

# Joint Fine-tuning and Conversion of Pretrained Speech and Language Models towards Linear Complexity

Mutian He<sup>1,2</sup>, Philip N. Garner<sup>1</sup>  
<sup>1</sup>Idiap Research Institute <sup>2</sup>EPFL On ICLR 2025

## Background

### Transformers are so expensive!

- $O(L^2)$  time complexity
- $O(L)$  KV cache
- ...especially when handling speech
- few words  $\approx$  1sec = 16K samples = 50 frames

### An ever-growing arsenal of transformer alternatives

- Low-rank attention: **Linformer**
- Restricted attention: Longformer, Big Bird, MoBA, Native Sparse Attention...
- RNNs (linear attention): RetNet, **Mamba** (2), DeltaNet ...
- ...still increasing!

## Motivation

### How to make use of these new archs?

- Pretrained parameters often unavailable, esp. on speech
- New models emerge rapidly
- ...redo the whole pretraining for each new arch?
- Computational costs & access to pretraining data

### Convert/fine-tune pretrained transformers into the target arch on the target downstream task

- Use only the downstream target task data, avoid re-pretraining
- Without performance degradation

## Methods: Cross-Arch Layerwise Distillation

### Knowledge transfer from original transformer

- **Unguided** (a): Parameter transfer
  - Replace attention layers with, e.g., Mamba layers, then fine-tuning
  - Other parameters (e.g., MLPs) are reused
- **Guided**: Behavior transfer
  - Reproduce the original behavior (hidden states) by layerwise distillation

$$\mathcal{L}_{CE}(\mathbf{y}^{(s)}, \mathbf{y}) = - \sum_i \mathbf{y}_i \log(\mathbf{y}_i^{(s)})$$

Cross-entropy loss for the target classification task

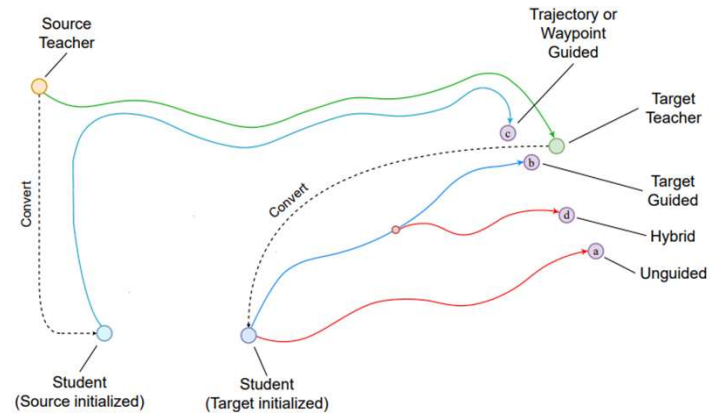
$$\mathcal{L}_{KD}(\mathbf{y}^{(s)}, \mathbf{y}^{(t)}) = \sum_i \left( \frac{\mathbf{y}_i^{(t)}}{\beta} \right) \log \left( \frac{\mathbf{y}_i^{(t)}}{\mathbf{y}_i^{(s)}/\beta} \right)$$

KL-div between student (ours) and teacher (transformer) outputs

$$\mathcal{L}_{LD}(\mathbf{H}^{(s)}, \mathbf{H}^{(t)}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{H}_i^{(s)} - \mathbf{H}_i^{(t)})^2$$

L2 loss between hidden states in each layer of the student and the teacher

$$\mathcal{L} = \alpha_{CE} \mathcal{L}_{CE} + \alpha_{KD} \mathcal{L}_{KD} + \alpha_{LD} \mathcal{L}_{LD}$$



### What should be the teacher?

- **Target-guided** (b)
  - Directly distill from the fine-tuned transformer (*target teacher*)
- **Trajectory/Waypoint guided** (c)
  - Original pretrained transformer (*source teacher*) carries important knowledge that leads to downstream capabilities
    - Essential to fine-tuning, but will be lost in the end
    - Preserving the knowledge helps, as found by e.g. L2-SP
- Can we reproduce the trajectory of transformer fine-tuning?
  - Simultaneously fine-tune transformer & target model
  - Distill from hidden states at each fine-tuning step
  - Approximation: distill from several checkpoints (*waypoints*) during transformer fine-tuning

### When should we distill?

- Distillation loss terms are like splints, stabilizes optimization in the early stage but restrains it later
- **Hybrid** (d): stop distillation later and set the target loss free

## Configuration

- RoBERTa  $\rightarrow$  Linformer for NLP: QNLI, QQP, SST2, IMDB
- Wav2Vec2  $\rightarrow$  Bidirectional Mamba2 for speech tasks: TEDLIUM (ASR), SLURP (IC), VoxCeleb1 (Speaker ID)
- Pythia-1B  $\rightarrow$  Mamba, on zero-shot LM tasks \*

\* An ancillary experiment, since there isn't a separate *target task* here; hence trajectory/waypoint guided conversion doesn't apply

## Results

- Guided conversion matches standard transformer results
- Trajectory/Waypoint guided mode helps in the NLP case
  - Waypoint approximation works well
  - Hybrid mode performs worse

	QNLI	QQP	SST2	IMDB	Average
Pretrained Linformer	91.2%	90.8%	93.1%	94.1%	92.3%
Std. RoBERTa	92.4%	91.8%	95.3%	95.7%	93.8% $+1.5$
Unguided	69.4%	84.3%	83.6%	82.6%	80.0% $-12.3$

	ASR WER $\downarrow$	IC Acc. $\uparrow$	SID Acc. $\uparrow$
Std. Wav2Vec2	6.24	91.70	96.09
Unguided	11.29	79.68	84.24

	ASR WER $\downarrow$	IC Acc. $\uparrow$	SID Acc. $\uparrow$
CALD			
- Target Guided	89.0%	91.8%	93.3%
- Trajectory Guided	89.0%	91.9%	94.0%
- Waypoint Guided	89.9%	91.9%	93.7%
- Hybrid	86.8%	90.8%	91.4%

### Works also on speech; hybrid mode works better

Model	Lambada	PIQA	Winog.	WSC	ArcE	ArcC	SciQ	LogiQA	Avg.
Pythia-1B	0.562	0.707	0.535	0.670	0.570	0.244	0.839	0.221	0.544

	ASR WER $\downarrow$	IC Acc. $\uparrow$	SID Acc. $\uparrow$
CALD			
- Target Guided	6.56	90.43	96.56
- Waypoint Guided	6.92	90.32	96.16
- Hybrid	6.41	91.23	96.41

### Similar results on LM, converted using 0.5%/1%/2% Pile

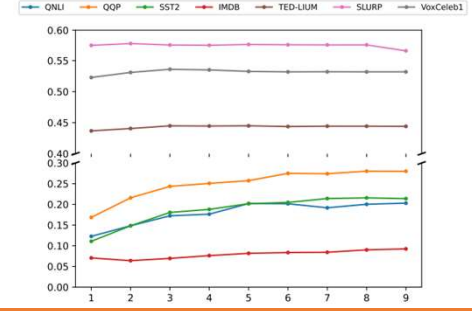
Model	Lambada	PIQA	Winog.	WSC	ArcE	ArcC	SciQ	LogiQA	Avg.
Pythia-1B	0.562	0.707	0.535	0.670	0.570	0.244	0.839	0.221	0.544

	ASR WER $\downarrow$	IC Acc. $\uparrow$	SID Acc. $\uparrow$
Unguided	0.394	0.671	0.493
Tgt. Gd.	0.410	0.686	0.504
Hybrid	0.432	0.683	0.507

	ASR WER $\downarrow$	IC Acc. $\uparrow$	SID Acc. $\uparrow$
Unguided	0.453	0.673	0.514
Tgt. Gd.	0.449	0.689	0.518
Hybrid	0.479	0.693	0.520

### Why difference between hybrid & trajectory modes?

- Hidden states shift gradually, epoch-by-epoch in NLP
- Significant shifts have occurred in a few steps in speech



## Takeaway

- Pretrained transformers can be converted to linear-complexity (Linformer, Mamba) downstream models
  - Guided by layerwise distillation only on the target task data
- Alternative distillation modes help per task
  - Trajectory guided models, using knowledge from pre-fine-tuned pretrained transformer that will be lost later in FT
  - Time-hybrid of guided and unguided fine-tuning