



北京航空航天大學  
BEIHANG UNIVERSITY

# 大语言模型

人工智能研究院

主讲教师 刘偲 沙磊 黄雷 郭晋阳



# 教学团队



**刘偲** 教授 国家优青  
新加坡国立大学博士后

## 课程负责人 具身智能大模型单元 授课教师

带领团队开展**无人机具身大模型**相关研究，建立第一视角下的具身认知系统，为课程内容提供了丰富案例和实践素材



**沙磊** 教授 海外优青  
苹果公司NLP组  
资深研究科学家

## 大语言模型单元 授课教师

带领团队研究**大模型架构优化**、安全性增强、数据增强技术、智能agents 开发等前沿课题



**黄雷** 副教授 青年基金  
阿联酋起源人工  
智能研究院研究员

## 多模态大模型单元 授课教师

带领团队训练并发布了小尺寸**多模态大模型 TinyLLaVA**，相关模型已有 **14 万余次**的下载量



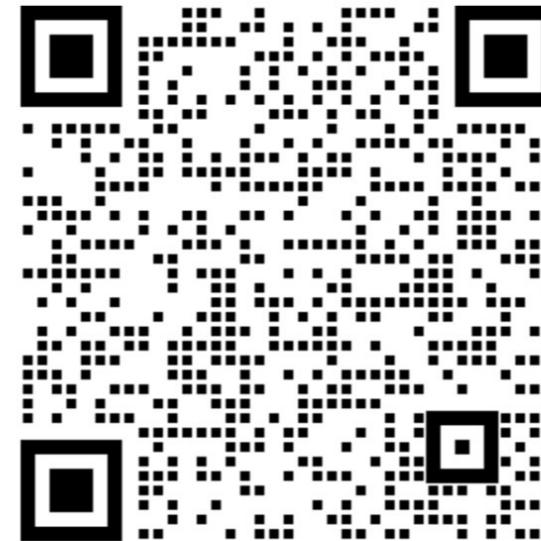
**郭晋阳** 讲师 QM  
悉尼大学博士

## 大模型部署单元 授课教师

带领团队研究**大模型轻量高效部署**，已实现将模型计算量降低 10 倍且性能无损

# 课程信息

- 主页：  
<https://shalei120.github.io/course/LLM/LLM.html>
- 课程群



该二维码 3月 4日前有效，重新进入将更新



# Transformers

# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

# Transformers: Is Attention All We Need?

Spoiler: Not Quite!

---

## Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
[avaswani@google.com](mailto:avaswani@google.com)

**Noam Shazeer\***  
Google Brain  
[noam@google.com](mailto:noam@google.com)

**Niki Parmar\***  
Google Research  
[nikip@google.com](mailto:nikip@google.com)

**Jakob Uszkoreit\***  
Google Research  
[usz@google.com](mailto:usz@google.com)

**Llion Jones\***  
Google Research  
[llion@google.com](mailto:llion@google.com)

**Aidan N. Gomez\* †**  
University of Toronto  
[aidan@cs.toronto.edu](mailto:aidan@cs.toronto.edu)

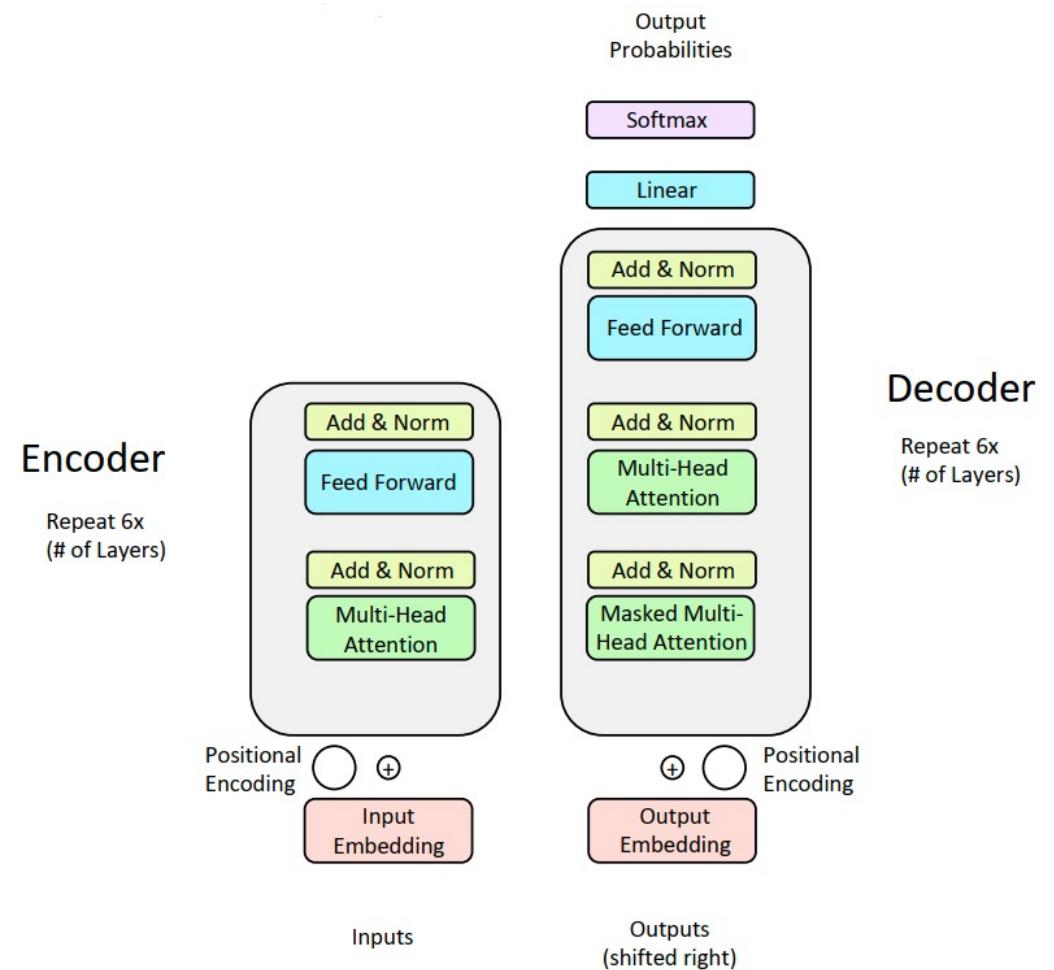
**Łukasz Kaiser\***  
Google Brain  
[lukaszkaiser@google.com](mailto:lukaszkaiser@google.com)

**Illia Polosukhin\* ‡**  
[illia.polosukhin@gmail.com](mailto:illia.polosukhin@gmail.com)

# Transformers Have Revolutionized the Field of NLP



Courtesy of Paramount Pictures



# Great Results with Transformers: Machine Translation

- First, Machine Translation results from the original Transformers paper!

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		<b><math>3.3 \cdot 10^{18}</math></b>
Transformer (big)	<b>28.4</b>	<b>41.8</b>		$2.3 \cdot 10^{19}$

# Great Results with Transformers: Document Generation

- Next, document generation!
- (For perplexity, lower is better; for ROUGE-L, higher is better.)

Model	Test perplexity	ROUGE-L
<i>seq2seq-attention, L = 500</i>	5.04952	12.7
<i>Transformer-ED, L = 500</i>	2.46645	34.2
<i>Transformer-D, L = 4000</i>	2.22216	33.6
<i>Transformer-DMCA, no MoE-layer, L = 11000</i>	2.05159	36.2
<i>Transformer-DMCA, MoE-128, L = 11000</i>	1.92871	37.9
<i>Transformer-DMCA, MoE-256, L = 7500</i>	1.90325	38.8

The old standard from last week!

Transformers dominating across the board.

# Preview: Great Results with (Pre-Trained) Transformers

- Transformers 可以用来预训练
- Transformer 的并行性使得预训练很高效

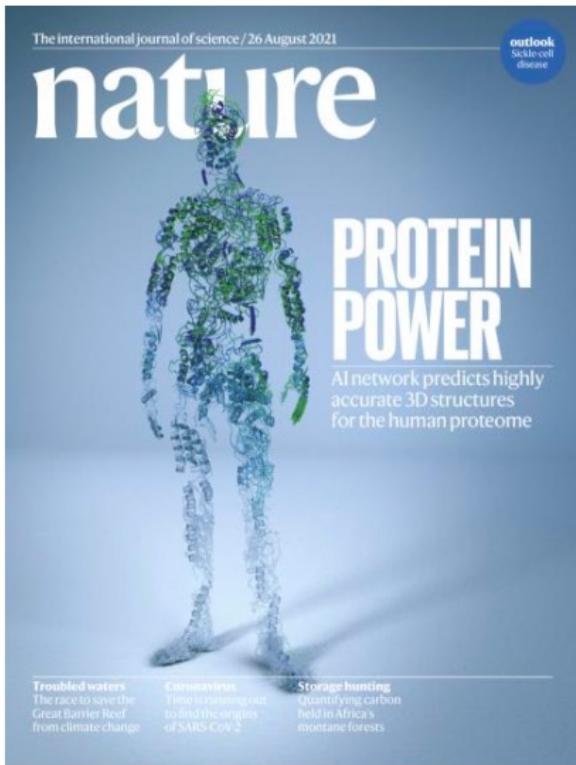


All top models are  
Transformer (and  
pretraining)-based

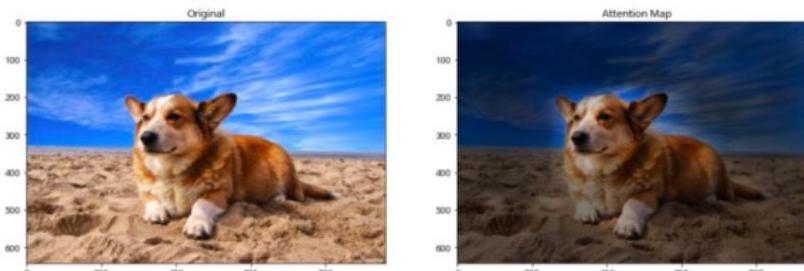
Rank	Name	Model	URL	Score
1	DeBERTa Team - Microsoft	DeBERTa / TuringNLRv4	<a href="#">↗</a>	90.8
2	HFL iFLYTEK	MacALBERT + DKM		90.7
3	+ Alibaba DAMO NLP	StructBERT + TAPT	<a href="#">↗</a>	90.6
4	+ PING-AN Omni-Sinitic	ALBERT + DAAF + NAS		90.6
5	ERNIE Team - Baidu	ERNIE	<a href="#">↗</a>	90.4
6	T5 Team - Google	T5	<a href="#">↗</a>	90.3

# Transformers Even Show Promise Outside of NLP

## Protein Folding



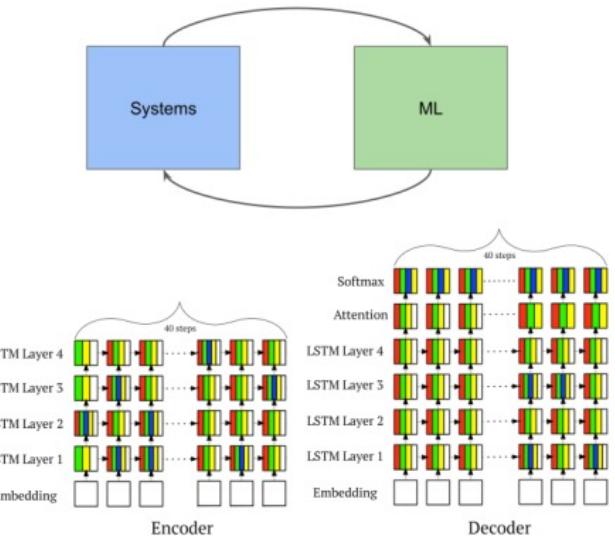
[Jumper et al. 2021] aka AlphaFold2!



## Image Classification

[Dosovitskiy et al. 2020]: Vision Transformer (ViT) outperforms ResNet-based baselines with substantially less compute.

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4 / 88.5*
ImageNet ReaL	<b>90.72</b> ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	<b>99.50</b> ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	<b>94.55</b> ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	<b>97.56</b> ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	<b>99.74</b> ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	<b>77.63</b> ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k



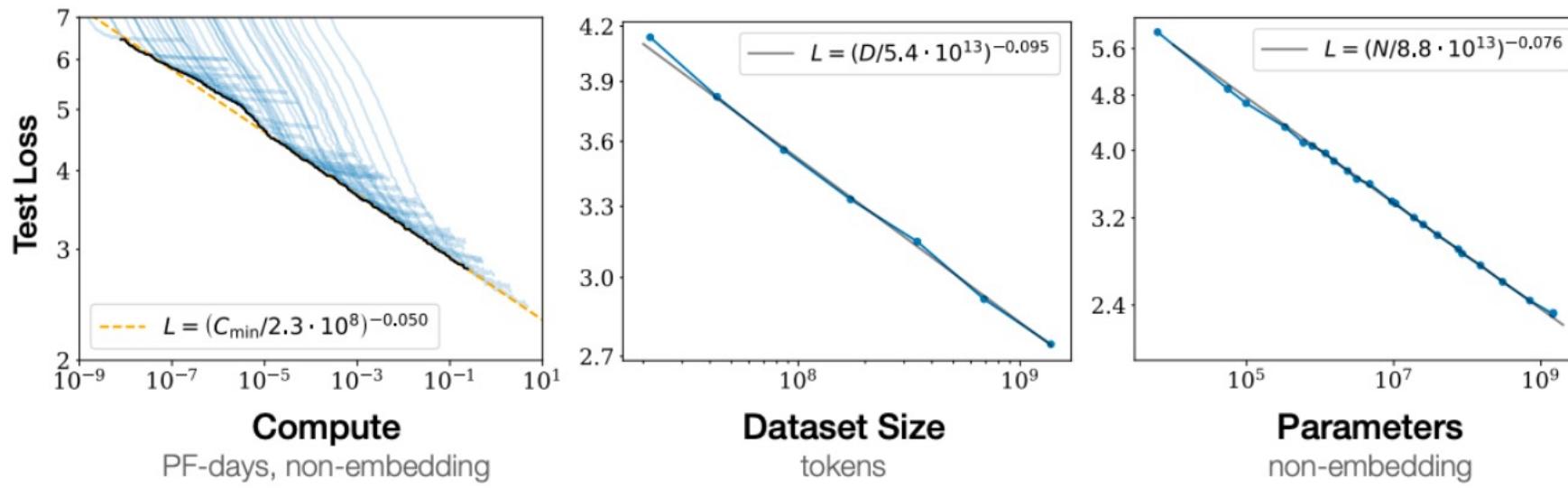
## ML for Systems

[Zhou et al. 2020]: A Transformer-based compiler model (GO-one) speeds up a Transformer model!

Model (#devices)	GO-one (s)	HP (s)	METIS (s)	HDP (s)	Run time speed up over HP / HDP	Search speed up over HDP
2-layer RNNLM (2)	0.173	0.192	0.355	0.191	9.9% / 9.4%	2.95x
4-layer RNNLM (4)	0.210	0.239	0.503	0.251	13.8% / 16.3%	1.76x
8-layer RNNLM (8)	0.320	0.332	0.644	0.764	3.8% / 58.1%	27.8x
2-layer GNMT (2)	0.301	0.384	0.344	0.327	27.6% / 14.3%	30x
4-layer GNMT (4)	0.350	0.469	0.466	0.432	34% / 23.4%	58.8x
8-layer GNMT (8)	0.440	0.562	0.600	0.693	21.7% / 36.5%	7.35x
2-layer Transformer-XL (2)	0.223	0.268	0.37	0.262	20.1% / 17.4%	40x
4-layer Transformer-XL (4)	0.230	0.27	0.40	0.259	17.4% / 12.6%	26.7x
8-layer Transformer-XL (8)	0.350	0.46	0.60	0.425	23.9% / 16.7%	16.7x
Inception (2) b64	0.229	0.312	0.600	0.301	26.6% / 23.9%	13.5x
Inception (2) b64	0.423	0.731	0.600	0.498	42.1% / 29.3%	21.0x
AmoebaNet (4)	0.394	0.44	0.426	0.418	26.1% / 6.1%	58.8x
2-stack 18-layer WaveNet (2)	0.317	0.376	0.600	0.354	18.6% / 11.7%	6.67x
4-stack 36-layer WaveNet (4)	0.659	0.988	0.600	0.721	50% / 9.4%	20x
GEOOMEAN	-	-	-	-	<b>20.5% / 18.2%</b>	<b>15x</b>

# Scaling Laws: Are Transformers All We Need?

- 有了Transformers，语言模型的效果可以跟随模型大小，训练数据规模以及计算资源量的增加而提升
- 如果我们在不改变架构的情况下持续增大这些模型，有无可能最终超越人类？

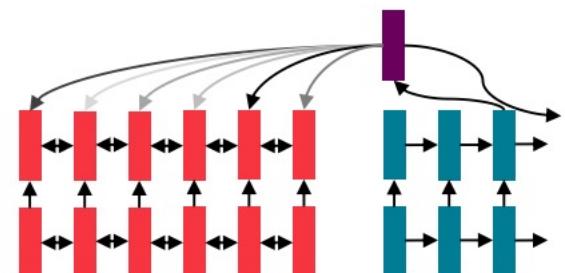
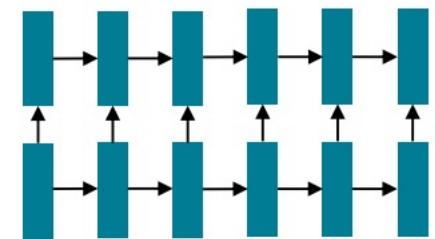
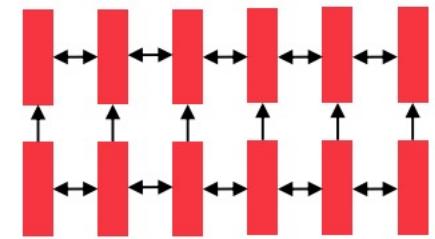


# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - **From Recurrence (RNNs) to Attention-Based NLP Models**
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

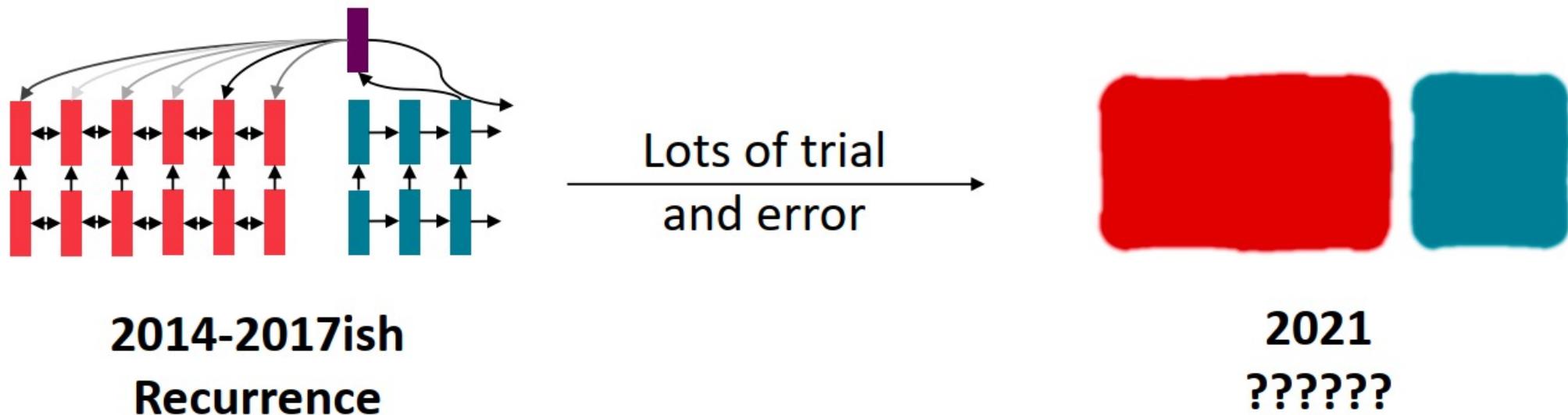
# Recurrent models for (most) NLP!

- Since 2016, the de facto strategy in NLP is to **encode** sentences with a bidirectional LSTM: (for example, the source sentence in a translation)
- Define your output (parse, sentence, summary) as a sequence, and use an LSTM to generate it.
- Use attention to allow flexible access to memory



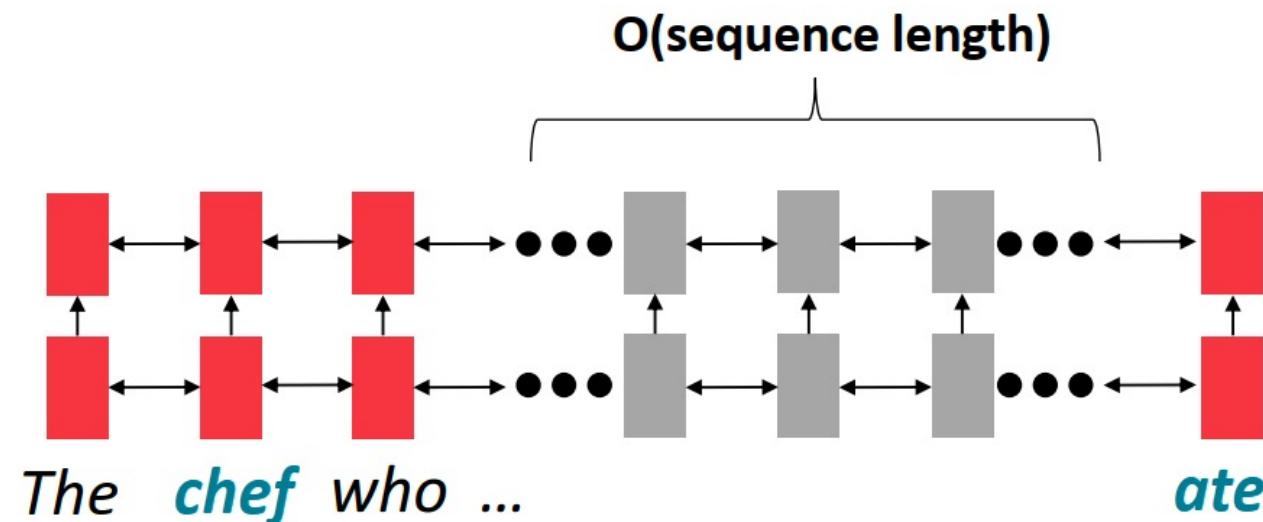
# Today: Same goals, different building blocks

- We have learned about sequence-to-sequence problems and encoder-decoder models.
- Today, we're not trying to motivate entirely new ways of looking at problems (like Machine Translation)
- Instead, we're trying to find the best building blocks to plug into our models and enable broad progress.



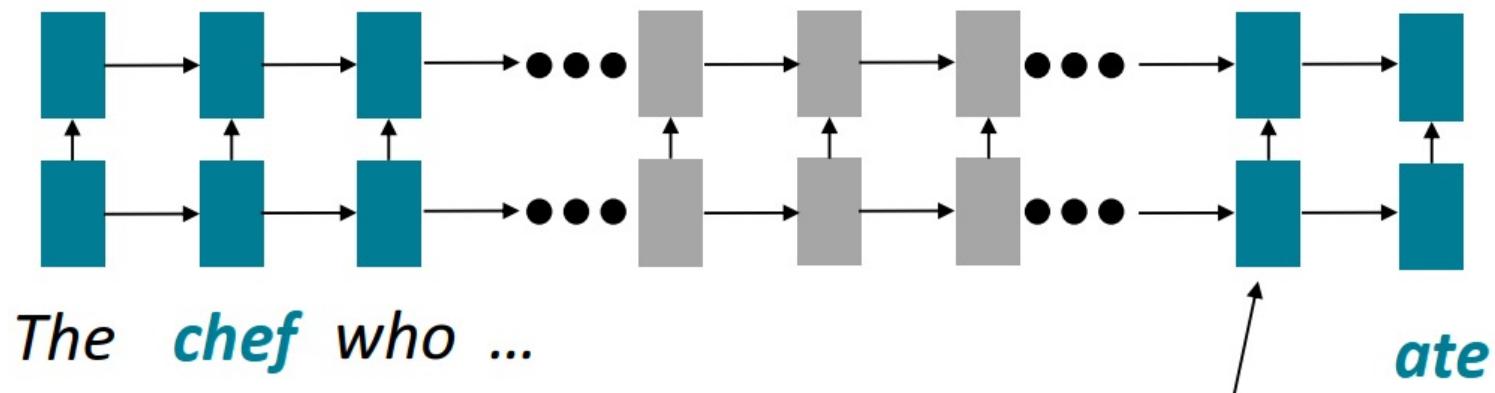
# Issues with recurrent models: Linear interaction distance

- RNNs are unrolled “left-to-right”.
- It encodes linear locality: a useful heuristic!
  - Nearby words often affect each other’s meanings
- Problem: RNNs take **O(sequence length)** steps for distant word pairs to interact.



# Issues with recurrent models: Linear interaction distance

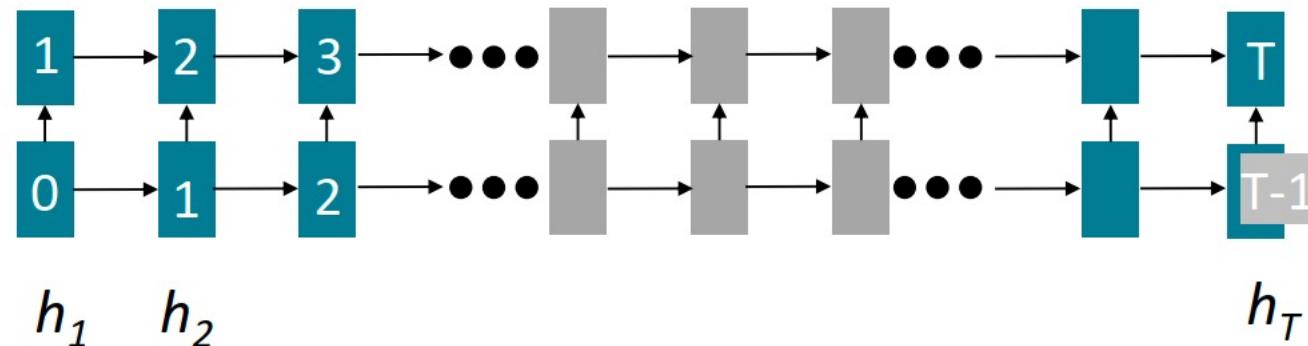
- 远距离的词的交互需要 $O(\text{sequence length})$  步， 意味着：
  - 长距离依存关系很难学习（梯度消失）
  - 词的线性顺序被模型限定



Info of *chef* has gone through  
 $O(\text{sequence length})$  many layers!

# Issues with recurrent models: Lack of parallelizability

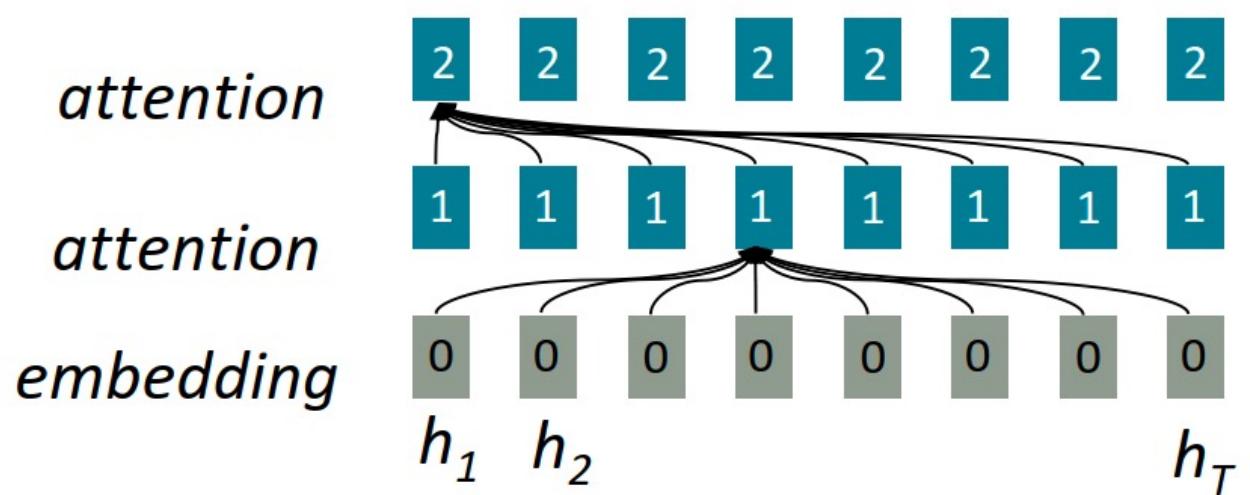
- 前向与后向计算均无法并行实现
  - GPU (或TPU) 可以同时做很多不相关的计算
  - 但未来的RNN隐状态无法在之前的RNN隐状态计算完成之前进行计算
  - 对于大数据集训练很不利
  - 序列长度越长问题越大。 Why? 因为batch就得相应调小



Numbers indicate min # of steps before a state can be computed

## If not recurrence, then what? How about (self) attention?

- 回忆：Attention：将每个词的表示看做query，从一组value中检索信息
  - 我们学过从decoder到encoder的attention
  - **Self-attention**是encoder-encoder或decoder-decoder的attention，每一个词都要关注到其他的词



All words attend  
to all words in  
previous layer;  
most arrows here  
are omitted

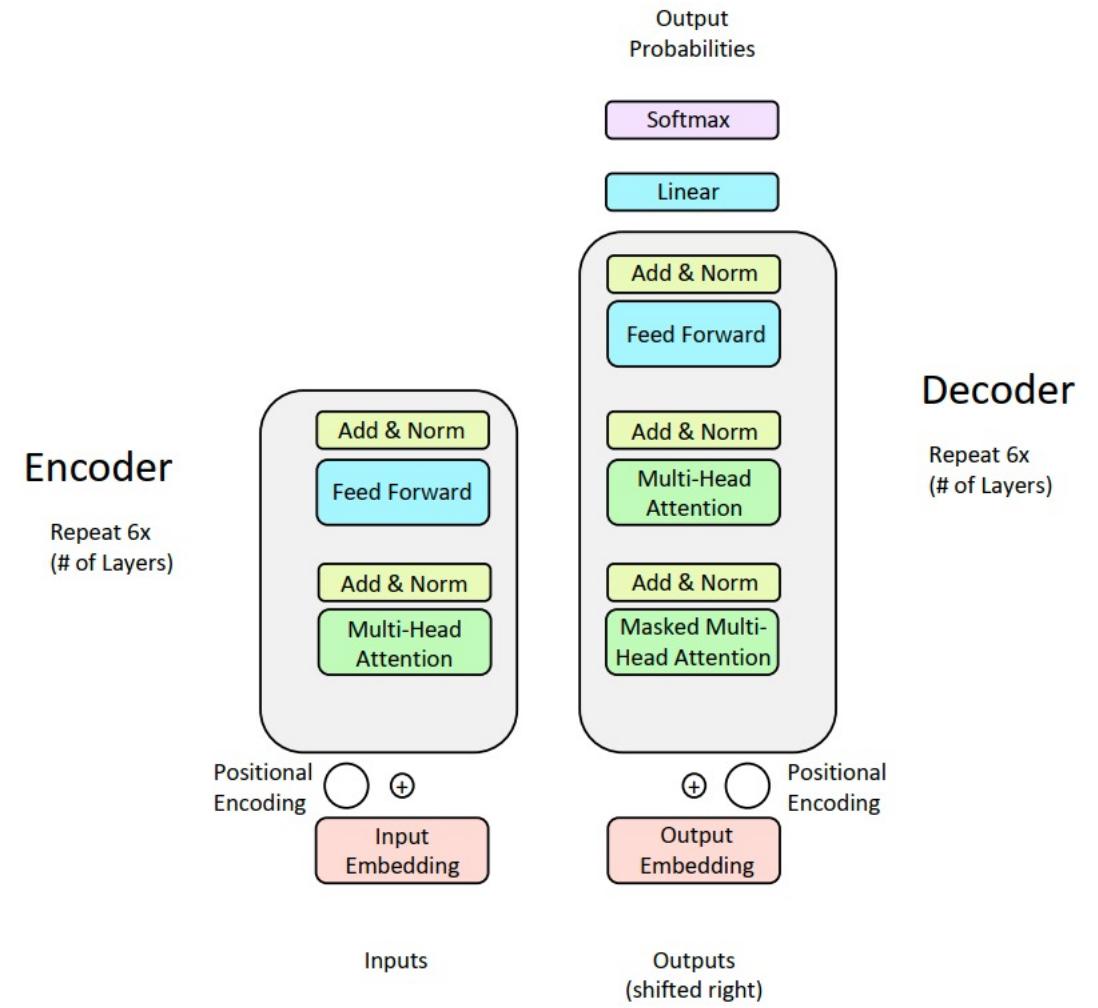
Maximum interact distance  $O(1)!$

# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - **Understanding the Transformer Model**
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

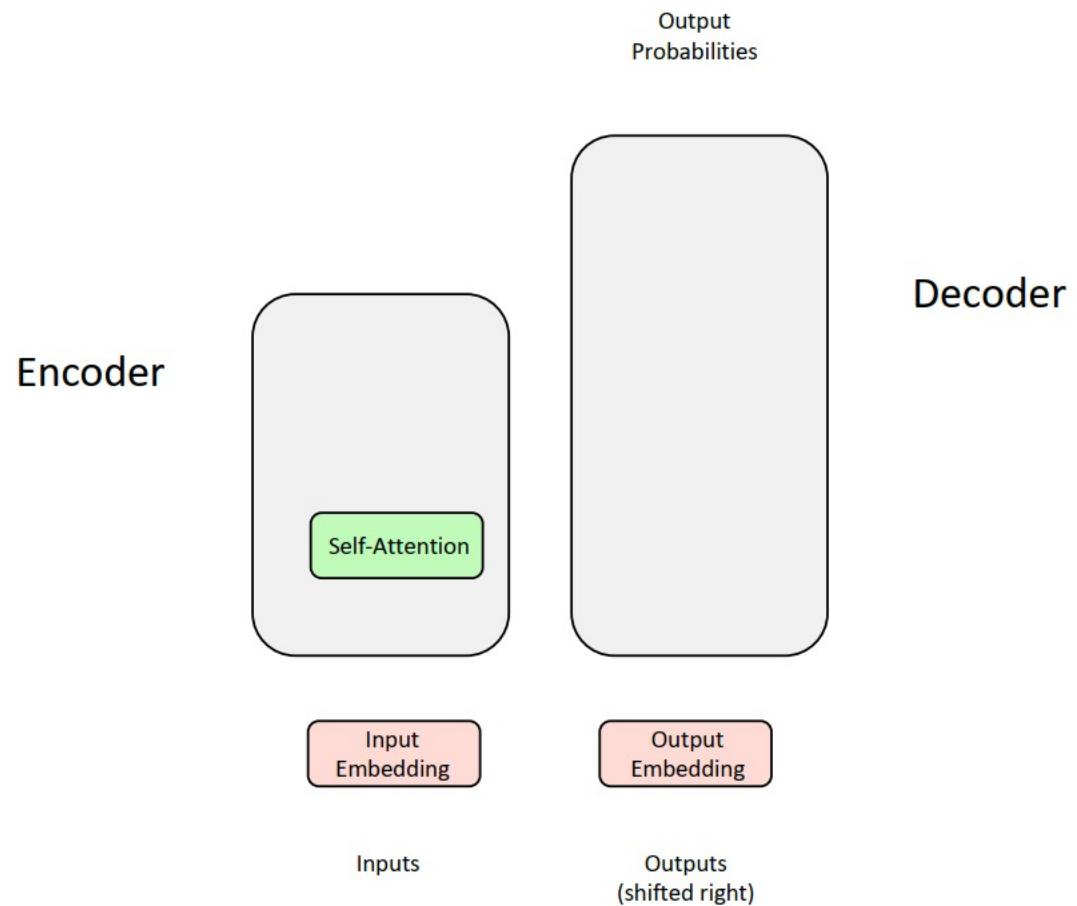
# The Transformer Encoder-Decoder [Vaswani et al., 2017]

- In this section, you will learn exactly how the Transformer architecture works:
  - First, we will talk about the Encoder!
  - Next, we will go through the Decoder (which is quite similar)!



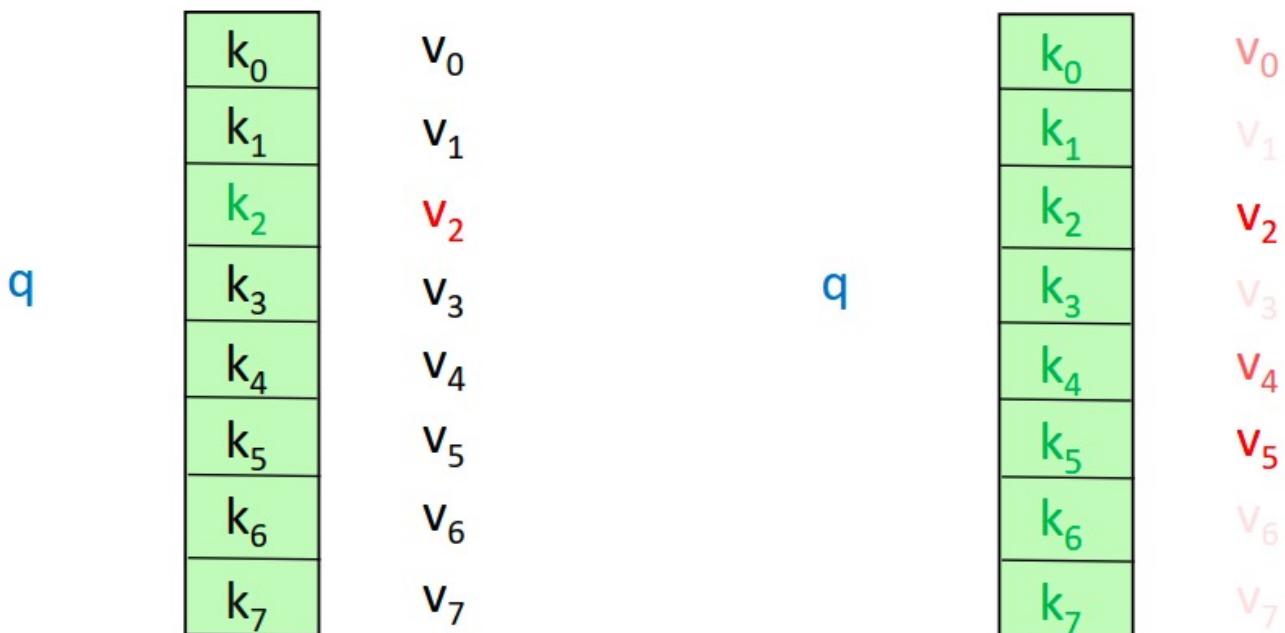
# Encoder: Self-Attention

- Self-Attention is the core building block of Transformer, so let's first focus on that!



# Intuition for Attention Mechanism

- Let's think of attention as a "fuzzy" or approximate hashtable:
  - To look up a **value**, we compare a **query** against **keys** in a table.
  - In a hashtable (shown on the bottom left):
    - Each **query** (hash) maps to exactly one **key-value** pair.
  - In (self-)attention (shown on the bottom right):
    - Each **query** matches each **key** to varying degrees.
    - We return a sum of **values** weighted by the **query-key** match



# Recipe for Self-Attention in the Transformer Encoder

- Step 1: For each word , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$

- Step 2: Calculate attention score between **query** and **keys**.

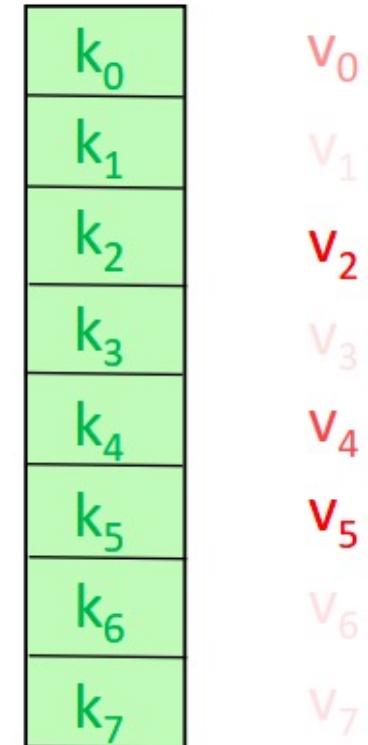
$$e_{ij} = q_i \cdot k_j$$

- Step 3: Take the softmax to normalize attention scores.

$$\alpha_{ij} = softmax(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

- Step 4: Take a weighted sum of **values**

$$Output_i = \sum_j \alpha_{ij} v_j$$



# Recipe for (Vectorized) Self-Attention in the Transformer Encoder

- Step 1: With embeddings stacked in  $X$ , calculate **queries**, **keys**, and **values**.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

- Step 2: Calculate attention scores between **query** and **keys**.

$$E = QK^T$$

- Step 3: Take the softmax to normalize attention scores.

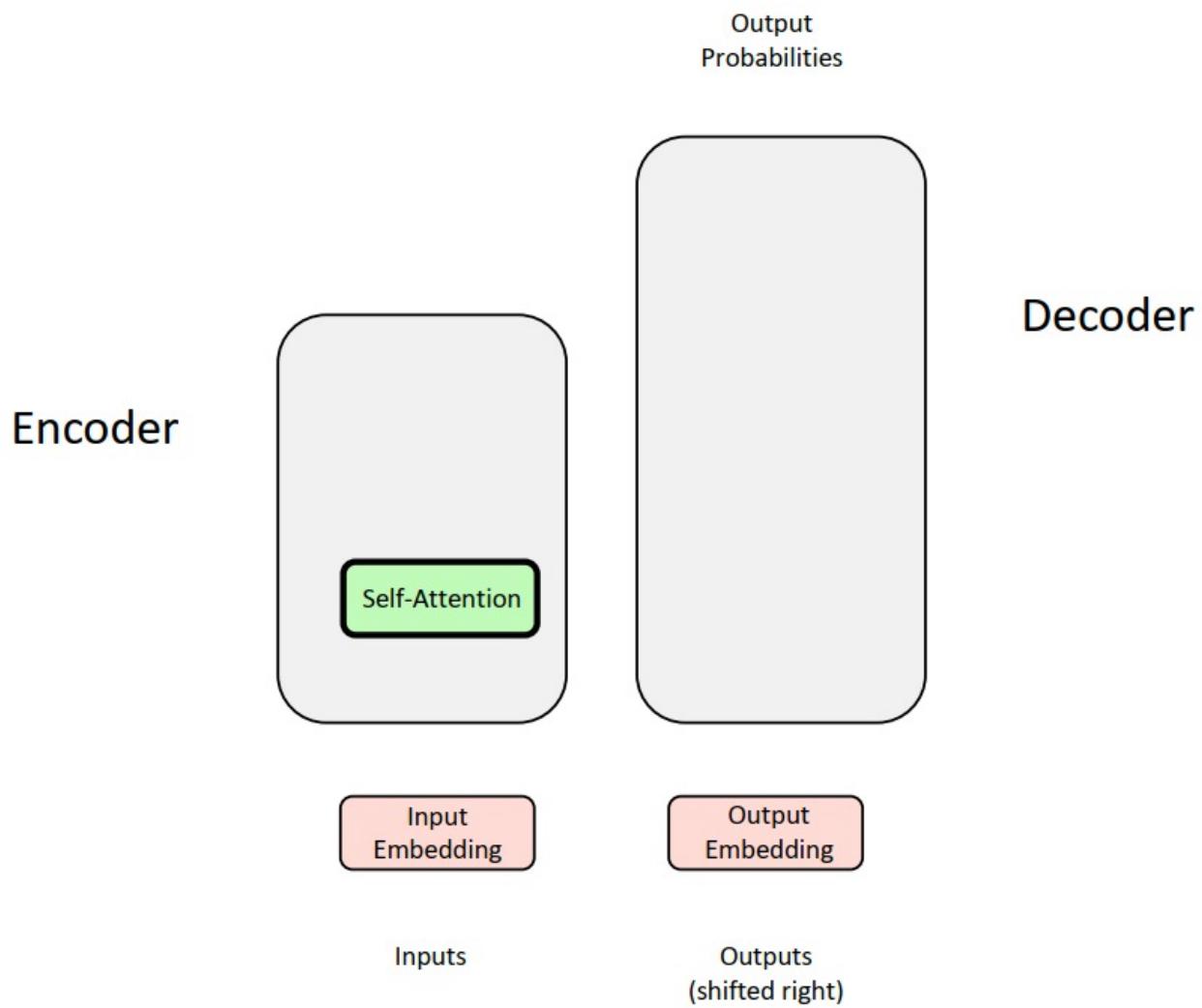
$$A = \text{softmax}(E)$$

- Step 4: Take a weighted sum of **values**

$$\text{Output} = AV$$

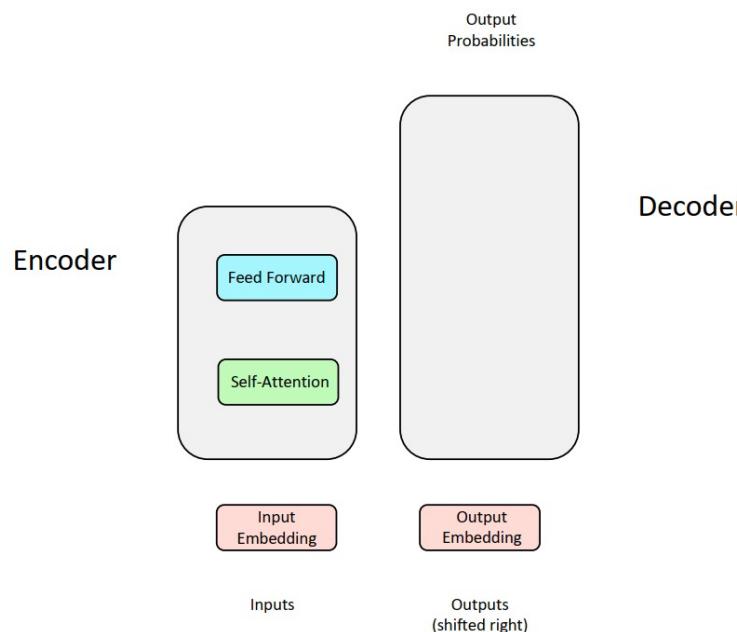
$$\text{Output} = \text{softmax}(QK^T)V$$

# What We Have So Far: (Encoder) Self-Attention!



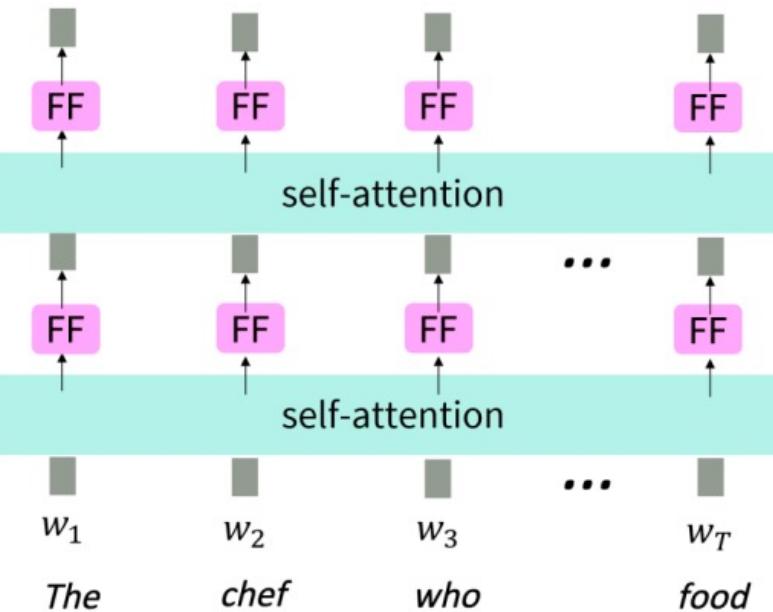
# But attention isn't quite all you need!

- 问题：
  - self-attention就是value向量的重新加权平均
  - 没有任何非线性层
- 简单的修复方式
  - Self-attention后加一个feedforward层，加入非线性激活函数，以及额外的表达能力

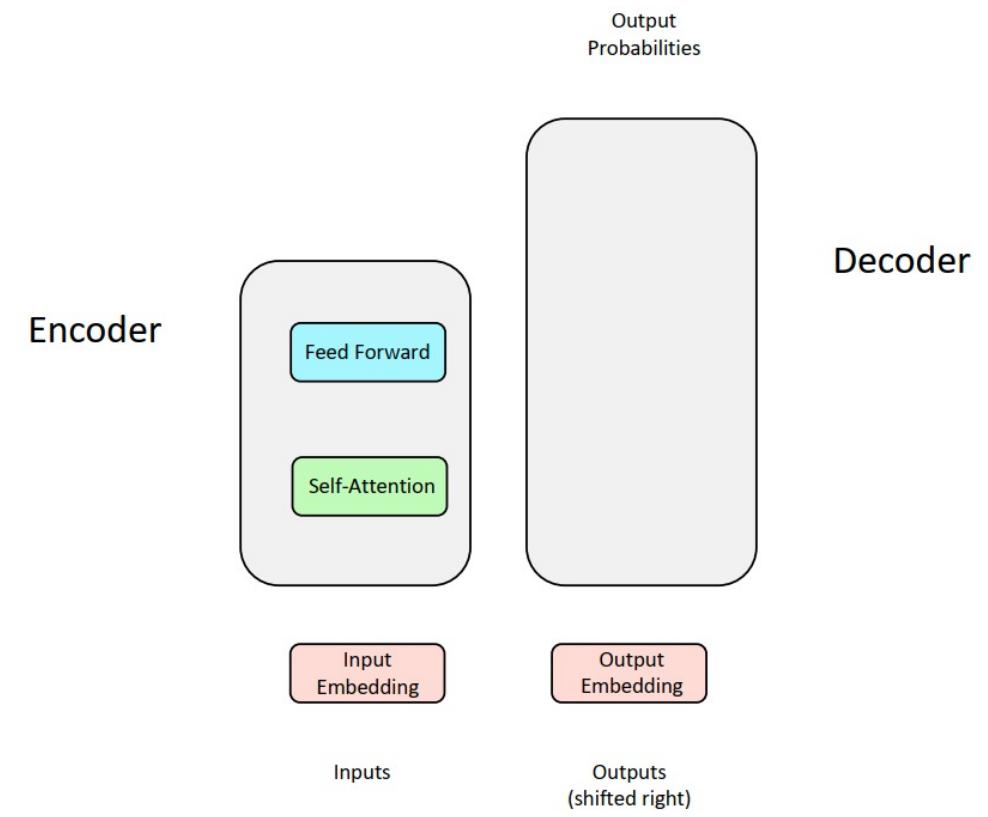
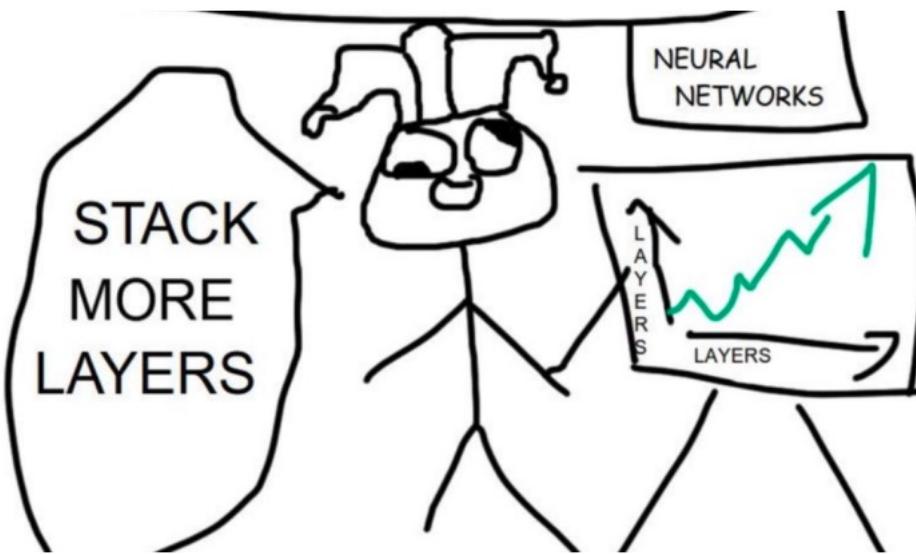


Equation for Feed Forward Layer

$$\begin{aligned}m_i &= \text{MLP}(\text{output}_i) \\&= W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2\end{aligned}$$



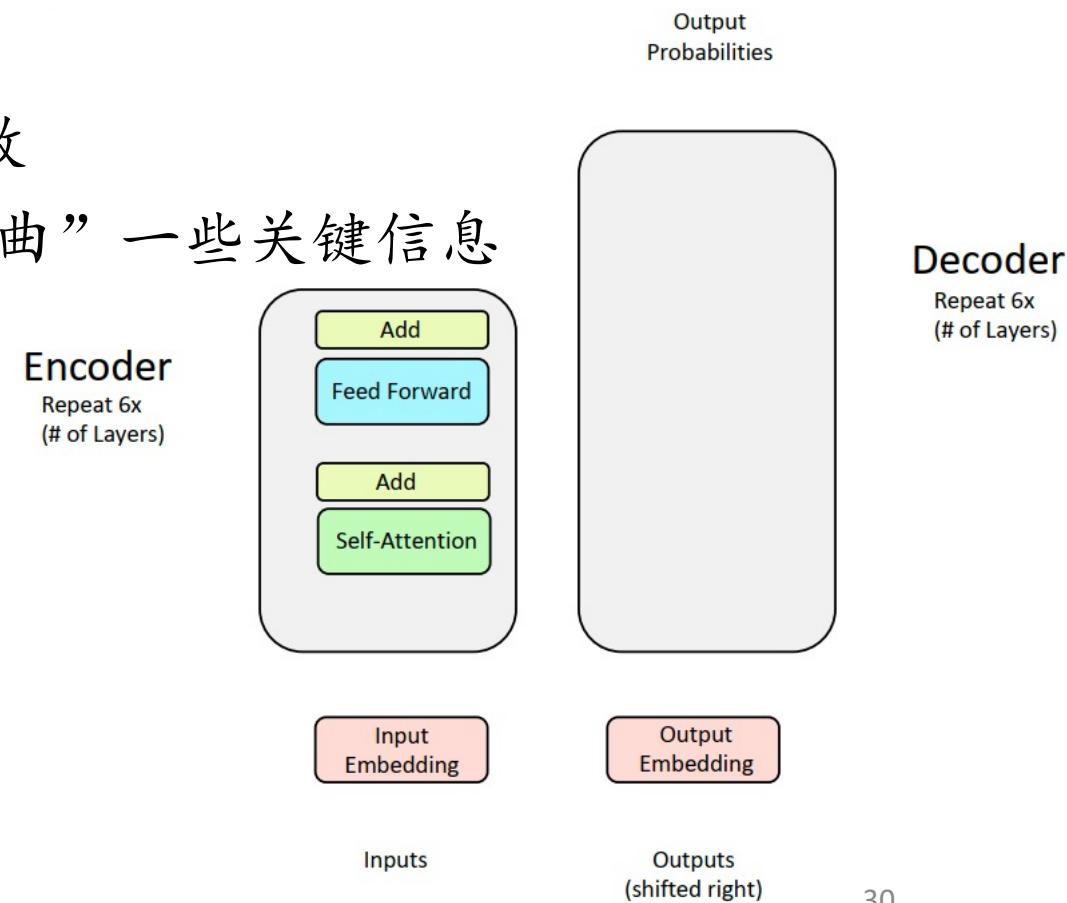
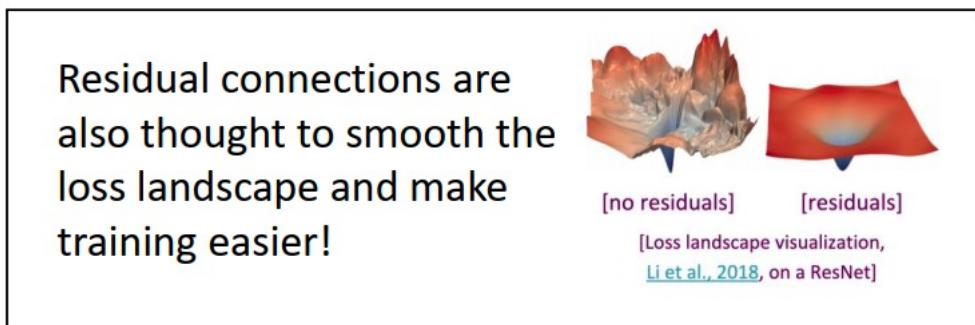
# But how do we make this work for deep networks?



- Training Trick #1: Residual Connections
- Training Trick #2: LayerNorm
- Training Trick #3: Scaled Dot Product Attention

# Training Trick #1: Residual Connections [He et al., 2016]

- Residual Connection: CV中的简单有效的方法
- 深度神经网络很难学习到identity function
- 直接将“raw” embeddings传到下一层比较有效
- 可以避免网络“遗忘”或经过很多层后“扭曲”一些关键信息

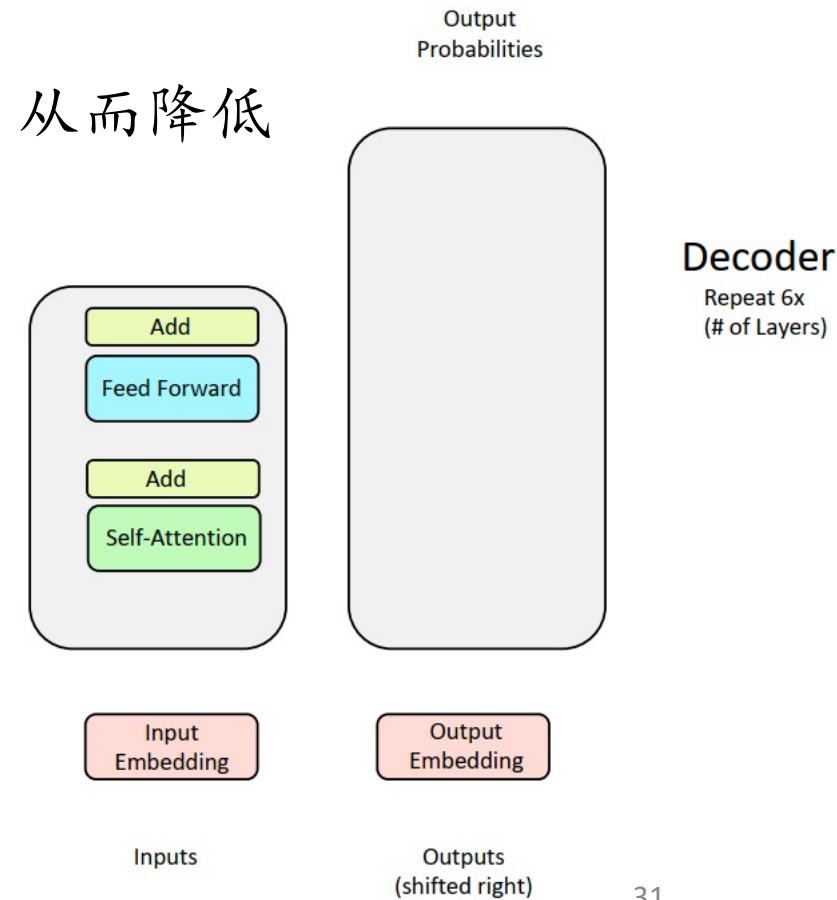


## Training Trick #2: Layer Normalization [Ba et al., 2016]

- 问题：由于下游层的输出会不断飘移，所以上层的参数会难以训练
- 解决方案：每层之内归一化为均值为0，标准差为1。从而降低非新信息导致的偏移。

$$\text{Mean: } \mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \text{Standard Deviation: } \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$x^{\ell'} = \frac{x^\ell - \mu^\ell}{\sigma^\ell + \epsilon}$$



# Training Trick #3: Scaled Dot Product Attention

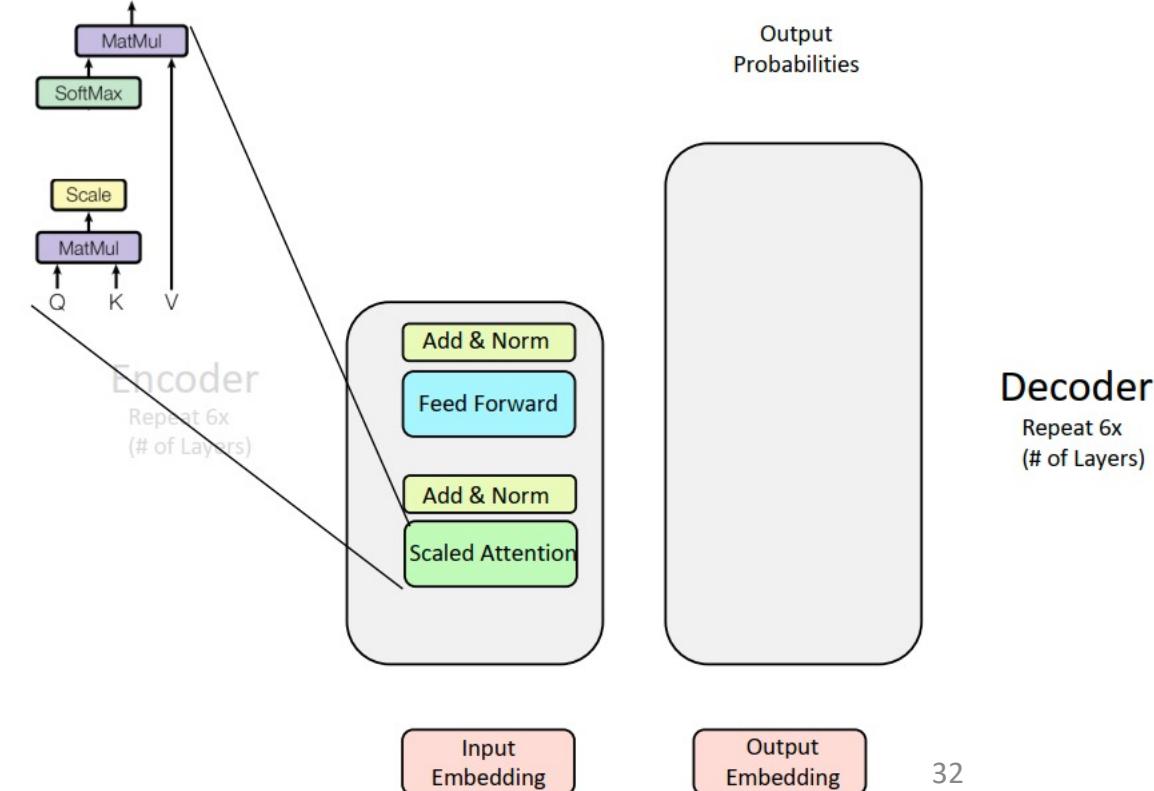
- LayerNorm之后，向量的均值方差被归为了0和1
- 然而，点积会给出极端值，向量越长，越容易出现。方差会因此

## Quick Statistics Review:

- Mean of sum = sum of means =  $d_k * 0 = 0$
- Variance of sum = sum of variances =  $d_k * 1 = d_k$
- To set the variance to 1, simply divide by  $\sqrt{d_k}$ !

## Updated Self-Attention Equation:

$$Output = \text{softmax}\left(QK^T / \sqrt{d_k}\right)V$$

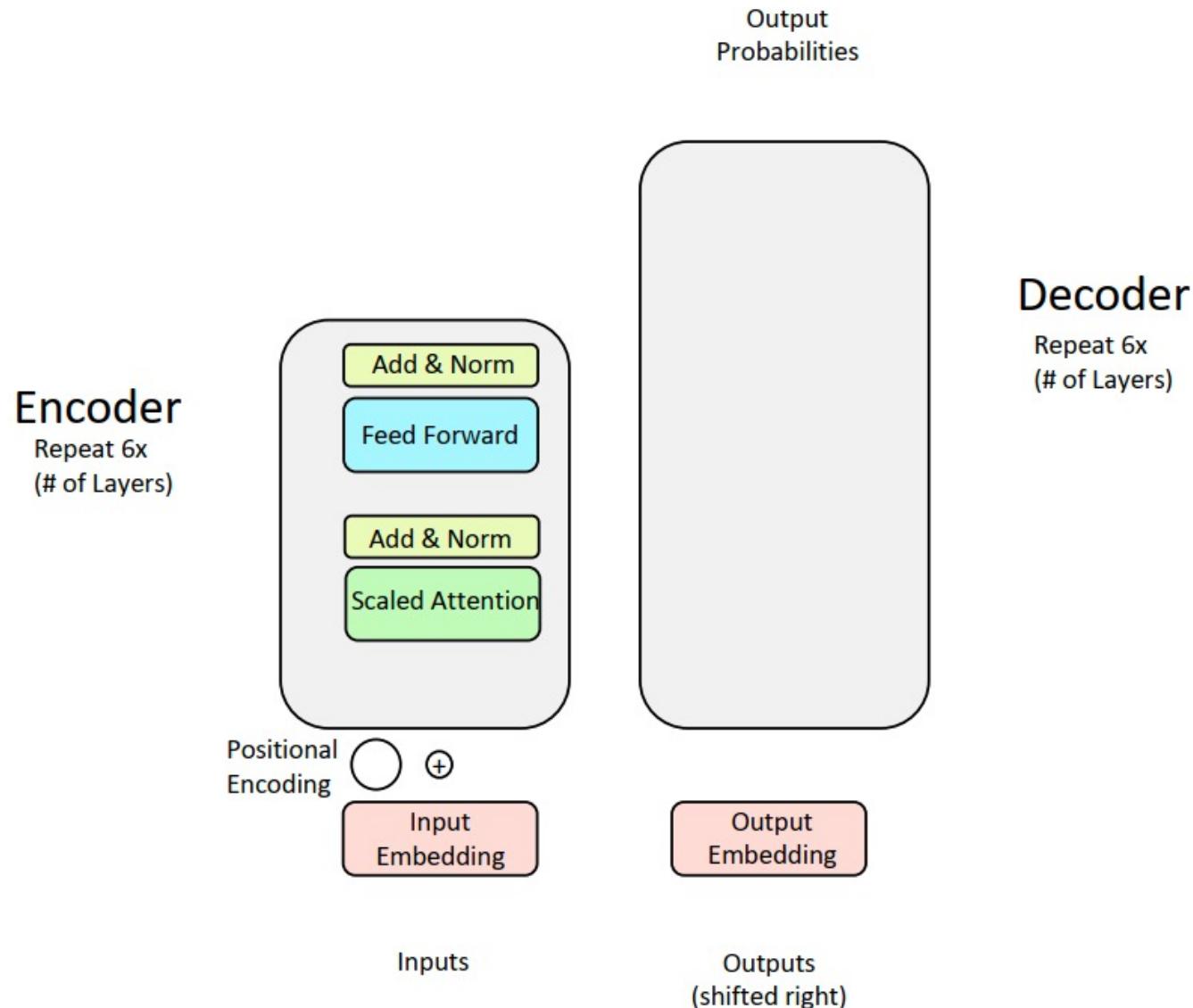


# Major issue!

- Encoder我们介绍的差不多了，但现在还有一个重要的问题，有人发现吗？
- 考虑这个句子
  - "Man eats small dinosaur."
- 用Transformer的话，词语顺序并不影响神经网络计算
- 词语顺序在很多语言中都很重要。所以此处问题很大。

$$Output = \text{softmax}\left(QK^T / \sqrt{d_k}\right)V$$

# Solution: Inject Order Information through Positional Encodings!



## Fixing the first self-attention problem: sequence order

- 由于self-attention无法利用到顺序信息，我们需要把顺序建模进 keys, queries, and values中
- 考虑把index表示成向量

$p_i \in \mathbb{R}^d$ , for  $i \in \{1, 2, \dots, T\}$  are position vectors

- 暂时先不考虑  $p_i$ 怎么算的
- 融合进self-attention很容易：简单的加在一起
- 令  $\tilde{v}_i, \tilde{k}_i, \tilde{q}_i$ 是旧的values, keys, and queries

$$v_i = \tilde{v}_i + p_i$$

$$q_i = \tilde{q}_i + p_i$$

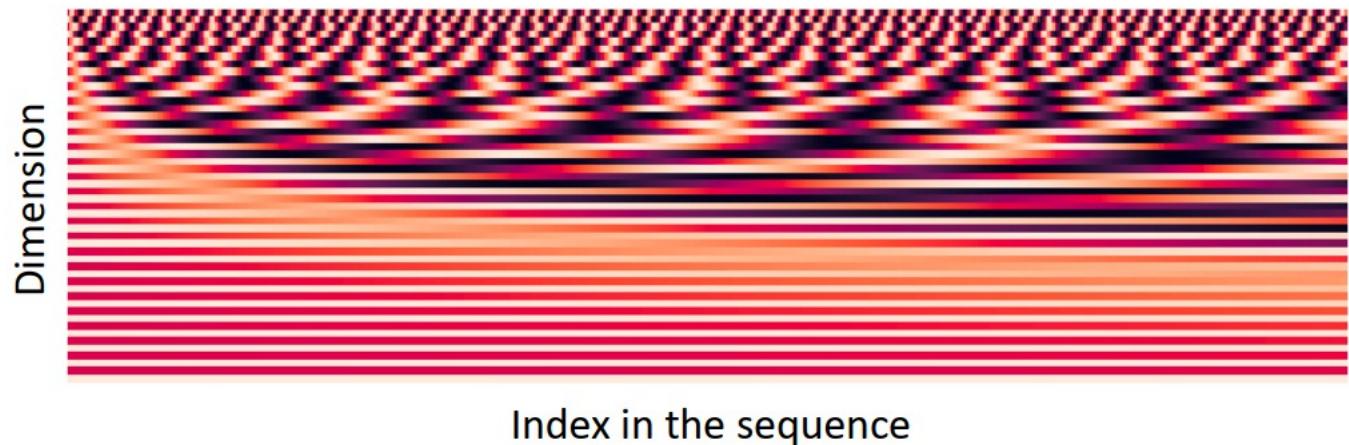
$$k_i = \tilde{k}_i + p_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...

# Position representation vectors through sinusoids

- Sinusoidal position representations: concatenate sinusoidal functions of varying periods:

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



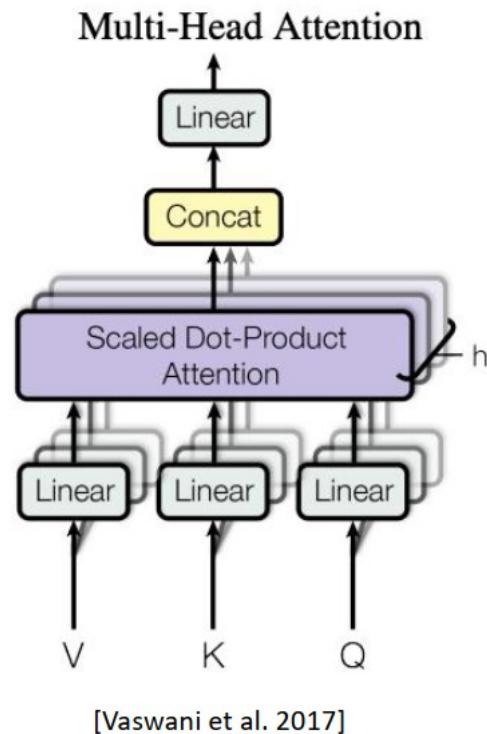
- 优点: (1) 周期性暗指绝对位置不是很重要 (2) 可以扩展到长序列
- 缺点: (1) 不可学习 (2) “扩展”不一定很有用

# Position representation vectors learned from scratch

- 学习一个绝对的位置表示：令所有的  $p_i$  可学习
  - 学习一个矩阵  $p \in \mathbb{R}^{d \times T}$ , 令  $p_i$  属于该矩阵的某一列
- 优点：
  - 灵活，每个位置的表示可以从数据中得到训练
- 缺点：
  - 无法扩展，超出  $1 \sim T$  的范围无法表示
  - 目前很多模型用的是这种方法
  - 其他灵活的位置向量方法
    - Relative linear position attention [Shaw et al., 2018]
    - Dependency syntax-based position [Wang et al., 2019]

# Multi-Headed Self-Attention: k heads are better than 1!

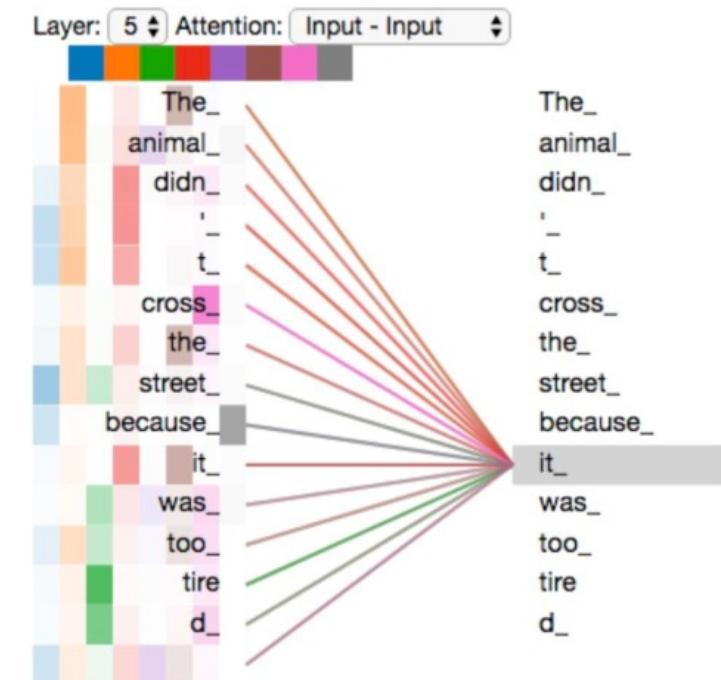
- High-Level Idea: Let's perform self-attention multiple times in parallel and combine the results



Wizards of the Coast, Artist: Todd Lockwood

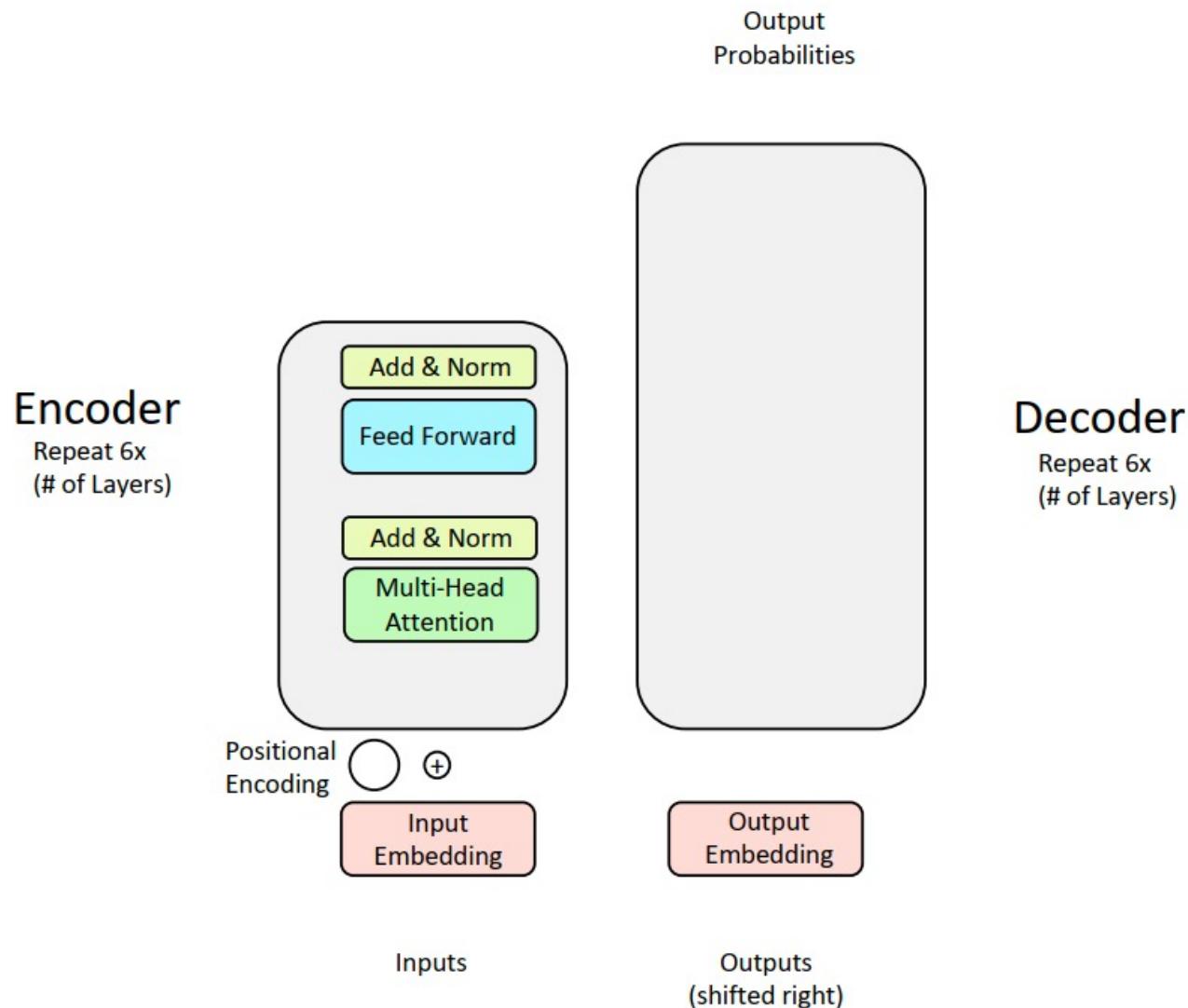
# The Transformer Encoder: Multi-headed Self-Attention

- 如果我们想同时关注句子中的多个位置呢?
  - 对于词*i*, self-attention关注的是  $x_i^T Q^T K x_j$  大的点。如果我们由于其他原因还想关注别的点呢?
- 通过多组Q,K,V来定义多个注意力头
- 令  $h$  是注意力头的数量,  $l=1\dots h$
- 每  $Q_l, K_l, V_l \in \mathbb{R}^{d \times \frac{d}{h}}$  计算过程都是一样的
- 最后将每个头的计算结果拼接起来



$$\text{output} = Y[\text{output}_1; \dots; \text{output}_h], \text{ where } Y \in \mathbb{R}^{d \times d}$$

# Yay, we've completed the Encoder! Time for the Decoder...



# Decoder: Masked Multi-Head Self-Attention

- 问题：
  - Self-attention 可以同时看到所有信息
  - Decode阶段会不会直接看到答案？
- 解答：
  - 将Multi-Head Attention 的未来信息mask掉

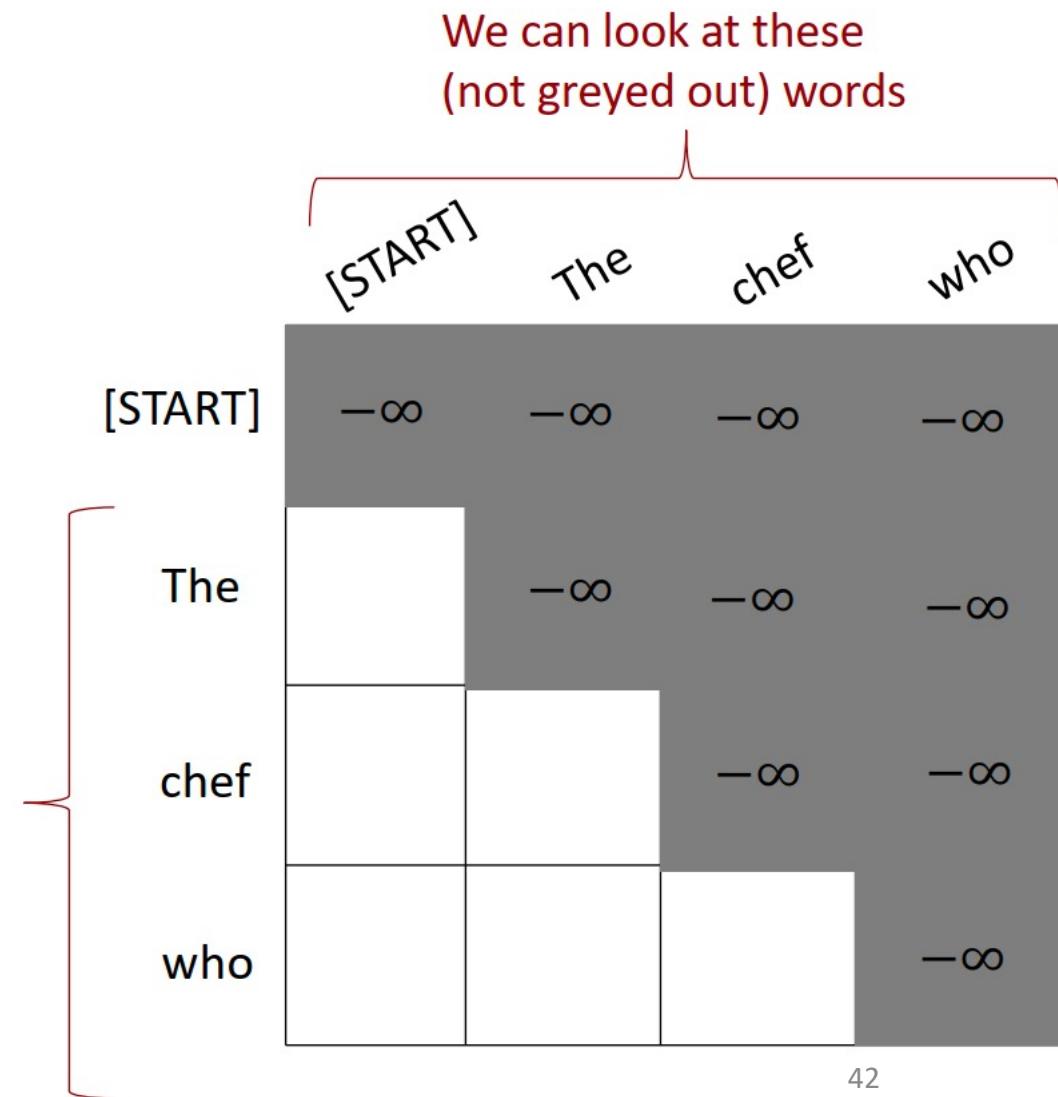
# Masking the future in self-attention

- To use self-attention in decoders, we need to ensure we can't peek at the future.
- To enable parallelization, we mask out attention to future words by setting attention scores to

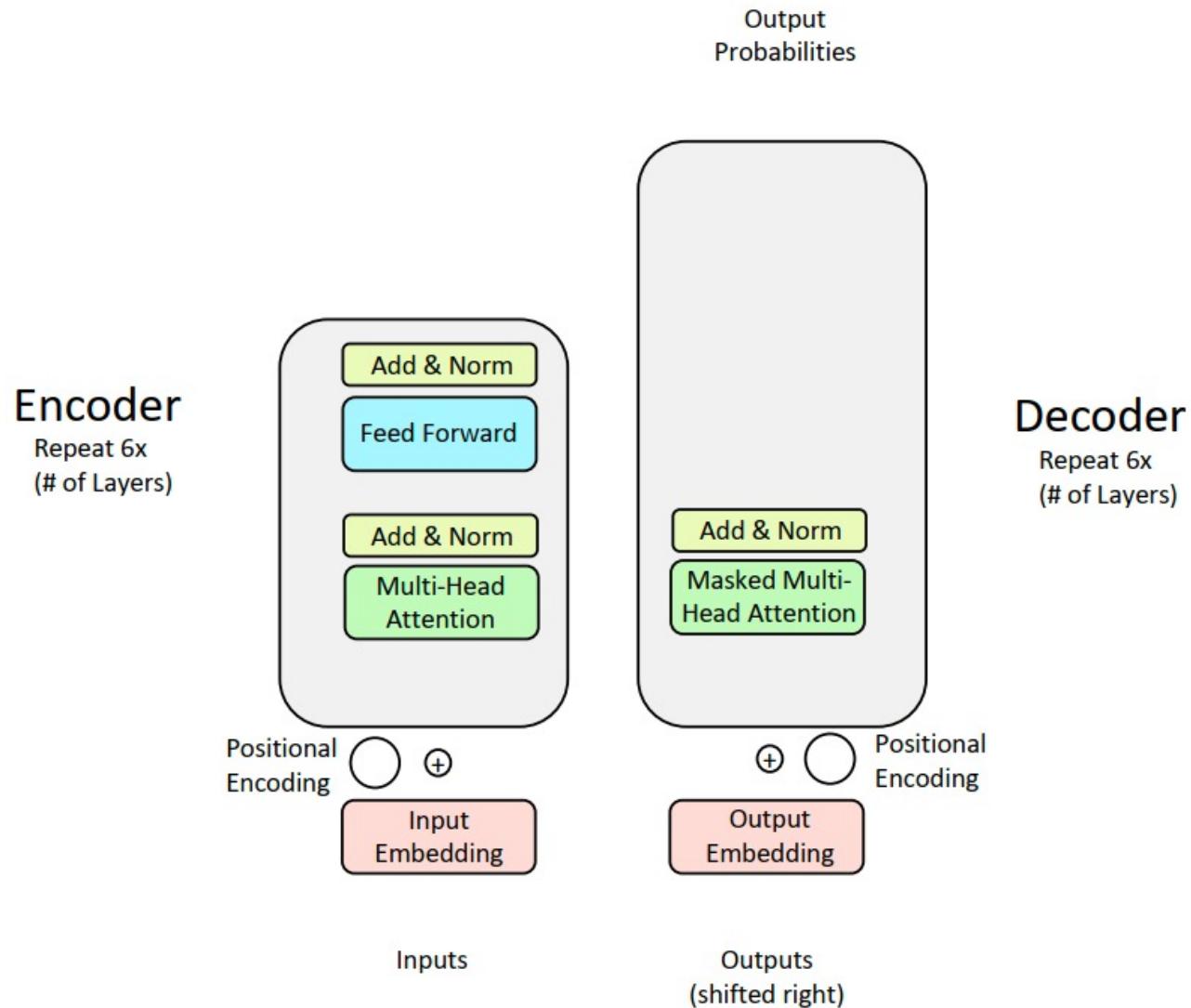
$$e_{ij} = \begin{cases} q_i^T k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

$-\infty$

For encoding  
these words



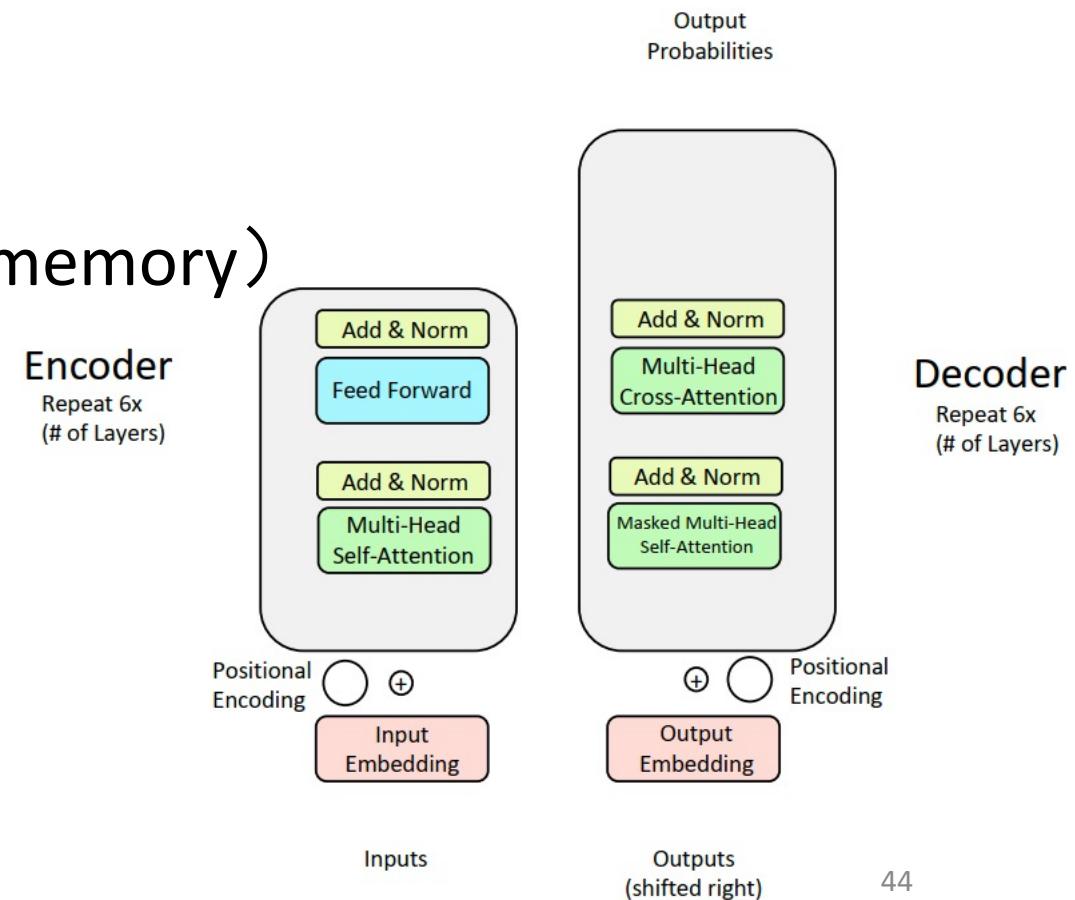
# Decoder: Masked Multi-Headed Self-Attention



# Encoder-Decoder Attention

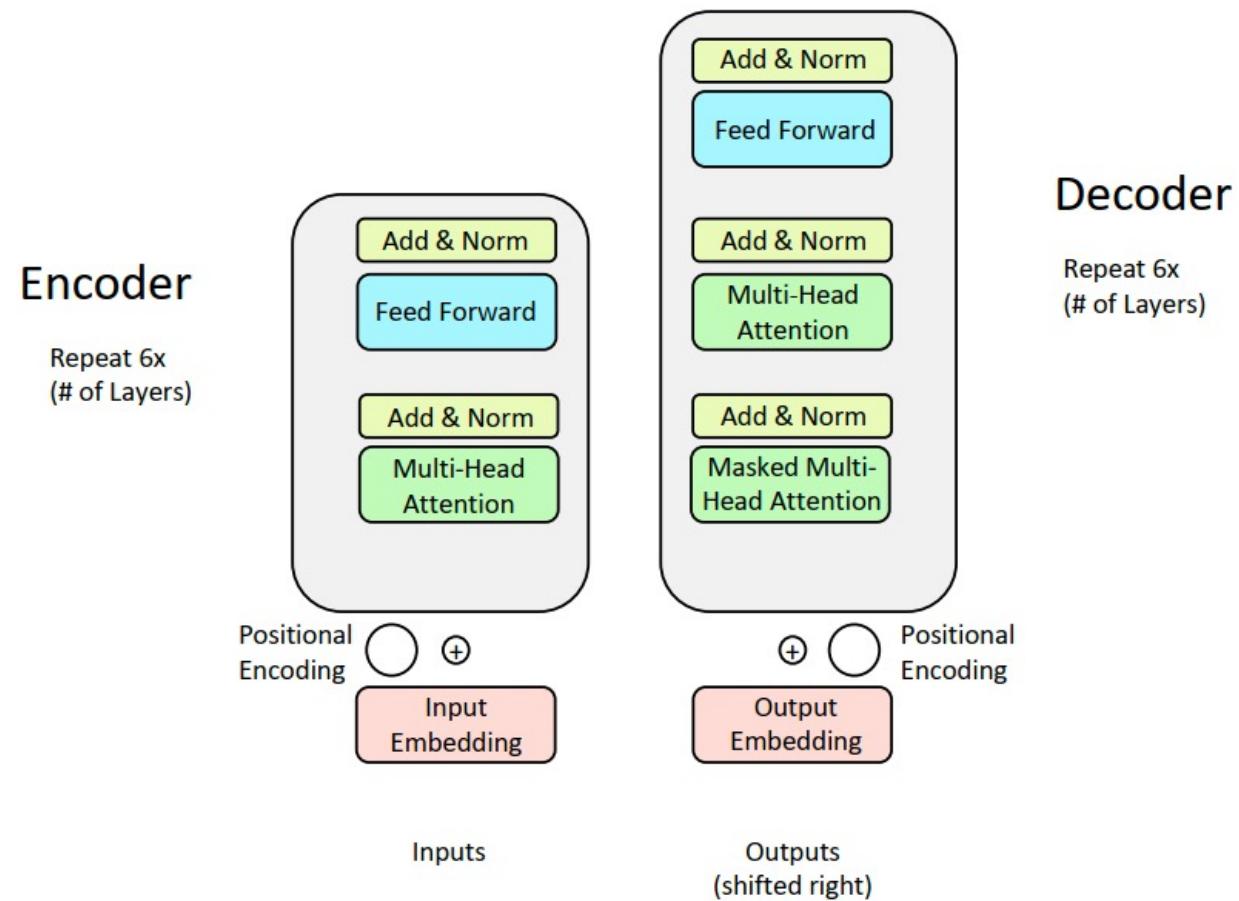
- Self-attention 中， keys, queries, values 都是同源的
- 令  $h_1, \dots, h_T$  encoder 的输出向量
- 令  $z_1, \dots, z_T$  decoder 的输入向量
- Key 和 value 是从 encoder 中提取的（像 memory）  
$$k_i = Kh_i, v_i = Vh_i$$

- Query 是从 decoder 中来的  
$$q_i = Qz_i$$



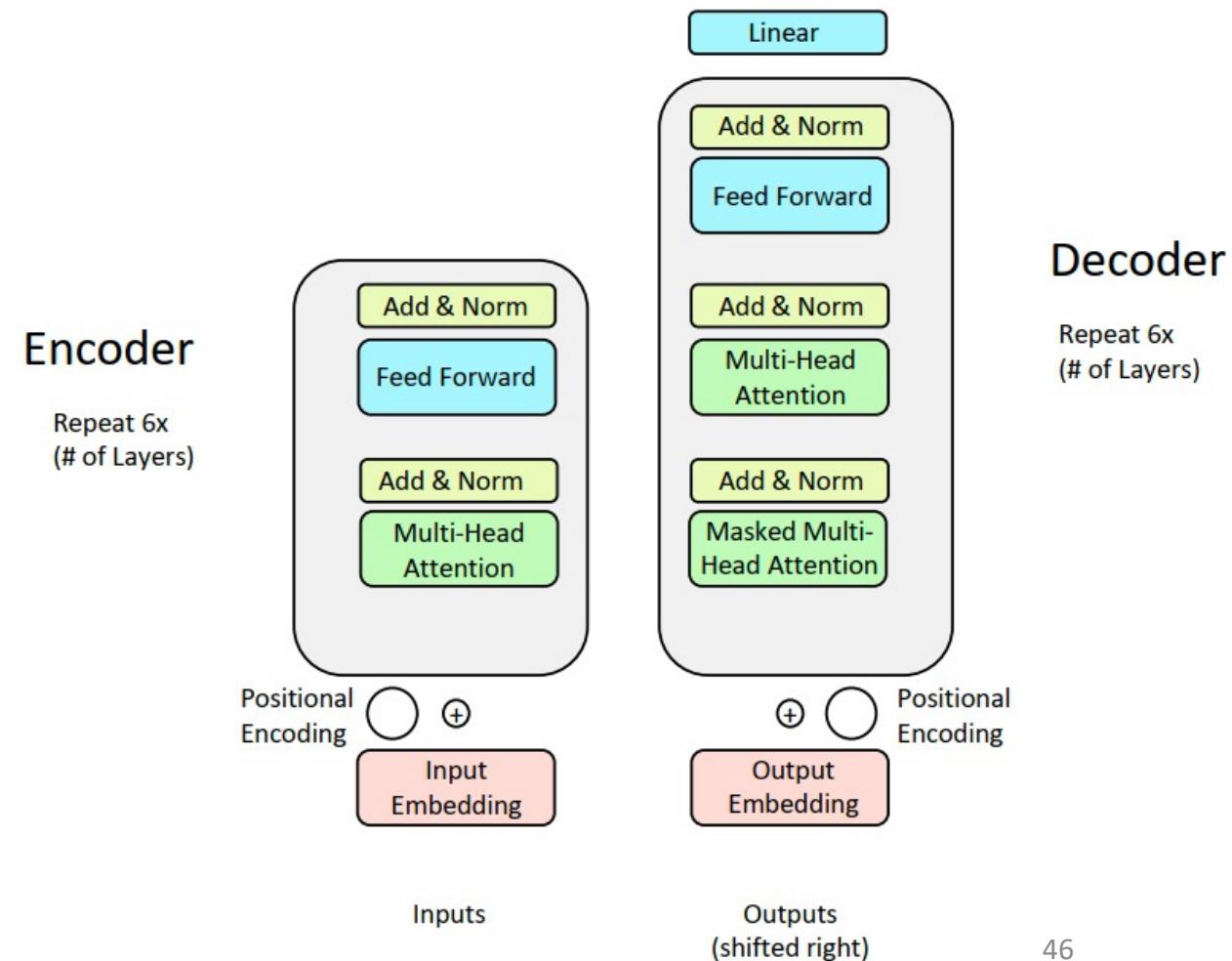
# Decoder: Finishing touches!

- Add a feed forward layer (with residual connections and layer norm)



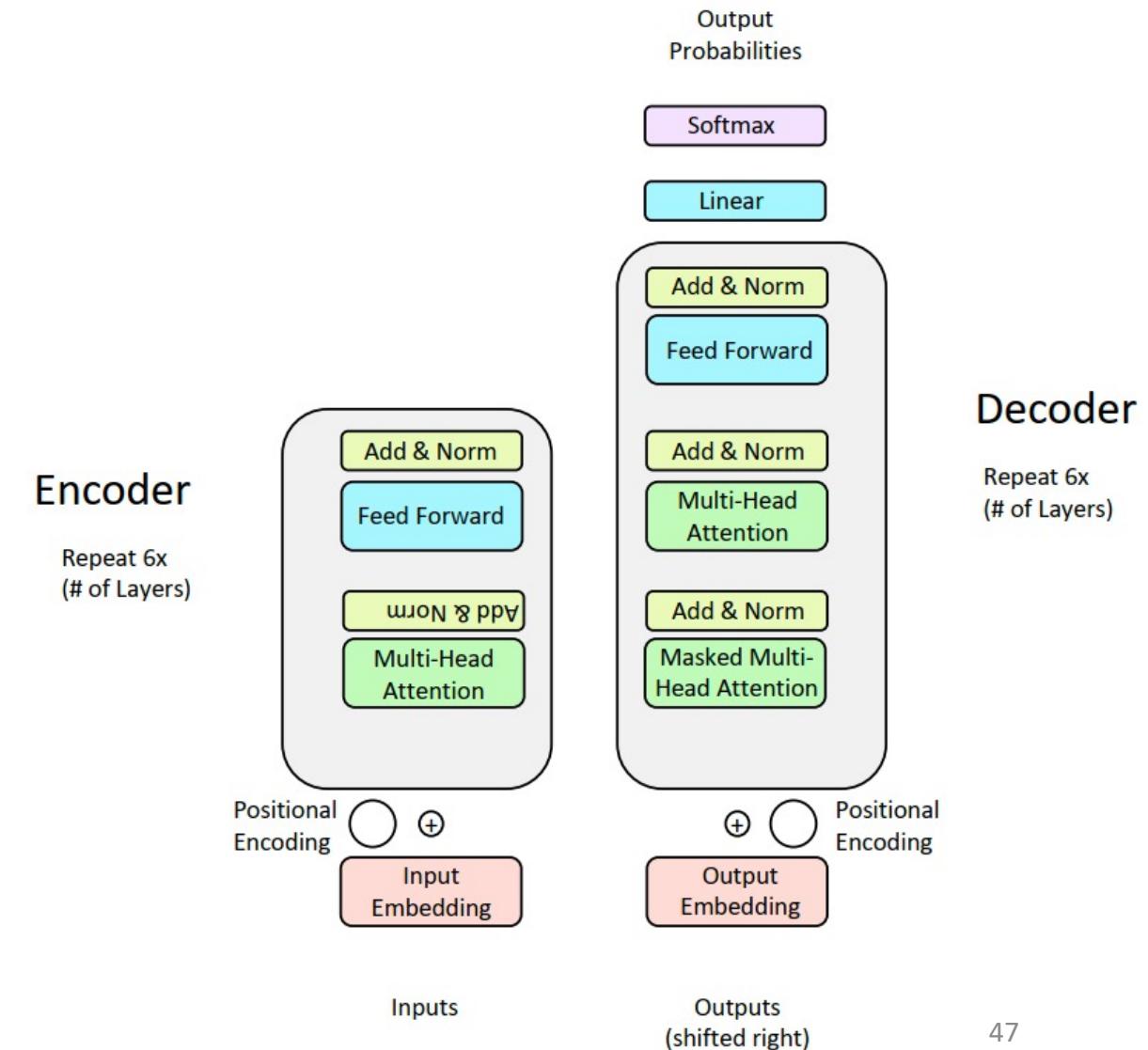
# Decoder: Finishing touches!

- Add a feed forward layer (with residual connections and layer norm)
- Add a final linear layer to project the embeddings into a much longer vector of length vocab size (logits)

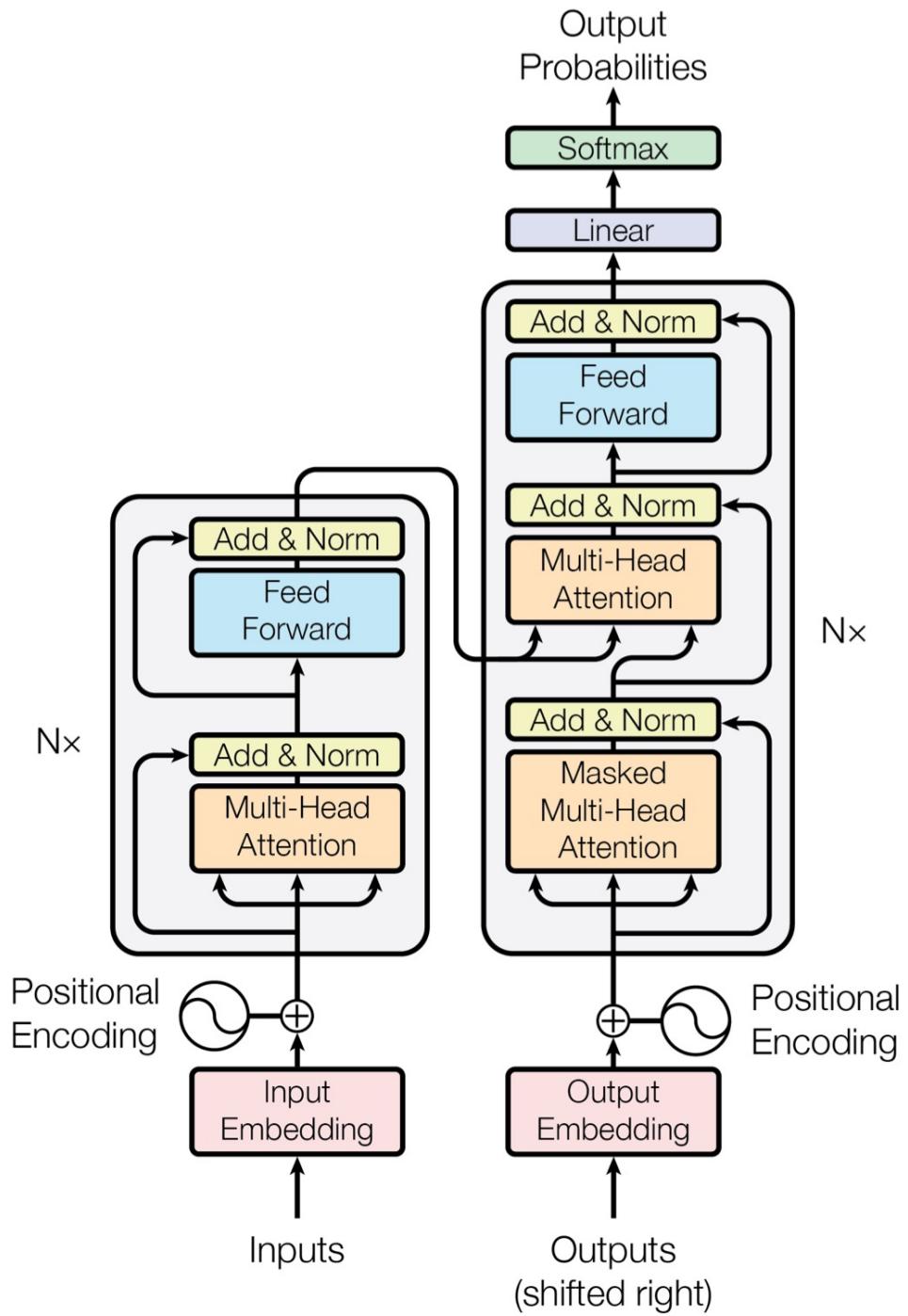


# Decoder: Finishing touches!

- Add a feed forward layer (with residual connections and layer norm)
- Add a final linear layer to project the embeddings into a much longer vector of length vocab size (logits)
- Add a final softmax to generate an probability distribution of possible next words!



# Recap of Transformer Architecture



# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

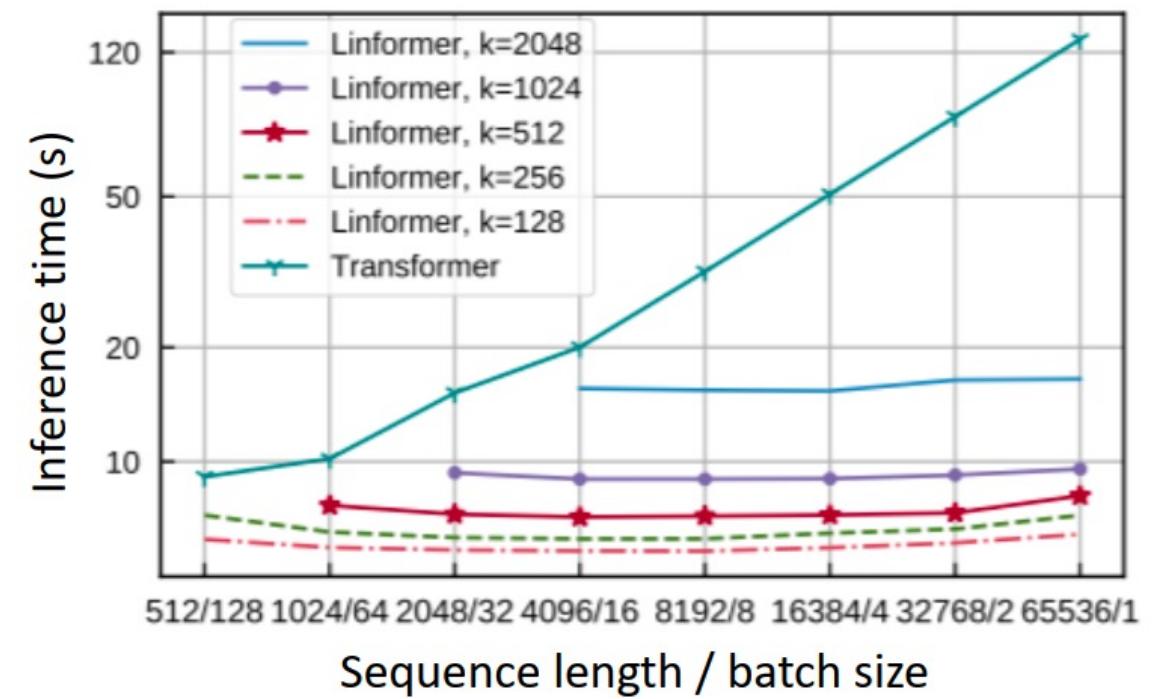
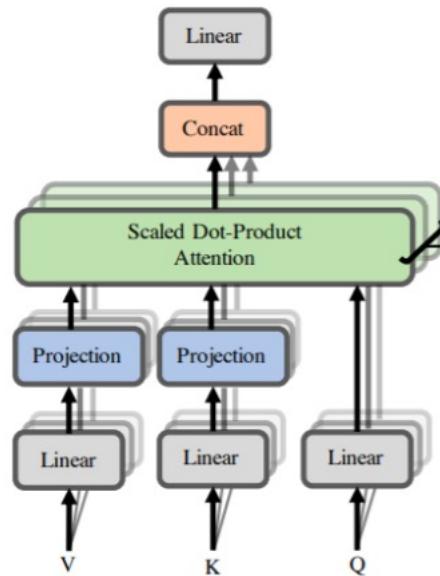
# What would we like to fix about the Transformer?

- Self-attention的计算是平方复杂度的
  - Recurrent模型仅仅是线性复杂度
- 位置表示
  - 绝对位置的可学习的向量是我们的最好选择吗？
  - Relative linear position attention [Shaw et al., 2018]
  - Dependency syntax-based position [Wang et al., 2019]

# Recent work on improving quadratic self-attention cost

- 很多目前的工作都要解决self-attention的 $O(T^2)$ 计算量问题
- 比如, Linformer [Wang et al., 2020]

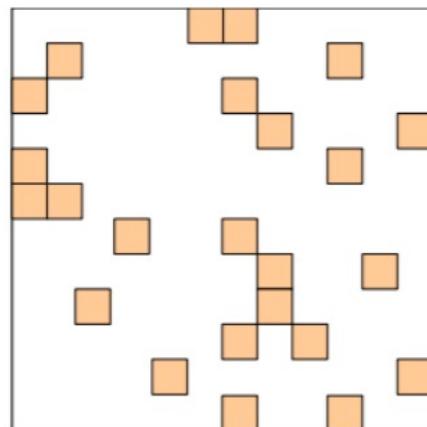
Key idea: map the sequence length dimension to a lower-dimensional space for values, keys



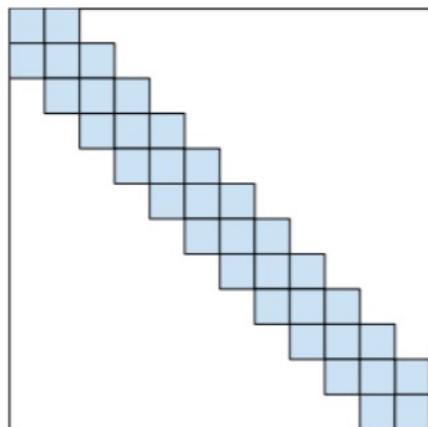
# Recent work on improving on quadratic self-attention cost

- 很多目前的工作都要解决self-attention的 $O(T^2)$ 计算量问题
- 比如, BigBird [Zaheer et al., 2021]

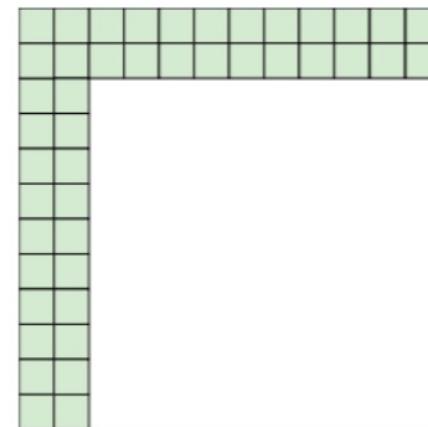
Key idea: replace all-pairs interactions with a family of other interactions, **like local windows, looking at everything, and random interactions.**



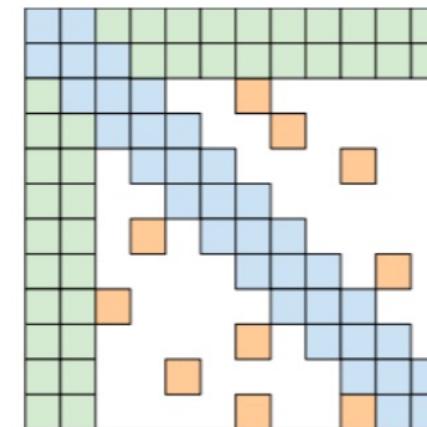
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD



# PLMs

# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

# Breaking (Transformer) News!

- AlphaCode (a pre-trained Transformer-based code generation model) achieved a top 54.3% rating on Codeforces programming competitions!

AlphaCode Attention Visualization

Hover over tokens in the solution to see which tokens the model attended to when generating the solution. Click a token to select it; clicking in empty space will deselect.

1548\_A. Web of Lies C++ incorrect Layer 18

Head 1 Head 2 Head 3 Head 4 Head 5 Head 6 Head 7 Head 8 Head 9 Head 10 Head 11 all none

Problem Description Solution

```
// RATING: 1200
// TAGS: dp,implementation
// LANGUAGE IS cpp
// CORRECT SOLUTION
// When you play the game of thrones, you win, or you die. There is no middle ground.
// Cersei Lannister, A Game of Thrones by George R. R. Martin
// There are n nobles, numbered from 1 to n. Noble i has a power of i. There are also m "friendships". A friendship between nobles a and b is always mutual.
// A noble is defined to be vulnerable if both of the following conditions are satisfied:
//   * the noble has at least one friend, and
//   * all of that noble's friends have a higher power.
// You will have to process the following three types of queries.
//   1. Add a friendship between nobles u and v.
//   2. Remove a friendship between nobles u and v.
//   3. Calculate the answer to the following process.
// 
```

```
#include <bits/stdc++.h>
using namespace std;
int n, m, q, u, v, t;
vector<int> p[200001];
int main() {
    scanf("%d%d", &n, &m);
    for (int i = 0; i < m; i++) {
        scanf("%d%d", &u, &v);
        p[u].push_back(v);
        p[v].push_back(u);
    }
    vector<int> a(n + 1, 0);
    for (int i = 1; i <= n; i++) {
        a[i] = i;
        for (int j = 0; j < p[i].size(); j++) {
            a[p[i][j]] = i;
        }
    }
    scanf("%d", &q);
    for (int i = 0; i < q; i++) {
        scanf("%d", &t);
        if (t == 1) {
            scanf("%d%d", &u, &v);
            while (a[u] != u) {
                u = a[u];
            }
            u = a[u];
            while (a[v] != v) {
                v = a[v];
            }
            a[v] = u;
        }
    }
}
```

# More Breaking (Transformer) News!

- Pre-Trained Transformer-Based theorem prover sets new state-of-the-art (41.2% vs. 29.3%) on a collection of challenging math Olympiad questions (miniF2F) [Polu et al., 2022]

## PROBLEM 1

Adapted from AMC12 2000 Problem 5

Prove that if  $|x - 2| = p$ , where  $x < 2$ , then  $x - p = 2 - 2p$ .

◊ FORMAL

INFORMAL

```
theorem amc12_2000_p5      -- ← theorem name
  (x p : ℝ)                  -- ← the statement we want
  (h₀ : x < 2)                --   to prove
  (h₁ : abs (x - 2) = p) :
  x - p = 2 - 2 * p :=       -- ← formal proof starts here
begin                         -- This first tactic requires that the prover invent
  -- the term: `abs (x - 2) = -(x - 2)`.
  have h₂ : abs (x - 2) = -(x - 2), {
    apply abs_of_neg,
    linarith,
  },
  rw h₁ at h₂,
  -- At this stage the remaining goal to prove is:
  -- `x - p = 2 - 2 * p` knowing that `p = -(x - 2)`.
  linarith,
end
```

# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

# Word structure and subword models

- 首先回忆一下我们对于语言词表的假设
- 我们首先从训练集中构造一个固定的词表，可能包含几万个词
- 测试集中的词如果无法在此表中找到，一律标为UNK

	word	vocab mapping	embedding
Common words	hat	→ hat	
	learn	→ learn	
Variations	taaaaaasty	→ UNK	
	laern	→ UNK	
novel items	Transformerify	→ UNK	

# Word structure and subword models

- 有限的词表在很多语言中不太合理
  - 很多语言都有复杂的形态学变化
  - 后果是更多的词语形态占据词表中很多位置，有些可能很低频
- 例子：斯瓦西里语的动词有上百种变形，每种对应不同的时态，情绪，确定性，否定，等等

Conjugation of -ambia												[less ▲]		
Polarity	Form Infinitive				Non-finite forms								[less ▲]	
	Positive form				Simple finite forms				Complex finite forms					
	Imperative Habitual				Singular ambia				huambia					
Polarity	Persons				Persons / Classes				M-mi	4	5 Ma	7 Ki-vi	Classes	
	1st Sg.	Pl.	2nd Sg.	Pl.	3rd / M-wa	Sg. / Pl. / 2			3	4	5	6	N	
Positive	niiambia naliambia	tuliambia tvaliambia	uliambia waliambia	mliambia mwaliambia	aliambia	wallambia	ullambia	illambia	yallambia	kiliambia	villiambia	illambia	zillambia uillambia	
Negative	sikuambia	hatukuambia	hukuambia	hamkuambia	hakuambia	a	haukuambia	haikuambia	halikuambia	hayakuambia	hakikuambia	havikuambia	haikuambia hazikuambia	
Present													[less ▲]	
Positive	ninaambia naambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	linambia	yanaambia	kinaambia	vinaambia	inaambia zinaambia	
Negative	siambii	hatuambii	huambii	hamambii	haambii	hawaambii	hauambii	haiambii	halambii	hayaambii	hakiambii	haviambi	halambii haziambi	
Future													[less ▲]	
Positive	nitaambia	tutaambia	utaambia	mtaambia	ataambia	wataambia	utaambia	itaambia	itaambia	yataambia	kitaambia	vitaambia	itaambia zitaambia	
Negative	sitaambia	hatutaambia	hutaambia	hamtaambia	hataambia	hawataambia	hautaambia	altaambia	altaambia	hayataambia	hakitaambia	havitaambia	altaambia hazitaambia	
Subjunctive													[less ▲]	
Positive	niambie	tuamvie	uamble	mambie	aambie	waambie	uamble	iambie	liambie	yaambie	kiambie	viambie	iambie ziambie	
Negative	nisiambie	tusiambie	usiambie	msiambie	asiambie	wasiambie	usiambie	isiambie	lisiambie	yasiambie	kisiambie	visiambie	isiambie zisiambie	
Present Conditional													[less ▲]	
Positive	ningeambia	tungeambia	ungeambia	mngeambia	angeambia	wangeambia	ungeambia	ingeambia	lingeambia	yangeambia	kingeambia	vingeambia	ingeambia zingeambia	
Negative	nisingeambi a	tusingeambi a	usingeambi a	msingeambi a	asingeambi a	wasingeambi a	usingeambi a	isingeambi a	lisingeambi a	yasingeambi a	kisingeambi a	visingeambi a	zingeambia zisingeambia	
Past Conditional													[less ▲]	
Positive	ningaliambia	tungaliambia	ungaliambia	mgngaliambia	angaliambia	wangaliambia	ungaliambia	ingaliambia	lingaliambia	yangaliambia	kingaliambia	vngaliambia	ngaliambia zingaliambia	
Negative	nisingaliamb ia	tusingaliamb ia	usingaliamb ia	msngaliamb ia	asingaliamb ia	wasingaliamb ia	usingaliamb ia	isingaliamb ia	lisngaliamb ia	yasingaliamb ia	kisingaliamb ia	visingaliamb ia	zingaliambia zisingaliambia	
Conditional Contrary to Fact													[less ▲]	
Positive	ningeliambia	tngeliambia	ngeliambia	mnngeliambia	ngeliambia	wngeliambia	ngeliambia	ingeliambia	lingeliambia	yngeliambia	kingeliambia	vngeliambia	ngeliambia zingeliambia	
Positive	naambia	twaambia	waambia	mwaambia	aambia	waambia	yaambia	laambia	yaambia	chaambia	vyaambia	yaambia	zaambia waambia	
Positive	naambia	twaambia	waambia	mwaambia	aambia	waambia	yaambia	laambia	yaambia	chaambia	vyaambia	yaambia	kwaambia paambia	
													[less ▲]	
													[less ▲]	

# The byte-pair encoding algorithm

- NLP中的Subword建模有很多方法
  - 目前最通用的方法是建立subword的词表
  - 在训练和测试时，每个词都被分为几个subword
- Byte-pair encoding是一种简单有效的方法
  - 最开始，我们只有一个只包含字符和“词结束符”的词表
  - 在语料库中，找到最常见的字符对（“如ab”），则将ab加入subword词表中
  - 不断重复此过程，直到词表大小符合要求
- 另一种常用方法，WordPiece

# Word structure and subword models

- 常见词最终都会变成subword词表的一部分，罕见词会被拆成几个部分
- 最坏情况下，单词会被拆成字符

	word	vocab mapping	embedding
Common words	hat	→ hat	
	learn	→ learn	
Variations	taaaaasty	→ taa## aaa## sty	
	laern	→ la## ern	
misspellings			
novel items	Transformerify	→ Transformer## ify	

# Contents

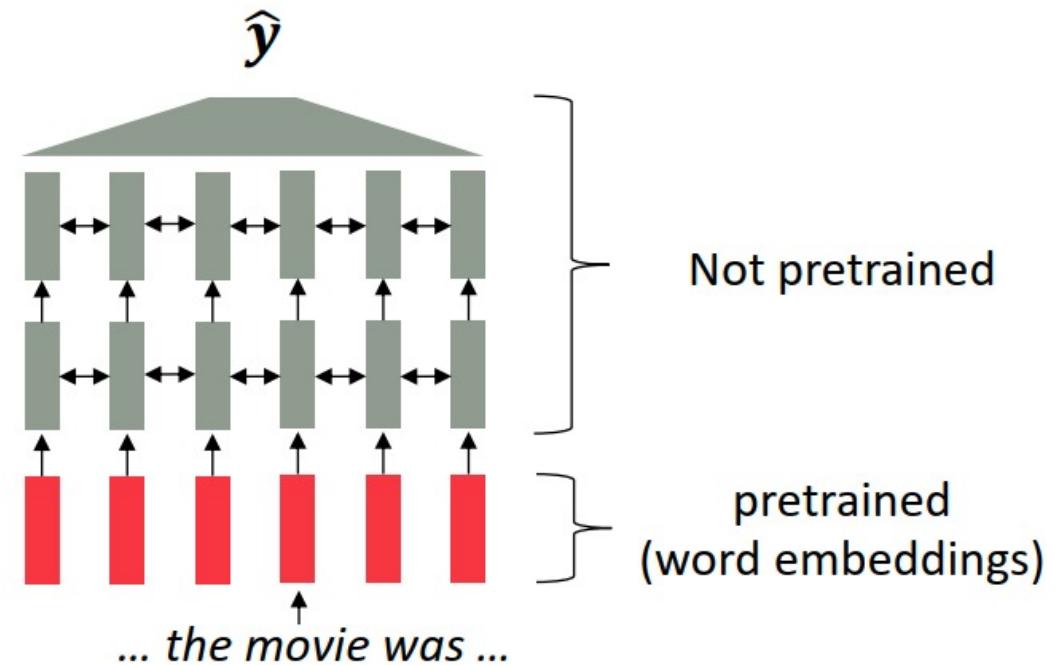
- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

# Motivating word meaning and context

- 回忆之前讲word embedding时候引用的一句话
  - “You shall know a word by the company it keeps” (J. R. Firth 1957: 11)
- 这句话总结了分布式语义，同时引发Word2vec的研究
  - “... the complete meaning of a word is always contextual, and no study of meaning apart from a complete context can be taken seriously.” (J. R. Firth 1935)
- 考虑 I record the record, 两个record有不同含义

# Where we were: pretrained word embeddings

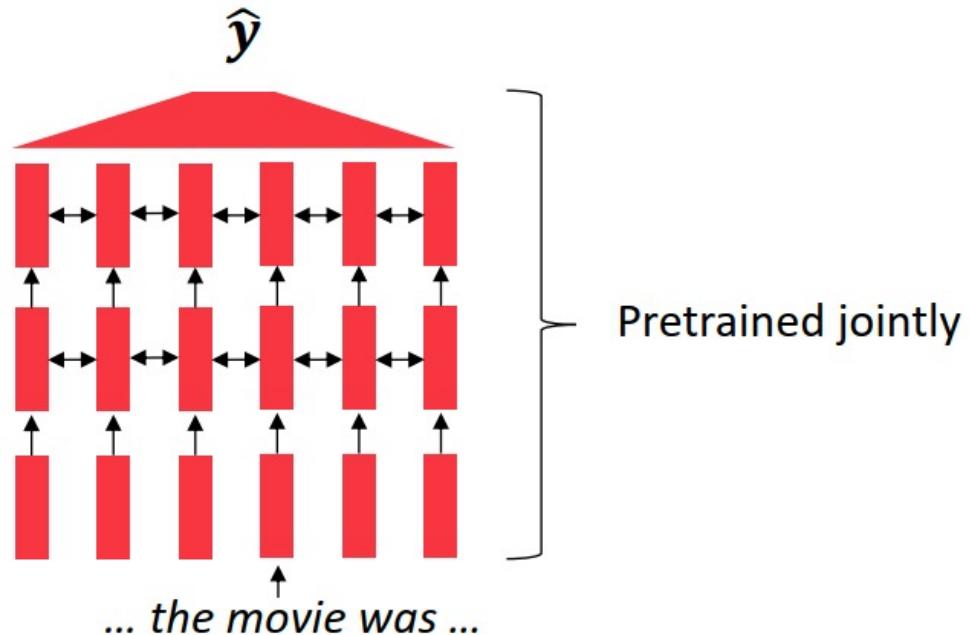
- Since 2017:
  - 仅用预训练的词向量
  - 在训练模型时考虑如何利用上下文信息
- 问题:
  - 下游任务的训练数据必须足以学习上下文信息
  - 大多数参数都是随机初始化的



[Recall, *movie* gets the same word embedding, no matter what sentence it shows up in]

# Where we're going: pretraining whole models

- 现代NLP:
  - NLP模型中几乎所有参数都用预训练来初始化
  - 预训练方法隐藏一部分输入，令模型重构这部分输入
- 在建立强大的以下三种模型时非常有效
  - 语言的表示
  - 对于强大NLP模型的参数初始化
  - 可以被用来采样的语言分布



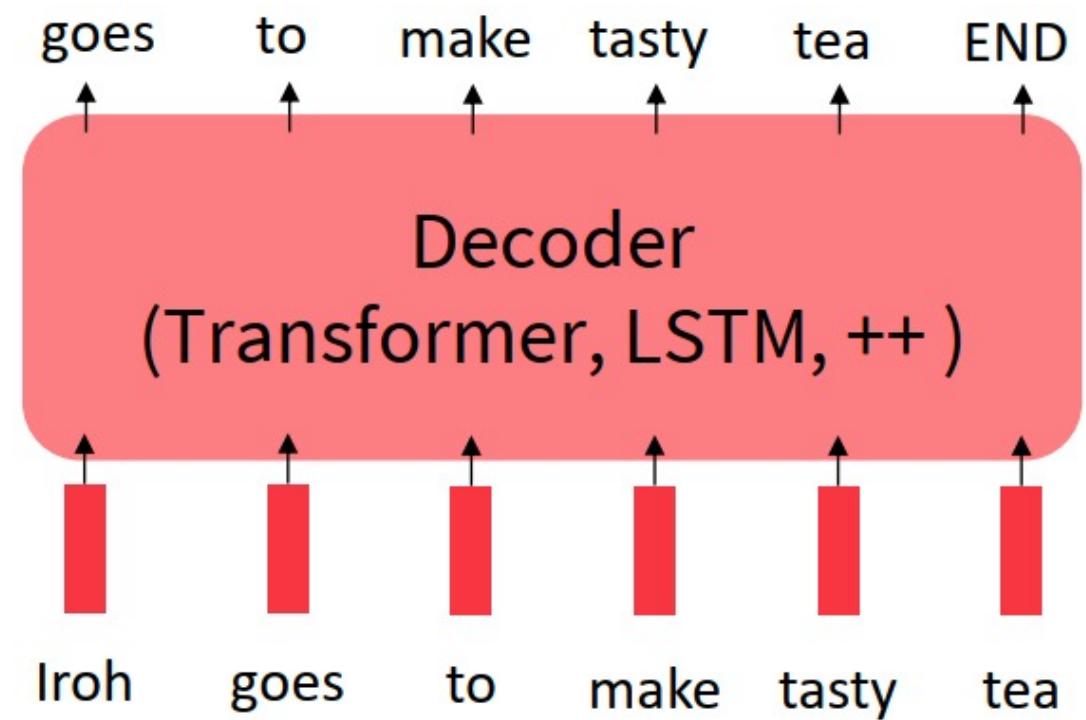
[This model has learned how to represent entire sentences through pretraining]

# What can we learn from reconstructing the input?

- The woman walked across the street, checking for traffic over \_\_\_\_ shoulder.
- I went to the ocean to see the fish, turtles, seals, and \_\_\_\_.
- Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was \_\_\_\_.
- I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_

# Pretraining through language modeling [Dai and Le, 2015]

- 回忆语言模型任务
  - 给定前面的若干词，预测后一个词  $p_\theta(w_t|w_{1:t-1})$
  - 这种训练数据很多
- 通过语言模型来预训练
  - 利用大量语料来训练神经语言模型
  - 存储训好的参数

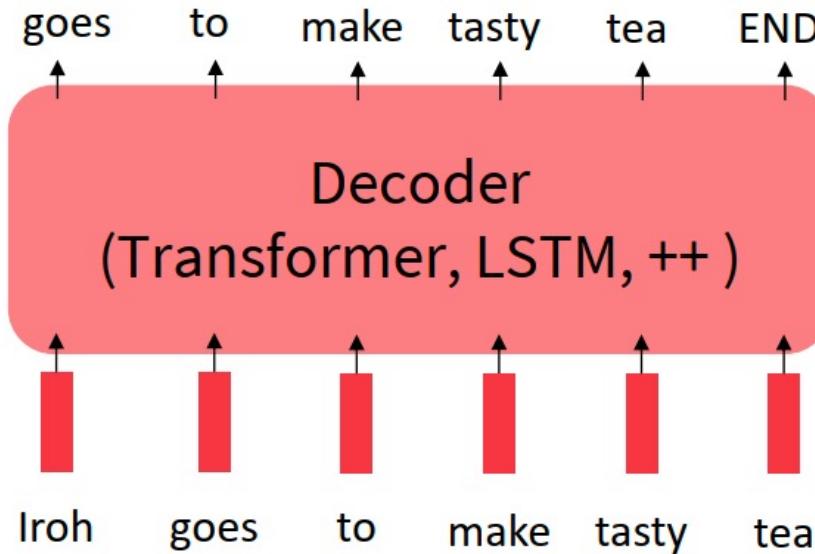


# The Pretraining / Finetuning Paradigm

- 预训练可以以参数初始化的方式来提升NLP模型的效果

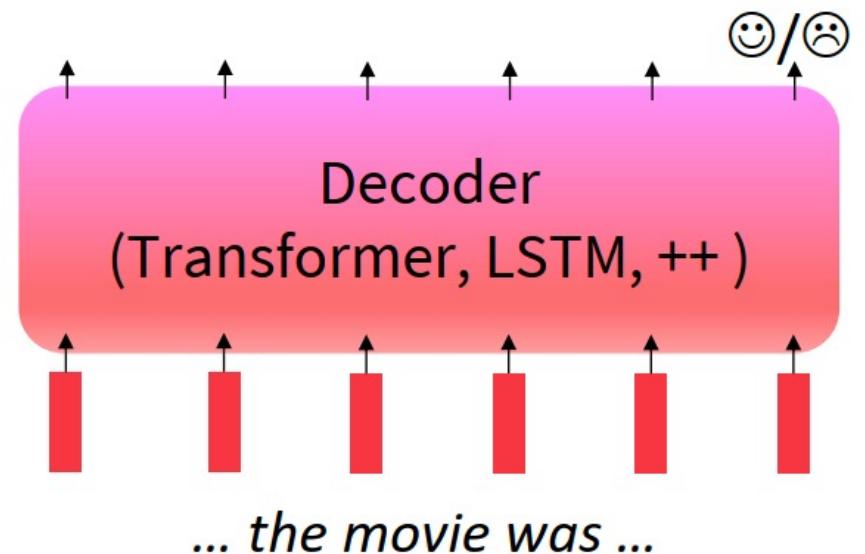
## Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



## Step 2: Finetune (on your task)

Not many labels; adapt to the task!

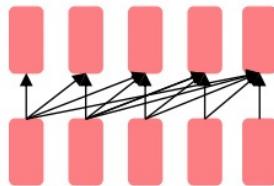


# Contents

- Transformers
  - Impact of Transformers on NLP (and ML more broadly)
  - From Recurrence (RNNs) to Attention-Based NLP Models
  - Understanding the Transformer Model
  - Drawbacks and Variants of Transformers
- Pretraining Language Models(PLMs)
  - Subword modeling
  - Motivating model pretraining from word embeddings
  - Model pretraining three ways
    - Decoders
    - Encoders
    - Encoder-Decoders
  - Very large models and in-context learning

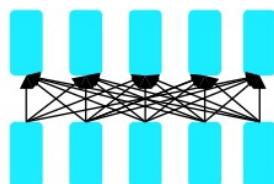
# Pretraining for three types of architectures

- The neural architecture influences the type of pretraining, and natural use cases.



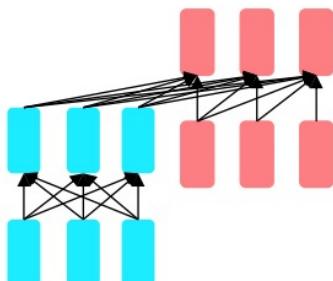
**Decoders**

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words
- Examples: GPT-2, GPT-3, LaMDA, LLaMa



**Encoders**

- Gets bidirectional context – can condition on future!
- Wait, how do we pretrain them?
- Examples: BERT and its many variants, e.g. RoBERTa

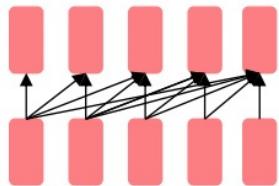


**Encoder-Decoders**

- Good parts of decoders and encoders?
- What's the best way to pretrain them?
- Examples: Transformer, T5, Meena

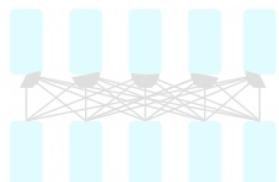
# Pretraining for three types of architectures

- The neural architecture influences the type of pretraining, and natural use cases.



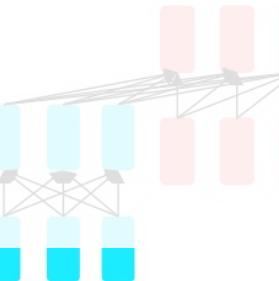
**Decoders**

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words
- Examples: GPT-2, GPT-3, LaMDA



**Encoders**

- Gets bidirectional context – can condition on future!
- Wait, how do we pretrain them?
- Examples: BERT and its many variants, e.g. RoBERTa



**Encoder-Decoders**

- Good parts of decoders and encoders?
- What's the best way to pretrain them?
- Examples: Transformer, T5, Meena

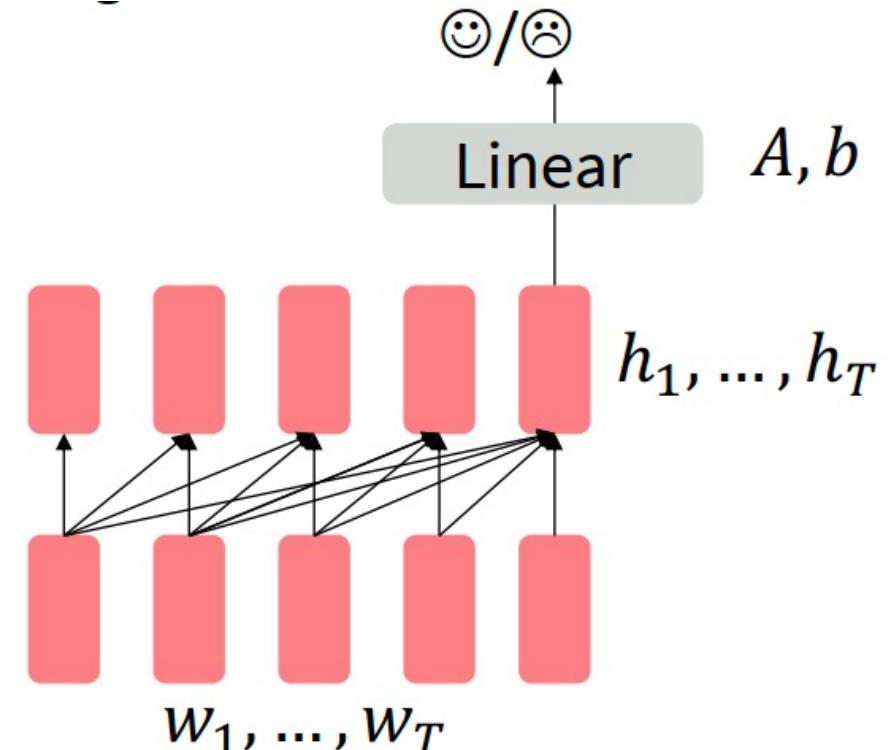
# Pretraining decoders

- 当我们使用利用LM预训练的decoder时，暂时忘记LM是为了建模  $p(w_t|w_{1:t-1})$

- 我们利用最后一个词的隐状态来训练一个分类器，从而微调这个预训练模型

$$\begin{aligned} h_1, \dots, h_T &= \text{Decoder}(w_1, \dots, w_T) \\ y &\sim Ah_T + b \end{aligned}$$

- $A, b$ 是随机初始化的参数，用下游模型来训练
- 梯度回传整个网络



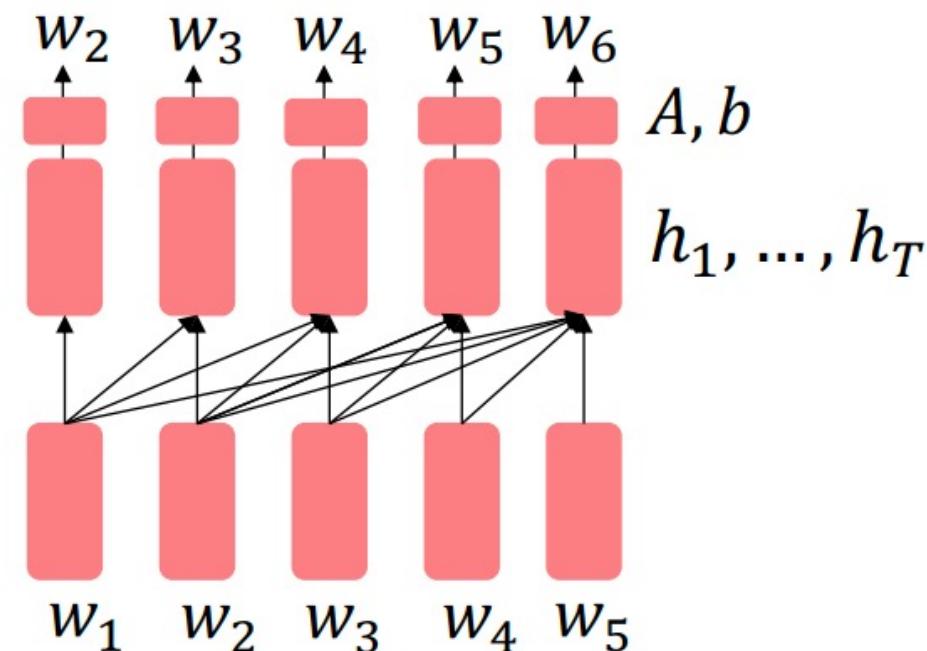
[Note how the linear layer hasn't been pretrained and must be learned from scratch.]

# Pretraining decoders

- 下游任务如果是生成任务，那么很自然的，我们要微调  $p(w_t|w_{1:t-1})$
- 如果任务的输出是个序列，而且词表与预训练数据高度重合，那么将会非常有效
  - 对话（上下文=对话历史）
  - 摘要（上下文=文档）

$$h_1, \dots, h_T = \text{Decoder}(w_1, \dots, w_T)$$
$$w_t \sim Ah_{t-1} + b$$

- $A, b$ 是随机初始化的参数，用下游模型来训练



[Note how the linear layer has been pretrained.]

# Generative Pretrained Transformer (GPT) [Radford et al., 2018]

2018's GPT was a big success in pretraining a decoder!

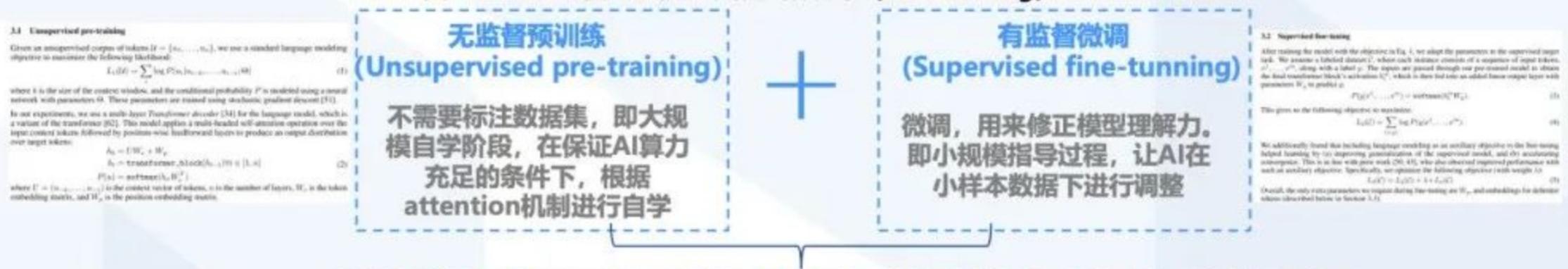
- Transformer decoder with 12 layers.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
- Trained on BooksCorpus: over 7000 unique books.
  - Contains long spans of contiguous text, for learning long-distance dependencies

# GPT-1（预训练+微调）

■ GPT-1模型基于Transformer解除了顺序关联和依赖性的前提，采用生成式模型方式，重点考虑了从原始文本中有效学习的能力，这对于减轻自然语言处理（NLP）中对监督学习的依赖至关重要

- ✓ GPT (Generative Pre-training Transformer) 于2018年6月由OpenAI首次提出。GPT模型考虑到在自然语言理解中有大量不同的任务，尽管大量的未标记文本语料库非常丰富，但用于学习这些特定任务的标记数据却很少，这使得经过区分训练的模型很难充分执行。同时，大多数深度学习方法需要大量手动标记的数据，这限制了它们在许多缺少注释资源的领域的适用性。
- ✓ 在考虑以上局限性的前提下，GPT论文中证明，通过对未标记文本的不同语料库进行语言模型的生成性预训练，然后对每个特定任务进行区分性微调，可以实现这些任务上的巨大收益。和之前方法不同，GPT在微调期间使用任务感知输入转换，以实现有效的传输，同时对模型架构的更改最小。

图29：GPT-1模型的核心手段是预训练（Pre-training）



Thank you