No Need To Talk 🕍: Asynchronous Mixture of Language Models



Anastasiia Filippova, Angelos Katharopoulos, David Grangier, Ronan Collobert ICLR 2025 · Apple

Motivation

Training & inference of LLMs requires high-bandwidth > & low-latency interconnect

But what if we don't have access to state-of-the-art interconnect hardware?

Prior works to reduce communication cost:

- require some form of synchronization [1, 2]
- relay on full corpus clustering [3]
- often under-perform or does not scale [4, 5]

How can we achieve efficient training and inference without

- relying on fast interconnect
- compromising model performance?

Background

Asynchronous training & sparse inference with Mixture of Experts [6]

$$\mathcal{L}(\mathbf{x}; \theta) = -\sum_{s=1}^{S-1} \log \sum_{e=1}^{E} p(x_{s+1} \mid \mathbf{x}_{1:s}, e ; \theta^e) p(e \mid \mathbf{x}_{1:s}; \theta^r),$$

Can be sparse by using hard assignments:

$$\mathcal{L}(\mathbf{x};\theta) \le -\sum_{s=1}^{s-1} \log p(x_{s+1} \mid \mathbf{x}_{1:s}, e^{\star}; \theta^{e^{\star}}) p(e^{\star} \mid \mathbf{x}_{1:s}; \theta^{r}),$$

Compared to Switch Transformer MoE [7]

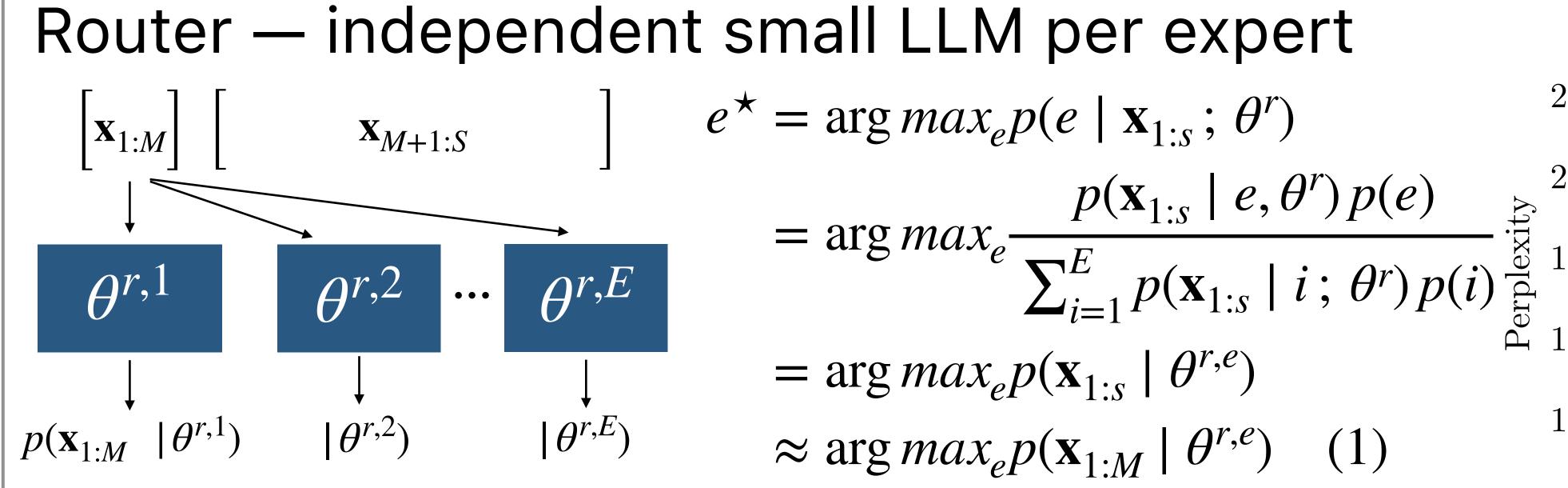
Routing is done per token ⇒

- all experts must be retained in RAM during both training and inference
- high communication overhead is required to synchronise gradients

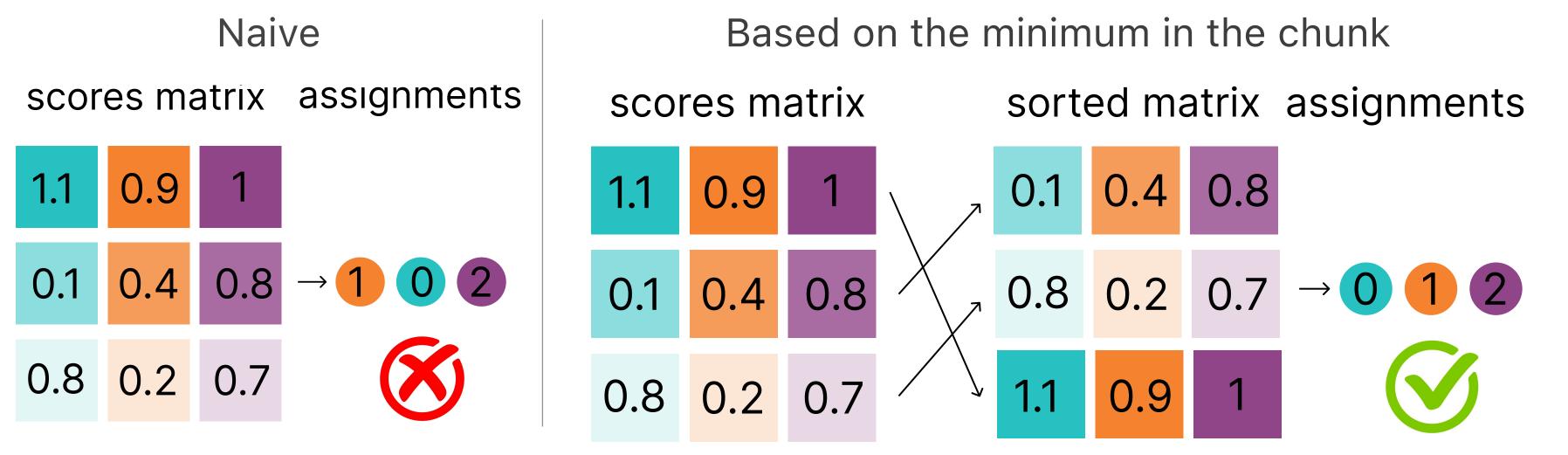
CITATIONSSS

SmallTalk LM

Routing With Independent Language Models



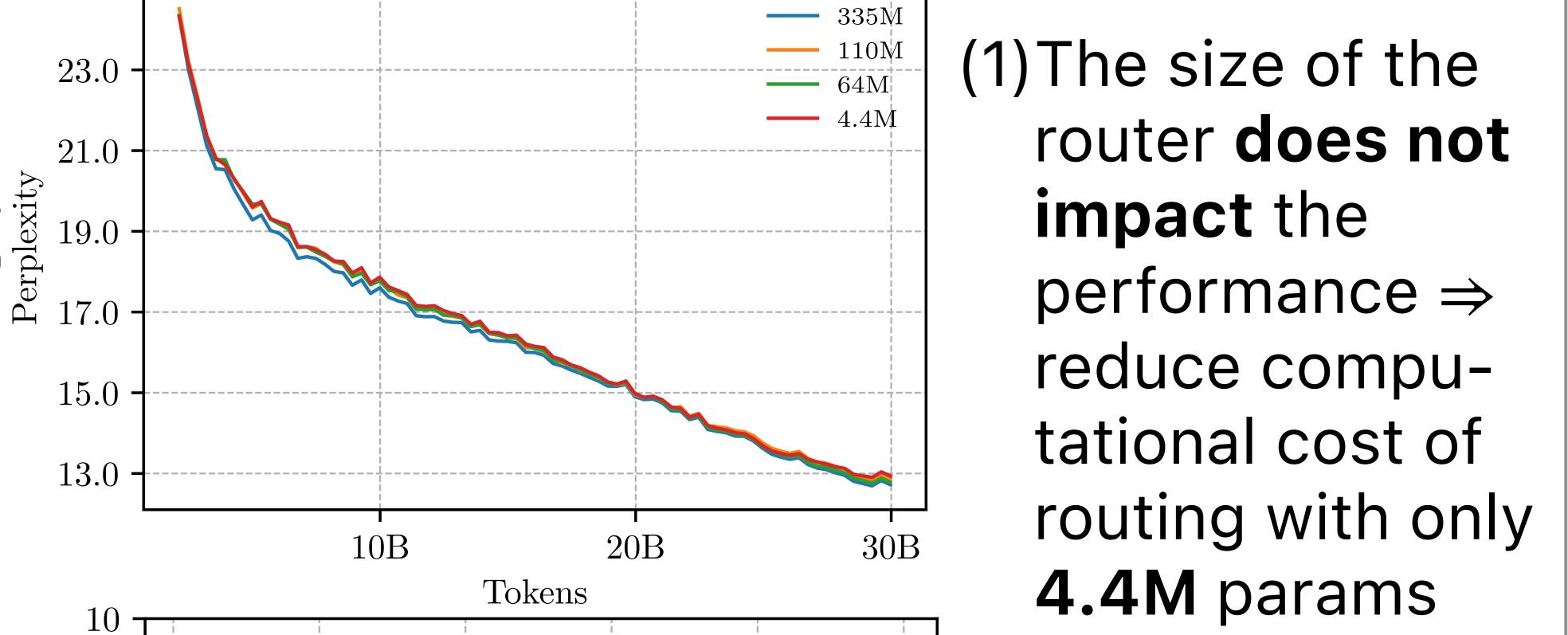
Assignments balancing



Training Procedure

- 1) Train the routers with EM without using the experts at all, alternating between log-likelihood optimization & assignments
- 2) Train the experts as E independent lms:
- 1: Train the routers
- 2: $\mathbf{X} \leftarrow N$ new sequences from the dataset
- 3: $\mathbf{X}_{1:E} = \text{random_assignments}(\mathbf{X})$
- 4: for i = 1...T do
- 5: for e = 1...E do
- $heta^{r,e} pprox rg \min_{ heta^{r,e}} \mathcal{L}(\mathbf{X}_e; heta^{r,e})$
- 7: end for
- $\mathbf{X} \leftarrow N$ new sequences from the dataset
- 9: $\mathbf{X}_{1:E} = \text{balanced_assignments}(\mathbf{X}, \theta^r)$
- 10: end for
- 11: Train the experts
- 12: $\mathbf{X} \leftarrow M$ new sequences from the dataset
- 13: $\mathbf{X}_{1:E} = \text{balanced_assignments}(\mathbf{X}, \theta^r)$
- 14: for e = 1...E do
- 15: $\theta^e \approx \arg\min_{\theta^e} \mathcal{L}(\mathbf{X}_e; \theta^e)$
- 16: end for

Computational cost & Routing Analysis



(2) With 32 tokens prefix length SmallTalk LM still outperforms dense baseline!

routers!

(3)Simpler methods such as TF-IDF encoding + K-Means under performs compared to routing with 4.4M language model.

15 14 13 13 11 10 8 16 32 64 128 256 Prefix length

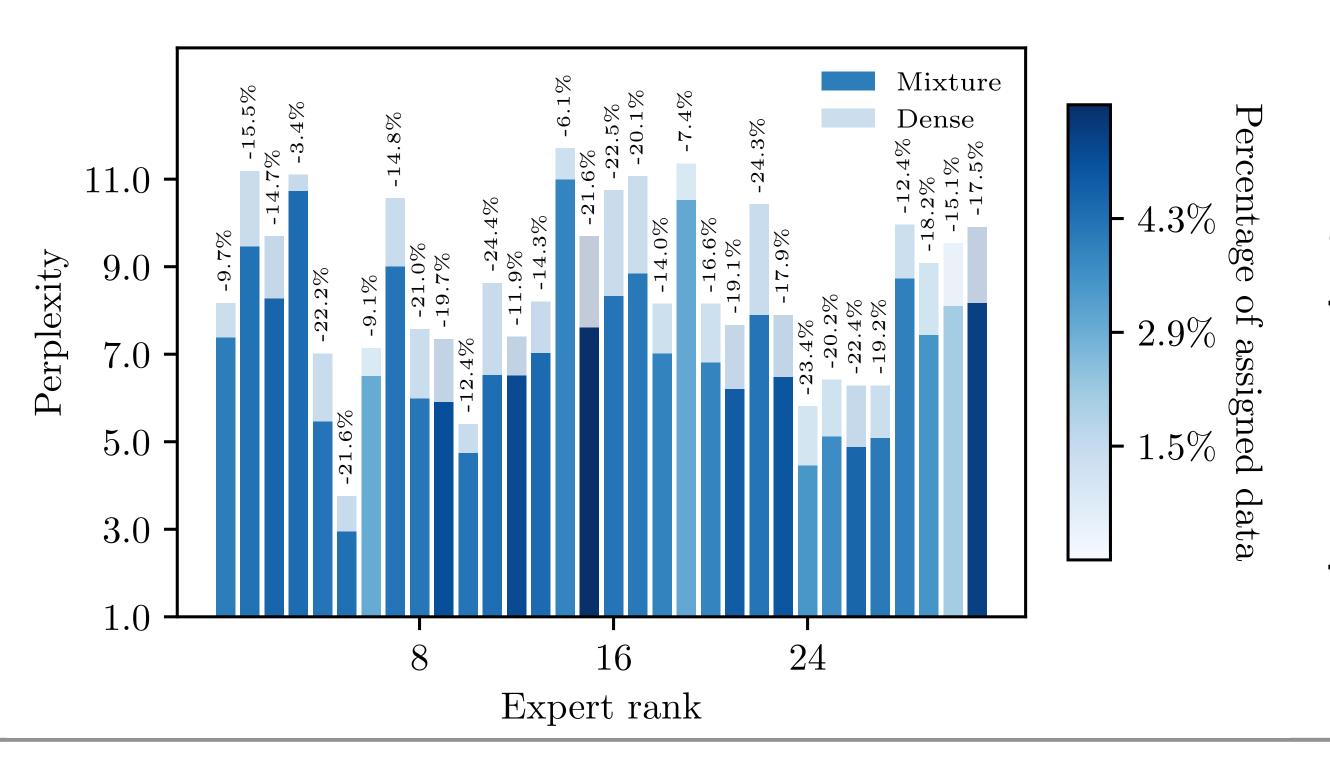
Results

4 experts, 266B

— 16 experts,

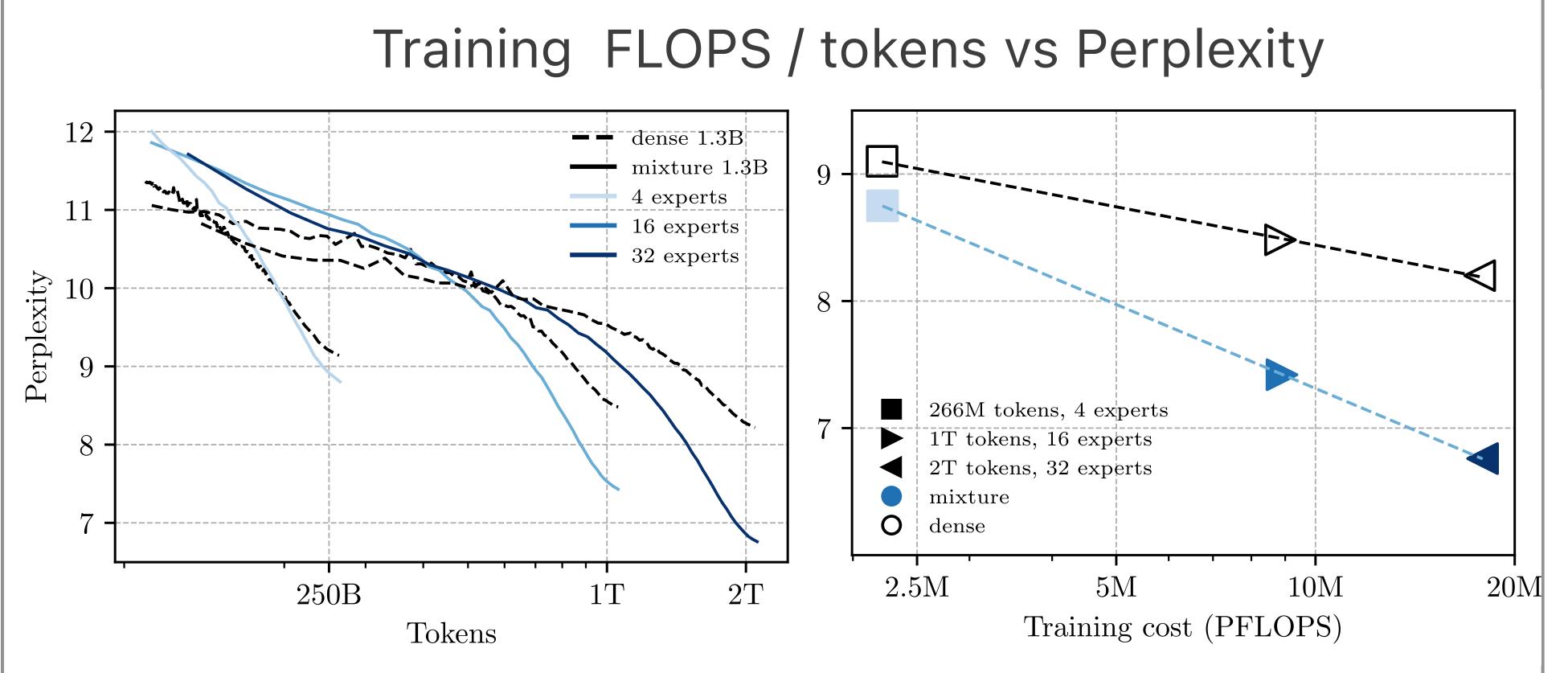
32 experts, 2T

+-- dense, 2T

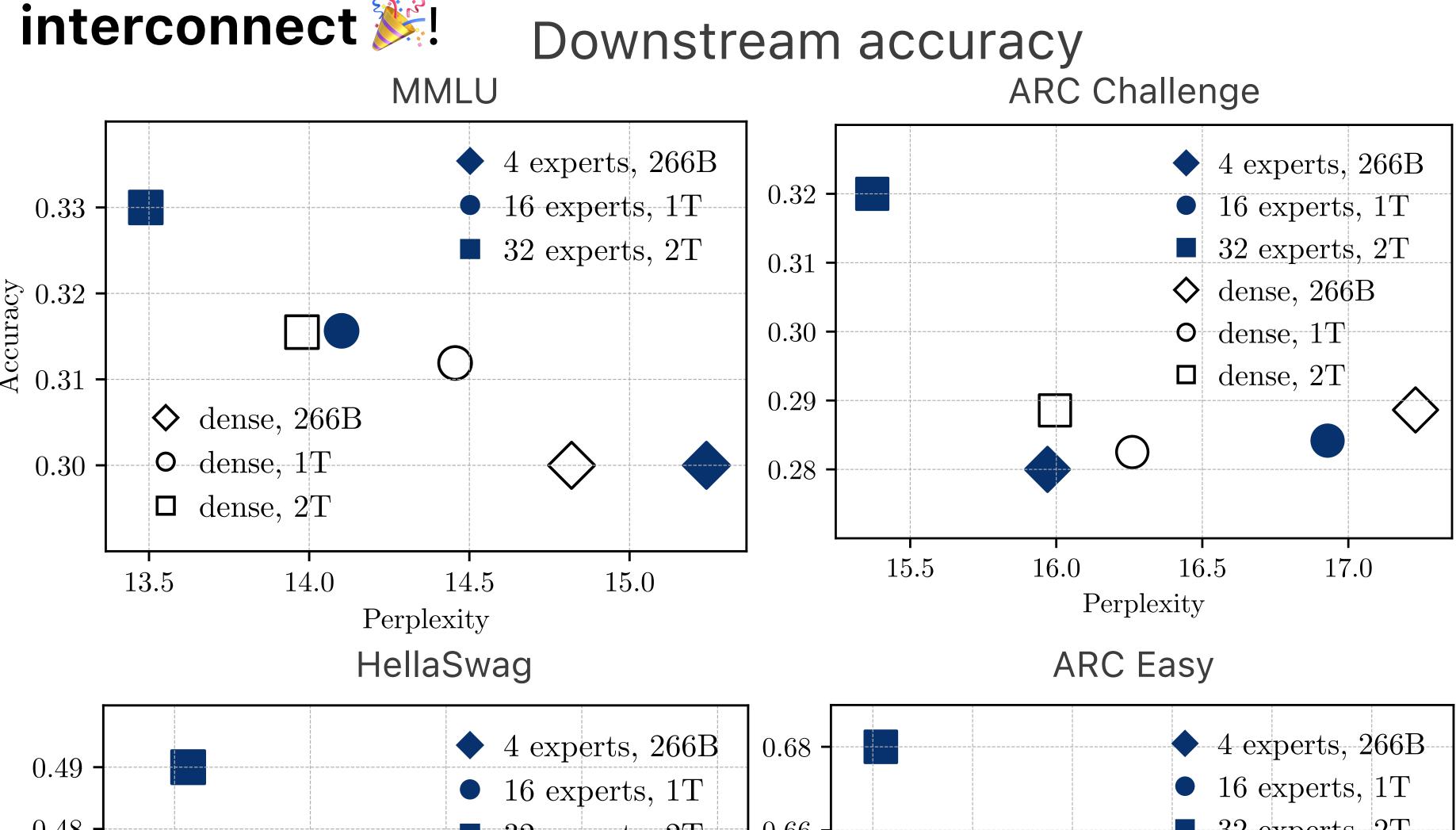


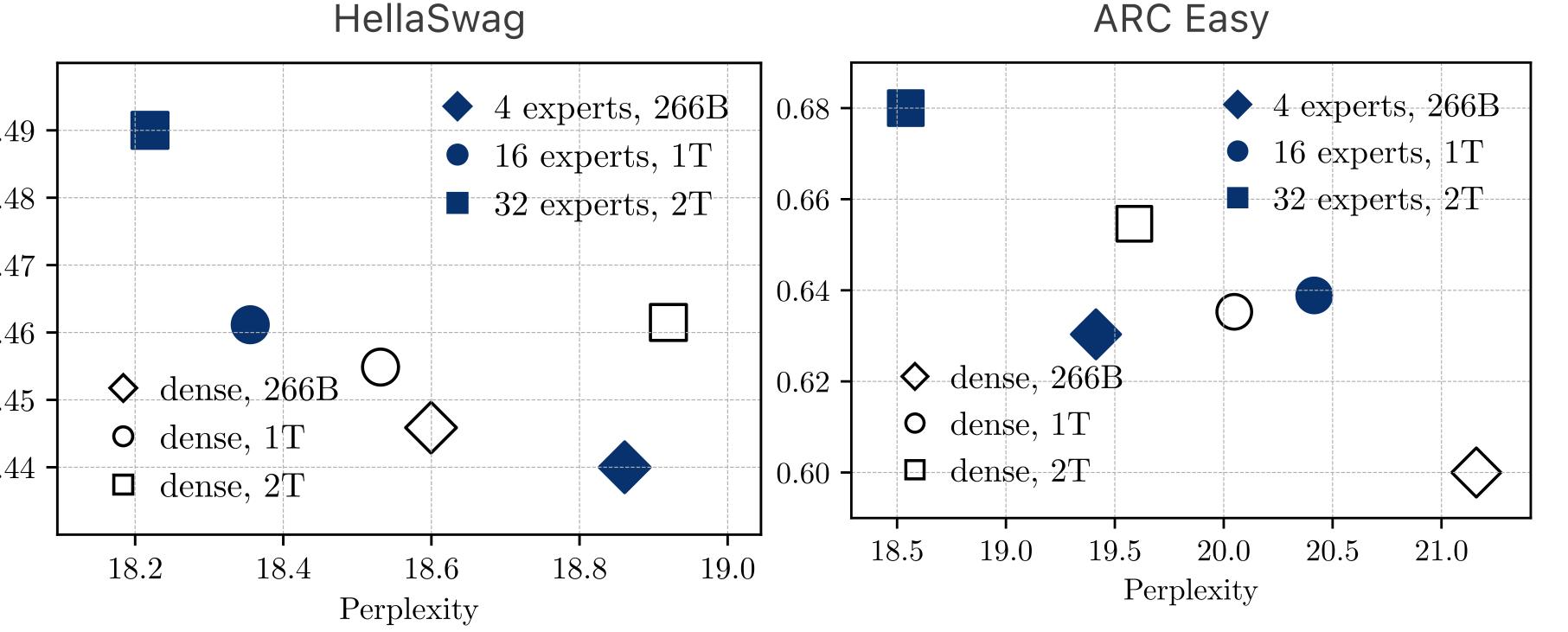
Experts do specialize and all contribute to the overall improvement instead of only a few dominant experts.

Results



Significant performance boost with near-identical FLOPS, identical data volume and **no reliance on fast**





Conclusion: No Need To Talk

Compared to dense baseline SmallTalk LM:

- reduces bandwidth requirement to zero
- archives significantly lower perplexity with near-identical training and inference FLOPs and identical training data volume
- achieves better accuracy on downstream tasks