# Modern LLMs are becoming larger and computational costly



PaLM (540B)

Nemotron-4 (340B)

GPT3 (175B)

LLaMa2 (70B)

100B

T5 (11B)

10B

GPT-J (6B)

1B

GPT2 (1.5B)

BERT (340M)

100M

ELMo(94M)

Data-center GPUs' capacities **cannot** scale their memory as fast as modern LLMs' sizes.

2018    2019    2020    2021    2022    2023    2024

# Current solutions to run LMs efficiently

| Category | Method | Training Time | Training Mem. | Inference Time | Inference Mem. |
|----------|--------|:---:|:---:|:---:|:---:|
| PEFT | Adapter | +++ | - | + | + |
| | LoRA | +++ | - | = | = |
| Pruning | MvP | +++ | + | - | - |
| | CoFi | +++ | ++ | --- | - |
| Combined | SPA | +++ | + | --- | - |
| | LRP | +++ | - | --- | - |

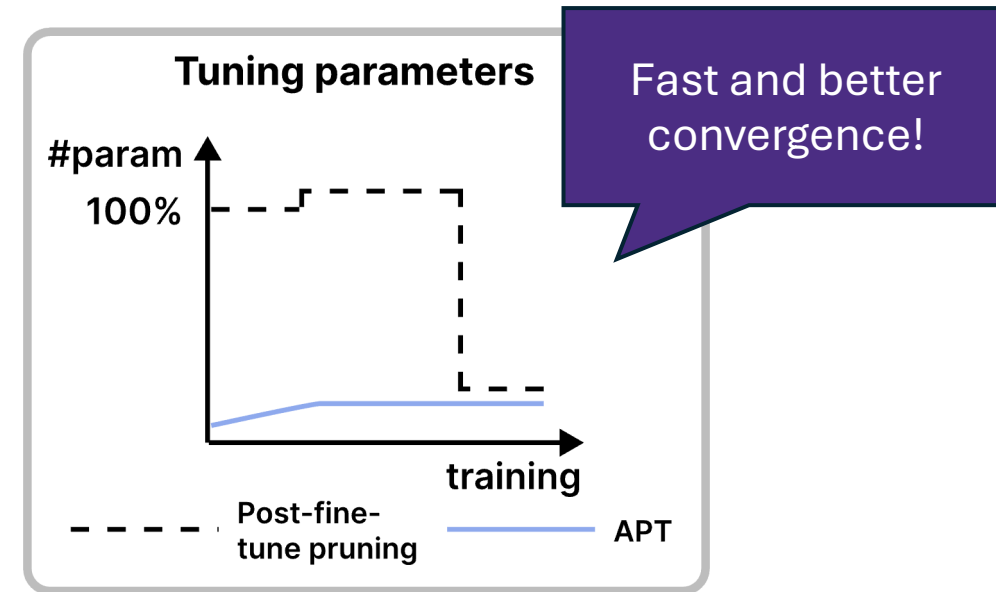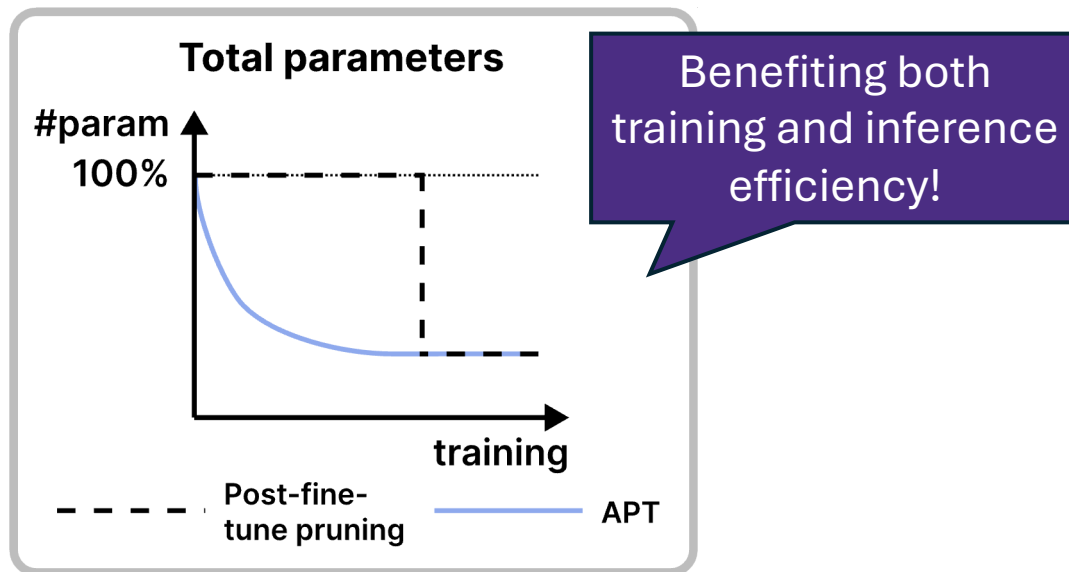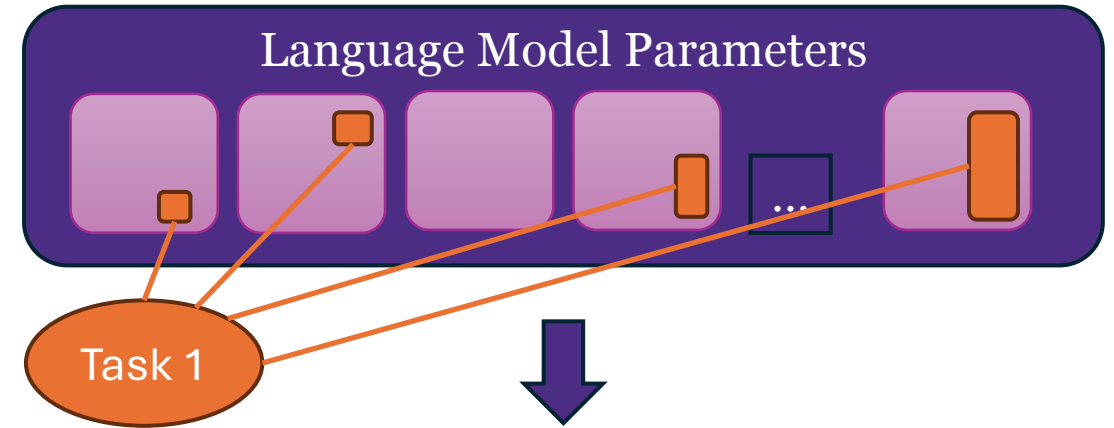Combined methods suffer from substantial end-task performance loss

Existing efficient methods often requires longer training time to converge the LM

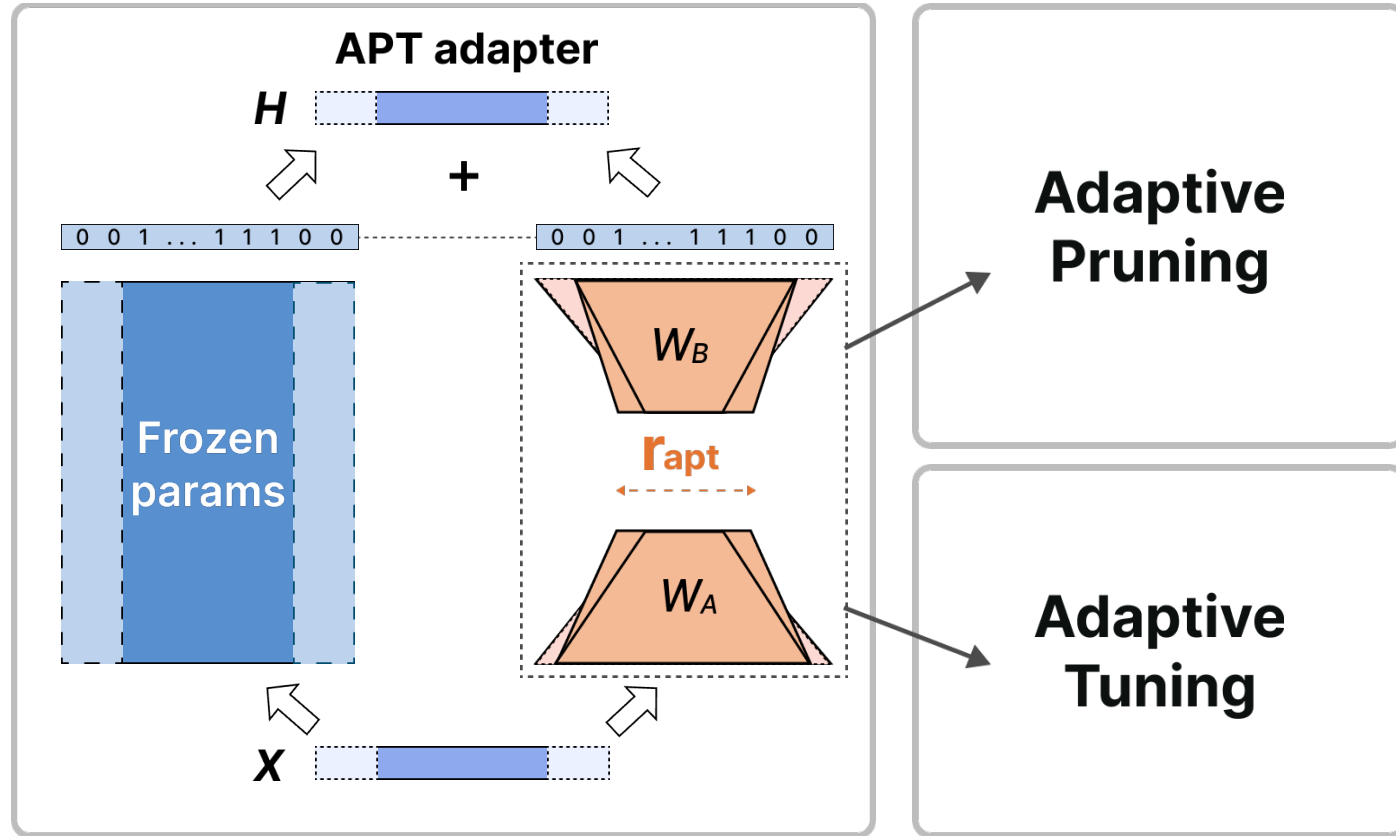Pruning methods tend to cost extra training memory due to knowledge distillation

# Improve training and inference efficiency

Question: can we combine the benefits of **PEFT and pruning** to improve both **training and inference efficiency** while maintaining **task performance**?
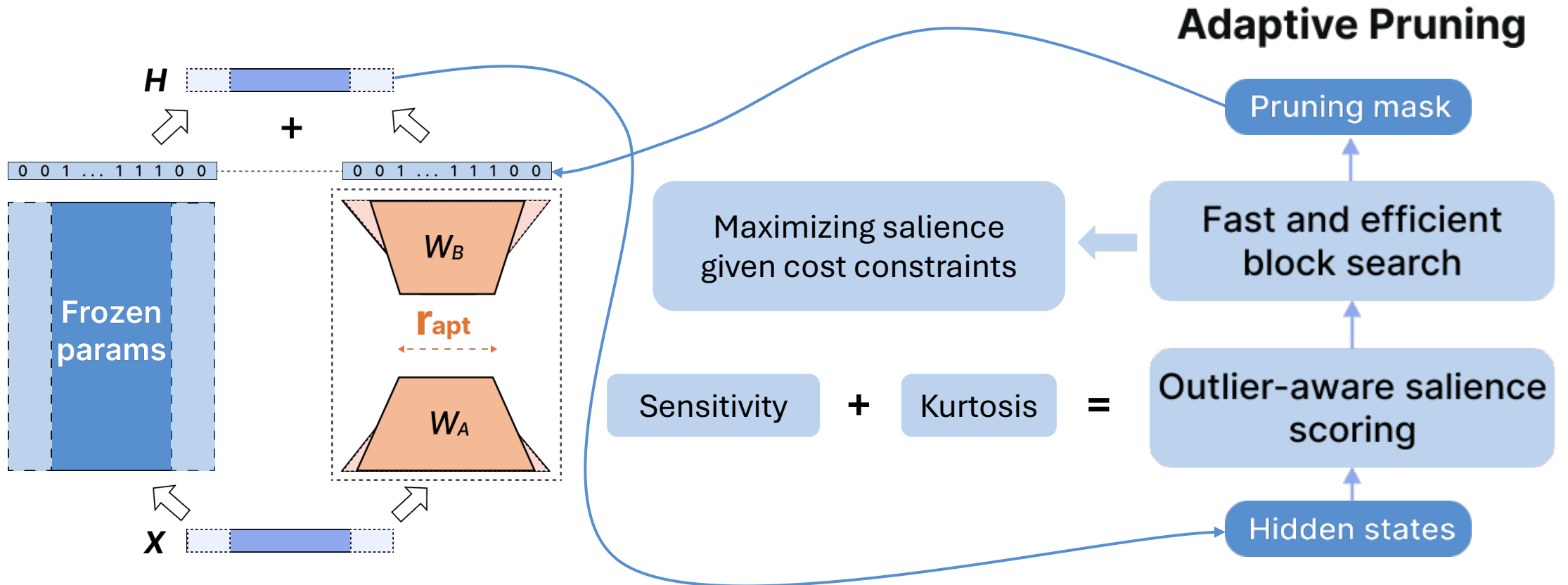
# Intuitions for improving LM efficiency

# Our solution – APT: pruning & tuning **adaptively**

# Low-cost adaptive pruning

# Efficient adaptive tuning

# Evaluation setup

- LM backbones and tasks:
  - Small-scale LMs:
    - BERT, RoBERTa: NLU tasks – GLUE, SQuAD
    - T5: NLU & NLG tasks – GLUE, CNN/DM
  - Large LMs:
    - LLaMa2 7B & 13B: standard few-shot tasks – ARC, HellaSwag, MMLU, TruthfulQA

- Metrics:
  - Task accuracy/F1/ROUGE score
  - Training efficiency: time to accuracy (seconds), training peak memory consumption (MB)
  - Inference efficiency: peak memory (MB), relative speedup

> TTA: training time to a percentage of the baseline (finetuning) accuracy

# Evaluation baselines

Direct baselines:

- Full-parameter finetuning
- LoRA

PEFT, pruning, and their combinations:

- LoRA+Prune: conducting post-training pruning (Mask-tuning; Kwon, et al., 2022) after LoRA-tuning
- Prune+Distill: structured pruning plus coarse-to-fine grained distillation (CoFi; Xia, et al., 2022)
- LoRA+Prune+Distill: using CoFi for pruning, but tuning LoRA only
- LLMPruner (Ma, et al., 2023): state-of-the-art structured pruning method on billion-level LLMs.

# APT speeds up small LMs pruning 8x faster compared to LoRA+Prune baseline

**Training convergence time comparison between APT and baselines**



1. Retraining-free pruning
2. Adaptive pruning: reduced training step time
3. Adaptive tuning: accelerate convergence

Relative Training Time to Accuracy

- FT: 100.00% (RoBERTa), 100.00% (T5)
- LoRA: 2137.00% (RoBERTa), 255.50% (T5)
- LoRA+Prune: 5128.30% (RoBERTa), 4523.50% (T5)
- Prune+Distill: 1495.30% (RoBERTa)
- LoRA+Prune+Distill: 6534.60% (RoBERTa)
- APT: 592.10% (RoBERTa), 484.70% (T5)

Methods

■ RoBERTa  ■ T5

# APT prunes LLMs with only 30% training memory consumption compared to LLMPruner

## Relative Training Time and Memory of LLaMa2-7B



Memory consumption reduced thanks to the pruning-before-finetuning scheduling
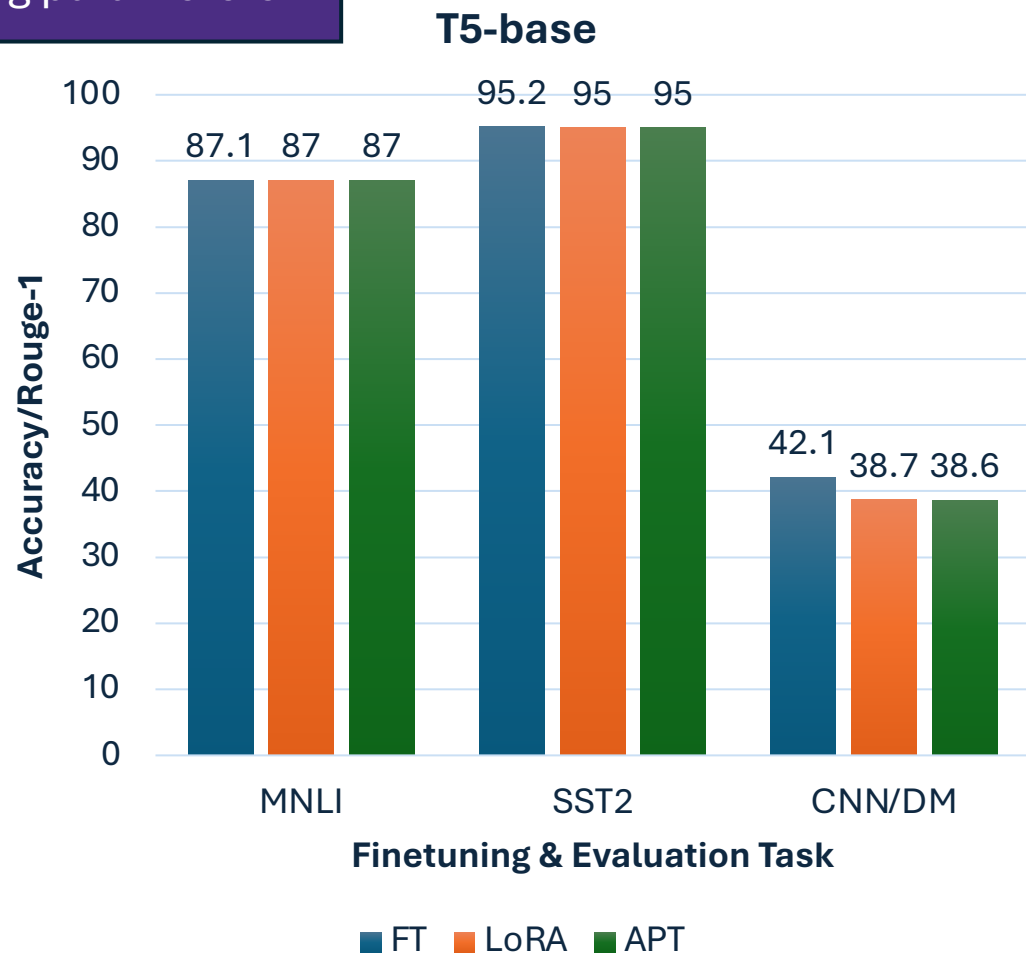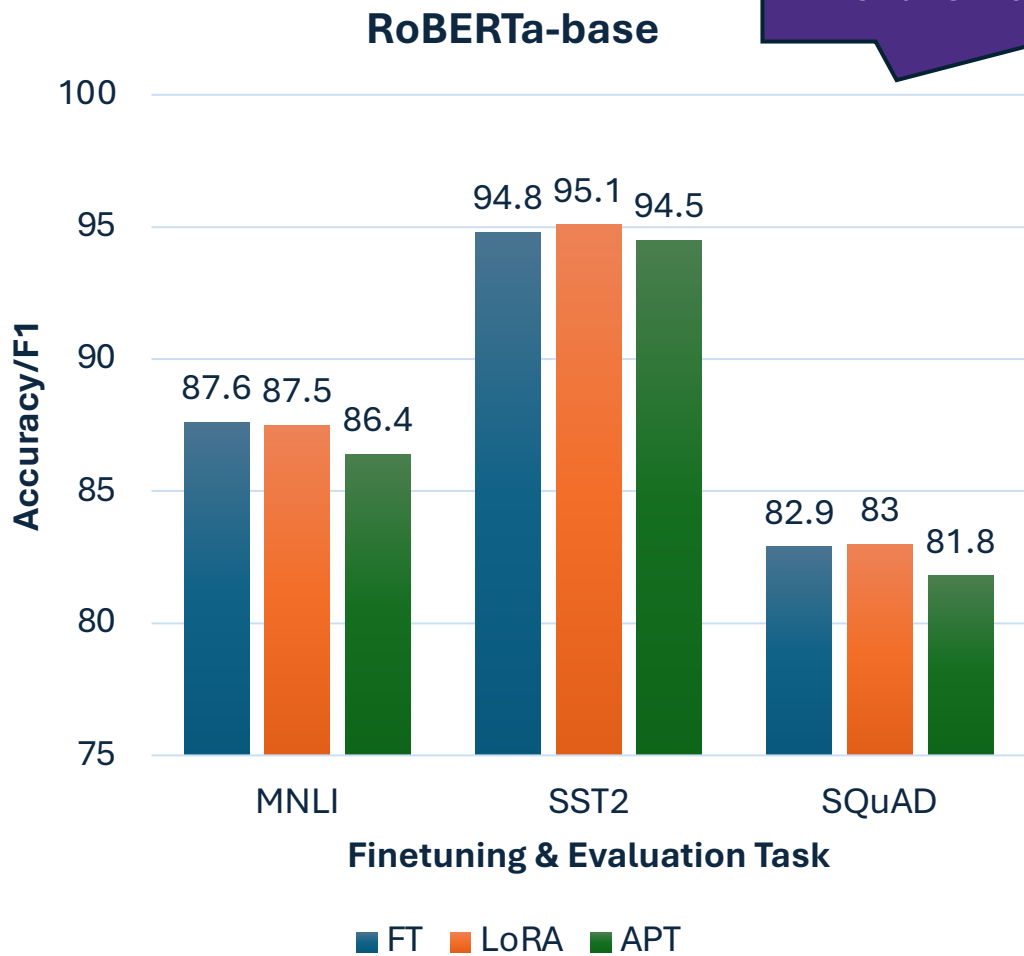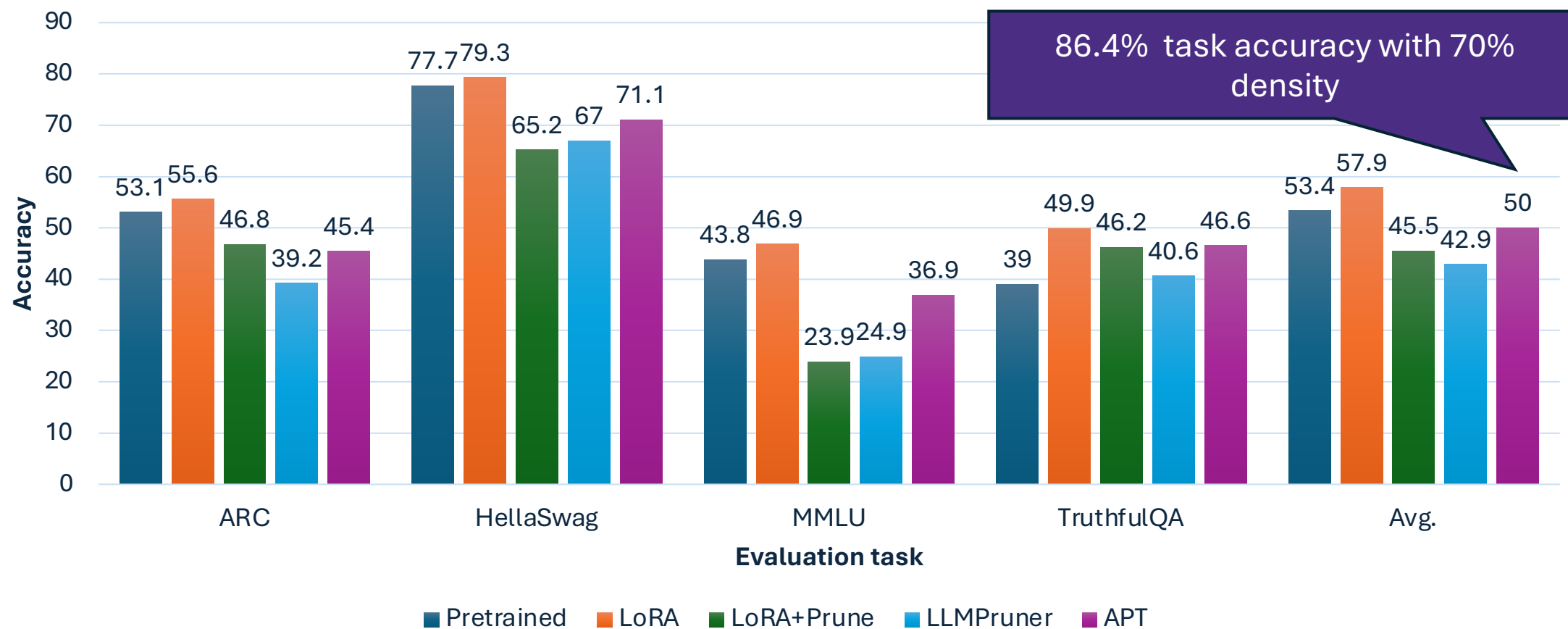
**Legend:** LoRA · LoRA+Prune · LLMPruner · APT

Y-axis: Relative Metric to FT (0.00% – 300.00%)
X-axis: Evaluation Metric (Train Mem., Train Time)

# APT recovers task accuracy for small and large LMs

Up to 98% performance with 40% remaining parameters



**RoBERTa-base**

**T5-base**

Legend: FT (blue), LoRA (orange), APT (green)

RoBERTa-base — Accuracy/F1 vs Finetuning & Evaluation Task:
- MNLI: FT 87.6, LoRA 87.5, APT 86.4
- SST2: FT 94.8, LoRA 95.1, APT 94.5
- SQuAD: FT 82.9, LoRA 83, APT 81.8

T5-base — Accuracy/Rouge-1 vs Finetuning & Evaluation Task:
- MNLI: FT 87.1, LoRA 87, APT 87
- SST2: FT 95.2, LoRA 95, APT 95
- CNN/DM: FT 42.1, LoRA 38.7, APT 38.6

# APT recovers task accuracy for small and large LMs



**LLaMa2 7B**

86.4%  task accuracy with 70% density

Accuracy

| Task | Pretrained | LoRA | LoRA+Prune | LLMPruner | APT |
|------|-----------|------|-----------|-----------|-----|
| ARC | 53.1 | 55.6 | 46.8 | 39.2 | 45.4 |
| HellaSwag | 77.7 | 79.3 | 65.2 | 67 | 71.1 |
| MMLU | 43.8 | 46.9 | 23.9 | 24.9 | 36.9 |
| TruthfulQA | 39 | 49.9 | 46.2 | 40.6 | 46.6 |
| Avg. | 53.4 | 57.9 | 45.5 | 42.9 | 50 |

Evaluation task

# APT achieves 2.5%-9.9% higher task performance than the LoRA+Prune baseline

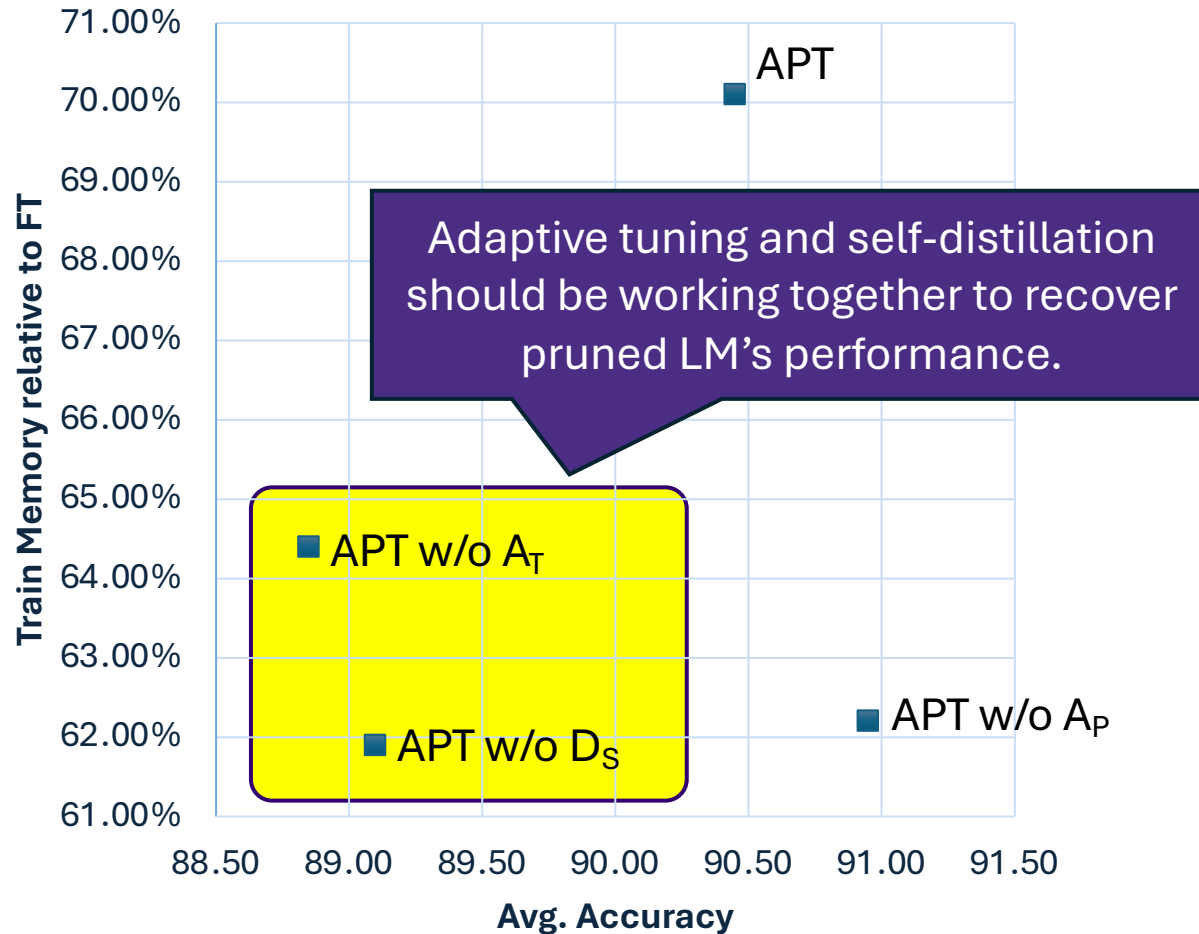**Performance comparison between APT and LoRA+Prune**

# APT reaches on-par performance with the Prune+Distill baseline but trains **2.5× faster** and costs **only 41.6% memory**.

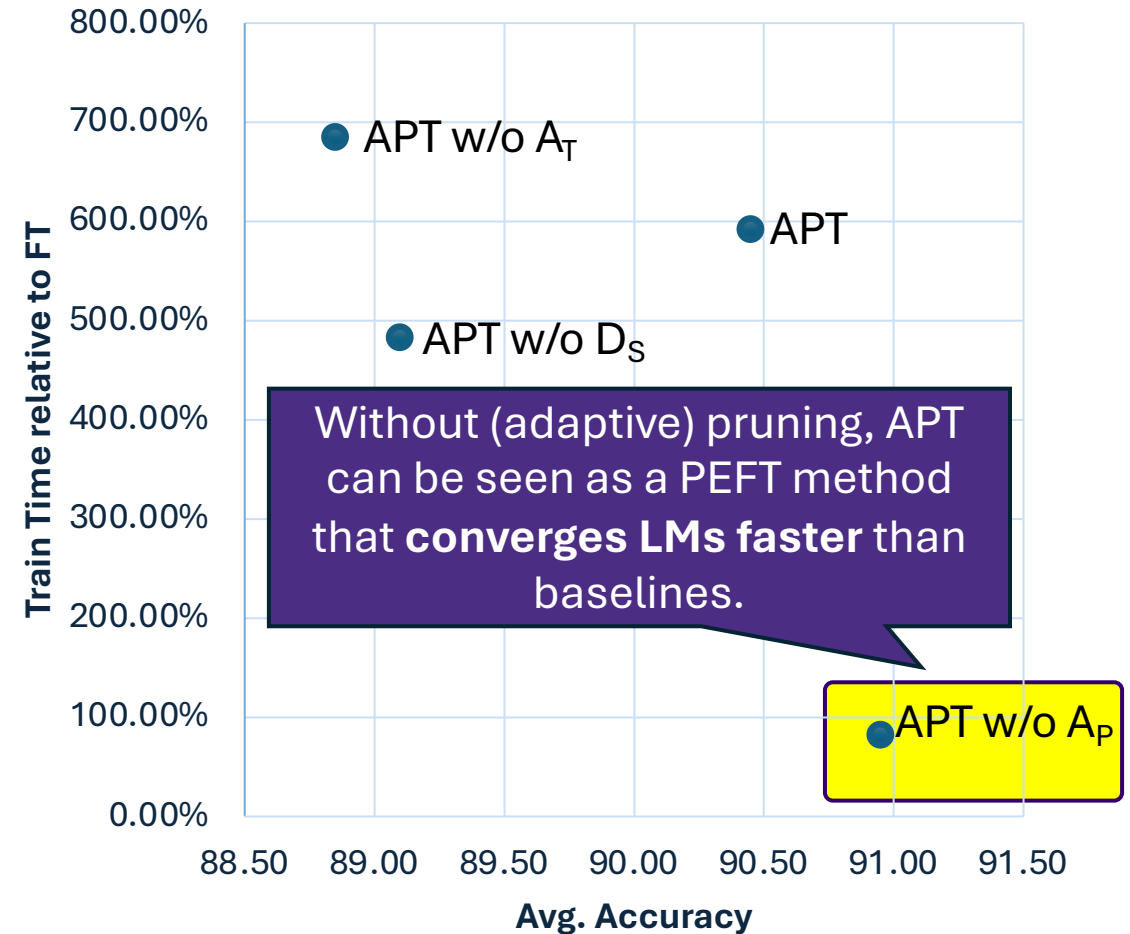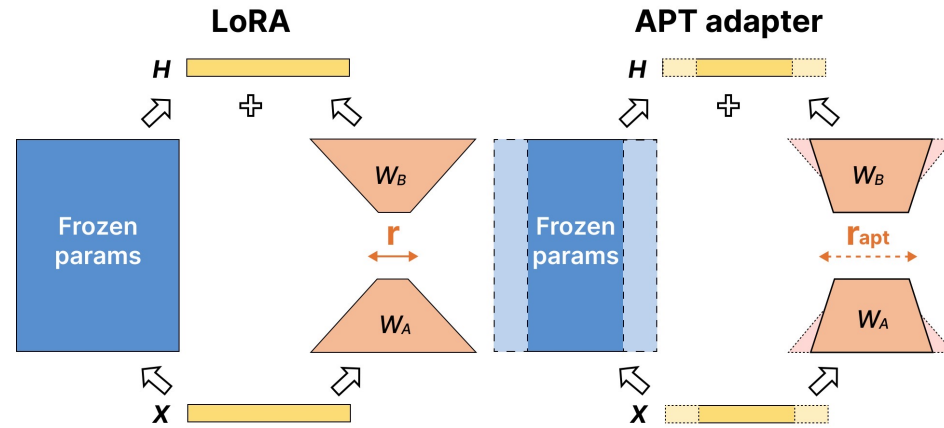**Performance and training efficiency of APT compared to baselines**



APT greatly reduces the training cost by **adaptive tuning + self-distillation**

Legend: ■ FT  ■ LoRA  ■ Prune+Distill  ■ APT

Y-axis: Relative Metric to FT
X-axis: Evaluation Metric

Avg. accuracy: 100%, 100.11%, 99.67%, 99.18%
Train Time: 100.00%, 2137.00%, 1495.30%, 592.10%
Train Mem.: 100.00%, 60.50%, 168.50%, 70.10%

APT, ICML2024

16

# Each component in APT is effective



**Accuracy - Train Memory Tradeoff**

Adaptive tuning and self-distillation should be working together to recover pruned LM's performance.

**Accuracy - Train Time Tradeoff**

Without (adaptive) pruning, APT can be seen as a PEFT method that **converges LMs faster** than baselines.

# APT key takeaways and impact



- We propose APT, a new adaptive paradigm to prune and tune LMs effectively, targeting both training and inference efficiency via APT adapters.

- APT dynamically adjusts (adds/reduces) APT adapter input/output dimensions and the rank ($r_{apt}$), thus accelerating LM training convergence and also reducing inference costs.

- APT preserves LM task performance while speeding up small-scaled LMs' fine-tuning by up to 8× and reducing large LMs' training memory footprint by up to 70%.

# Future work

- Even though APT proposes an efficient way to prune and tune LMs, it is definitely not always the optimal method for all LMs

- We hope that future work will focus on:

  - Adopting APT to a wider variety of PEFT backbones, e.g., prefix-tuning, prompt-tuning, parallel-adapter, VeRA, DoRA, etc.

  - Aiming at accurate, efficient, retraining-free pruning and distillation methods of large, billion-level LMs

  - Adapting APT with other efficient methods together for further inference efficiency gains, such as quantization, MoEfication, etc.