



# Training Transformers

Soumadeep Saha

# Language Modelling - Autoregressive



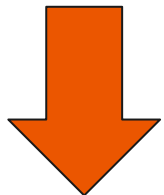
$$\begin{aligned} & [f(w_{T-1}, \dots, w_1)]_{w_T} \\ & \approx P(w_T | w_{T-1}, w_{T-2}, \dots, w_1) \end{aligned}$$

# Auto-encoder (Masked Language Modelling)



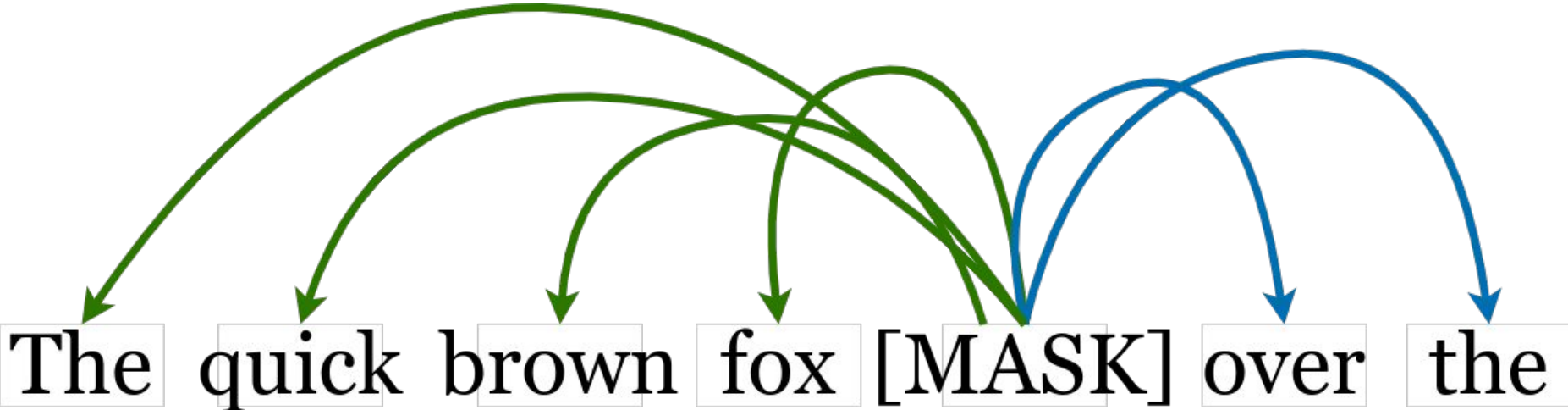
$$P(w_M | \dots, w_{M-2}, w_{M-1}, w_{M+1}, \dots)$$

The quick brown fox jumps over the lazy dog.



The quick brown fox [MASK] over the lazy dog.

# Are they the same?



# Is it all good?



$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}) = \sum_{t=1}^T \log \frac{\exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x_t))}{\sum_{x'} \exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x'))},$$

$$\max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} | \hat{\mathbf{x}}) \approx \sum_{t=1}^T m_t \log p_{\theta}(x_t | \hat{\mathbf{x}}) = \sum_{t=1}^T m_t \log \frac{\exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x_t))}{\sum_{x'} \exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x'))},$$

# Is it all good?



$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}) = \sum_{t=1}^T \log \frac{\exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x_t))}{\sum_{x'} \exp(h_{\theta}(\mathbf{x}_{1:t-1})^{\top} e(x'))},$$

$$\max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} | \hat{\mathbf{x}}) \approx \sum_{t=1}^T m_t \log p_{\theta}(x_t | \hat{\mathbf{x}}) = \sum_{t=1}^T m_t \log \frac{\exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x_t))}{\sum_{x'} \exp(H_{\theta}(\hat{\mathbf{x}})_t^{\top} e(x'))},$$

$$\hat{\mathbf{x}}) \approx \sum_{t=1}^T m_t$$

# Pros and Cons...



- 15% of [MASK] tokens...
- Independence assumption.

## **BUT...**

- We get bi-directional context

# Training pipeline

Clean

Tokenize

Batch

Q,K,V

+ Position  
Embedding

Token Embedding

$N^2$  attention

...

Loss



Next Batch

Update Weights

Back propagate gradients



# Walkthrough

The quick brown fox jumps over the lazy dog.   敏捷的棕色狐狸跳过了懒狗 `<a href="www.google.com">click`

Multilingual model, cased/uncased, etc

the quick brown fox jumps over the lazy dog

Special tokens

`[[CLS], the, quick, brown, fox, jumps, over, the, lazy, dog, [SEP]]`

# Walkthrough



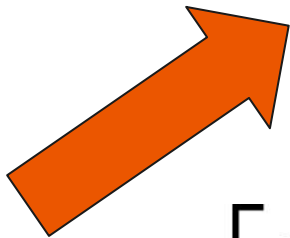
**[**[CLS], the, quick, br#, #own, fox, jumps, over, the, lazy, dog, **]** **]** Collect several in a batch [tokenized]

**[**[CLS], the, ..., [SEP], [PAD], [PAD], ..., [PAD] **]**

**[**[CLS], the, ..., [SEP], [PAD], [PAD], ..., [PAD]**]**

$$\left[ \begin{array}{c} a_{11} \\ a_{12} \\ \dots \\ a_{1d} \end{array} \right], \begin{array}{c} a_{21} \\ a_{22} \\ \dots \\ a_{2d} \end{array}, \dots, \begin{array}{c} a_{k1} \\ a_{k2} \\ \dots \\ a_{kd} \end{array} \right]$$

Batch



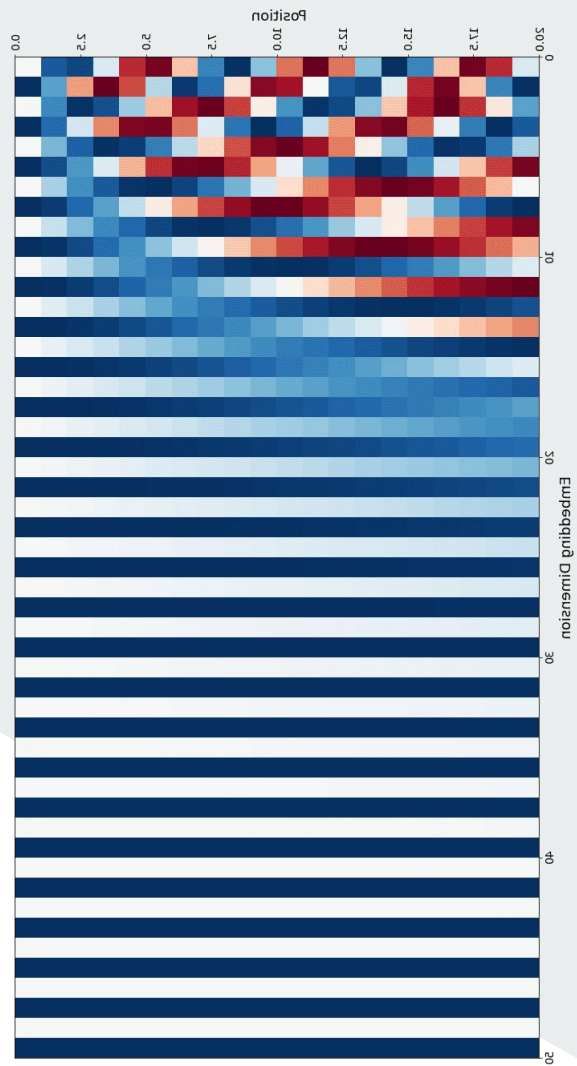
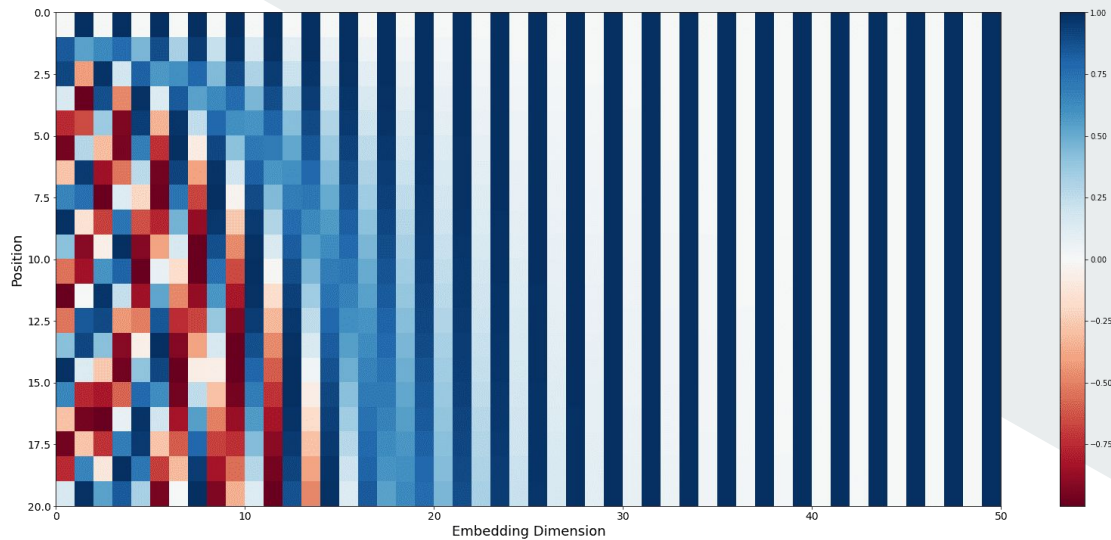
$$\begin{bmatrix} \begin{bmatrix} a_{11} \\ a_{12} \\ \dots \\ a_{1d} \end{bmatrix}, \begin{bmatrix} a_{21} \\ a_{22} \\ \dots \\ a_{2d} \end{bmatrix}, \dots, \begin{bmatrix} a_{k1} \\ a_{k2} \\ \dots \\ a_{kd} \end{bmatrix}, \begin{bmatrix} a_{11} \\ a_{21} \\ \dots \\ a_{k1} \end{bmatrix}, \begin{bmatrix} a_{12} \\ a_{22} \\ \dots \\ a_{k2} \end{bmatrix}, \dots, \begin{bmatrix} a_{1d} \\ a_{2d} \\ \dots \\ a_{kd} \end{bmatrix} \end{bmatrix}$$

+ Positional embeddings

# Positional Embeddings

$$\left[ \begin{array}{c} a_{11} \\ a_{12} \\ \dots \\ a_{1d} \end{array} \right], \begin{array}{c} a_{21} \\ a_{22} \\ \dots \\ a_{2d} \end{array}, \dots, \begin{array}{c} a_{k1} \\ a_{k2} \\ \dots \\ a_{kd} \end{array} \right]$$

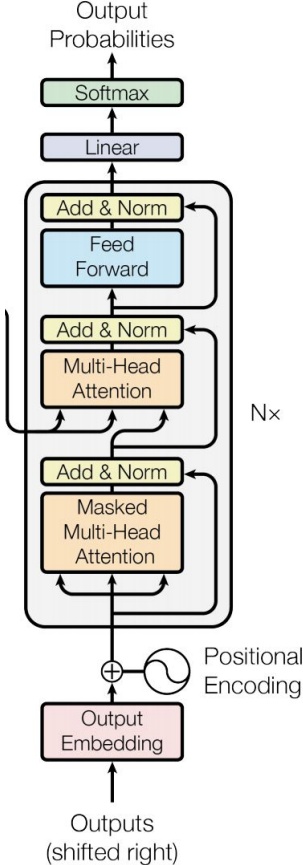
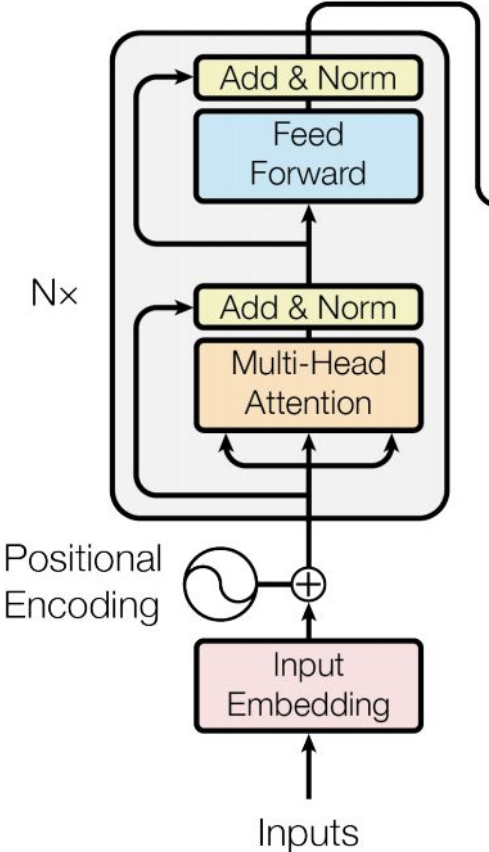
$$\left[ \begin{array}{c} 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{array} \right], \begin{array}{c} 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{array}, \dots, \begin{array}{c} 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{array}, \dots \right]$$



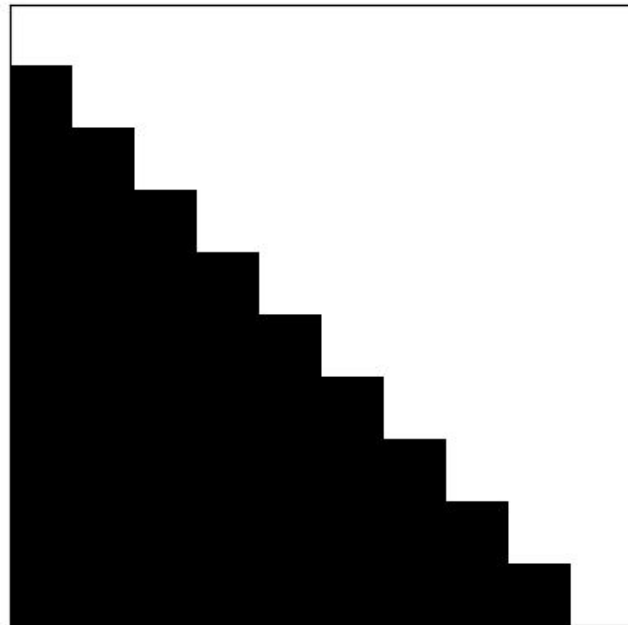
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

# BERT vs GPT



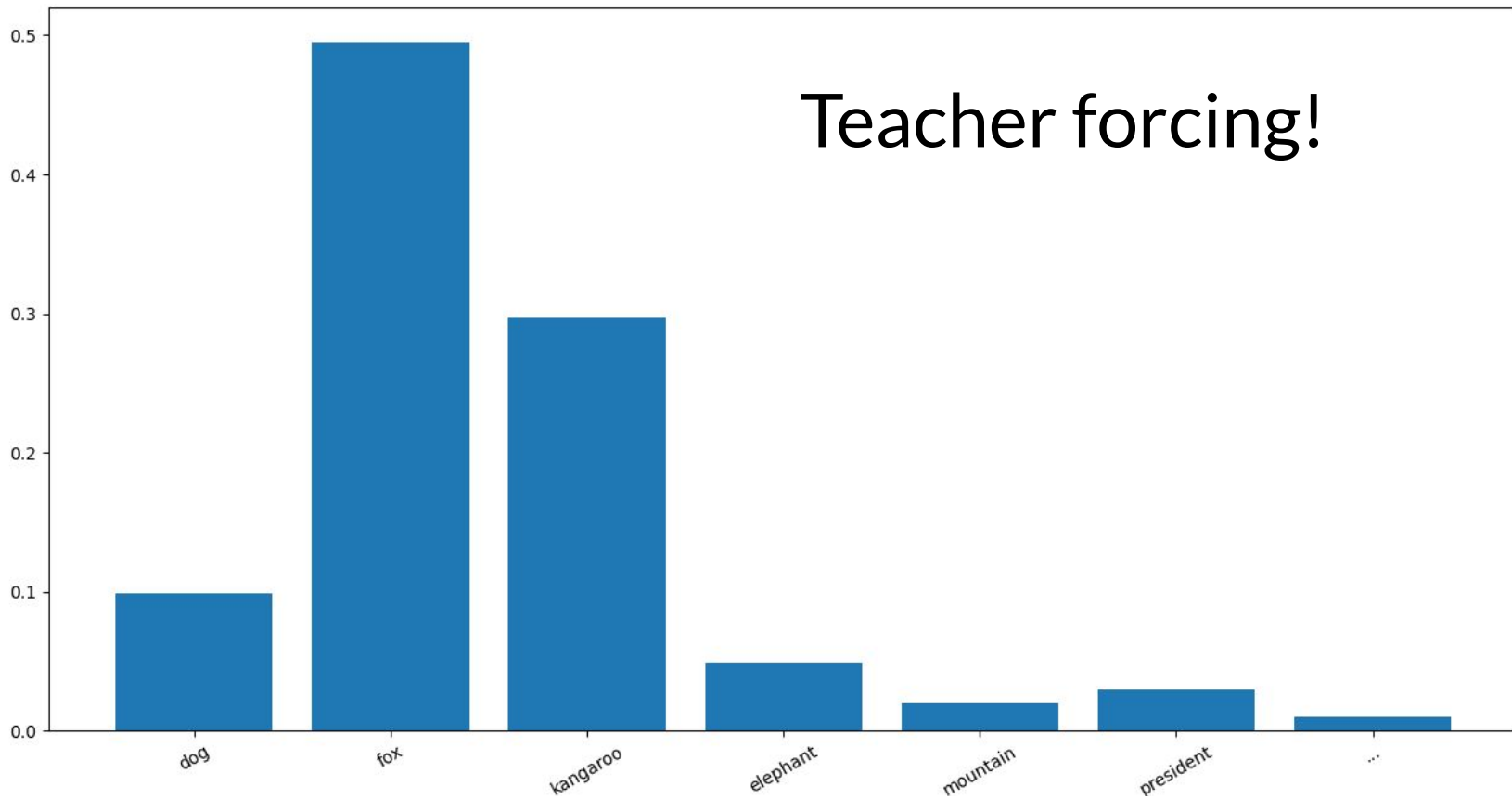
# BERT vs GPT



The quick brown fox jumps over the lazy dog.



The quick brown fox jumps over the lazy dog.



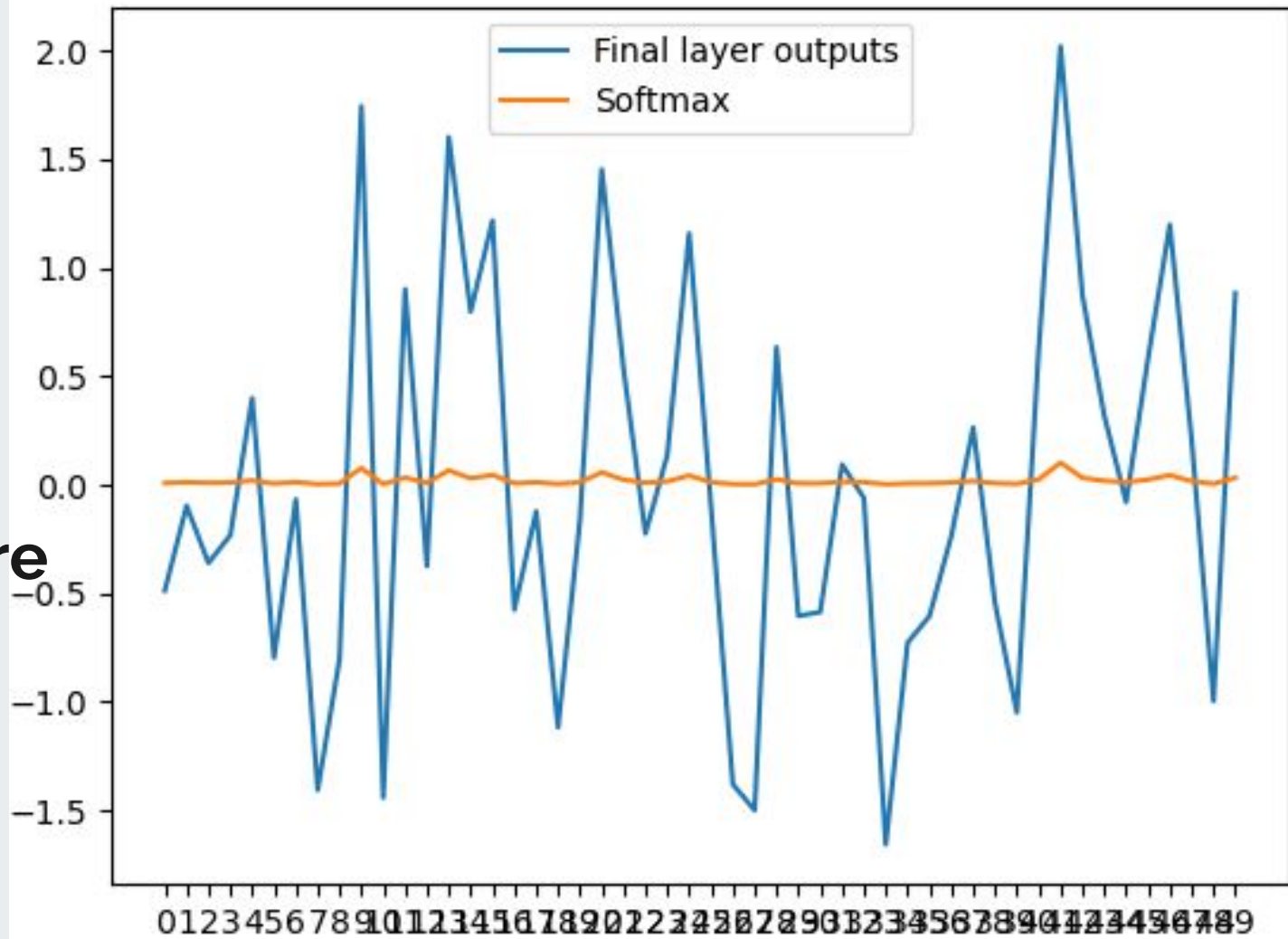
Teacher forcing!



# Inference?

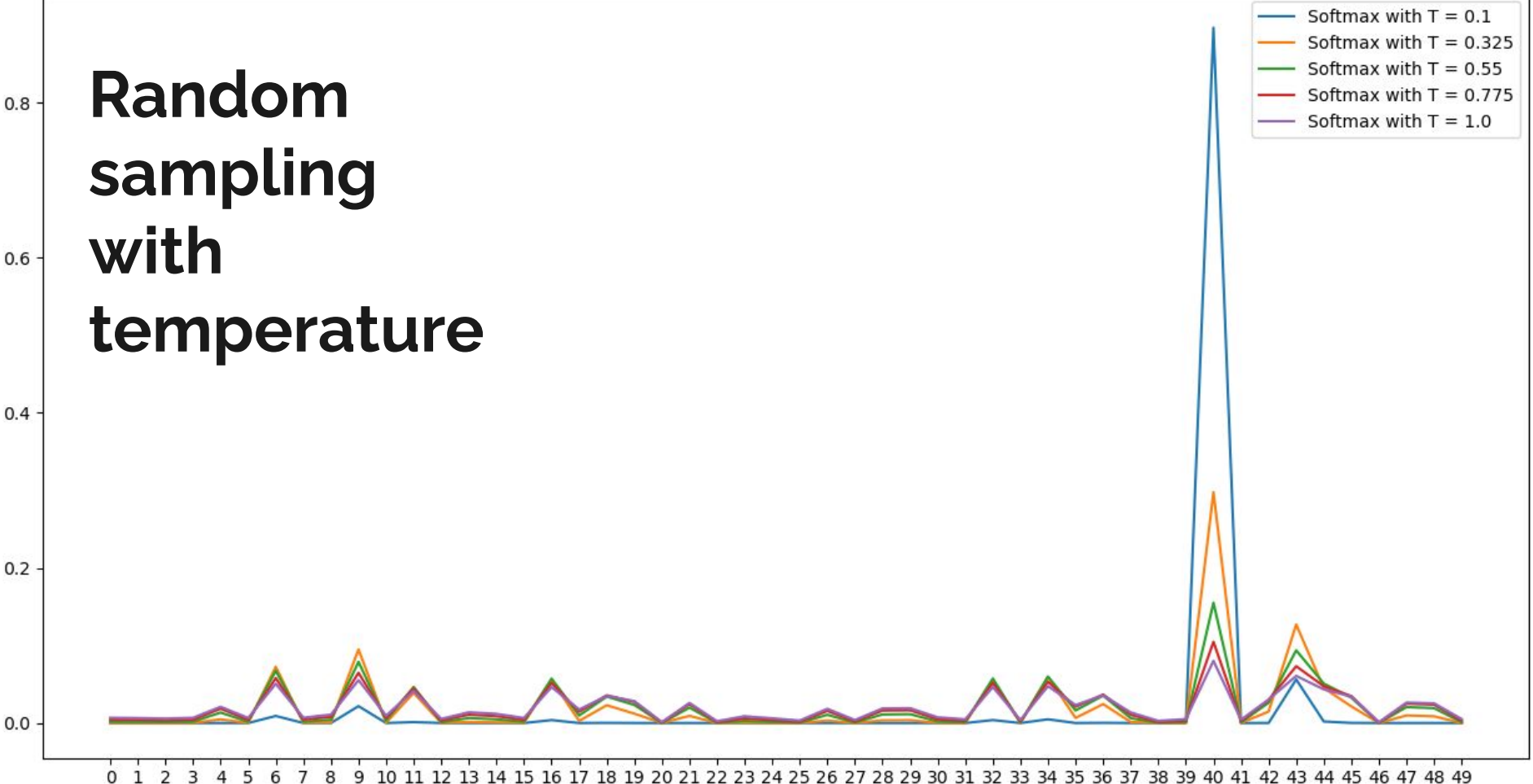
- Greedy - pick best answer.
- Random (with Temperature)
- Beam search
- Nucleus Sampling

Random  
sampling  
with  
temperature

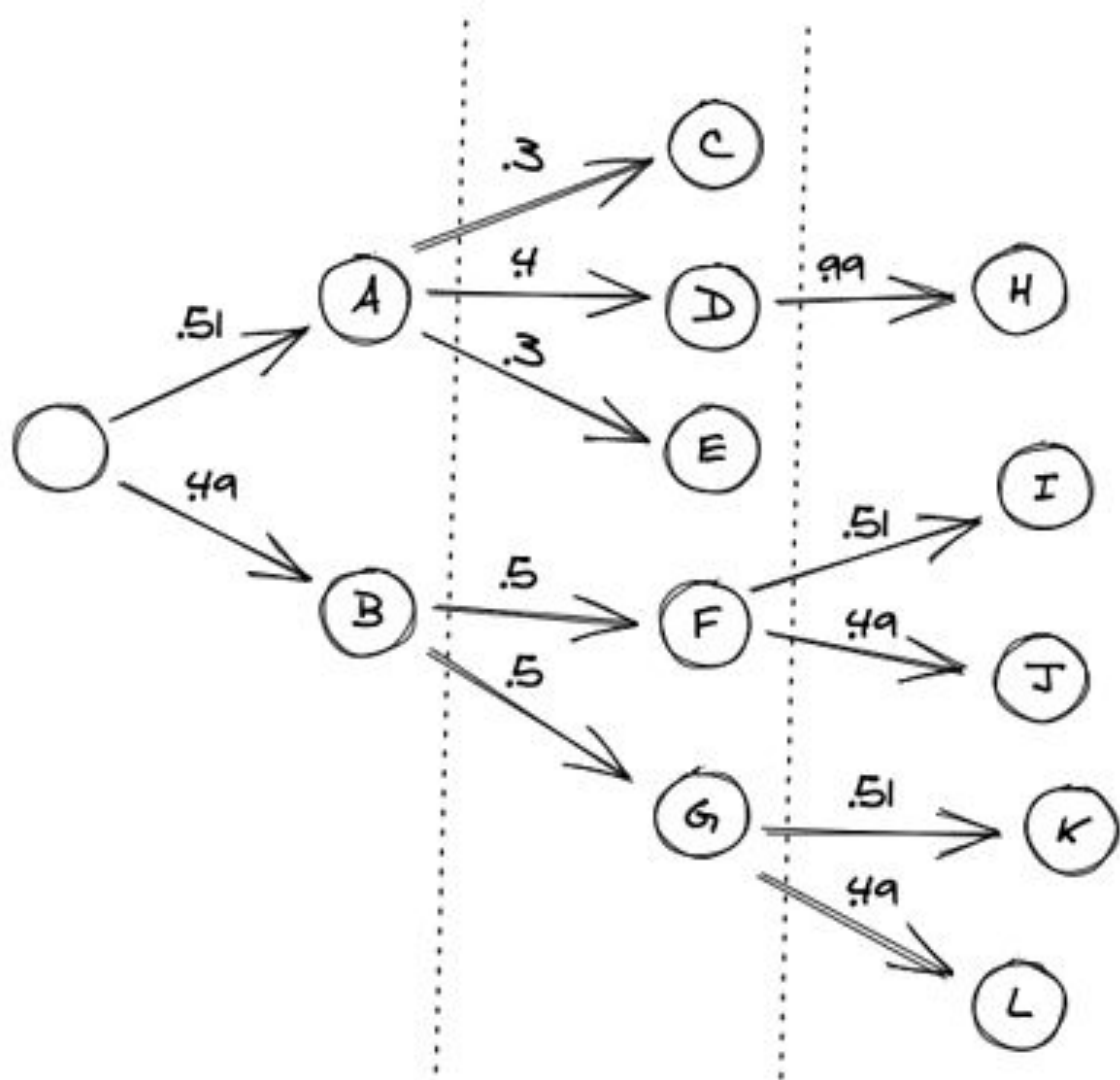


# Random sampling with temperature

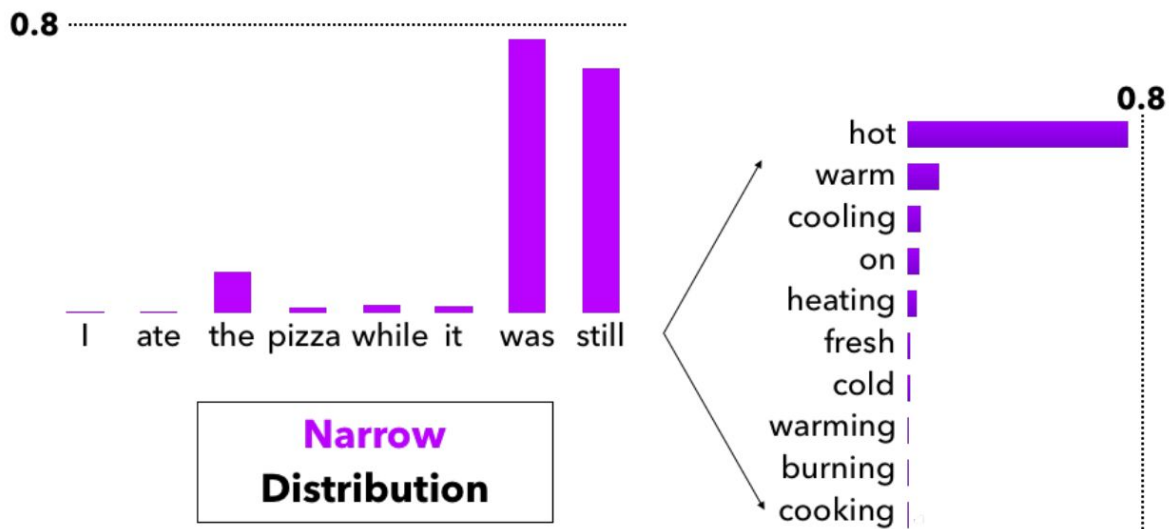
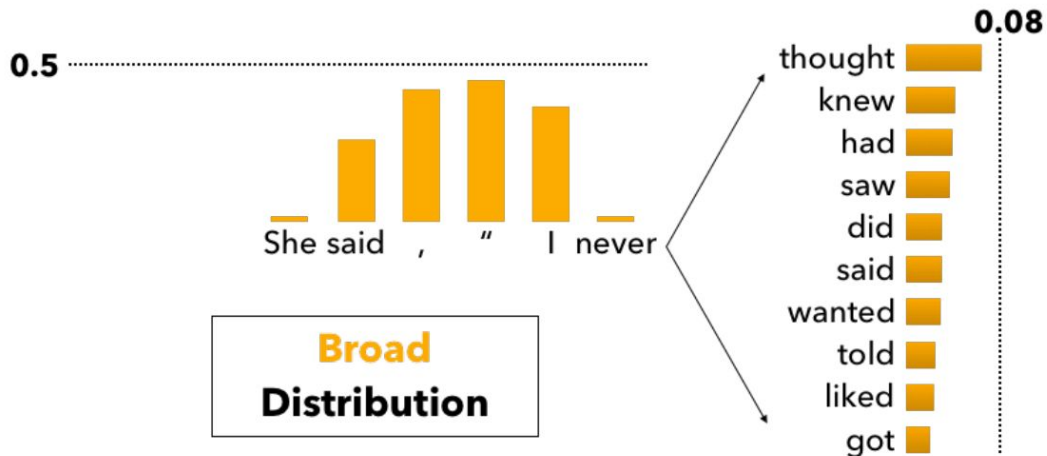
- Softmax with  $T = 0.1$
- Softmax with  $T = 0.325$
- Softmax with  $T = 0.55$
- Softmax with  $T = 0.775$
- Softmax with  $T = 1.0$



# Beam Search



# Nucleus Sampling



# Large Foundation Models - Conclusions

---

- Spoilt for choice [BERT, LLaMA, Falcon, etc]
- Size of model - bigger is not always better
  - 100M - 500B parameters
- Pre-training distribution
- AE vs AR

We'll cover more details in the session on scaling.