



# STAT 453: Introduction to Deep Learning and Generative Models

---

Ben Lengerich

Lecture 21: GPT Architectures

November 17, 2025

Reading: See course homepage

# From RNN...

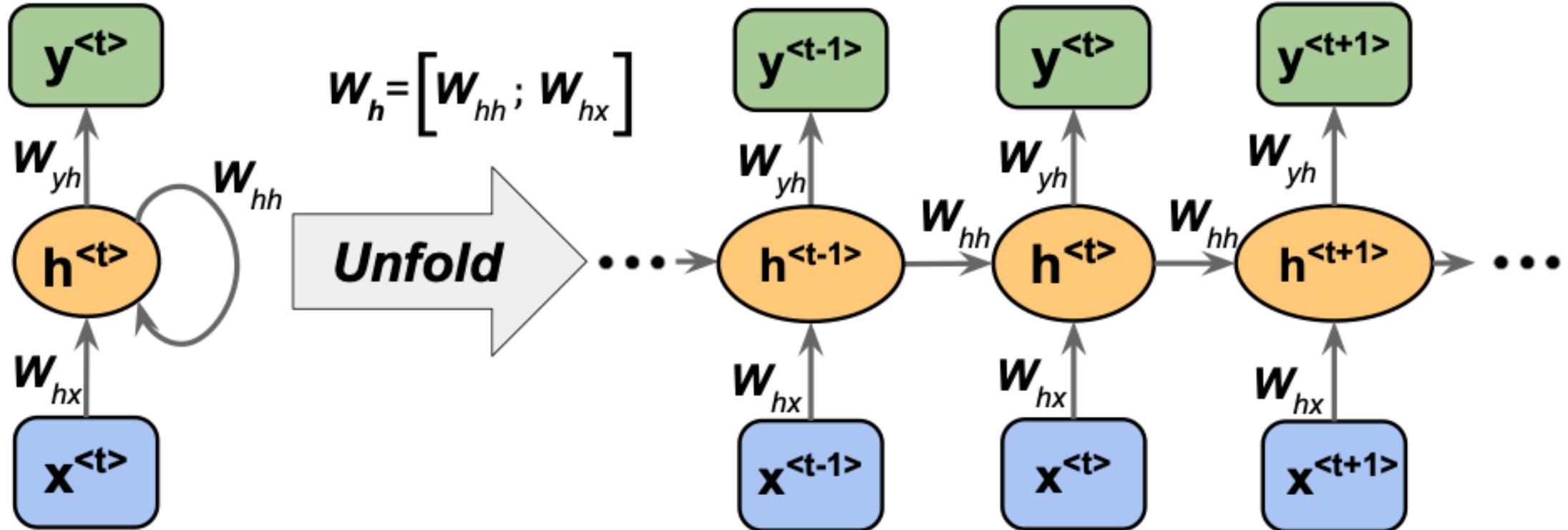


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

# From RNN...to GPT

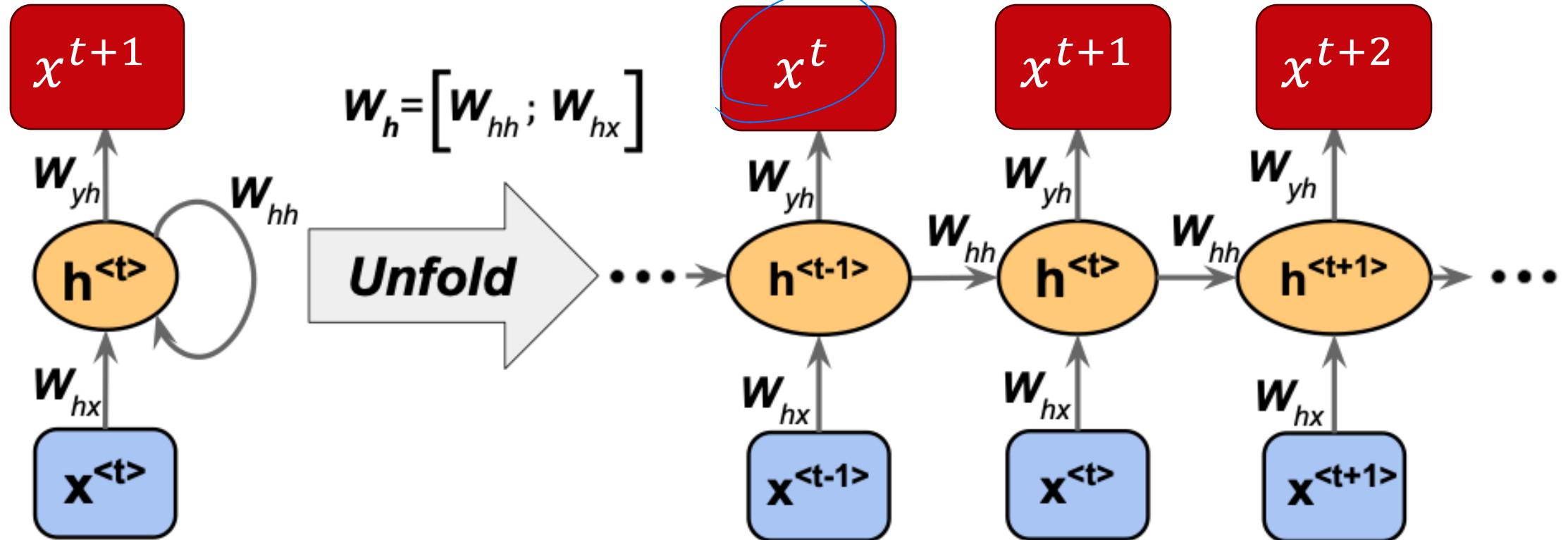
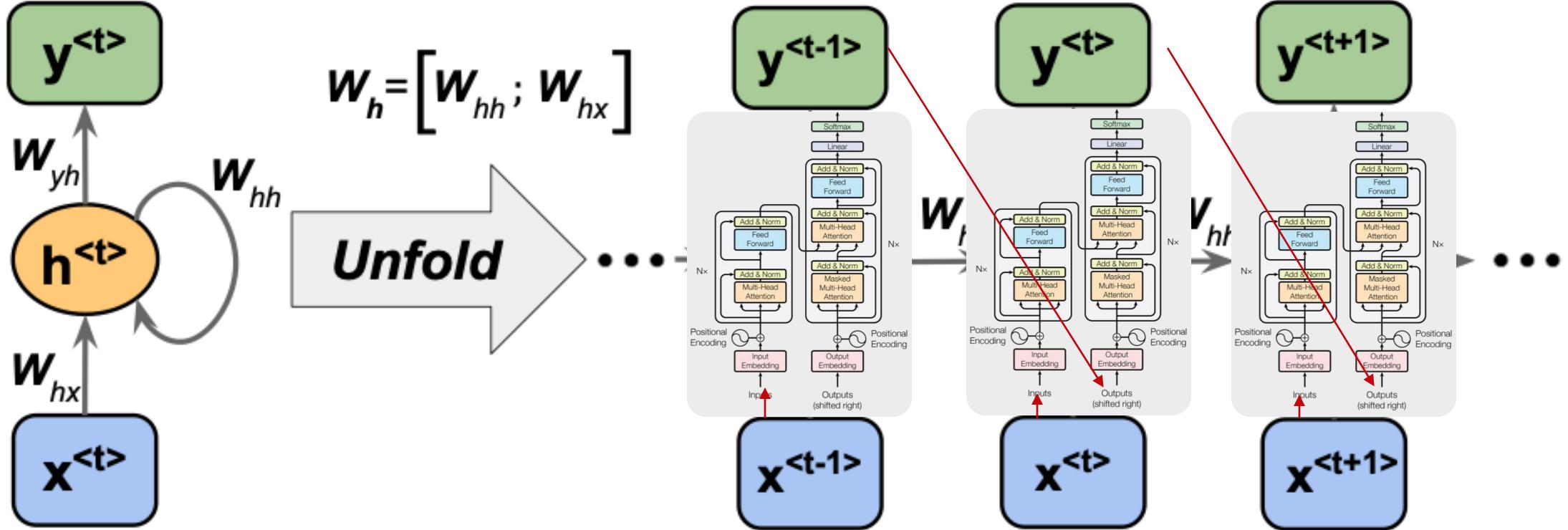


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019



# From RNN...to GPT...by Transformers

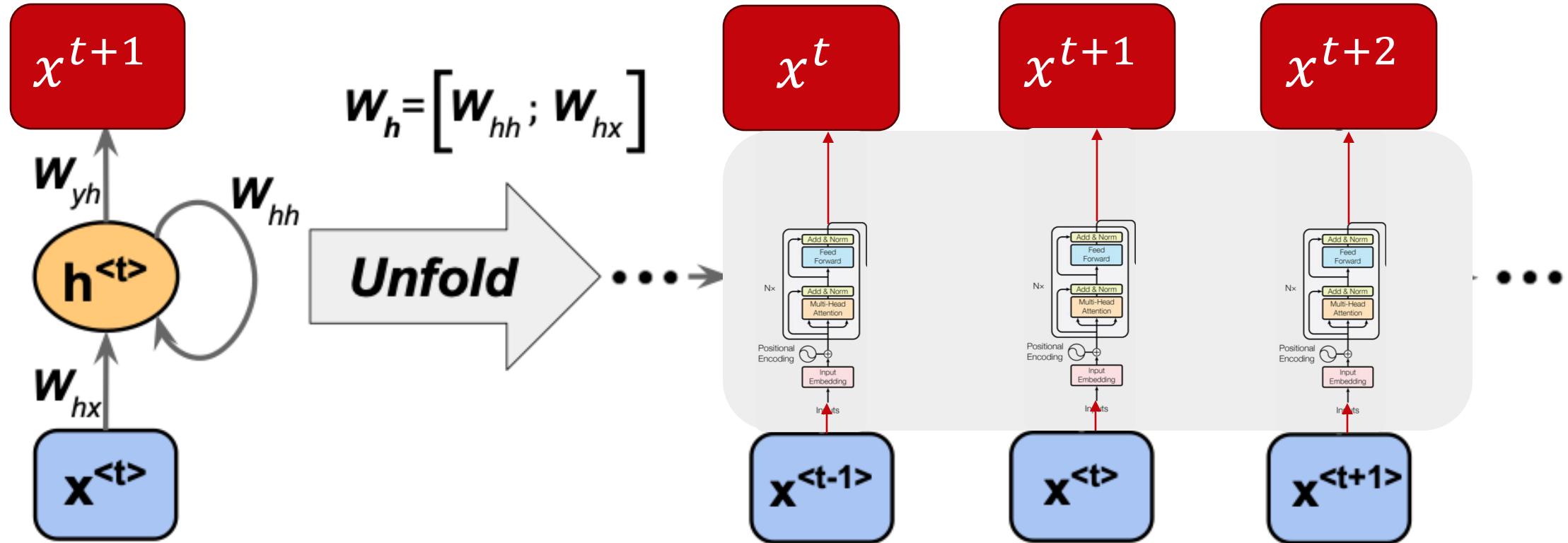


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.

Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019



# From RNN...to GPT...by Transformers



Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.

Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

# Self-attention (very basic form)

No learnable parameters?

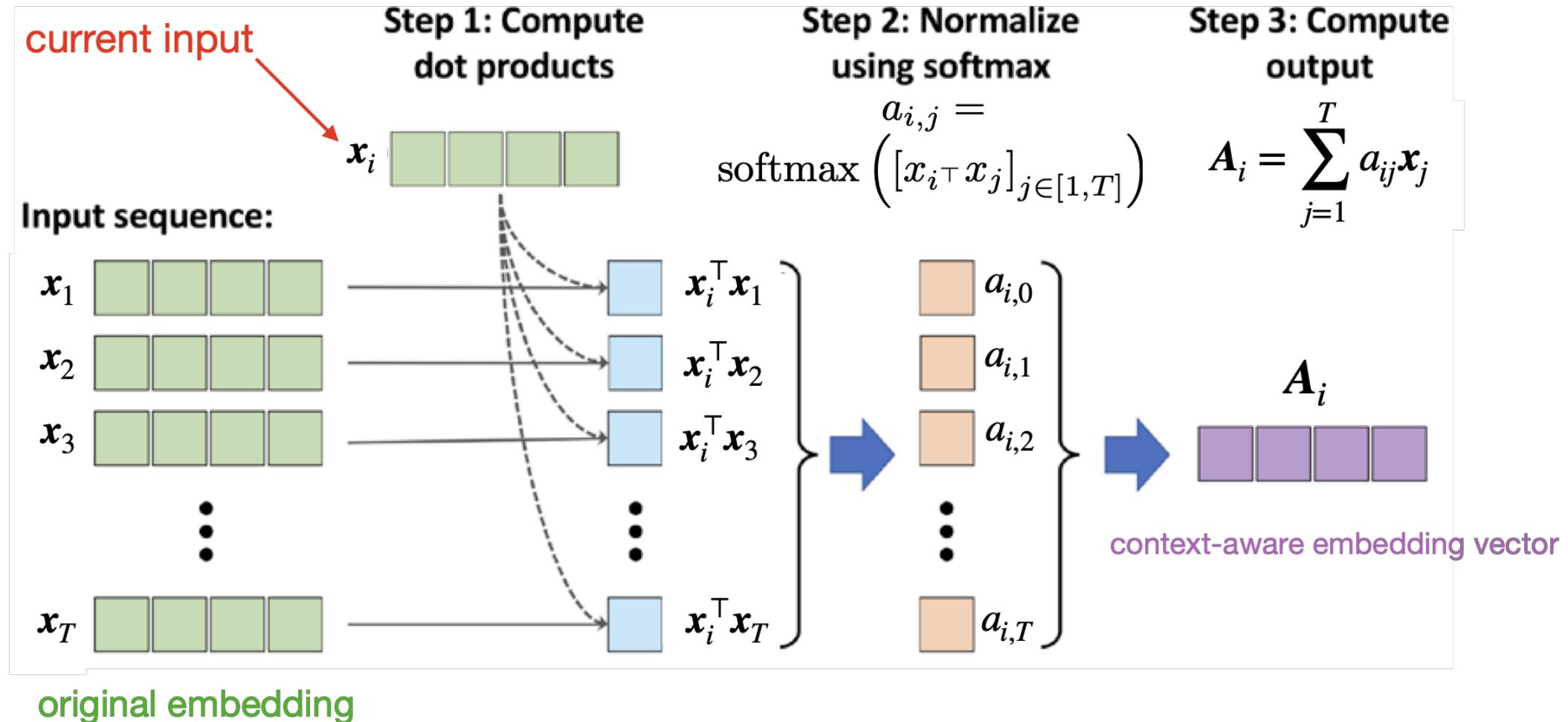


Image source: Raschka & Mirjalili 2019. Python Machine Learning, 3rd edition



# Learnable Self-attention

---

- Previous basic version did not involve any learnable parameters, so not very useful for learning a language model
- We are now adding 3 trainable weight matrices that are multiplied with the input sequence embeddings

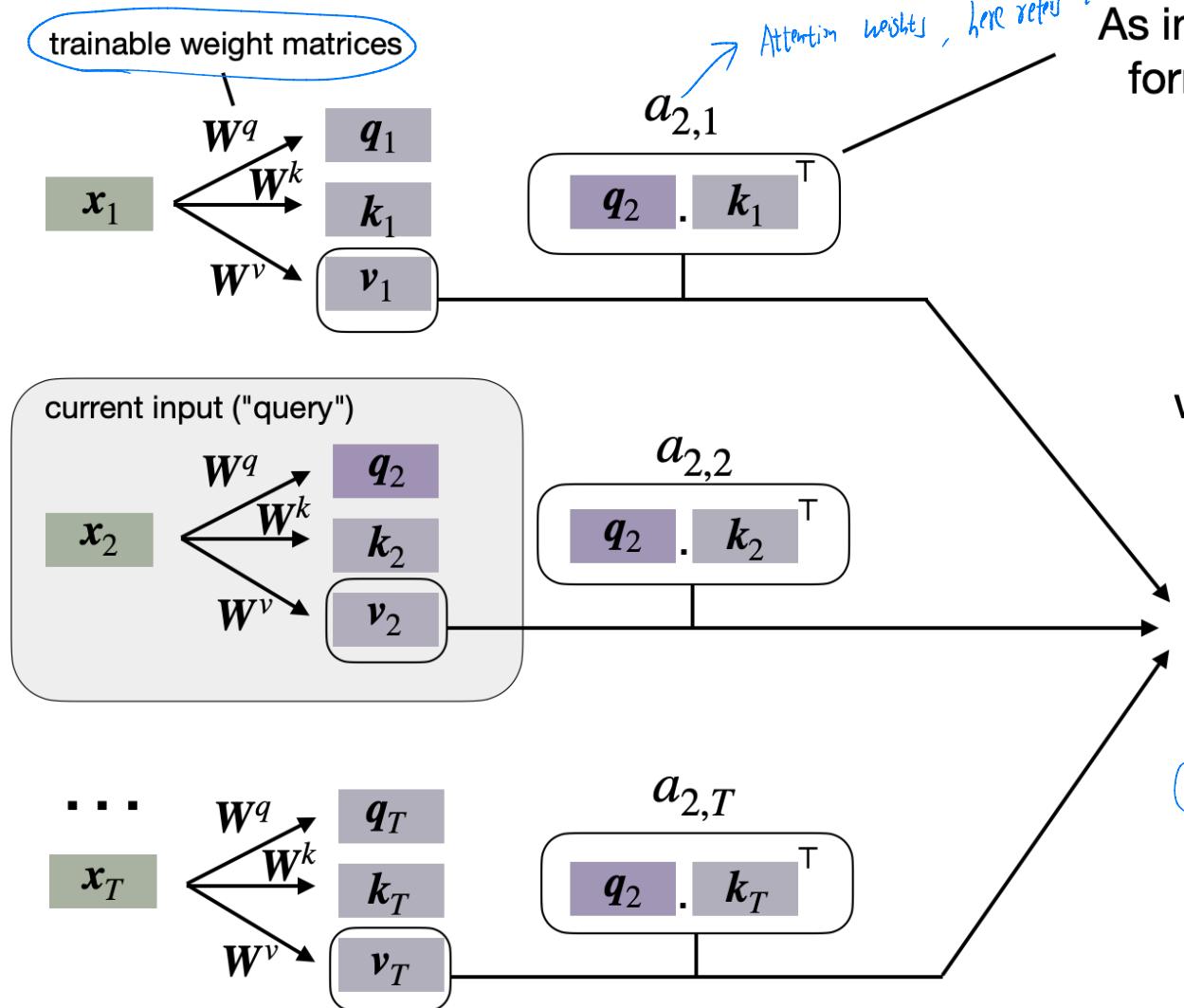
$$\text{query} = W^q \mathbf{x}_i$$

$$\text{key} = W^k \mathbf{x}_i$$

$$\text{value} = W^v \mathbf{x}_i$$



# Learnable Self-attention



to high match attention  $x_2$  should pay to  $k_1(x_1)$

As in the simplified version, this is a form of similarity or compatibility measure ("multiplicative attention")

For each query, model learns which **key-value** input it should attend to

$A_2$

$$A(q_2, K, V) = \sum_{i=1}^T \left[ \frac{\exp(q_2 \cdot k_i^T)}{\sum_j \exp(q_2 \cdot k_j^T)} \times v_i \right]$$

(信息, 被重要性  
挑选出用精华信息)

softmax

weighted sum: values weighted by attention weight (softmax score)

# At the end of the day...

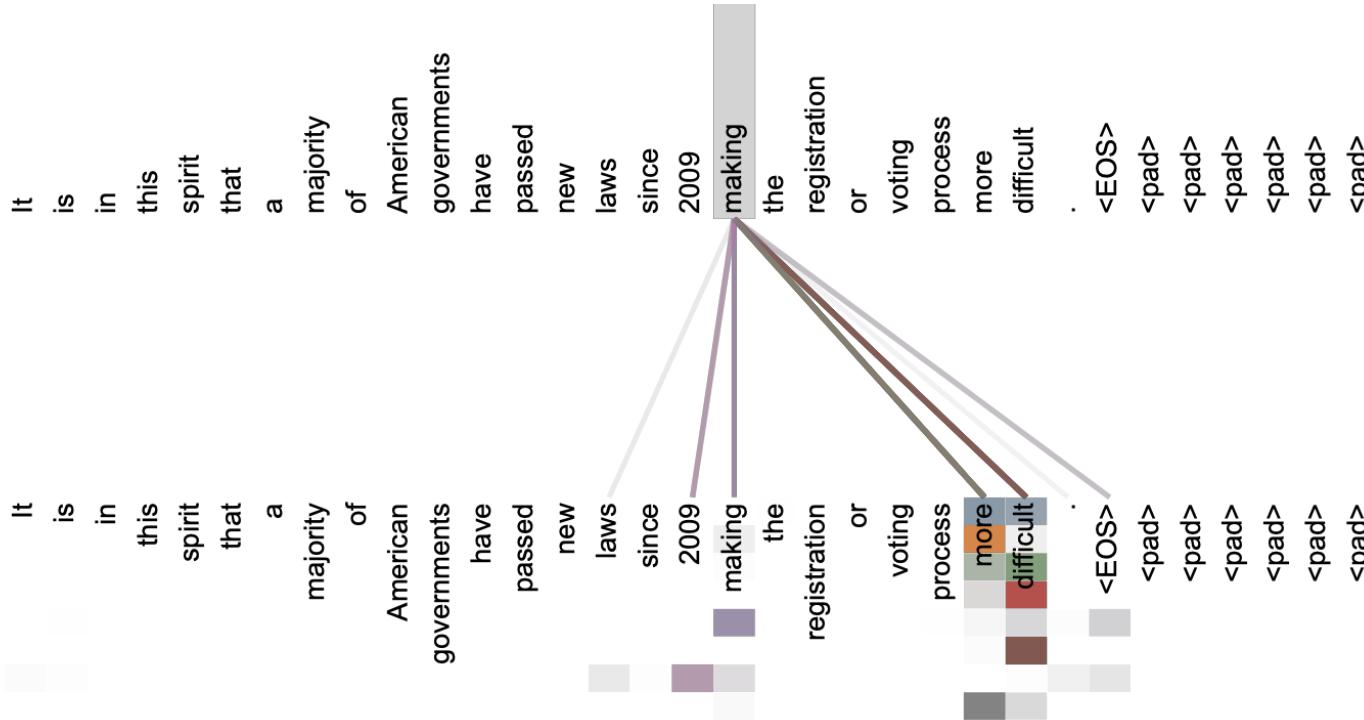


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb ‘making’, completing the phrase ‘making...more difficult’. Attentions here shown only for the word ‘making’. Different colors represent different heads. Best viewed in color.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

# The "Transformer"

## Attention Is All You Need

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaiser@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

### Attention is all you need

[A Vaswani, N Shazeer, N Parmar... - Advances in neural ...](#), 2017 - [proceedings.neurips.cc](#)

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent  
... **We implement** this inside of scaled dot-product **attention** by masking out (setting to  $-\infty$ ) ...

**☆ Save** **99 Cite** **Cited by 174852** **Related articles** **All 73 versions** **»**

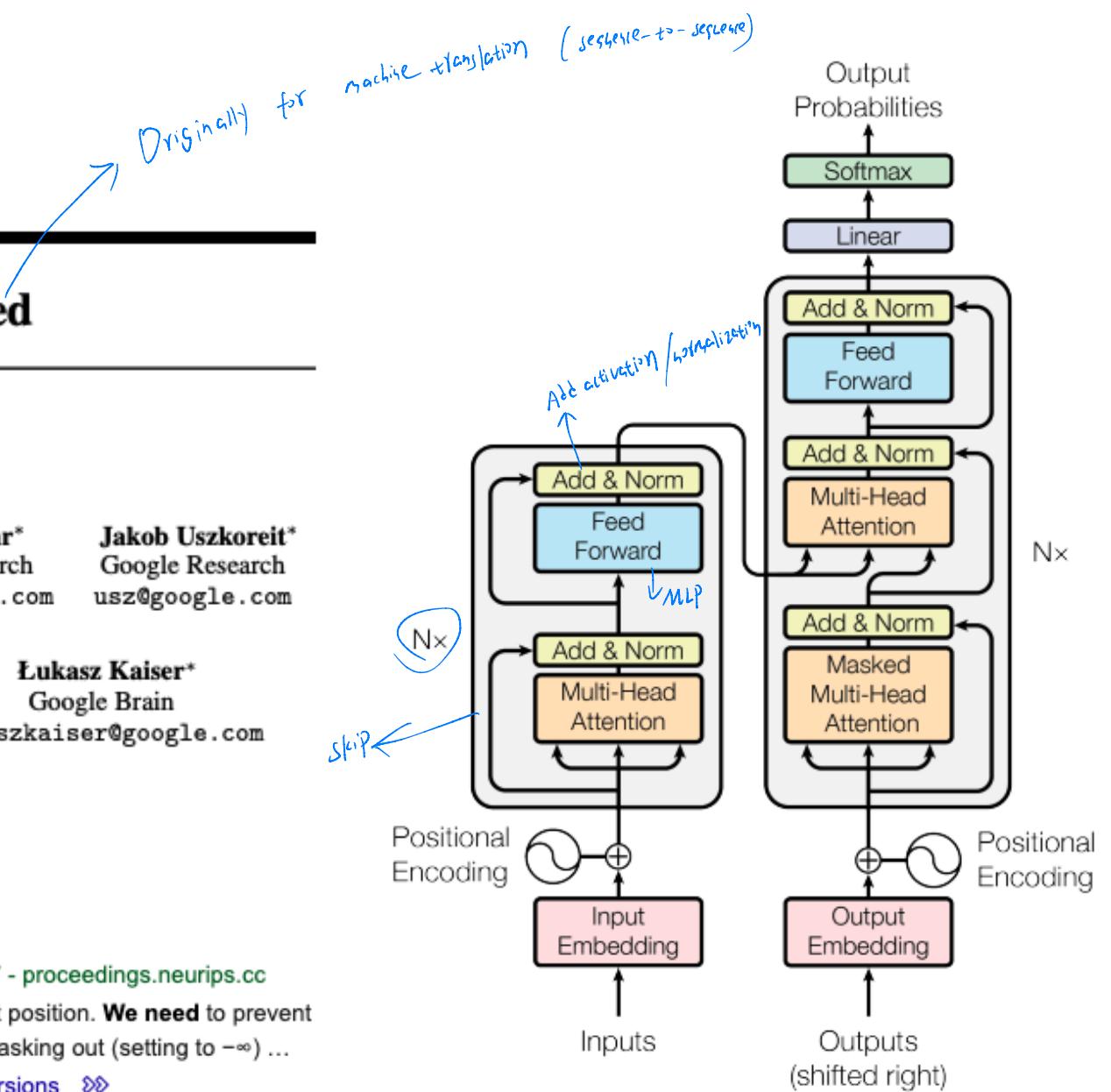
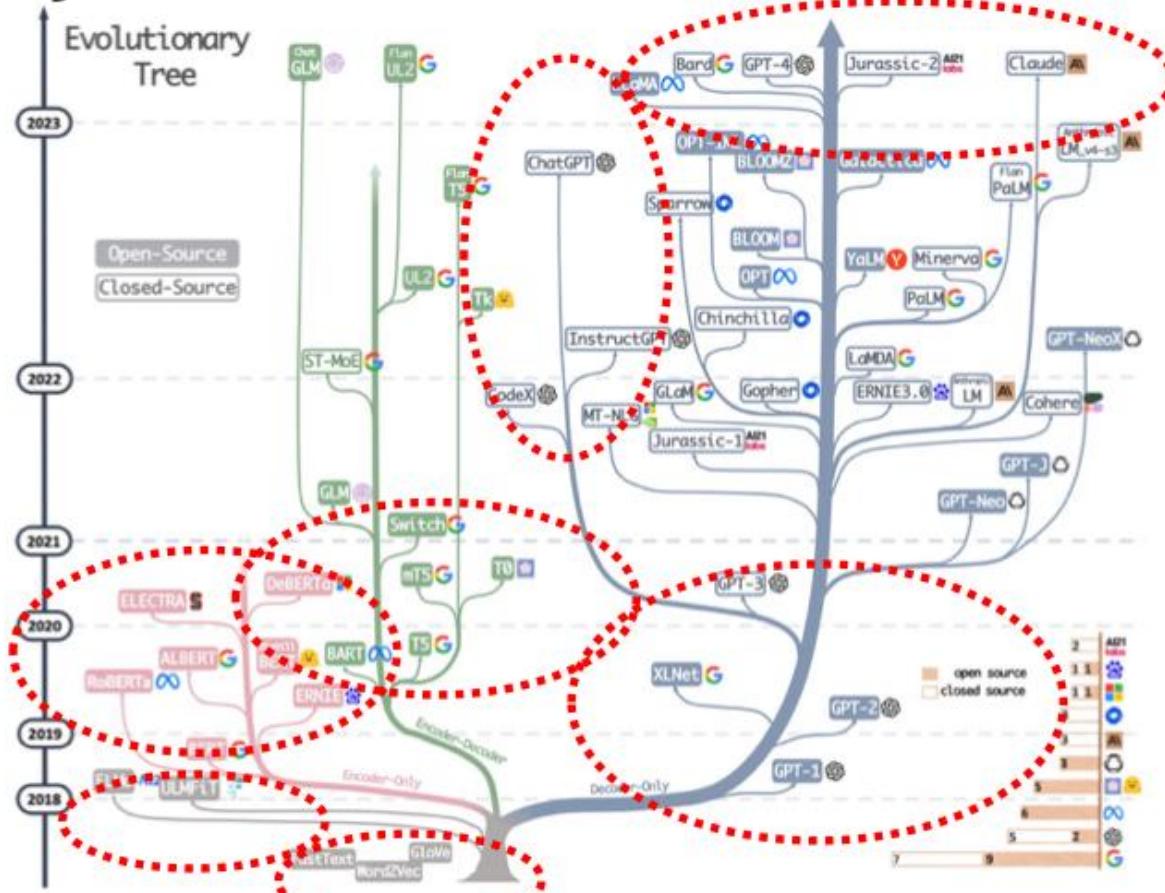


Figure 1: The Transformer - model architecture.

# Many variants of transformer architectures

Many variants: encoder, decoder, or both



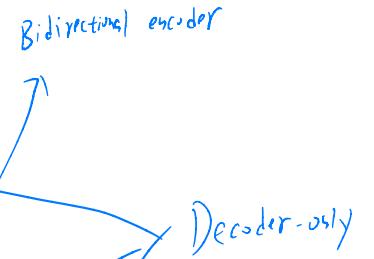
Yang et al., 2023, arxiv: 2304.13712



# Many variants of transformer architectures

Some early transformer architectures

- **GPT-v1:** Generative Pre-Trained Transformer
- **BERT:** Bidirectional Encoder Representations from Transformers
- **GPT-v2:** Language Models are Unsupervised Multitask Learners
- **GPT-v3:** Language Models are Few-Shot Learners
- **BART:** Combining Bidirectional and Auto-Regressive Transformers
- Closing Words -- The Recent Growth of Language Transformers





# Today: GPT

---

## Generative Pre-Trained Transformer



# From Transformer to GPT



# Recall the "Transformer"

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

Łukasz Kaiser\*  
Google Brain  
lukaszkaiser@google.com

Illia Polosukhin\* ‡  
illia.polosukhin@gmail.com

### Attention is all you need

[A Vaswani, N Shazeer, N Parmar... - Advances in neural ...](#), 2017 - [proceedings.neurips.cc](#)

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent  
... **We implement** this inside of scaled dot-product **attention** by masking out (setting to  $-\infty$ ) ...

☆ Save 99 Cite Cited by 174852 Related articles All 73 versions ➔

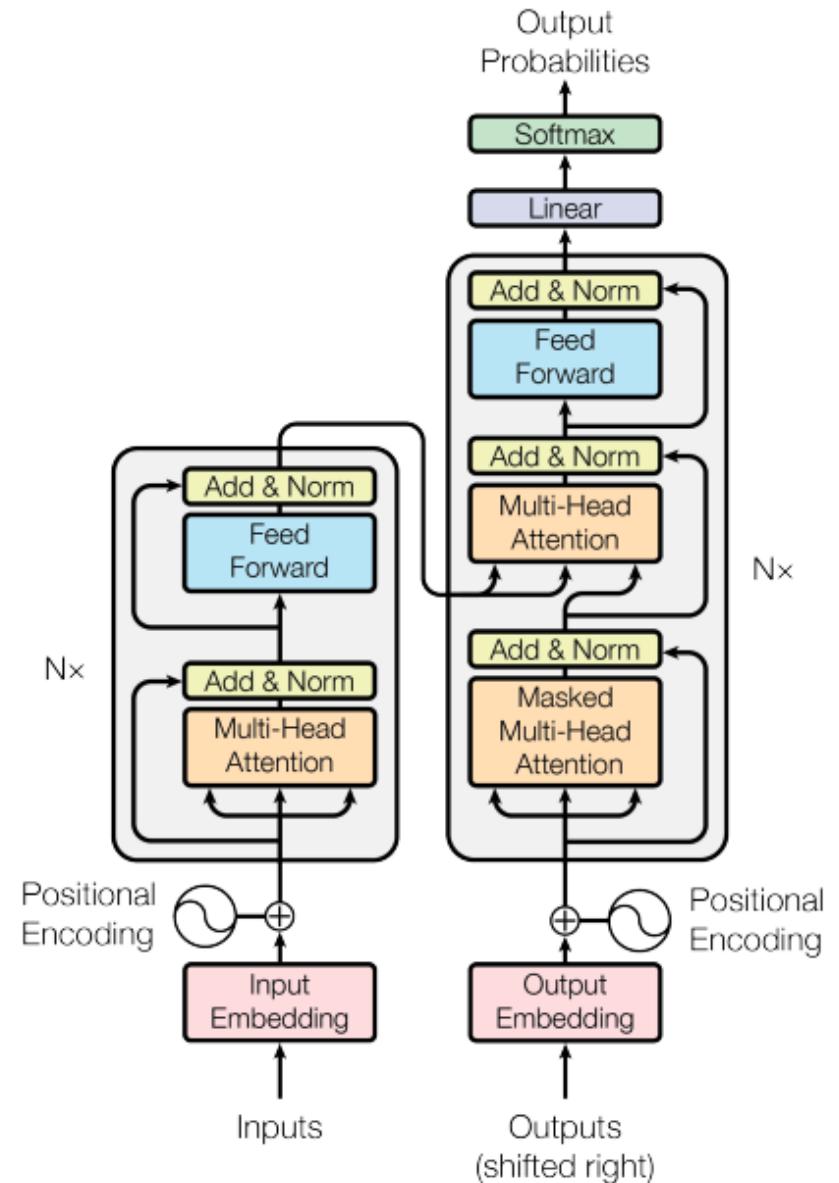


Figure 1: The Transformer - model architecture.

# Recall the "Transformer"

- **Original Transformer** (Vaswani et al., 2017):
  - Encoder-decoder architecture for sequence-to-sequence tasks
  - Parallelizable self-attention instead of recurrence
  - Positional encodings enable order sensitivity
- **Encoder:** Processes input sequence
- **Decoder:** Generates output sequence using masked attention + encoder output
- Inspired by machine translation (observe full input sequence, predict full output sequence)

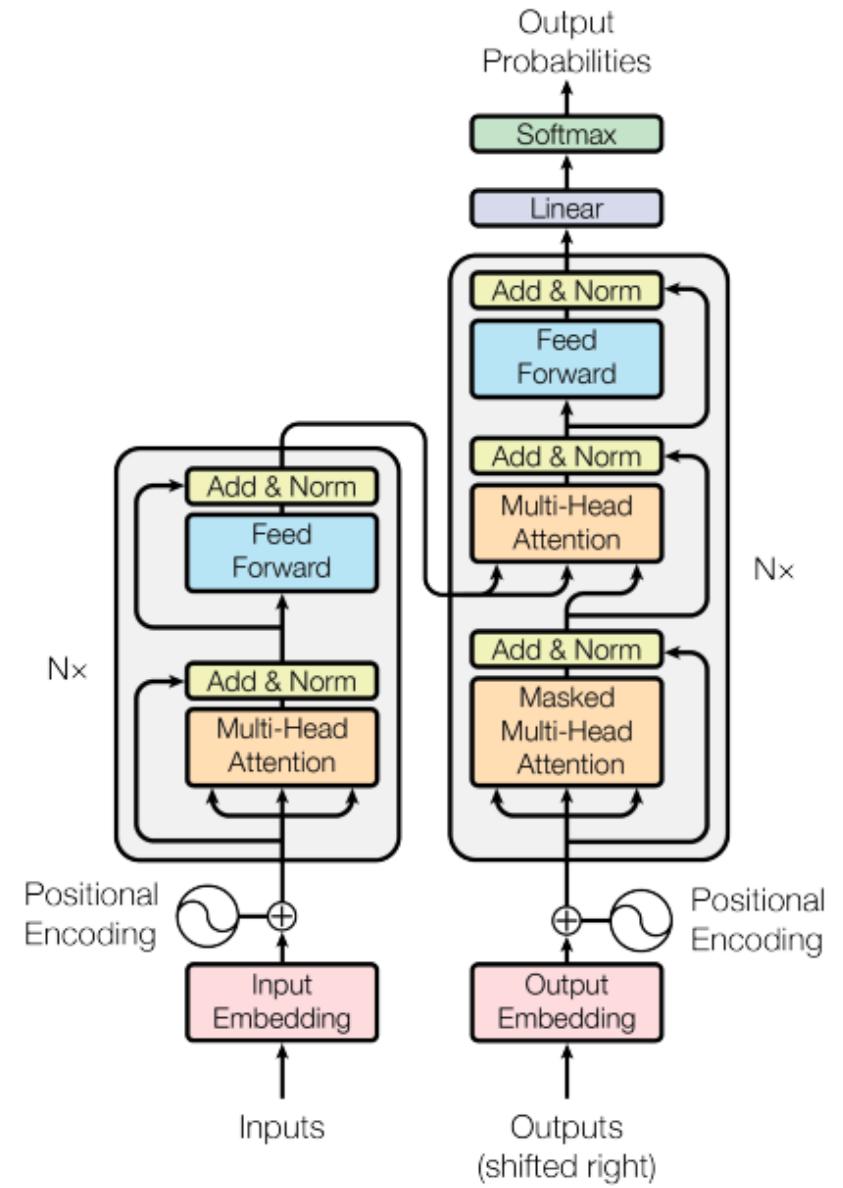


Figure 1: The Transformer - model architecture.



# From Sequence Transduction to Sequence Modeling

- Original Transformer (Vaswani et al., 2017):

$$P(Y | X) = \prod_t P(Y_t | Y_{<t}, X)$$

(翻译任务:  $X \rightarrow Y$ )  
从预训练阶段外部输入又操作  
(需要 Encoder)

- Conditional sequence model for tasks like translation (input → output)
- Generative Pretrained Transformer (GPT) Models:

$$P(X) = \prod_t P(X_t | X_{<t})$$

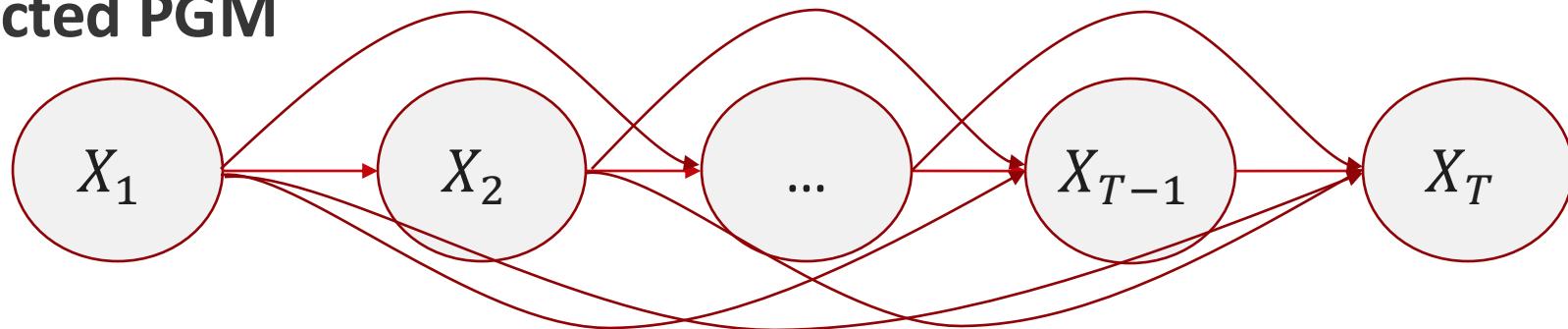
All the past  $X$

- Unconditional generative model over raw text
- Architectural consequence: no encoder, only a decoder with causal structure

Masking is done

# GPT = Probabilistic Model + Transformer Decoder

- Directed PGM



$$P_\theta(X) = \prod_i \prod_t P_\theta(X_{i,t} | X_{i,<t})$$

第 i 行里的第七个词  
→ *i* is the *i*th sample

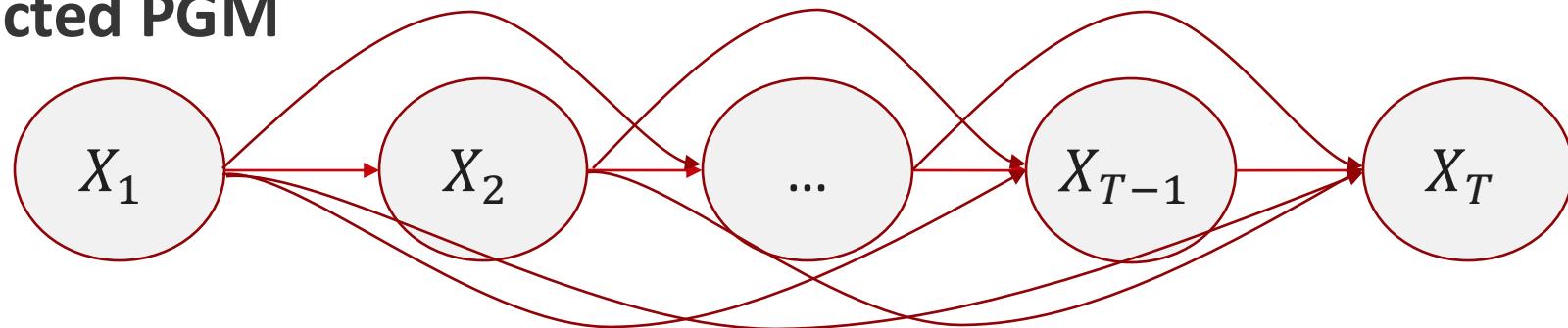
- **Probabilistic objective:** Max log-likelihood of observed seqs

$$\max_{\theta} \sum_i \sum_t \log P_\theta(X_{i,t} | X_{i,<t})$$

[Radford et al., [Improving Language Understanding by Generative Pre-Training](#)]

# GPT = Probabilistic Model + Transformer Decoder

- Directed PGM



$$P_{\theta}(X) = \prod_i \prod_t P_{\theta}(X_{i,t} \mid X_{i,<t})$$

- Model structure:

- Input: token embeddings + positional encodings
- Masked multi-head attention: Enforces “causality”
- Decoder stack: Learns  $P(X_t \mid X_{<t})$
- Output: softmax over vocabulary

We should tell the model where each input is in the sequence (because we do not have pun)  
But this is unnecessary

[Radford et al., [Improving Language Understanding by Generative Pre-Training](#)]

## 1. Masked (掩码): 防止作弊的“挡板”

核心目的: 强制遵守 “因果律” (Causality)。通俗解释: 考试时不准偷看答案。

- 问题出在哪? 回顾你第一张图 (Self-attention),  $x_2$  在计算时, 是可以看到  $x_T$  (句子结尾) 的信息的。
  - 在 **BERT/Encoder** 里, 这没问题, 因为 BERT 的任务是“理解”整句话, 它当然可以看全貌。
  - 但在 **GPT/Decoder** 里, 这是严重的作弊。GPT 的任务是“预测下一个词”。如果我在写  $x_2$  的时候, 已经偷看到了后面还没写出来的  $x_3, x_4$ , 那我就不用学逻辑了, 直接抄答案就行。这在训练时虽然爽, 但在实际应用时 (推理阶段), 后面是空白的, 模型就废了。
- Mask 是怎么工作的? 它是加在 Softmax 步骤之前的一个操作。还记得计算权重的公式吗?

$$\text{Score} = q_2 \cdot k_j^T$$

Mask 机制会强行把所有 “未来位置” 的分数修改为 负无穷大 ( $-\infty$ )。

- 没有 Mask:  $x_2$  对  $x_3$  的关注分可能是 0.8 (作弊)。
- 加上 Mask:  $x_2$  对  $x_3$  的关注分被强制改为  $-\infty$ 。
- 结果: 当这些分数进入 Softmax 公式时:

$$\text{Softmax}(-\infty) = 0$$

这意味着  $x_2$  对  $x_3$  的注意力权重变成了 0。物理意义: 模型“看不见”未来的词, 只能利用  $x_1$  和  $x_2$  自己的信息来推测  $x_3$ 。

视觉上的样子: 这就是为什么在相关论文里, 你会看到一个 “下三角矩阵” (Lower Triangular Matrix)。右上角全是黑的 (被 Mask 掉了)。



# Let's write a GPT



# Let's write a GPT

---

- [EasyGPT repo](#)
  - Adapted from Andrej Karpathy's [nanoGPT](#)



# From our “GPT” to GPT-4



# From our “GPT-0” to GPT-1

- Architecture:
  - Tokenizer: Characters → “Byte-Pair Encoder” tokenizer  
*(For GPT-0, each token is each character)*
  - Activation: ReLU → GELU
  - Weight sharing for embedding / output
  - Scale (117M params):
    - Layers: 4 → 12
    - Attention Heads: 4 → 12
    - Block size (max context): 32 → 512
    - Vocab: 65 → 40000 BPE tokens
    - Embedding dim: 64 → 768
- Training:
  - Dataset: TinyShakespeare (1MB) → BookCorpus (5GB)
  - Initialization & normalization: Default → Carefully tuned
  - Optimizer: Vanilla Adam → Adam + learning rate warmup + weight decay
- Inference:
  - Sampling: Greedy → Top-k



# From GPT-1 to GPT-2

GPT-2: [Radford et al., [Language Models are Unsupervised Multitask Learners](#)]

- Architecture:
  - Scale: Variety of options, with biggest (1.5B params):
    - Layers: 12 → 48
    - Attention Heads: 12 → 25
    - Embedding Dim: 768 → 1600
    - Block size (max context): 512 → 1024
    - Vocab: 40k → 50k tokens
- Training:
  - Dataset: BookCorpus (5GB) → WebText (40GB)

You can reproduce GPT-2 yourself:

<https://github.com/karpathy/nanoGPT>

(takes 4 days to train on an 8xA100 machine)



# From GPT-2 to GPT-3

---

- Architecture:
  - Scale (1.5B → 175B params):
    - Block size (max context): 1024 → 2048
    - Layers: 48 → 96
    - Embedding Dim: 1600 → 12,288
    - Attention Heads: 25 → 96
- Training:
  - Dataset: WebText (40GB) → Common Crawl + books, Wikipedia, code, etc. (~570GB)

GPT-3: [Brown et al., [Language Models are Few-Shot Learners](#)]

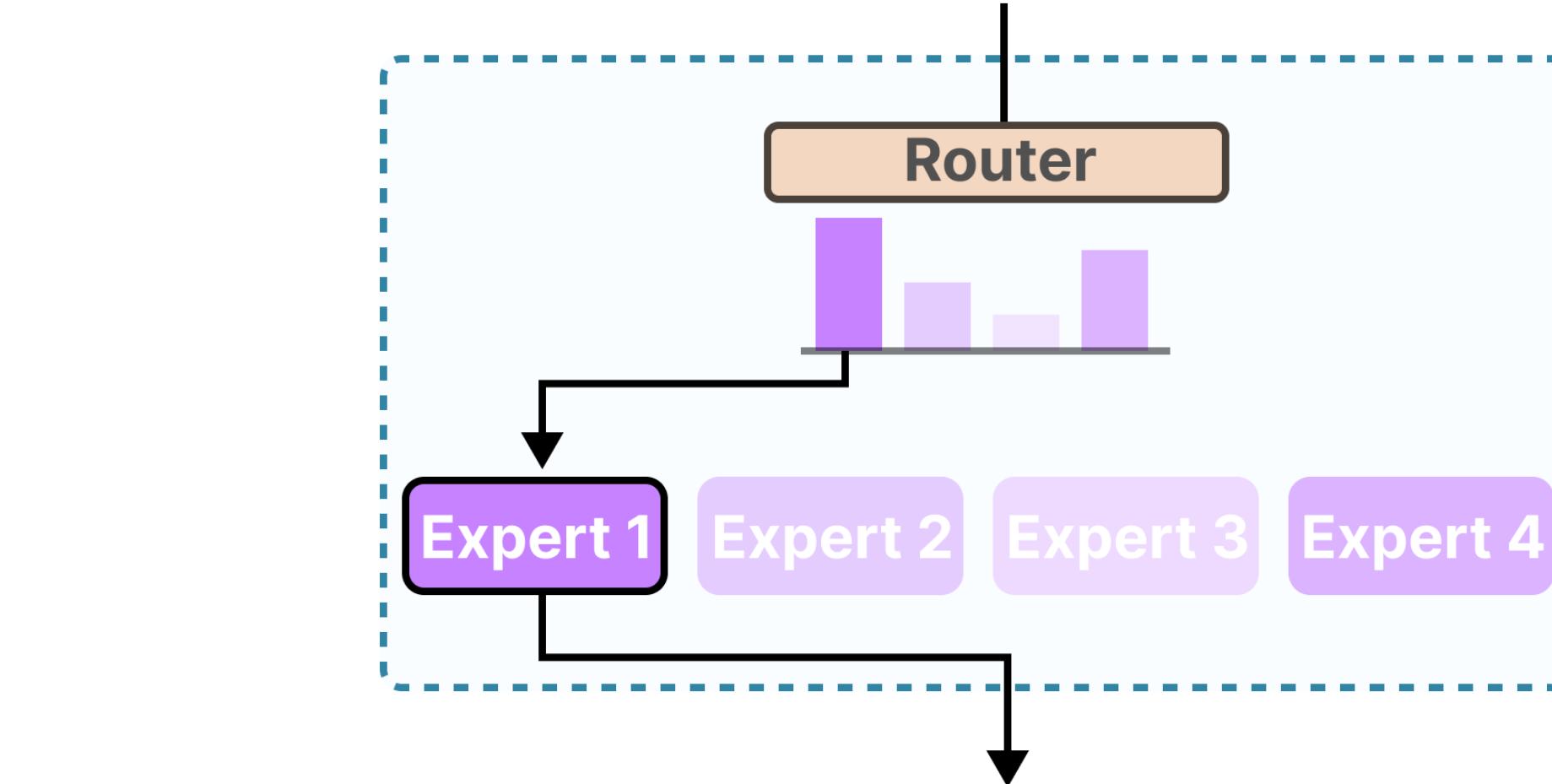


# From GPT-3 to GPT-4?

---

- Architecture:
  - Likely includes MoE
  - Tokenizer: Includes image patches for multi-modal
  - Scale:
    - Total parameters: 175B → Likely >1T
    - Block size (max context): 2048 → 128k
- Training:
  - Dataset: Common Crawl + books, Wikipedia, code, etc. (~570GB) → Larger, undisclosed data training. Reported 13T tokens (~50TB)
  - Alignment: Reinforcement learning + human feedback + system-level “safety”

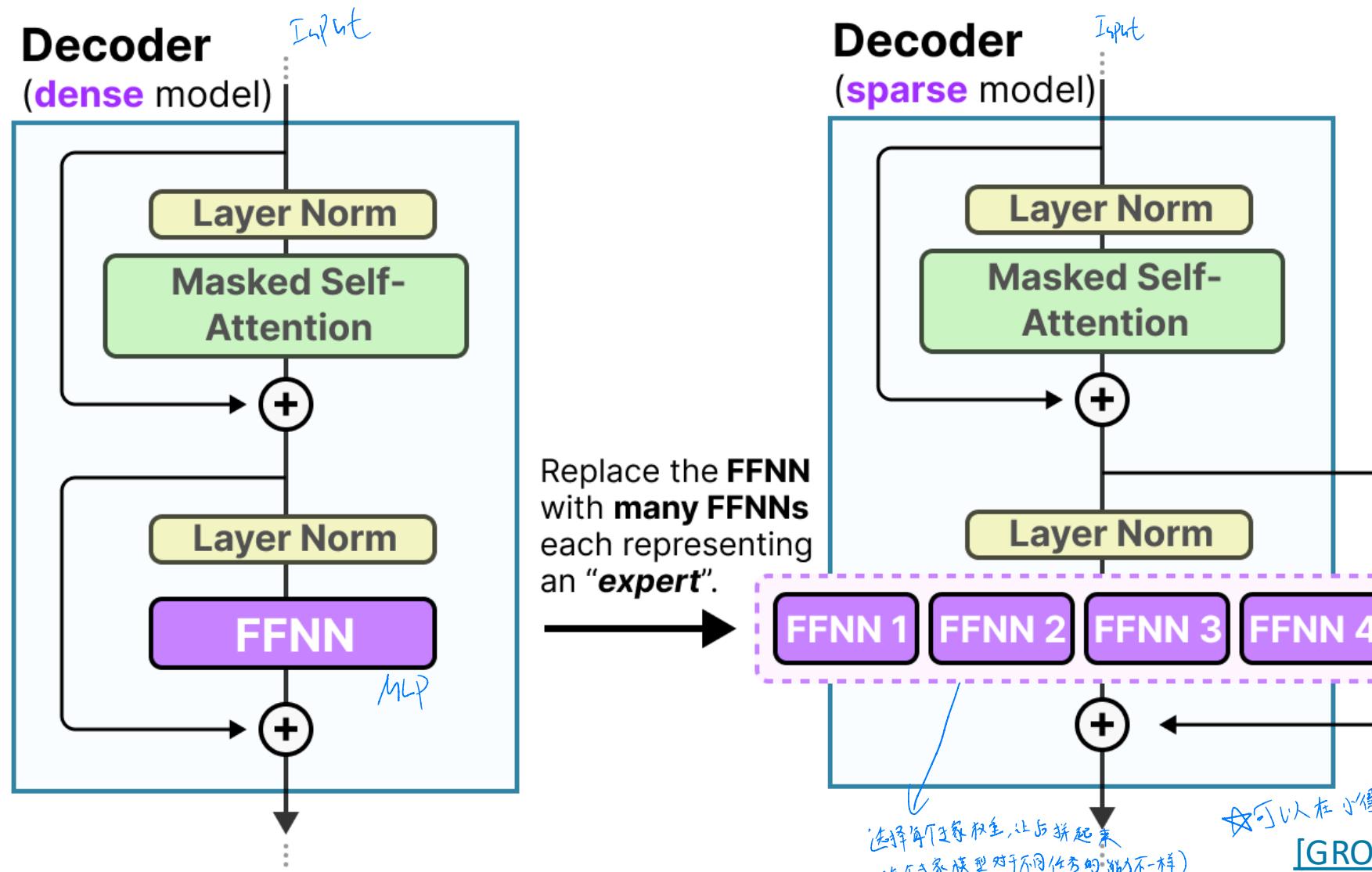
# Mixture of Experts



[GROOTENDORST]



# Mixture of Experts inside Transformer Decoder



# Mixture of Experts: A probabilistic idea

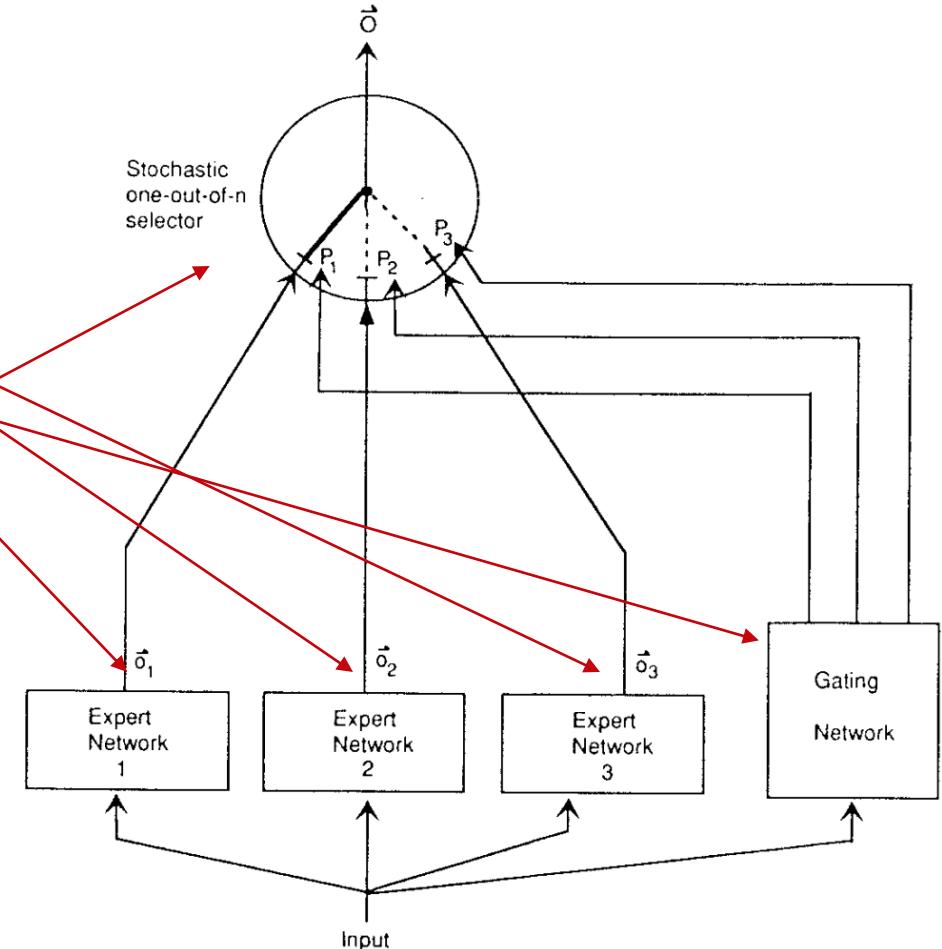
- Let

$$P(Y | X) = \sum_m g_m(X) \cdot P_m(Y | X)$$

gate  
 (Sum of all experts' predictions with their weight)  
 Prediction from expert

- Constrain  $\sum_m g_m(X) = 1$  and  
 $g_m(X) \geq 0 \forall m, X.$

(certain weight) between 0 - 1  
 (softmax)



How would you estimate these parameters?

[[Hierarchical mixtures of experts and the EM algorithm](#), 1993]

["Adaptive Mixtures of Local Experts"](#)  
 Jacobs et al 1991



# Mixture of Experts: Unifies Several Approaches

Let

$$P(Y | X) = \sum_m g_m(X) \cdot P_m(Y | X)$$

- **Mixture of Experts** [[Jacobs et al 1991](#)]:  $g_m(X)$  is a learned gating function.
- **Bagging** [[Breiman 1996](#)]:  $g_m(X) = \frac{1}{M}$  is a constant, uniform weighting.
- **Boosting** [[Freund & Schapire 1997](#)]:  $g_m(X) = \alpha_m$  is a constant per-expert weight.

→ Using  $M$

↓ Just averaging here

↙ (Learned)



# Mixture of Experts: Some analysis

Let

$$P(Y | X) = \sum_m g_m(X) \cdot P_m(Y | X)$$

- Let's examine the mean functions. Define:

$$\bar{f}(x) := E[Y | X] = \sum_m g_m(X) E_m[Y | X] = \sum_m g_m(X) f_m(x)$$

- Let's compare:

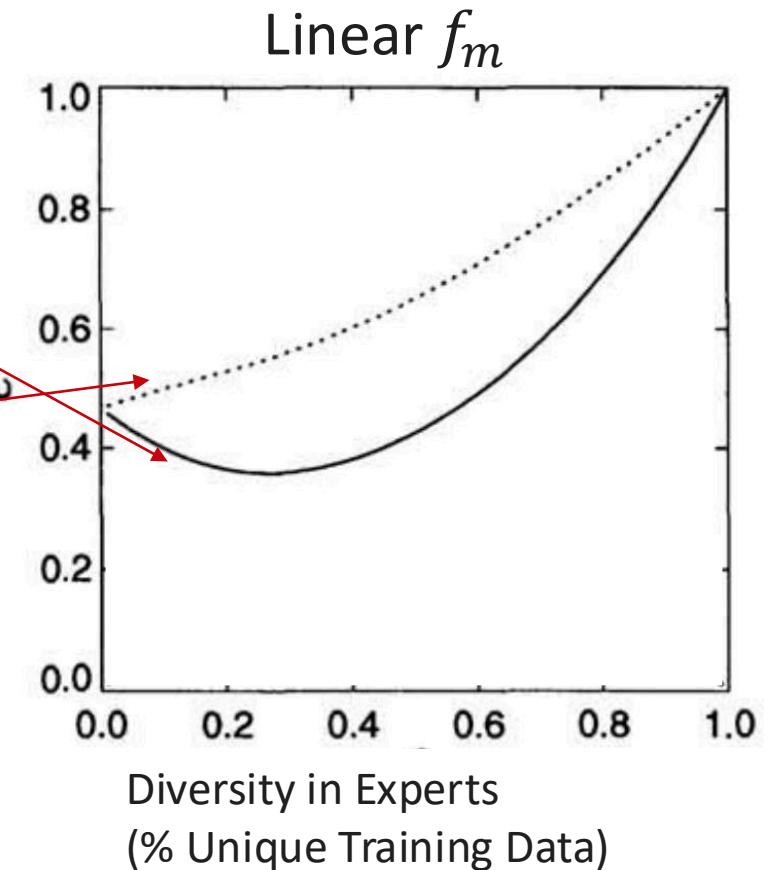
$$\epsilon(x) := (Y - \bar{f}(x))^2 \quad \text{“Ensemble error”}$$

$$\bar{\epsilon}(x) := \frac{1}{M} \sum_m (Y - f_m(x))^2 \quad \text{“Average expert error”}$$

**Will these be minimized for the same ensemble?**

# Mixture of Experts: Some analysis

- $\epsilon(x) := (Y - \bar{f}(x))^2$  ← “Ensemble error”
- $\bar{\epsilon}(x) := \frac{1}{M} \sum_m (Y - f_m(x))^2$  ← “Average expert error”



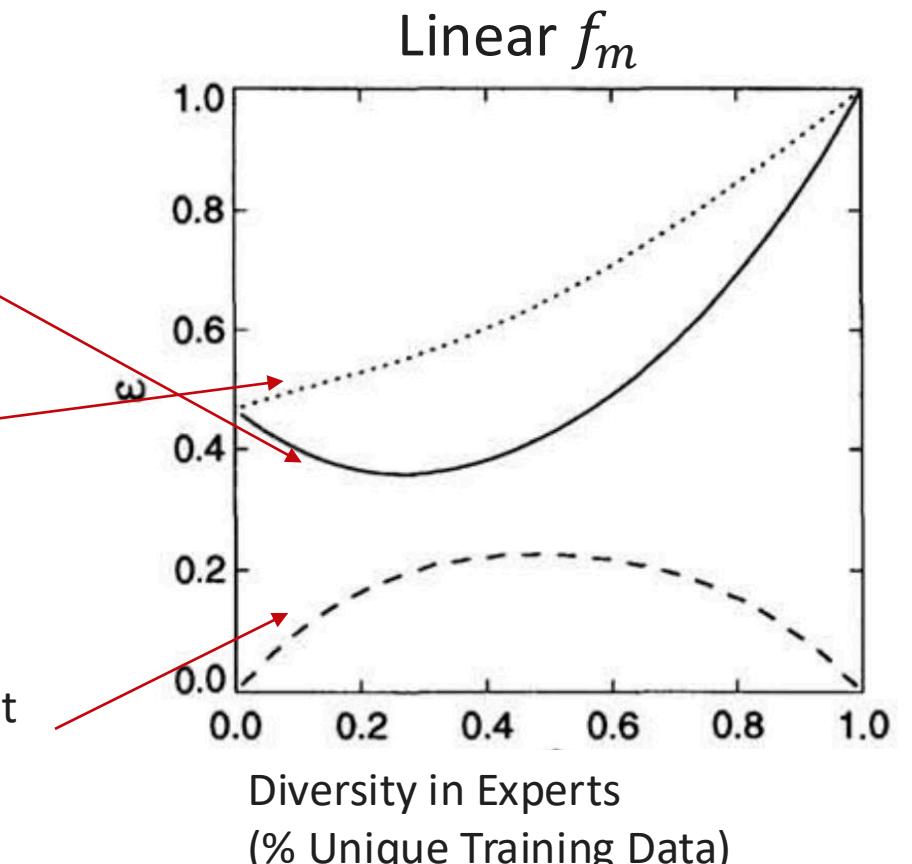
[Sollich & Krogh 1995]

# Mixture of Experts: Some analysis

- $\epsilon(x) := (Y - \bar{f}(x))^2$  ← “Ensemble error”
- $\bar{\epsilon}(x) := \frac{1}{M} \sum_m (Y - f_m(x))^2$  ← “Average expert error”

**Intuition:** Expert errors are wrong in individualized ways and cancel out through consensus.

“Disagreement between experts”



→ Slight “overfitting” of experts helps!

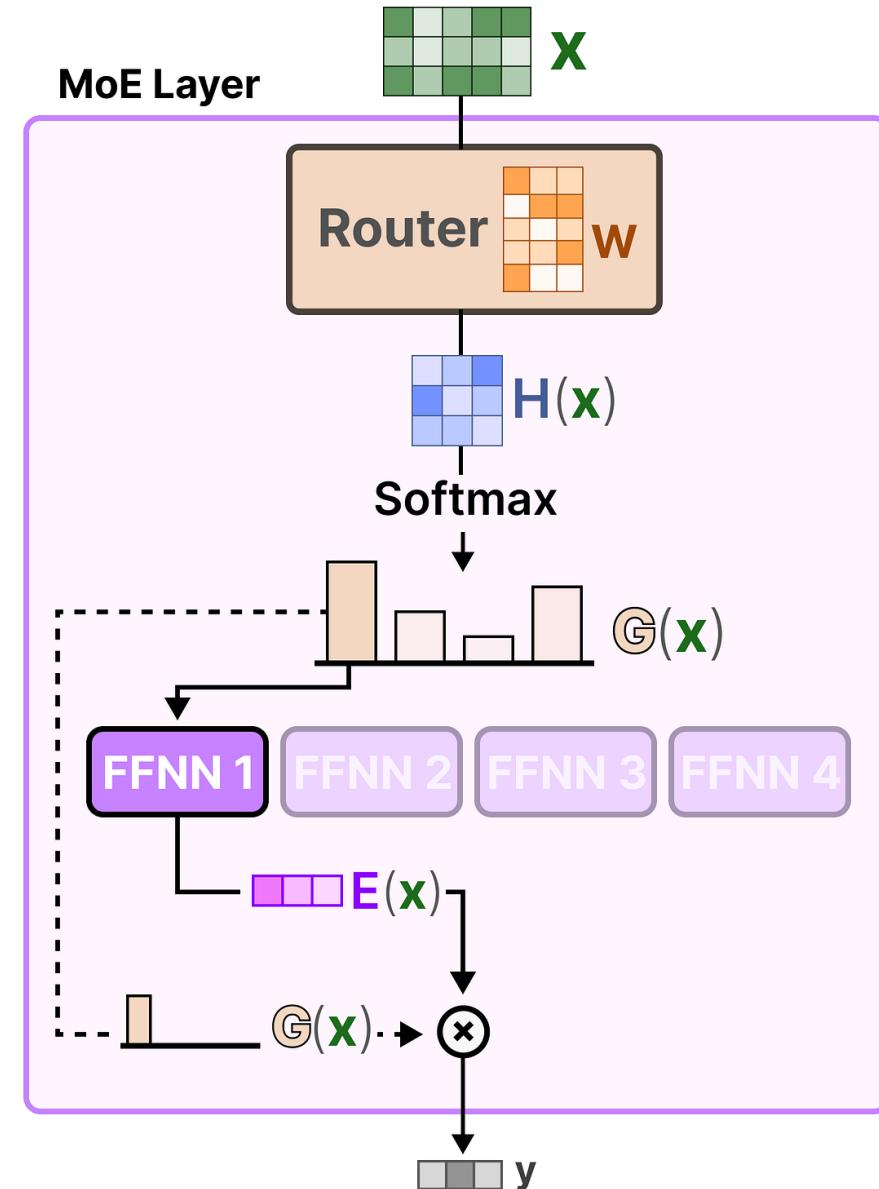
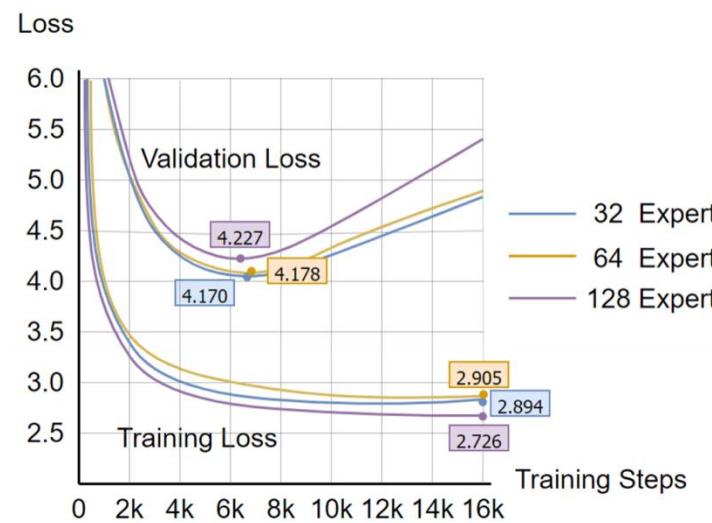
[Sollich & Krogh 1995]



# Mixture of Experts: In LLMs

- Implications for serving?
- Implications for training?

Easy to have experts over-specialize



[GROOTENDORST]



# Summary: From Transformer to GPT

Component	Transformer	GPT
Architecture	Encoder-decoder (full)	<u>Decoder-only</u> <i>Auto-regressive response generation</i>
Attention	Full self-attention	<u>Masked (causal) self-attention</u> <i>↓ Don't want to see the future words</i>
Positional encoding	Sinusoidal (original)	Learned positional embeddings
Output	<u>Task-specific</u>	Next-token prediction
Training objective	Flexible (e.g., <u>translation</u> )	<u>Language modeling (autoregressive)</u>
Inference	<u>Depends on task</u>	Greedy / sampling for text gen



# Summary: From GPT-1 to GPT-4

---

- **Architecture:**

- **Scale:** Variety of options, with biggest (1.5B params → >1T params):
  - Block size (max context): 512 → 128k
  - Layers: 12 → >96
  - Attention Heads: 12 → >96
  - Embedding Dim: 768 → >12,288
  - Vocab: 40k → >50k tokens
- Tokenizer: Includes image patches for multimodal
- **Mixture-of-Experts**

- **Training:**

- Dataset: BookCorpus (5GB) → Private 13T tokens (~50TB)
- Reinforcement learning for alignment

Questions?

