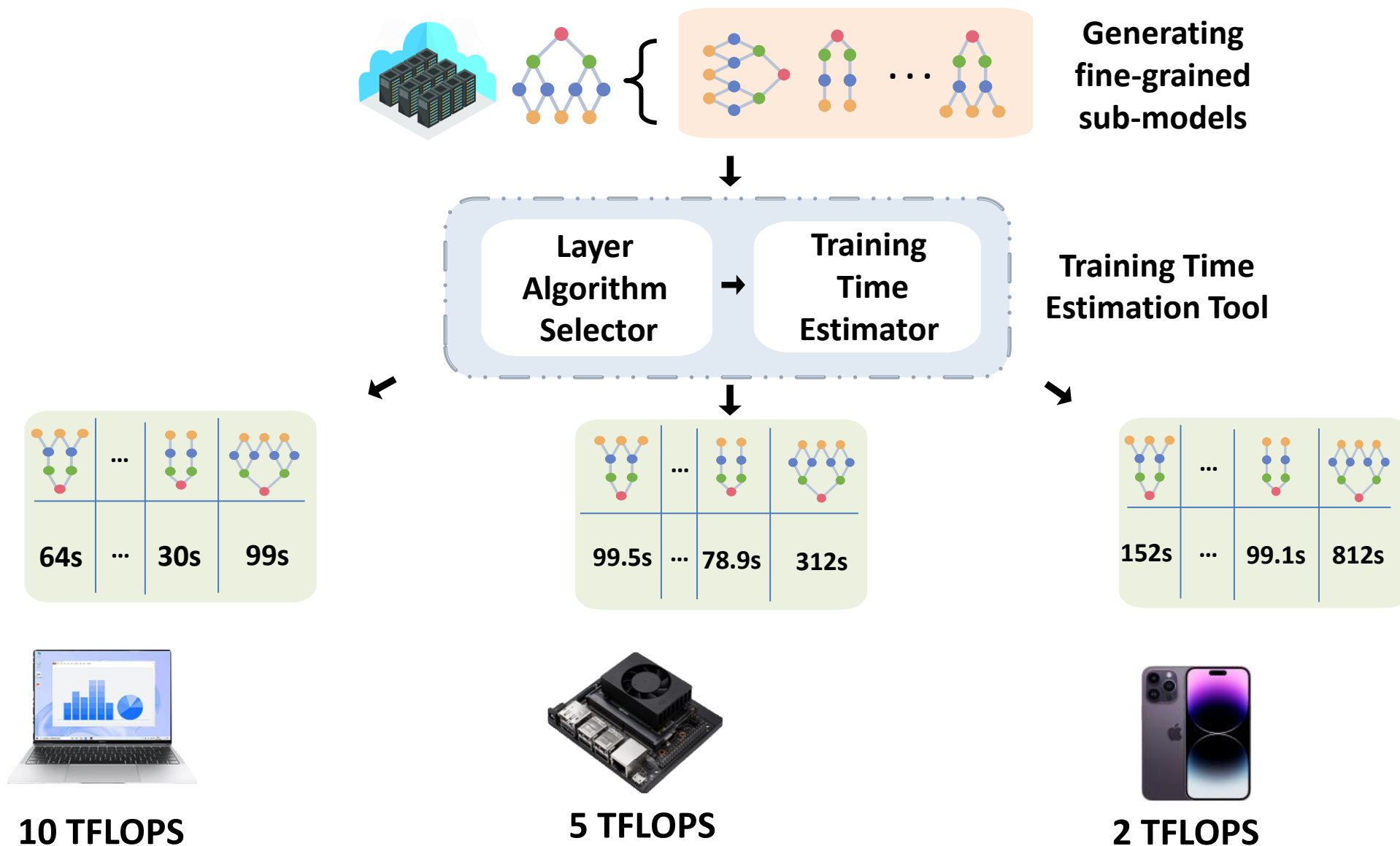
- Training Time Estimation by profiling



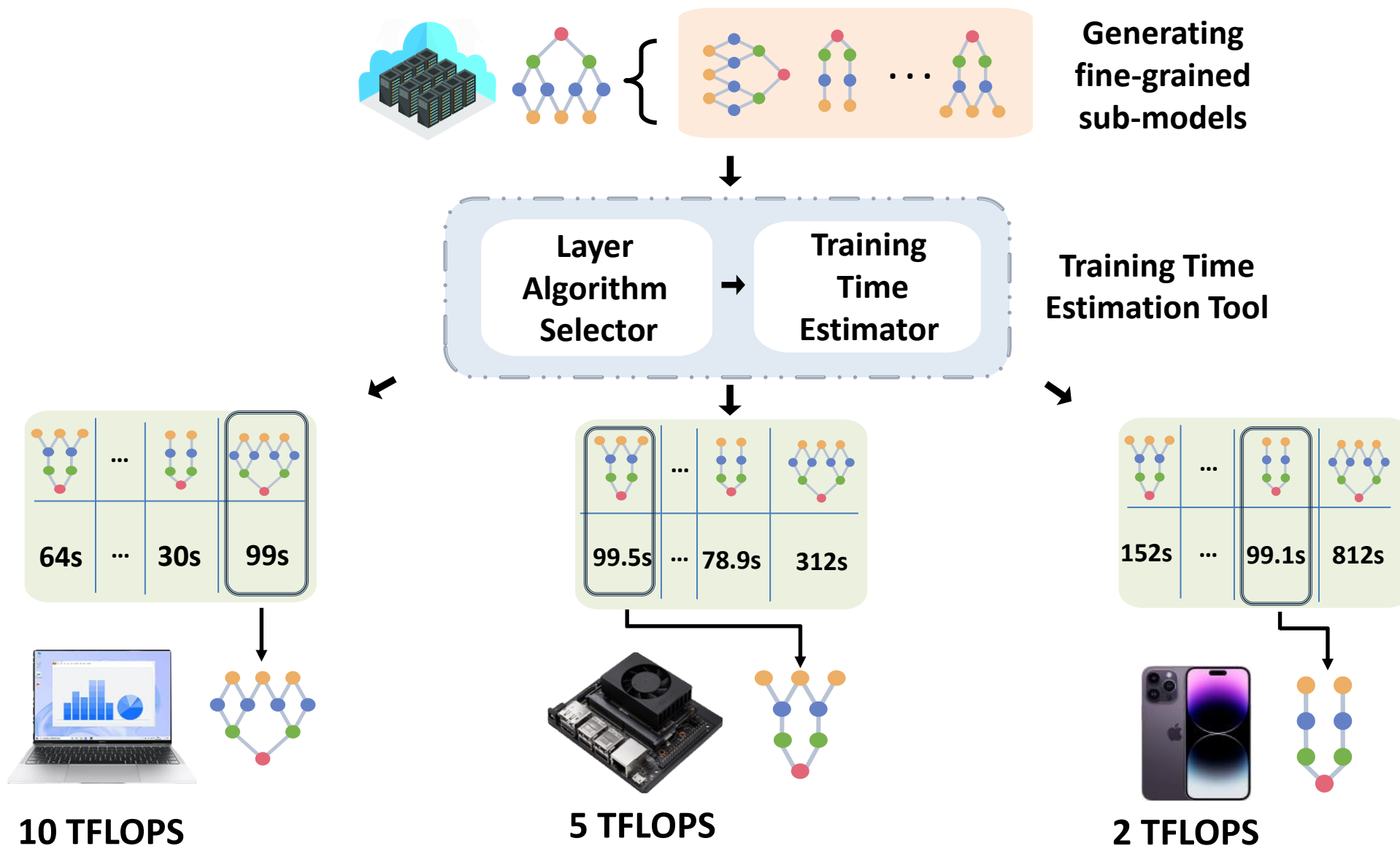
Put all pieces together



Put all pieces together

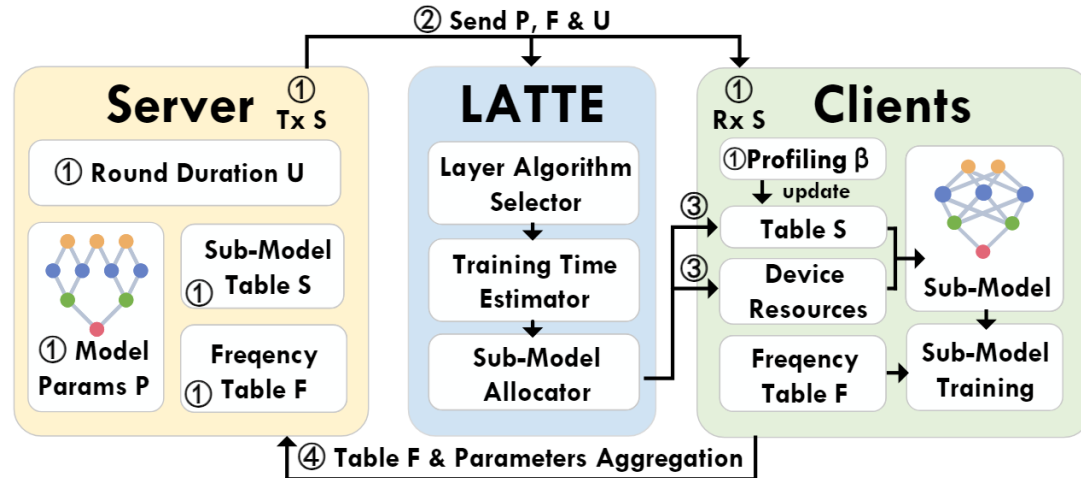


Put all pieces together

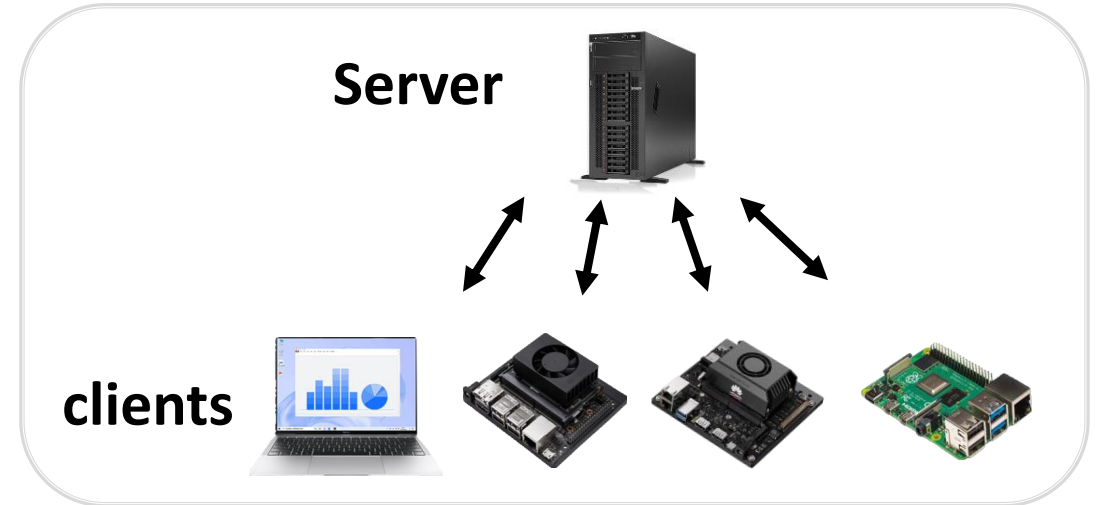


Implementation

Proposed LATTE Framework



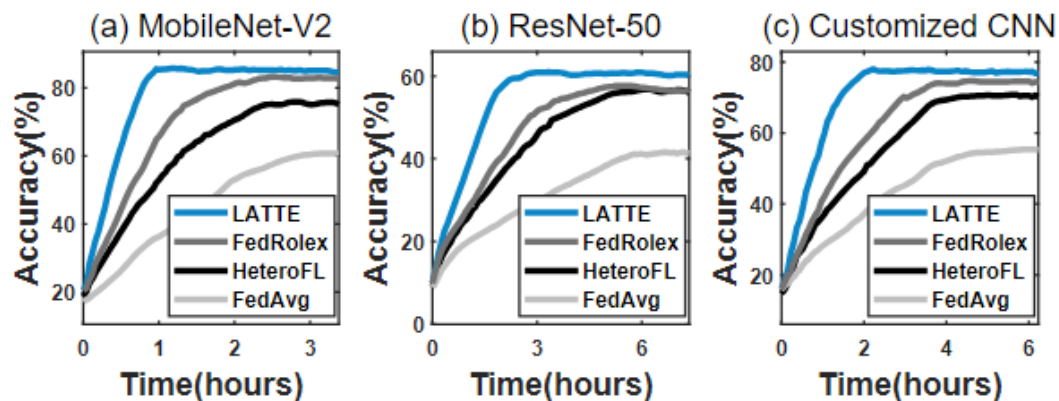
Test-bed



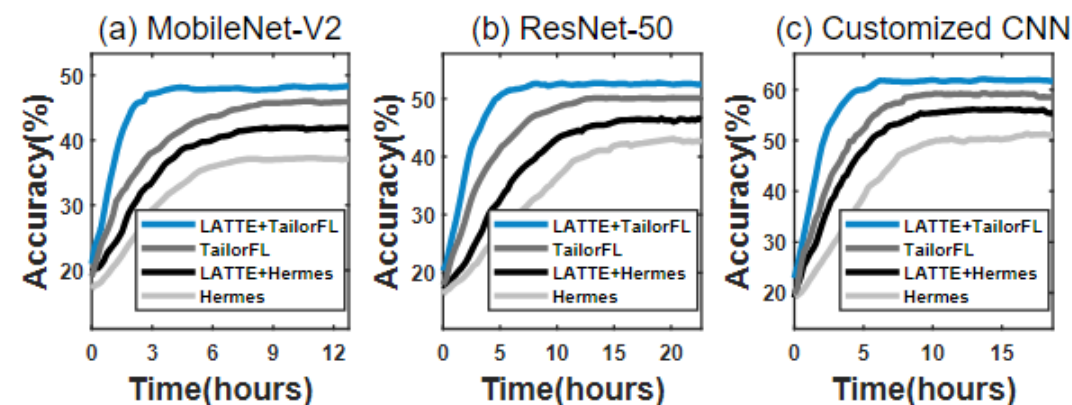
- 1. Time-to-Convergence:** Represent whether our system can accelerate converge speed.
- 2. Layer Algorithm Selector's accuracy:** Represent the Layer Algorithm Selector's performance.
- 3. Training Time Estimator's Precision:** Represent the Training Time Estimator's performance.

Metrics

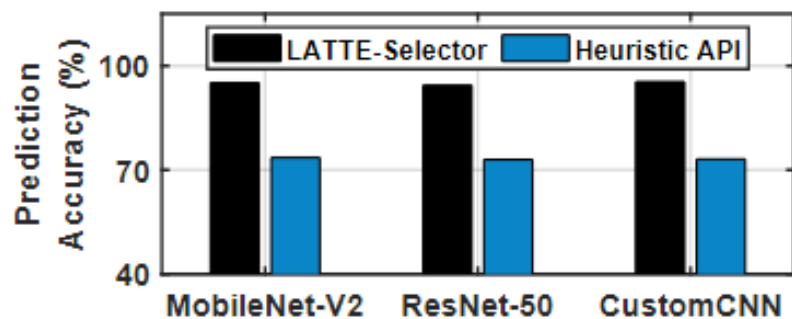
Evaluation



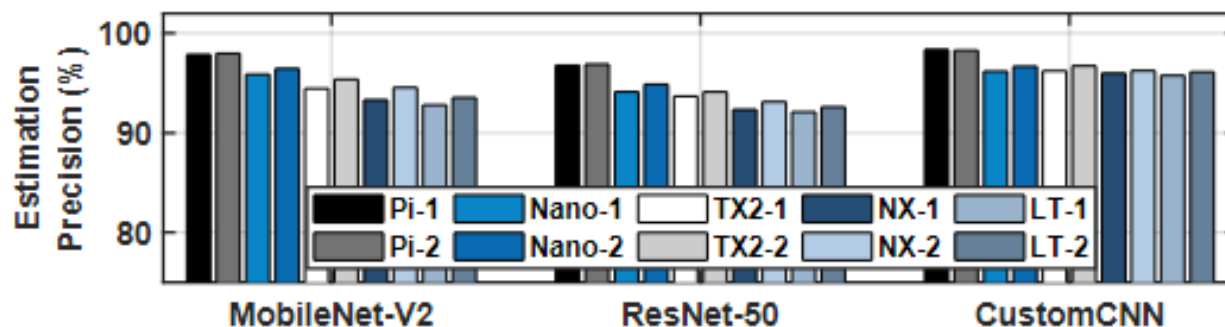
Compared with SOTA methods



Evaluating each components



Evaluating Selector's Performance



Evaluating Estimator's Performance

Conclusion

- We reveal the problem of **development-chain diversity** in federated learning systems and identify **diverse layer algorithms** as the key to explain the variability in training time. Based on this, accurate estimation of model training time can be achieved without complex operator or kernel-level modeling.
- We devise LATTE, with a novel **layer algorithm selector** and **training time estimator**, to accurately estimate the single-pass (forward/backward) propagation latency of a model given its architecture, expected hardware and runtime memory. We further showcase its usability in a client-side sub-model selection for HFL
- We conduct extensive experiments to evaluate LATTE in five typical HFL scenarios. The results show **significant improvements** in performance compared to seven classical or state-of-the-art methods.