# Switch Transformers

## Scaling to Trillion Parameter Models With Simple and Efficient Sparsity
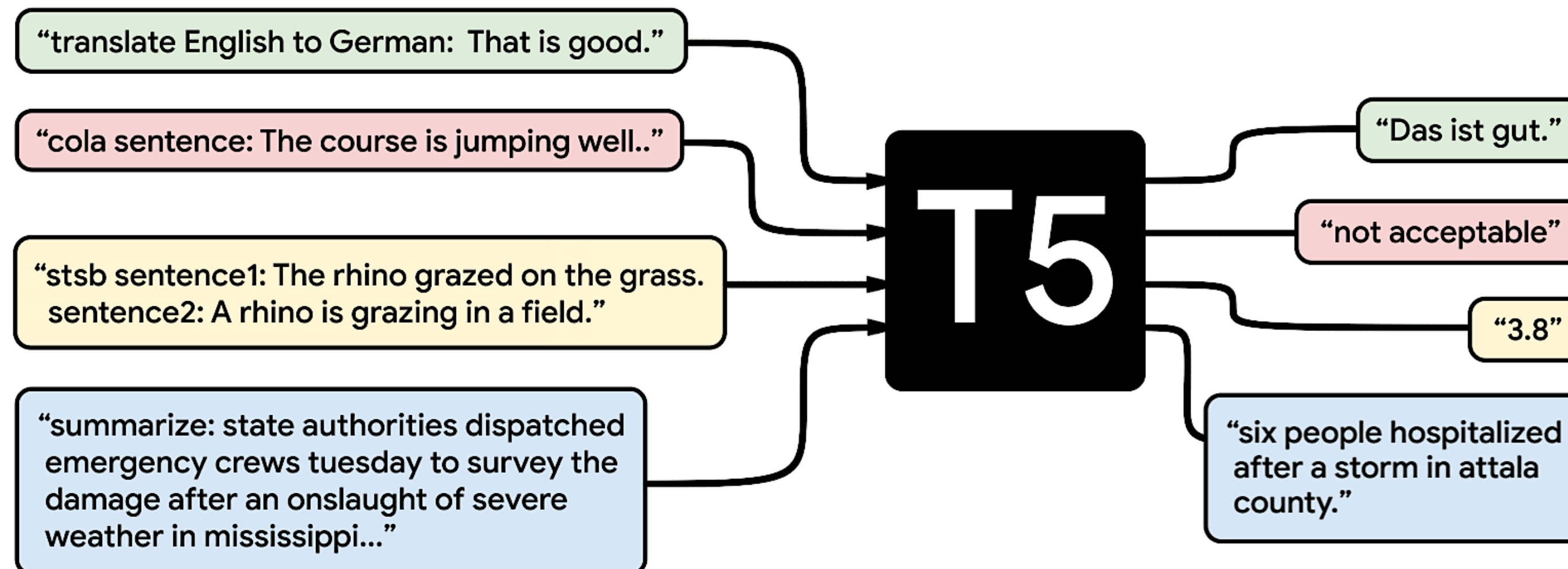
**Sam Foreman**
**June, 2021**

Argonne
NATIONAL LABORATORY
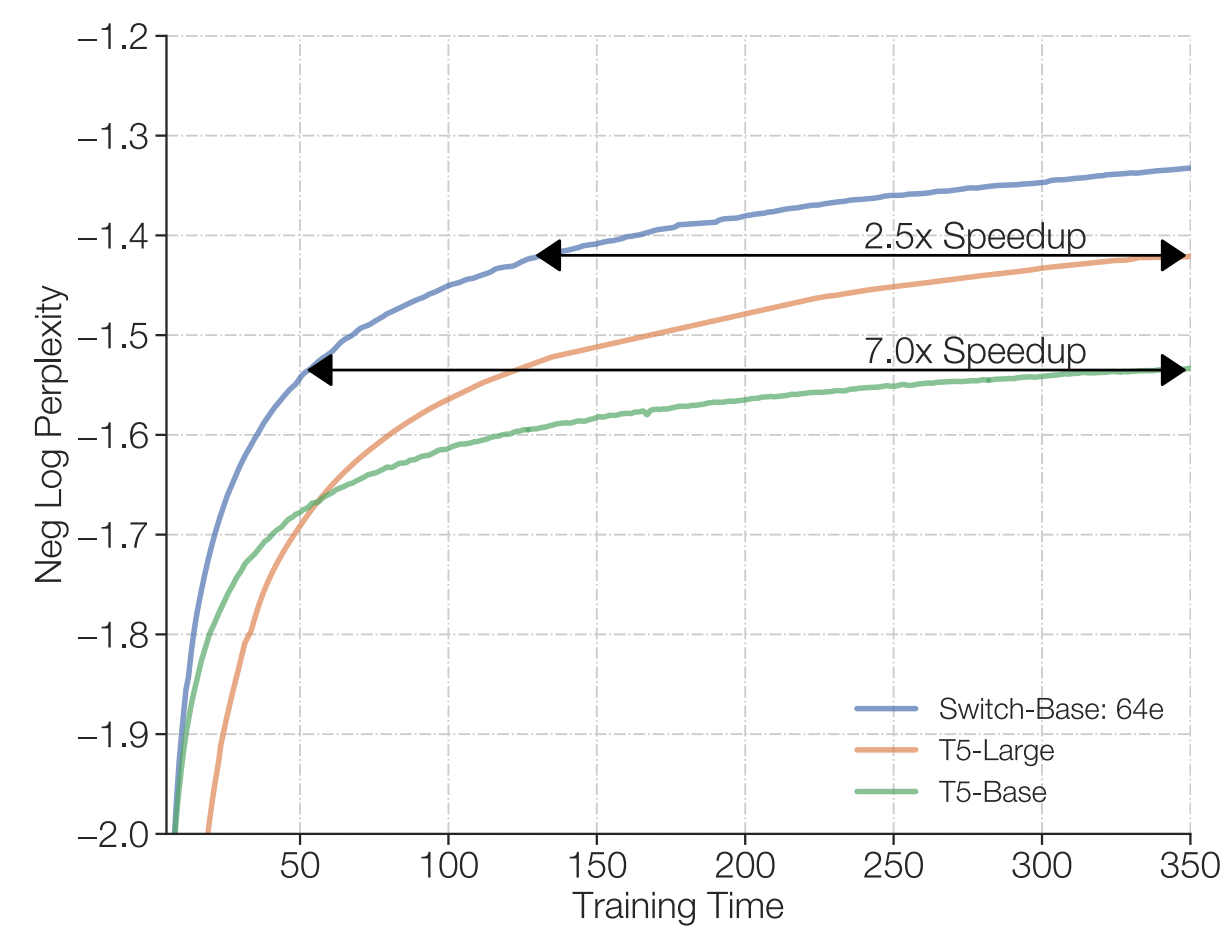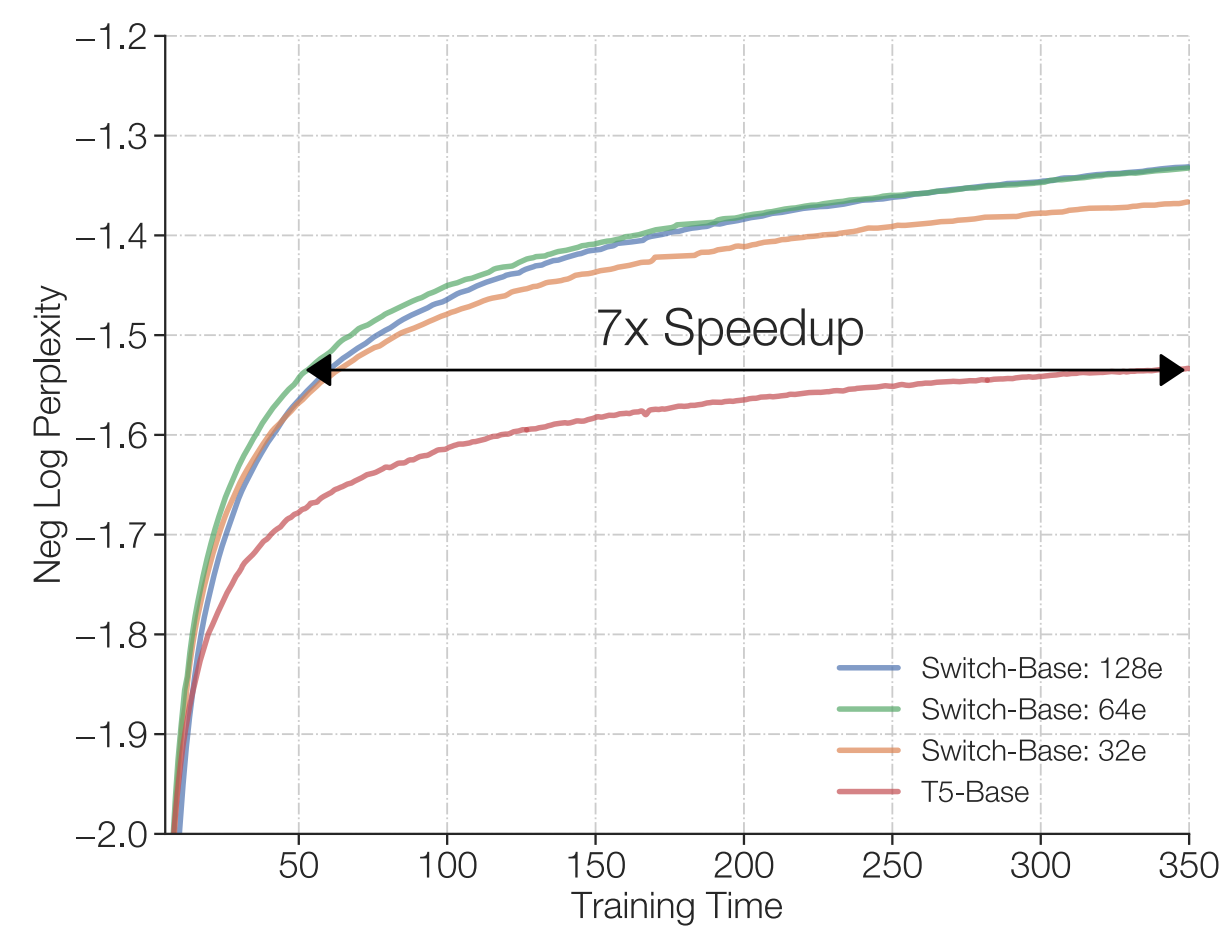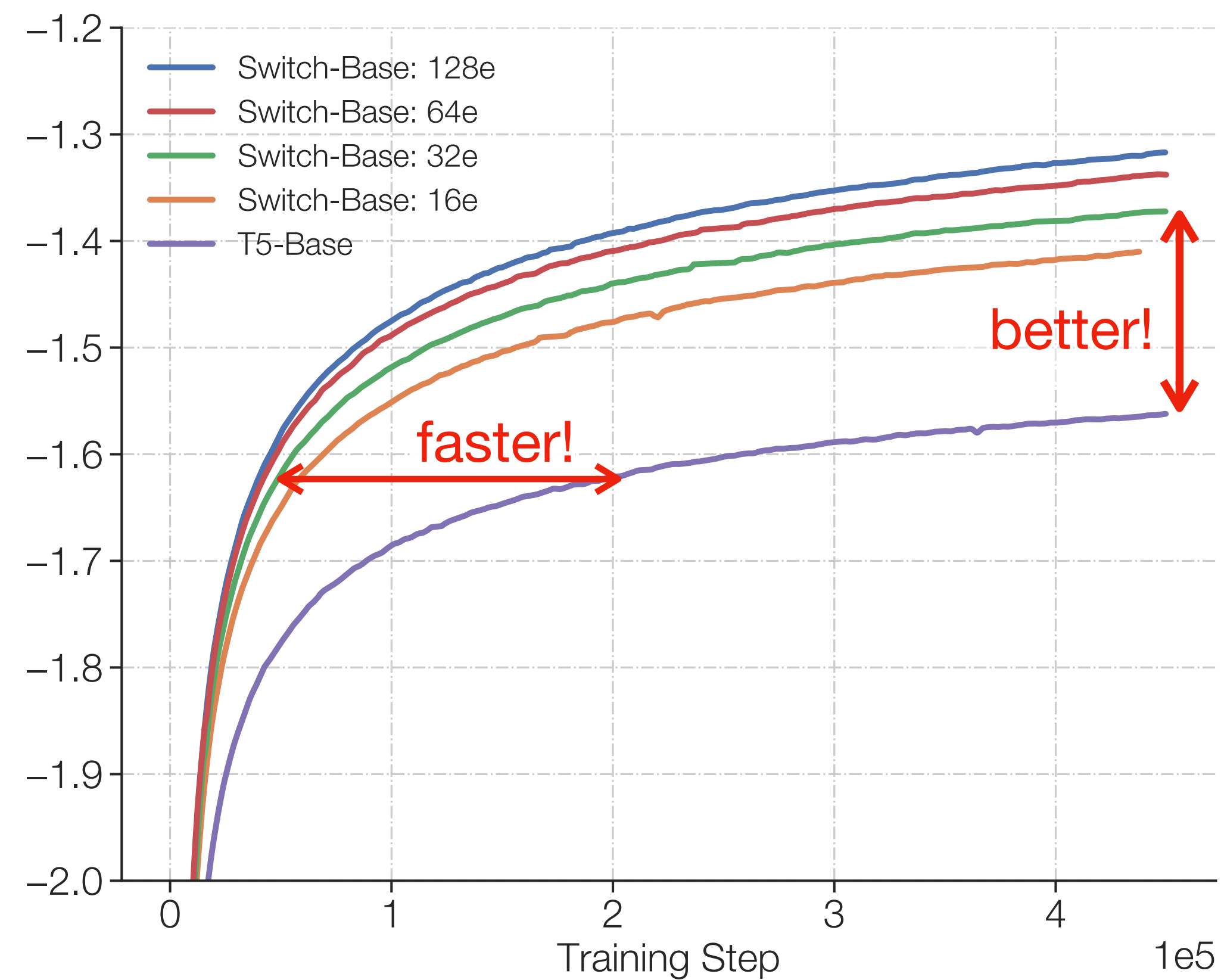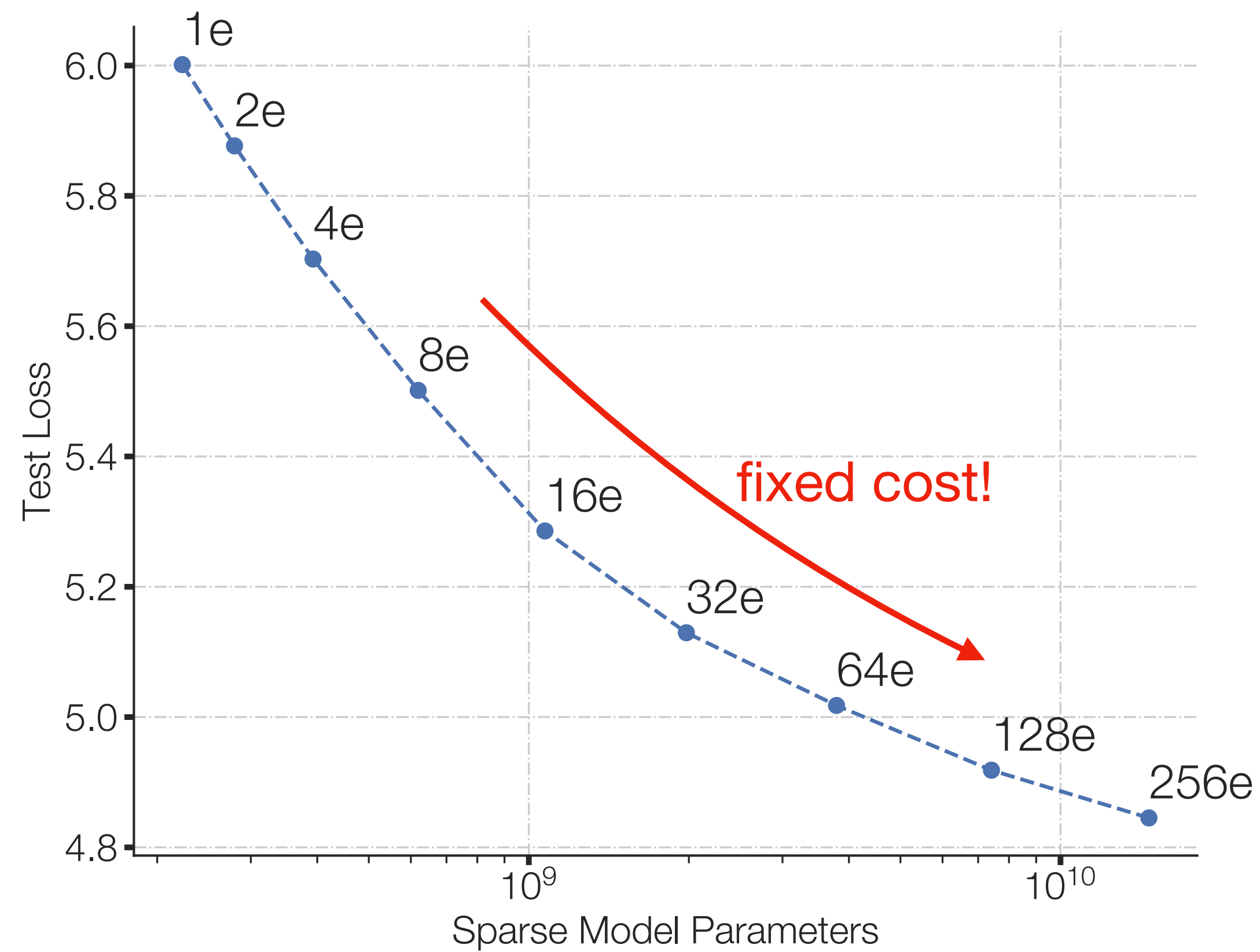
# Switch Transformers

More parameters, same cost!

- **Big Idea:** We don't *always* need to know *everything*.

- **Mixture-of-Experts** (MOE): provides an efficient mechanism for scaling up the number of parameters while keeping the total training cost fixed.

  ▸ Clever engineering avoids training instability

- **1.6 trillion parameters** (most to date), **7x speedup** over T5-BASE [1.]

[1.] Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer  arXiv:1910.10683
[2.] https://www.youtube.com/watch?v=iAR8LkkMMIM

Test Loss vs Sparse Model Parameters (top left): data points labeled 1e, 2e, 4e, 8e, 16e, 32e, 64e, 128e, 256e with "fixed cost!" arrow.

Neg Log Perplexity vs Training Step (top right): Switch-Base: 128e, Switch-Base: 64e, Switch-Base: 32e, Switch-Base: 16e, T5-Base, with "faster!" and "better!" annotations.

Neg Log Perplexity vs Training Time (bottom left): Switch-Base: 128e, Switch-Base: 64e, Switch-Base: 32e, T5-Base, with 7x Speedup.

Neg Log Perplexity vs Training Time (bottom right): Switch-Base: 64e, T5-Large, T5-Base, with 2.5x Speedup and 7.0x Speedup.
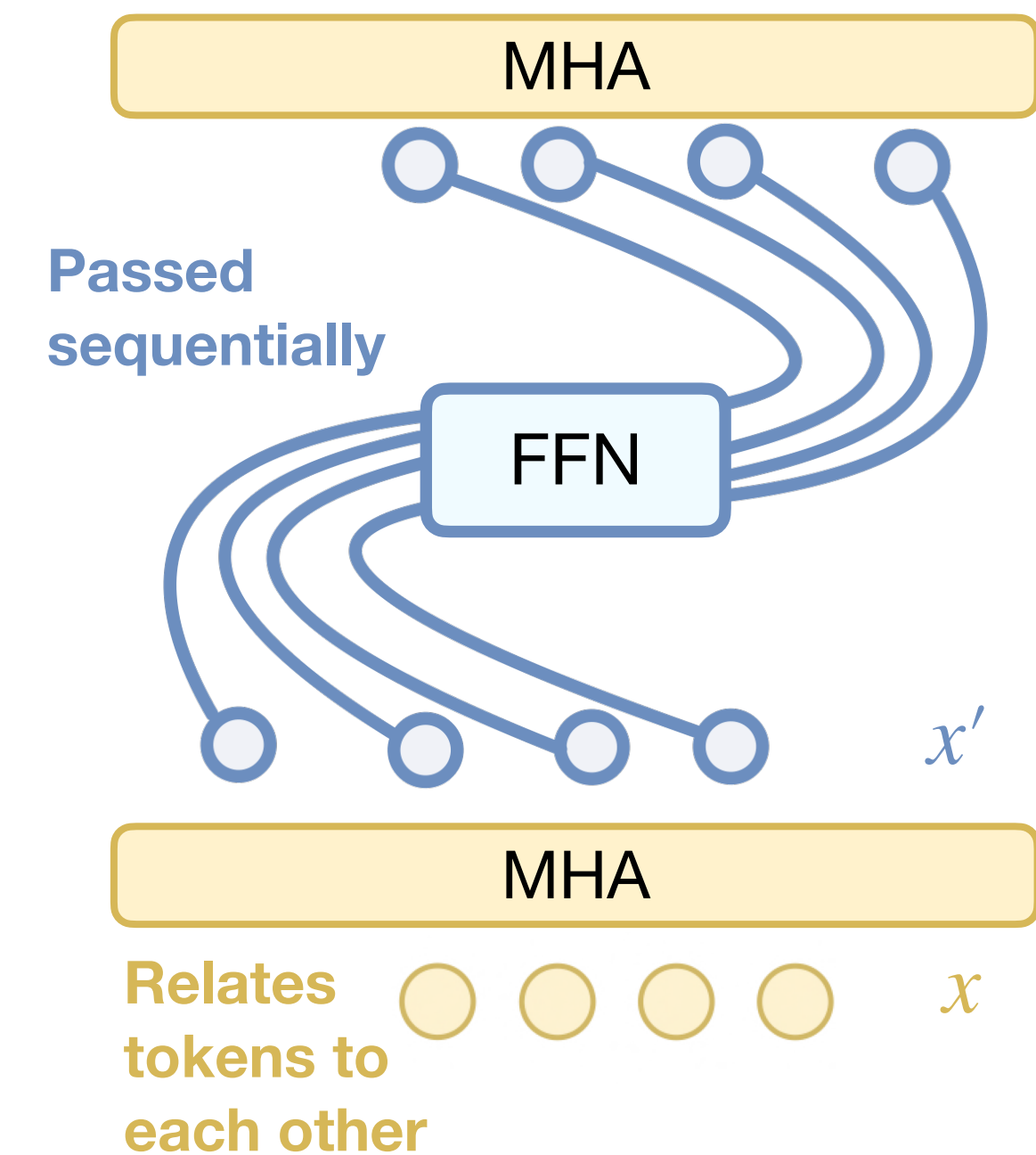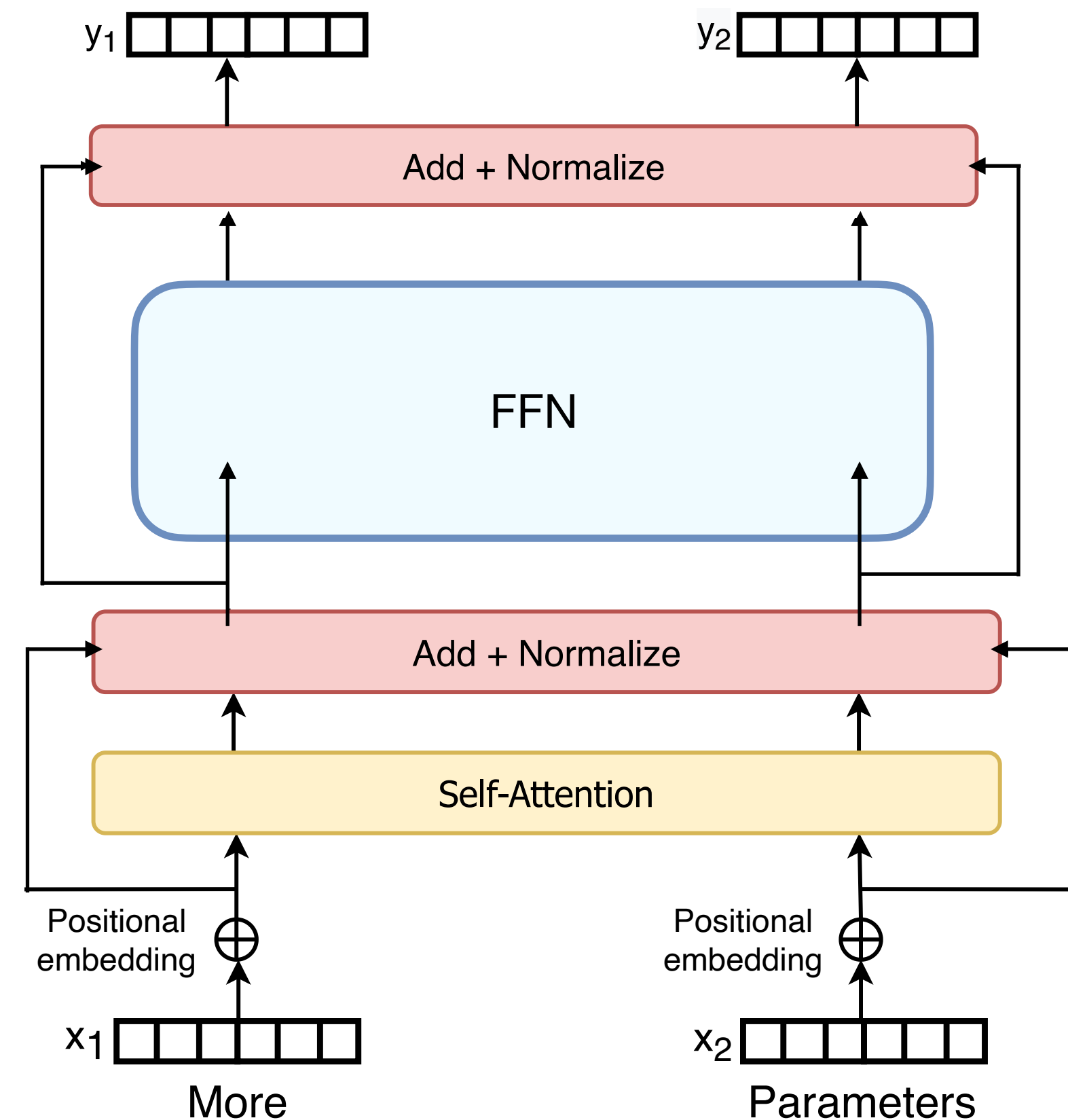
# Transformer Architecture

**Multi-Head Attention (MHA):**
- Aggregates information from sequences
- Relates tokens to each other

**Feed Forward (FFN):**
- Aggregates outputs from multiple heads
- Tokens in sequence are passed sequentially
- For a given token and its representation in this layer, what is the best representation in the next layer?
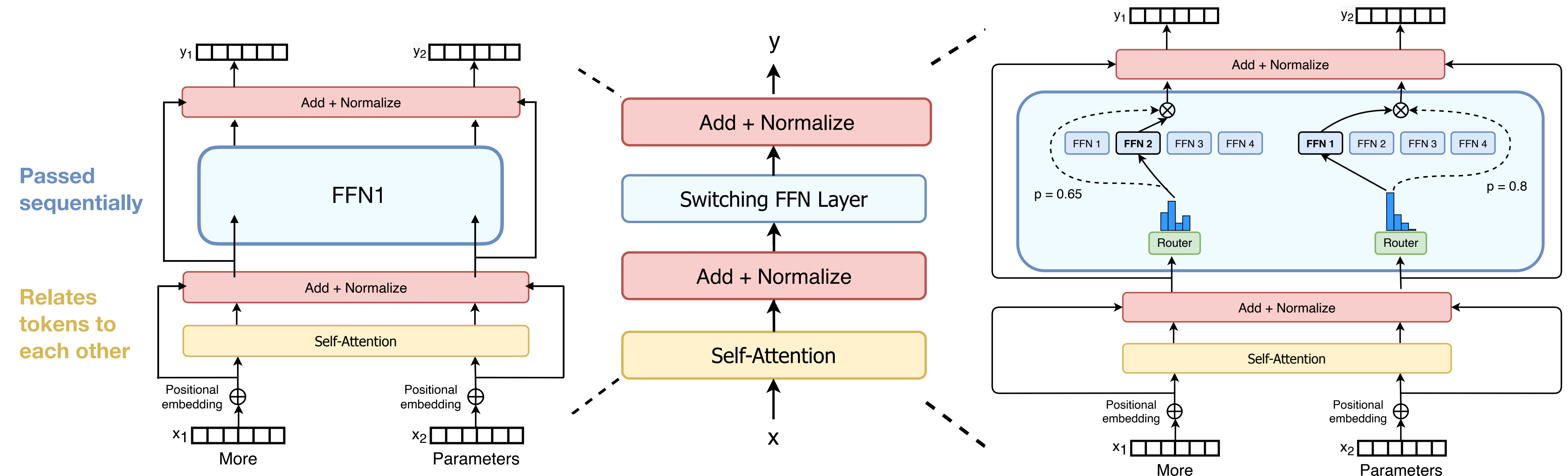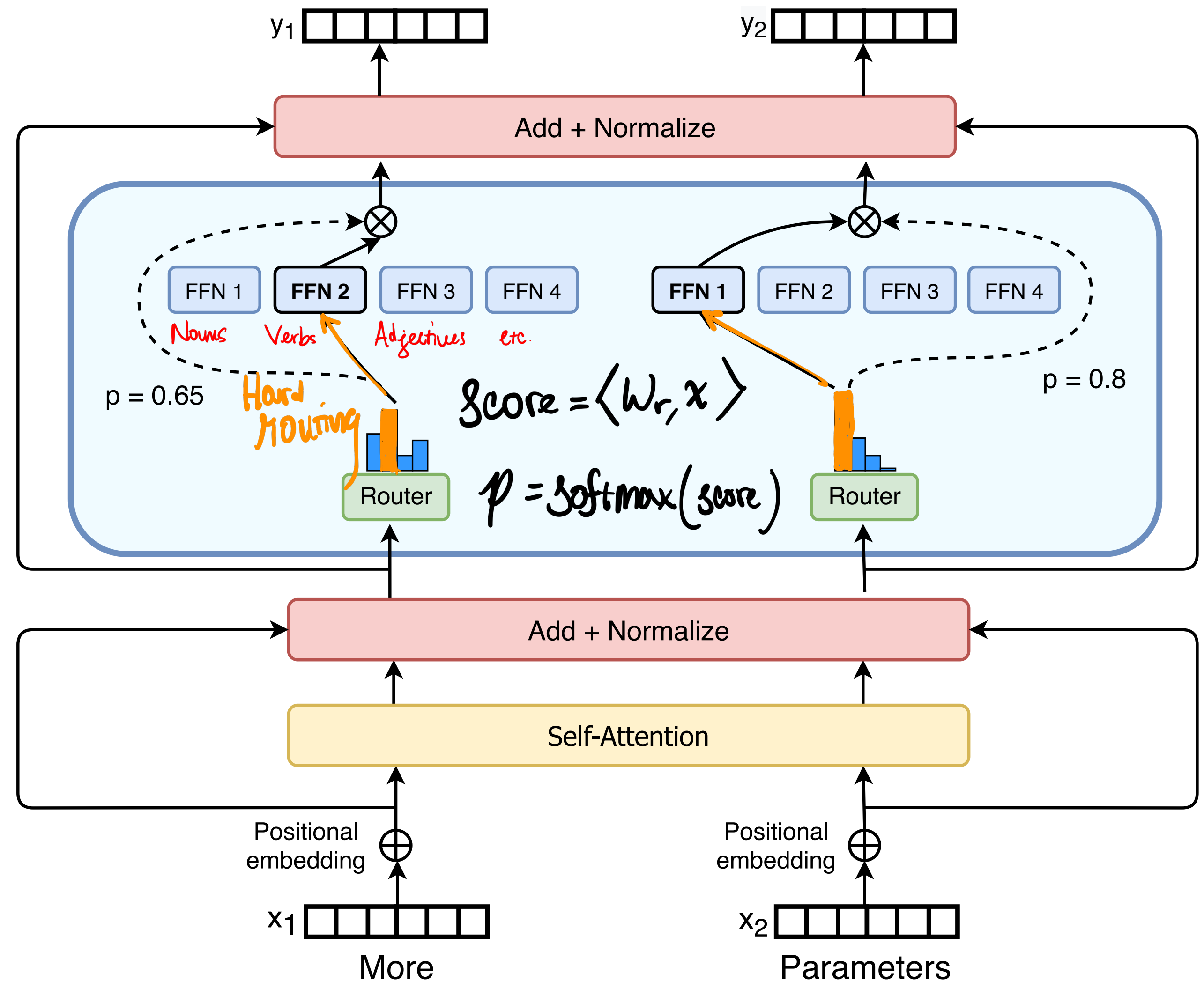
# Transformer Architecture
Using multiple FFNs

**Multi-Head Attention (MHA):**
- Aggregates information from sequences
- Relates tokens to each other

**Feed Forward (FFN):**
- Tokens in sequence are passed sequentially
- For a given token and its representation in this layer, what is the best representation in the next layer?
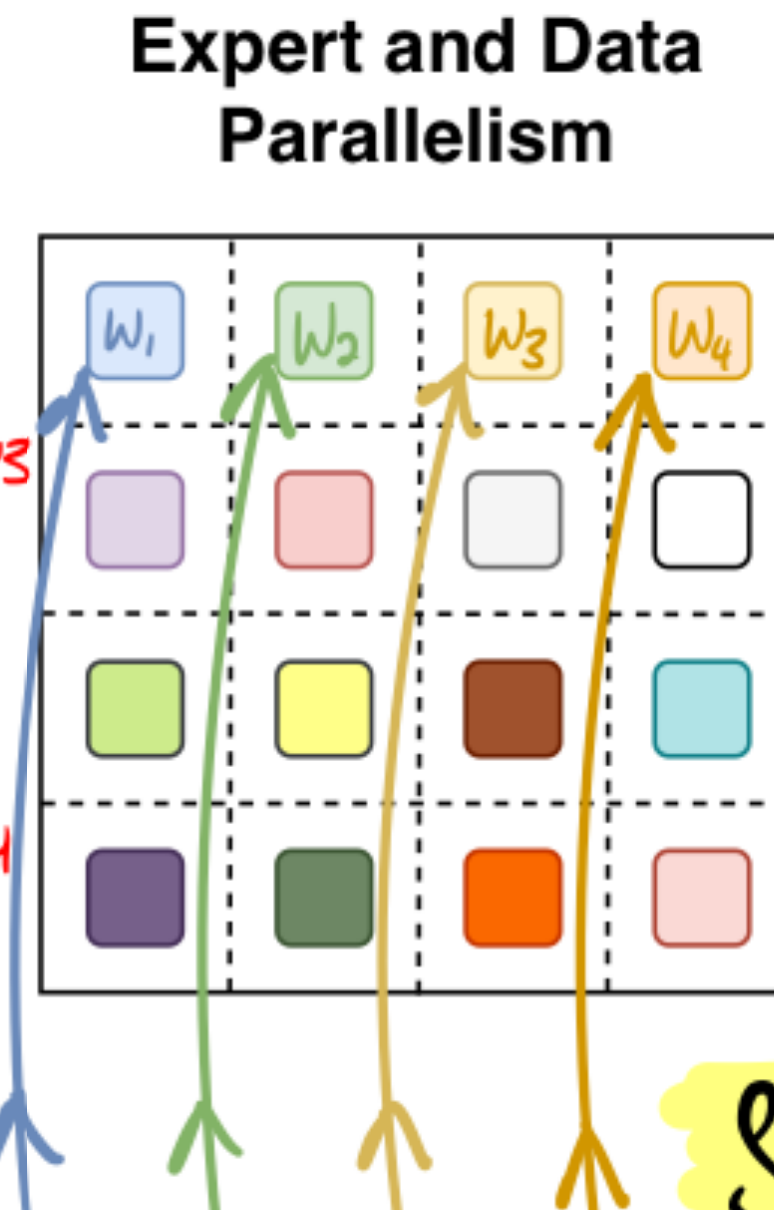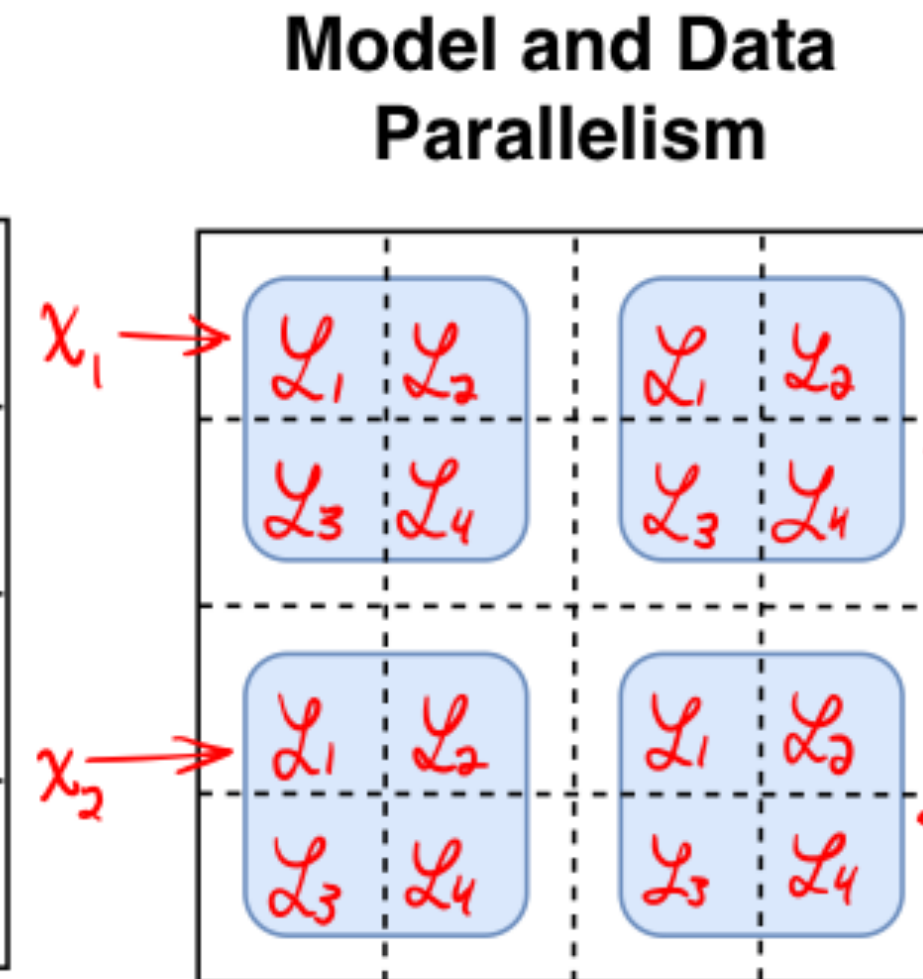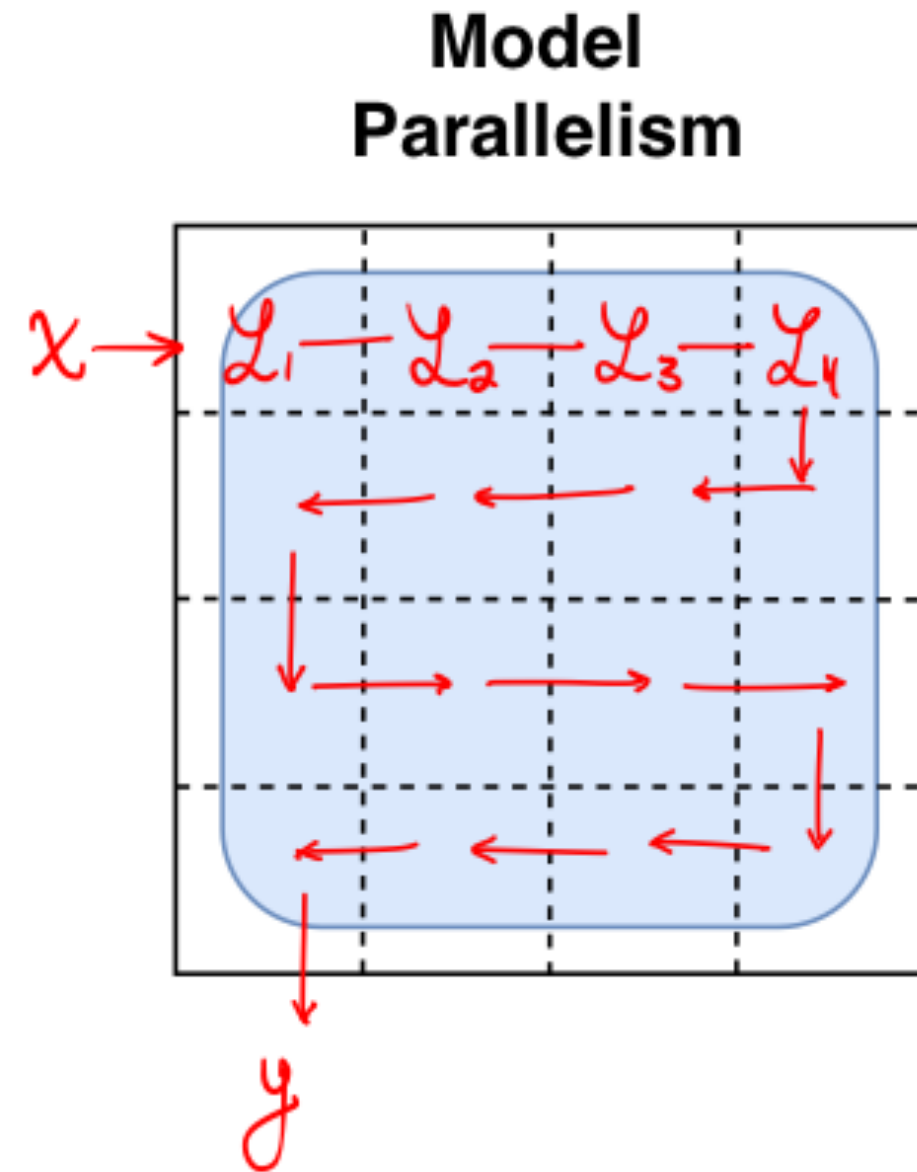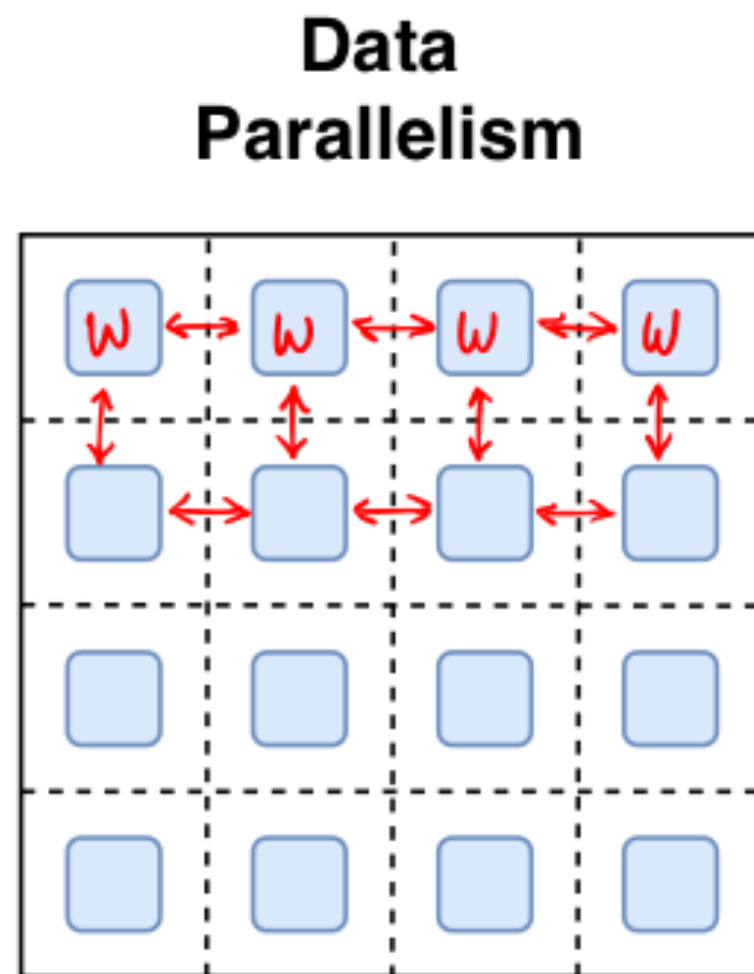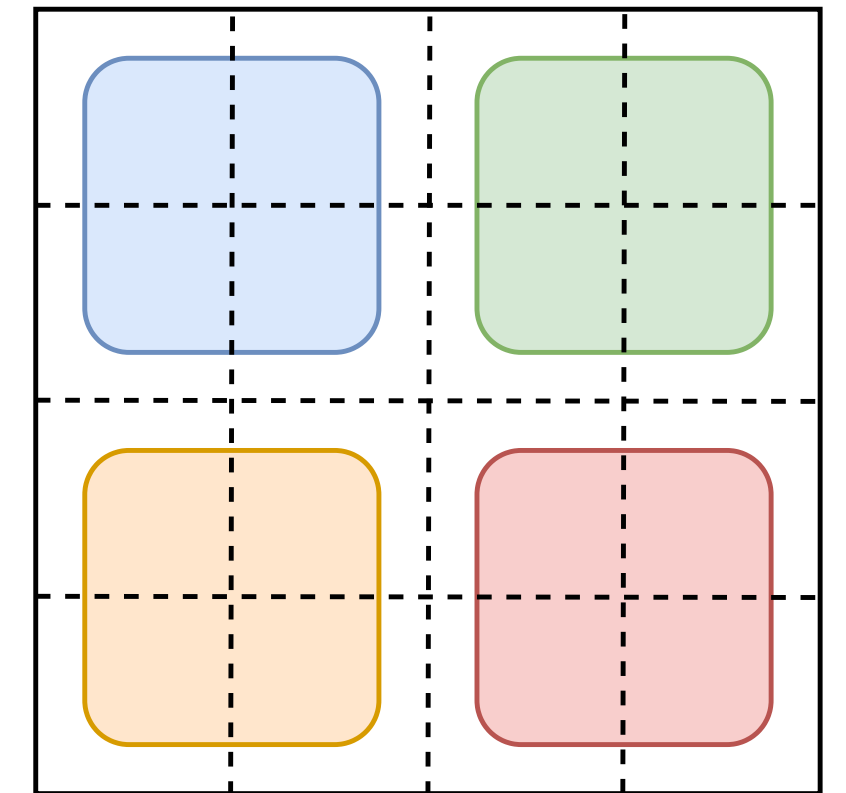
# Switch Layer

## Learning to Route inputs

- We don't need *all* the model's information for a particular input
  - ‣ **Use different parameters for different inputs**

- Replace *single* FFN with *multiple* FFNs, referred to as *experts.*
  - ‣ Each token is passed through exactly one FFN
  - ‣ Requires more memory
  - ‣ Constant FLOPs

- **Switch Layer:** Learns how to route each token to the most suitable expert.
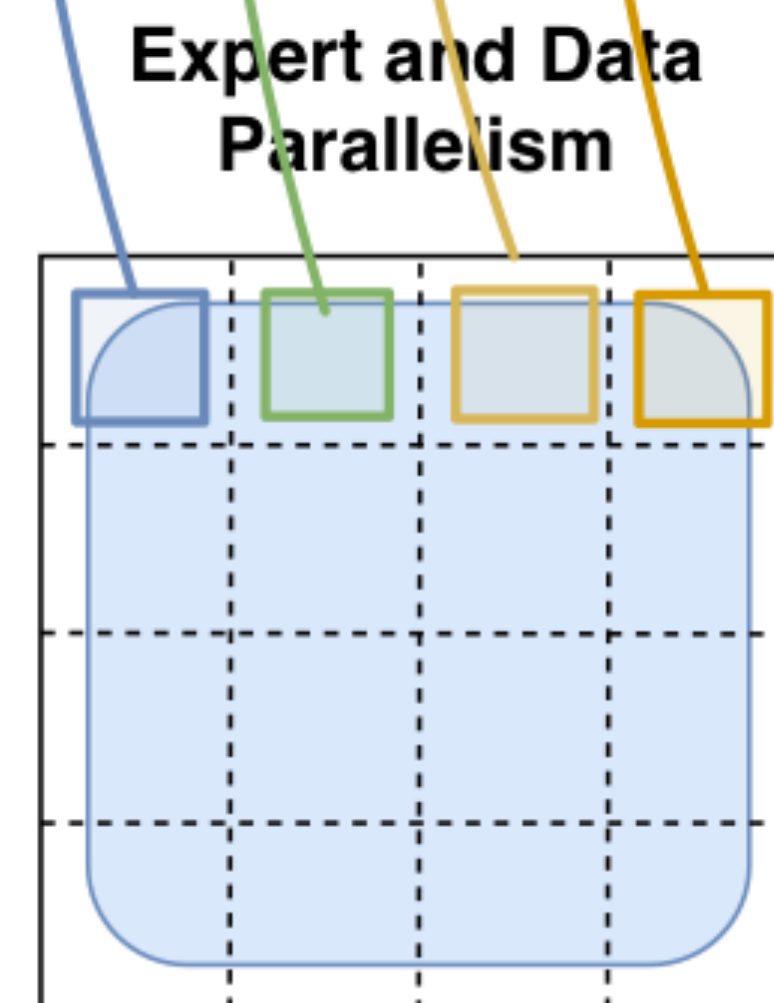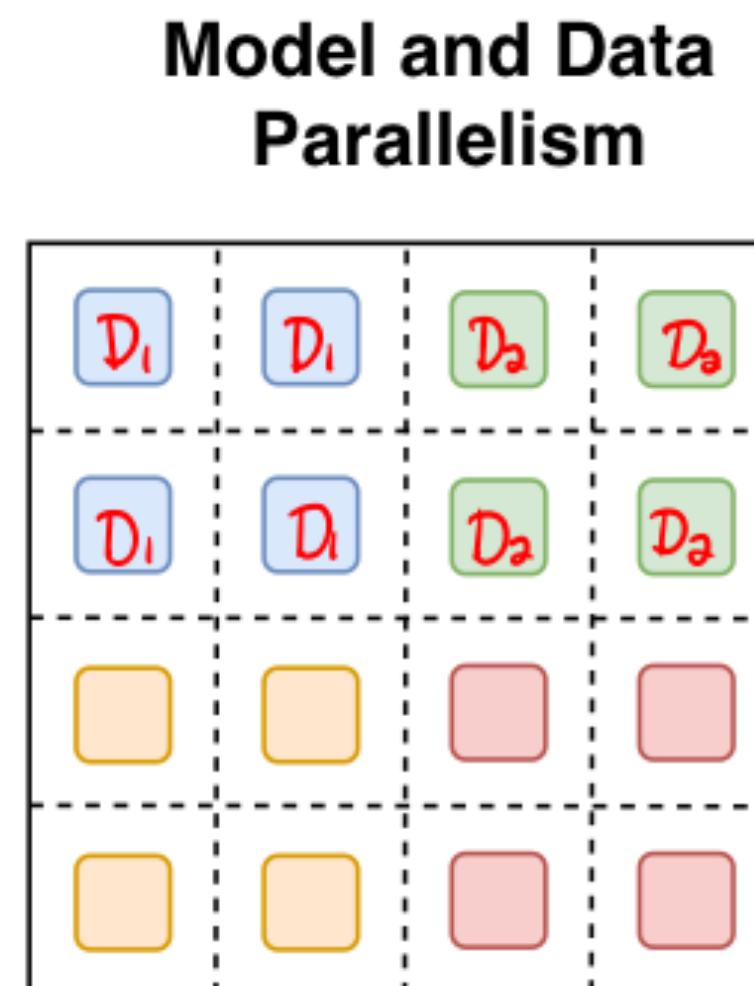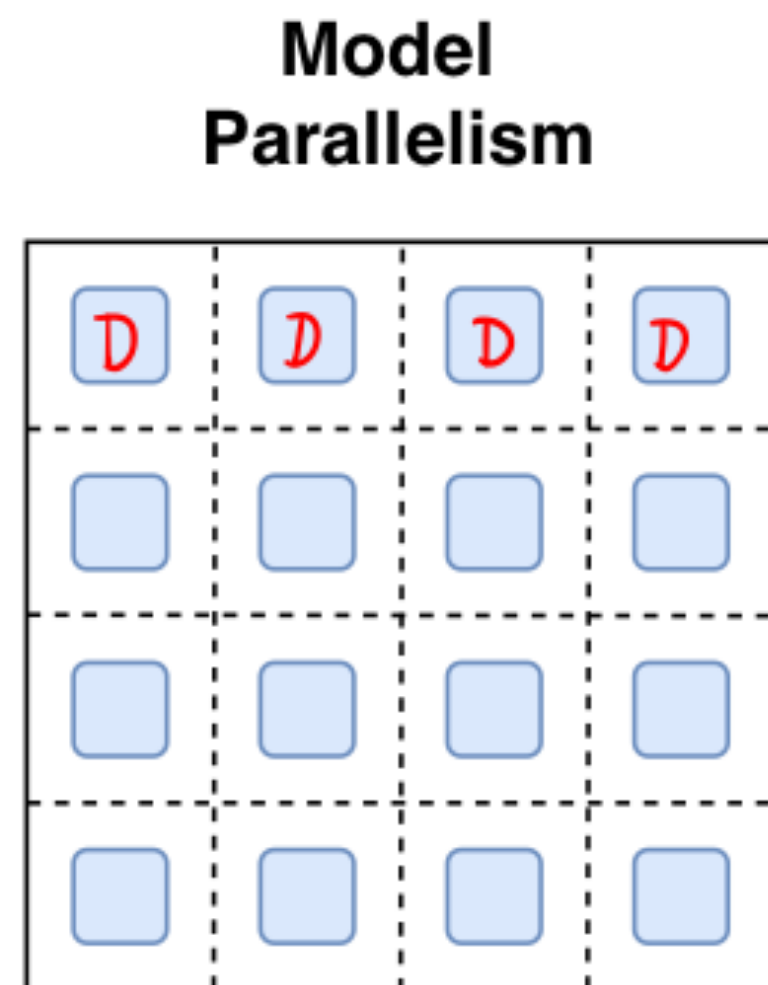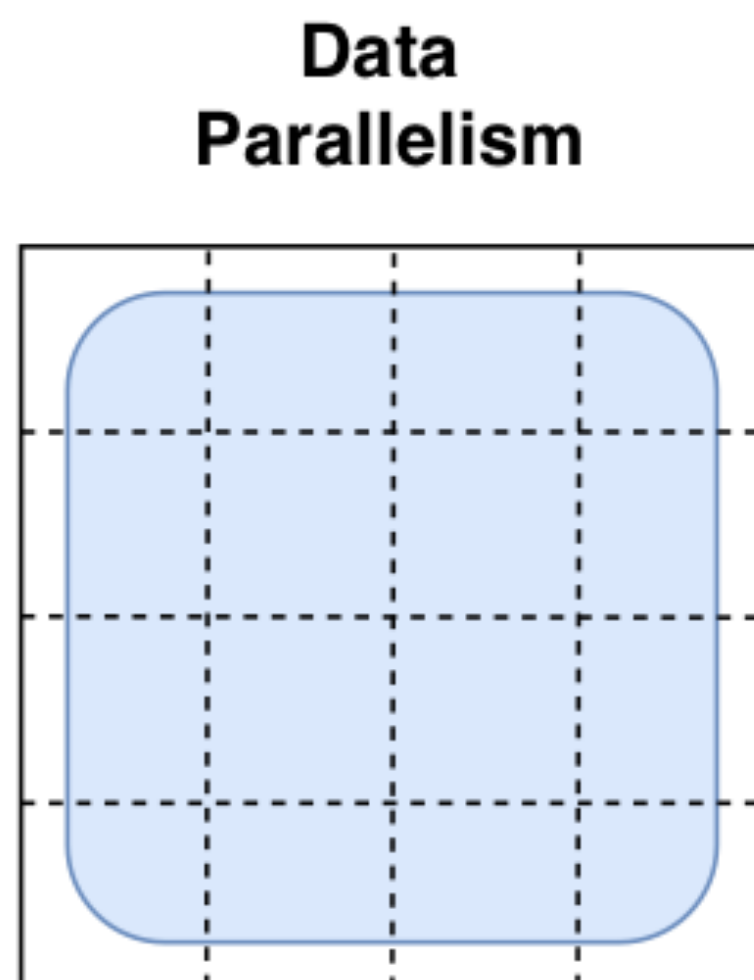
- Hard routing introduces sparsity!

# How the *model weights* are split over cores



**Data Parallelism**    **Model Parallelism**    **Model and Data Parallelism**    **Expert and Data Parallelism**    **Expert, Model and Data Parallelism**

# How the *data* is split over cores



**Data Parallelism**    **Model Parallelism**    **Model and Data Parallelism**    **Expert and Data Parallelism**    **Expert, Model and Data Parallelism**

# Stabilizing Training

- **Selective precision with large sparse models.**
  - ‣ Cast input to "switch" router up to float32 precision.
  - ‣ Cast output back to bfloat16
- **Smaller parameter initialization for stability.**
  - ‣ Initialize weight matrices from truncated normal with $\mu = 0$, $\sigma = \sqrt{s/n}$, and $s$ is a hyperparameter
    - – Recommend scaling default value ($s = 1.0$) by $\sim 1/10$.
- **Regularizing large sparse models**
  - ‣ Significantly increase dropout probability inside the experts

# Downstream Results

Multilingual Learning

- Outperforms T5 Base model *on all 101 languages!*