# Transformers

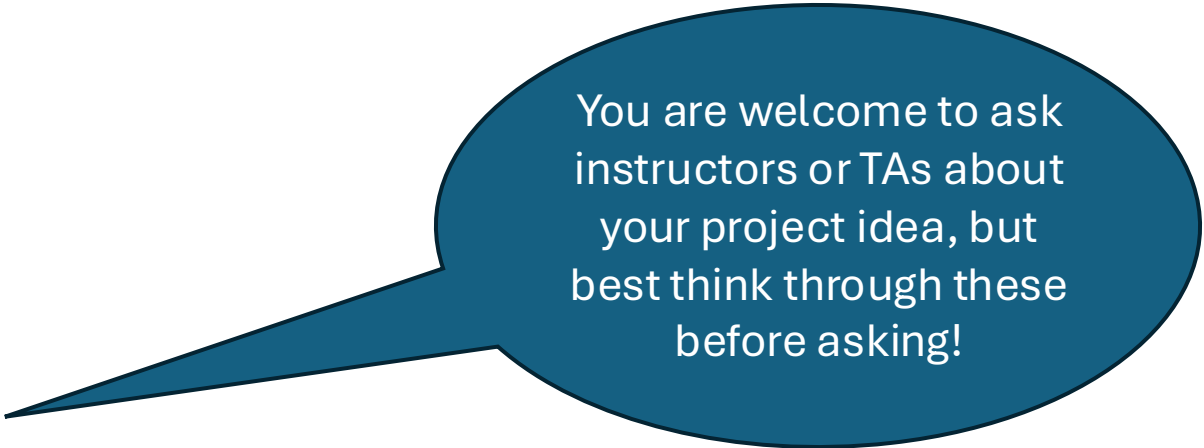CS 6120 Natural Language Processing

Northeastern University

Si Wu

# Logistics

- Coding assignment 2 due tonight.
- Coding assignment 1 grade will be released tonight.


- Today: the architecture of transformers and self-attention
- Next lecture: more in-depth discussion about different kinds of transformer models and masked language models

# More on the class project

Eventually in your final report, there will be these components at least:

- Problem statement,

- Prior /related work

- Data

- Models

- Experiment

- Analysis

- Conclusion

Make sure you pinpoint the exact problem in your problem statement

So your project needs to be sensible and have depth:

- Sensible/reasonable: could be finished before December, not a full-on research project.

- Interesting: it could be an open research project (e.g. poetry generation) that's difficult to evaluate; or, if you work on a very traditional, well-researched area, you need to provide some insight: comparing different approach, testing your own novel method, come up with new models or efficient implementation, or better evaluation that you come up with

- In short, you should have new insight from working on this project, that other people don't have.

You are welcome to ask instructors or TAs about your project idea, but best think through these before asking!

# Introduction

- Some of you already raised concerns in the last lecture about the "efficiency" of RNN:
  - vanishing gradient problem, memory, **parallelizability**.
- And we discussed these drawbacks of RNN, and said that's why we keep inventing new LMs!
  - Parallelizability especially at scale
    - Sure, it's not completely unparallelizable, but due to its sequential nature
  - Vanishing gradient: problematic for long-term dependency

- Today: a much better LM → transformer!
  - Mainly about autoregressive LMs, more will unfold over the next few lectures, e.g. masked language models and bidirectional transformers.

# Components of a transformer

- 3 components:
  - Input encoding components
  - Transformer blocks (each consists of multi-head attention layer, FFNN, layer normalization steps)
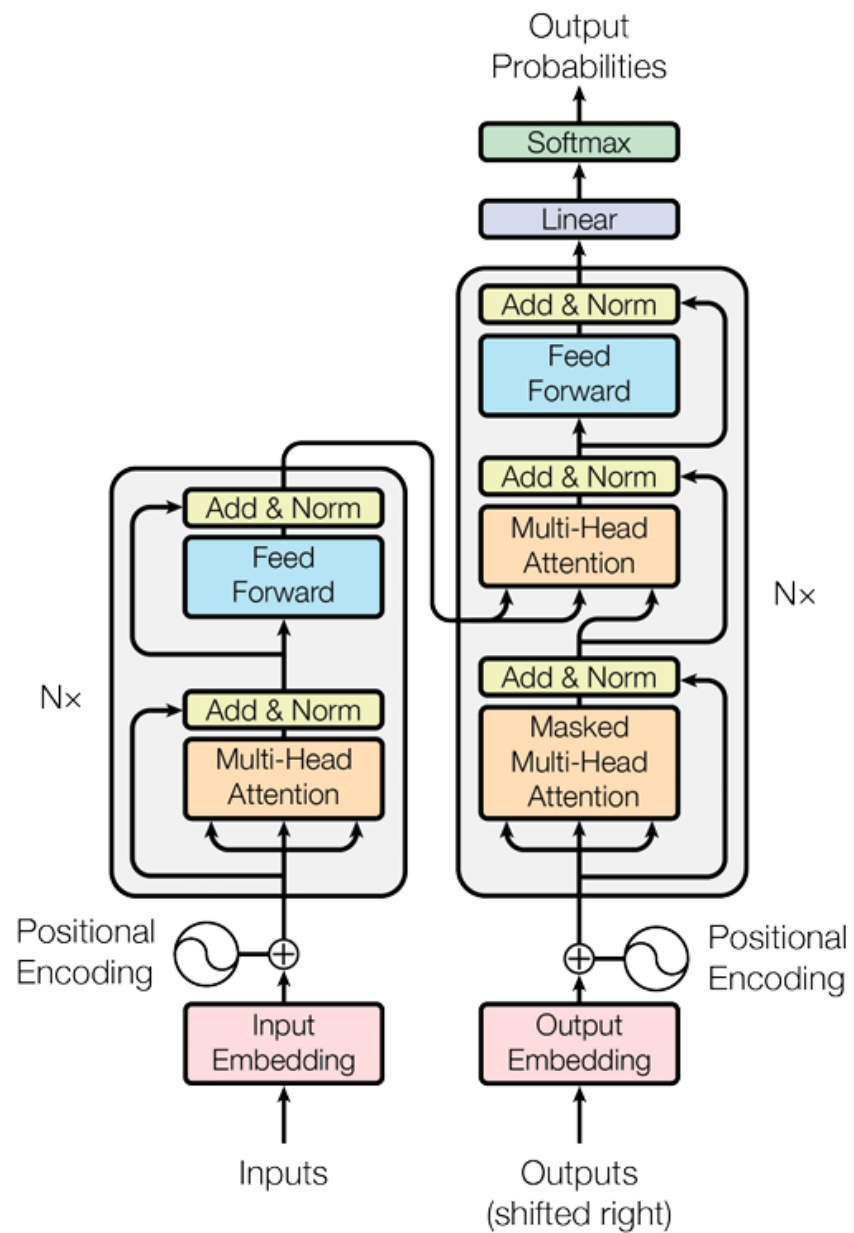  - Language modeling heads

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Feed
Forward

Nx

Add & Norm

Add & Norm

Multi-Head
Attention

Masked
Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

Figure from the original transformer paper Attention Is All You Need

# Attention

# Examples

Generating from left to right, use **are** instead of **is** to refer to keys. (grammar)

The keys to the cabinet are on the table

Imagine if we use "they" to refer to a group of people after a whole paragraph describing a party.

I walked along the pond, and noticed one of the tress along the bank

Recall lecture on linguistics where we talked about ambiguity, dependencies, polysemy. But also long-term dependency from RNN lecture.

Which word sense of bank?

# Many lecture ago

When talking about embeddings, we mentioned those were **static embeddings**, and that we would talk about **contextual embeddings**

Autoregressive generation (left to right)

The chicken didn't cross the road because **it ...**

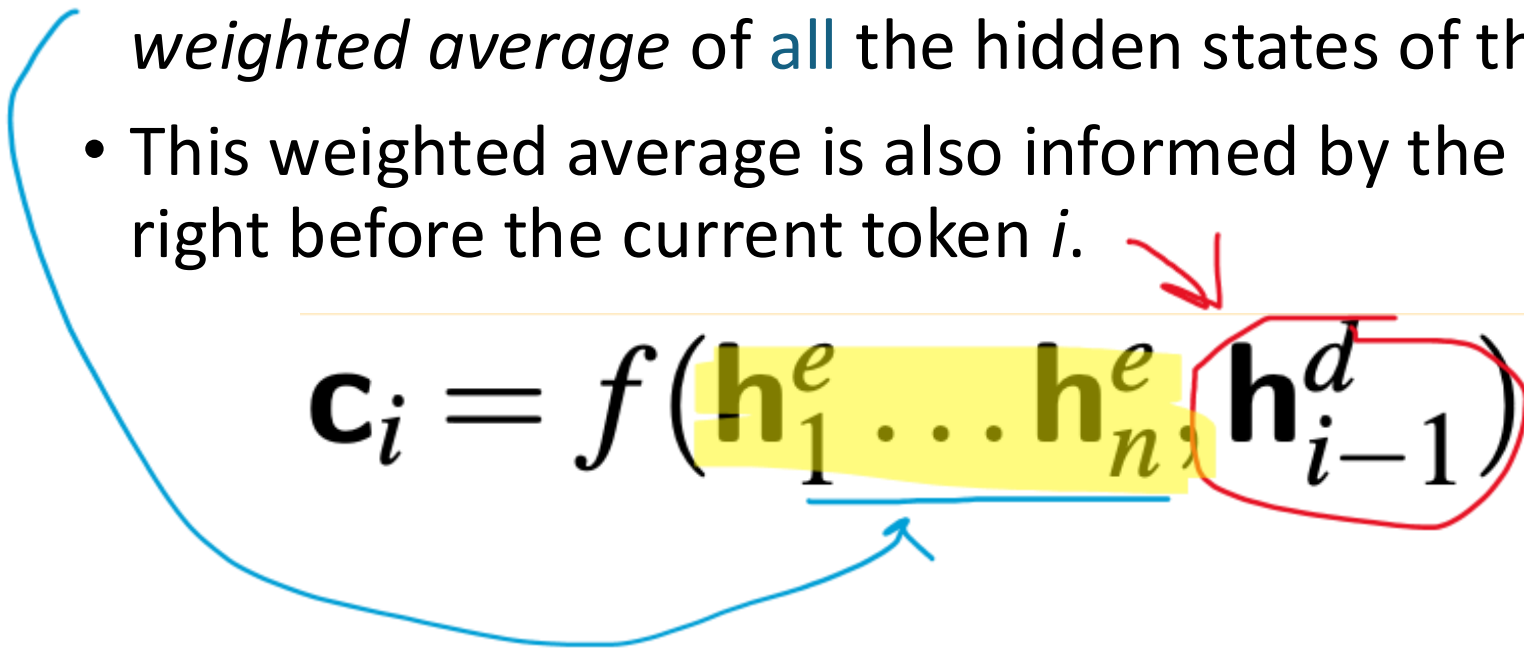The chicken didn't cross the road because **it** was too tired

The chicken didn't cross the road because **it** was too wide

# Attention from last lecture

In the RNN lecture, we introduced attention

# Solution: Attention!

- Instead of being taken from the last hidden state, the **context** it's a *weighted average* of all the hidden states of the encoder.

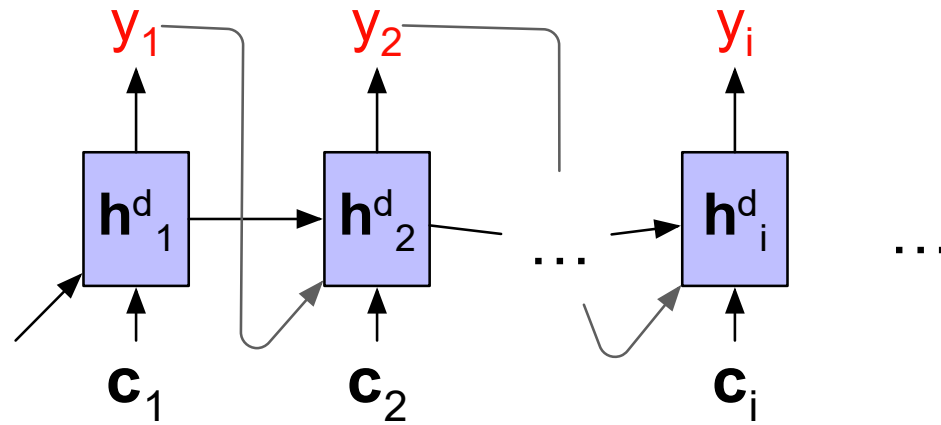- This weighted average is also informed by the state of the decoder right before the current token $i$.

$$c_i = f\left(\mathbf{h}_1^e \ldots \mathbf{h}_n^e, \mathbf{h}_{i-1}^d\right)$$

"weighted average" meaning, $c_i$ can attend to a particular part of the input text that is relevant to token I, which is what the decoder is trying to produce

# Attention

$$h_i^d = g(\hat{y}_{i-1}, h_{i-1}^d, c_i)$$

# How to compute $c_i$ ? How to decide what to pay attention to?

- One way is similarity!
  - Using similarity as a scoring function between last decoder state and each encoder hidden state
  - Simplest such score is **dot-product atten**tion.

For each encoder state j

$$\text{score}(h^d_{i-1}, h^e_j) \ = \ h^d_{i-1} \cdot h^e_j$$

# How to compute $c_i$ ? How to decide what

- We'll normalize these similarity scores of each encoder hidden states with a softmax to create weights $\alpha_{ij}$, that tell us the relevance of encoder hidden state $j$ to hidden decoder state, $h^d_{i-1}$
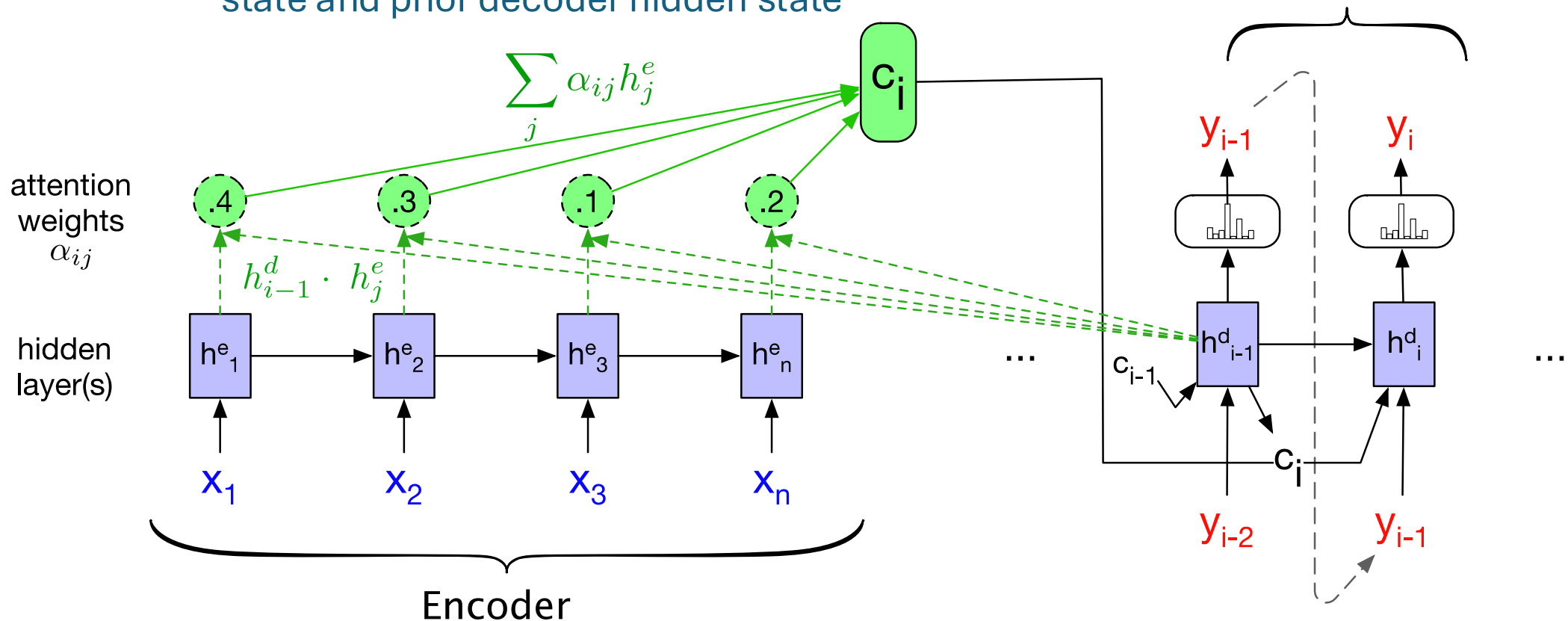
$$a_{ij} = \text{softmax}(\text{score}(h^d_{i-1}, h^e_j))$$

- And then use this to help create a weighted average of all the encoder hidden states:

$$\mathbf{c}_i = \sum_j \alpha_{ij} \, \mathbf{h}^e_j$$

# Encoder-decoder with attention, focusing on the computation of c

# Now, enter

## 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**[*]
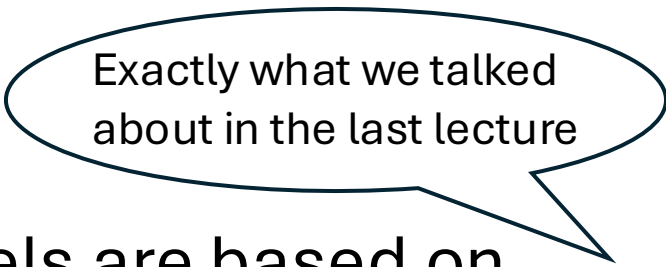Google Research
usz@google.com

**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

At the very beginning of the paper abstract:

Exactly what we talked about in the last lecture

"The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism.

We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train."

# Attention head

Recall the simple intuition of attention: begin by computing the similarity between current word and all previous words using their dot product. The similarity will weight the importance of that previous word to the current word.

The chicken didn't cross the road because **it**

$$x_i$$

Compute similarity with each previous word:

$$score(x_i, x_j) = x_i \cdot x_j$$

$\alpha_{ij}$ is the similarity between current word and a previous word at position j, it will become the weight that signifies its importance/relatedness → how much should we attend to this

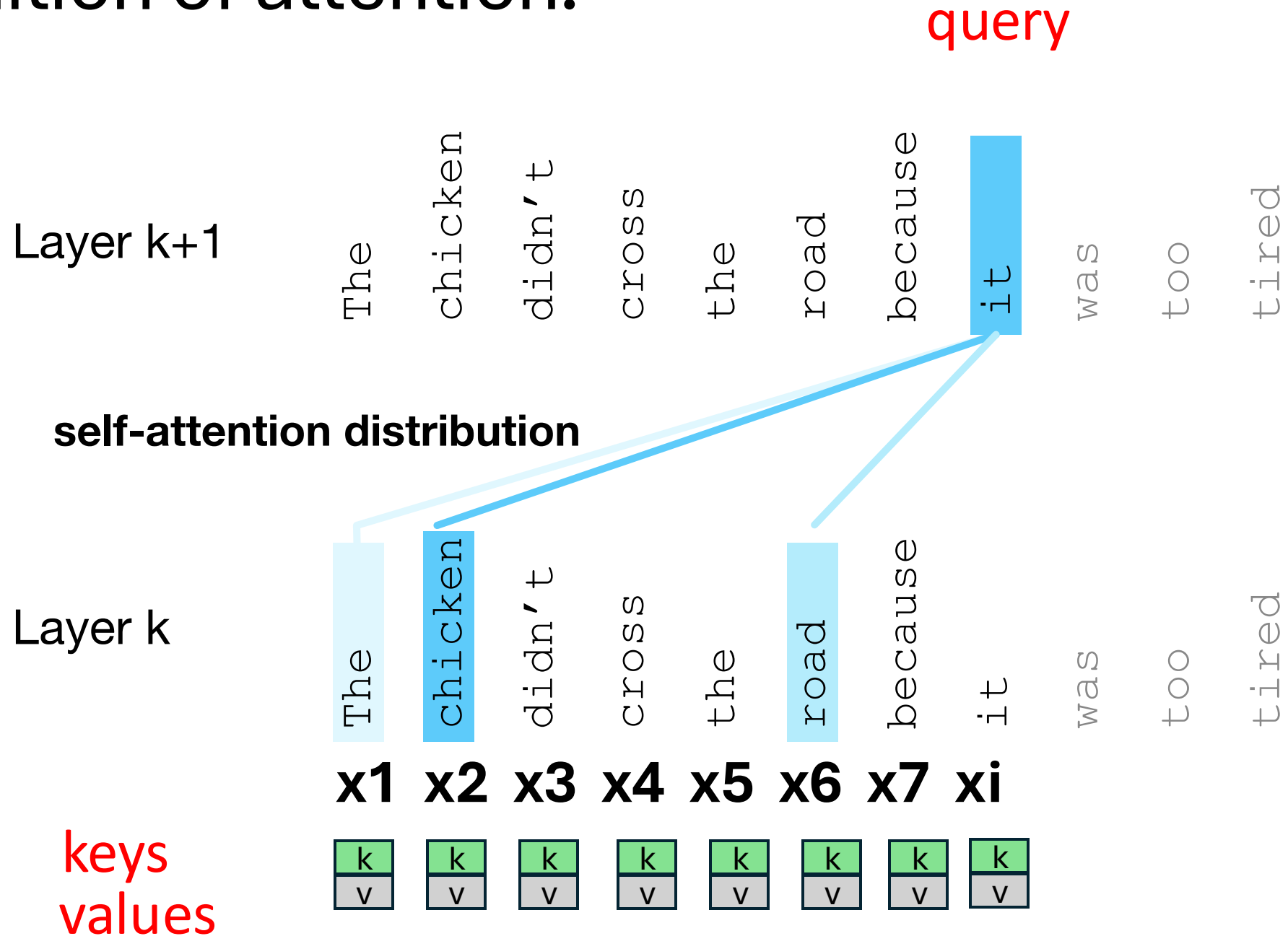$$a_{ij} = softmax(score(x_i, x_j)) \; 8j \leq i$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{x}_j$$    Weight sum

# Attention head

- In transformer, we have a slightly more complicated version of this
  - Attention head
    - "Head" in transformers: refer to specific structured layers
- A single attention head use **query, key, value** matrices.
  - 3 different roles each vector $x_i$ can play
    - Intuitively: current word $x_i$ (query), and key and value of preceding word $x_j$
  - Query: current element
  - Key: a preceding input that's being compared to (for similarity)
  - Value: value of a preceding element (that gets weighted and summed), the actual information, semantic contribution

# Intuition of attention:

# Single-head attention

- A single attention head uses these three matrices:
    - Query $W^Q$
    - Key $W^K$
    - Value $W^V$

    We'll use matrices to project each vector $\mathbf{x}_i$ into a representation of its role as query, key, value:

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W^Q}; \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W^K}; \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W^V}$$

# An Actual Attention Head: slightly more complicated

- Given these 3 representation of $\mathbf{x}_i$

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K; \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

- To compute similarity of current element $\mathbf{x}_i$ with some prior element $\mathbf{x}_j$

- We'll use dot product between $\mathbf{q}_i$ and $\mathbf{k}_j$.

- And instead of summing up $\mathbf{x}_j$, we'll sum up $\mathbf{v}_j$

# Final equations for one attention head

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_j = \mathbf{x}_j \mathbf{W}^K; \quad \mathbf{v}_j = \mathbf{x}_j \mathbf{W}^V$$

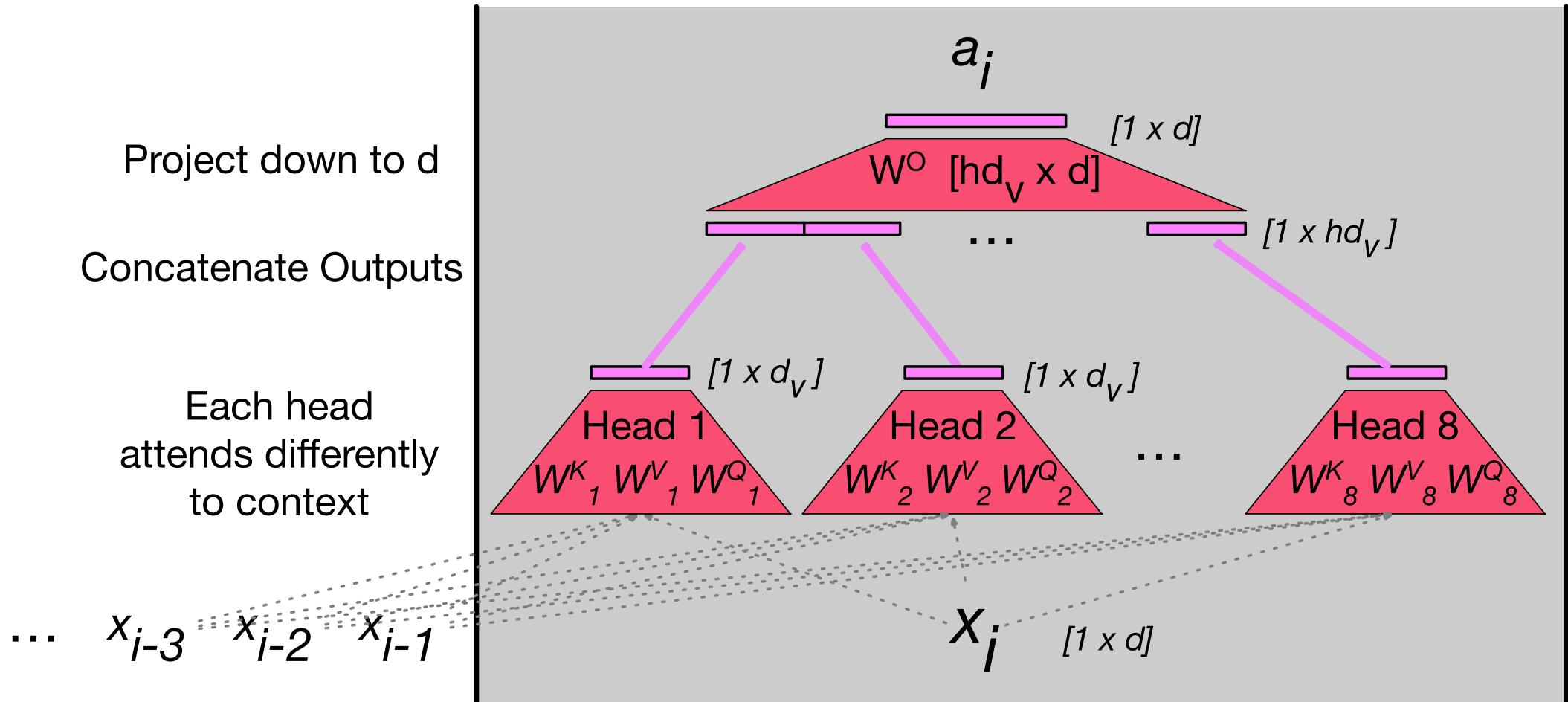$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \; \forall j \leq i$$

$$\mathbf{head}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

$$\mathbf{a}_i = \mathbf{head}_i \mathbf{W}^O$$

# Actual Attention: slightly more complicated

- Instead of one attention head, we'll have lots of them!

- Intuition: each head might be attending to the context for different purposes

  - Different linguistic relationships or patterns in the context

$$\mathbf{q}_i^c = \mathbf{x}_i \mathbf{W^{Qc}}; \quad \mathbf{k}_j^c = \mathbf{x}_j \mathbf{W^{Kc}}; \quad \mathbf{v}_j^c = \mathbf{x}_j \mathbf{W^{Vc}}; \quad \forall\, c \;\; 1 \le c \le h$$

$$\mathrm{score}^c(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i^c \cdot \mathbf{k}_j^c}{\sqrt{d_k}}$$

$$\alpha_{ij}^c = \mathrm{softmax}(\mathrm{score}^c(\mathbf{x}_i, \mathbf{x}_j)) \;\; \forall\, j \le i$$

$$\mathbf{head}_i^c = \sum_{j \le i} \alpha_{ij}^c \mathbf{v}_j^c$$

$$\mathbf{a}_i = (\mathbf{head}^1 \oplus \mathbf{head}^2 \ldots \oplus \mathbf{head}^h) \mathbf{W}^O$$

$$\mathrm{MultiHeadAttention}(\mathbf{x}_i, [\mathbf{x}_1, \cdots, \mathbf{x}_N]) = \mathbf{a}_i$$
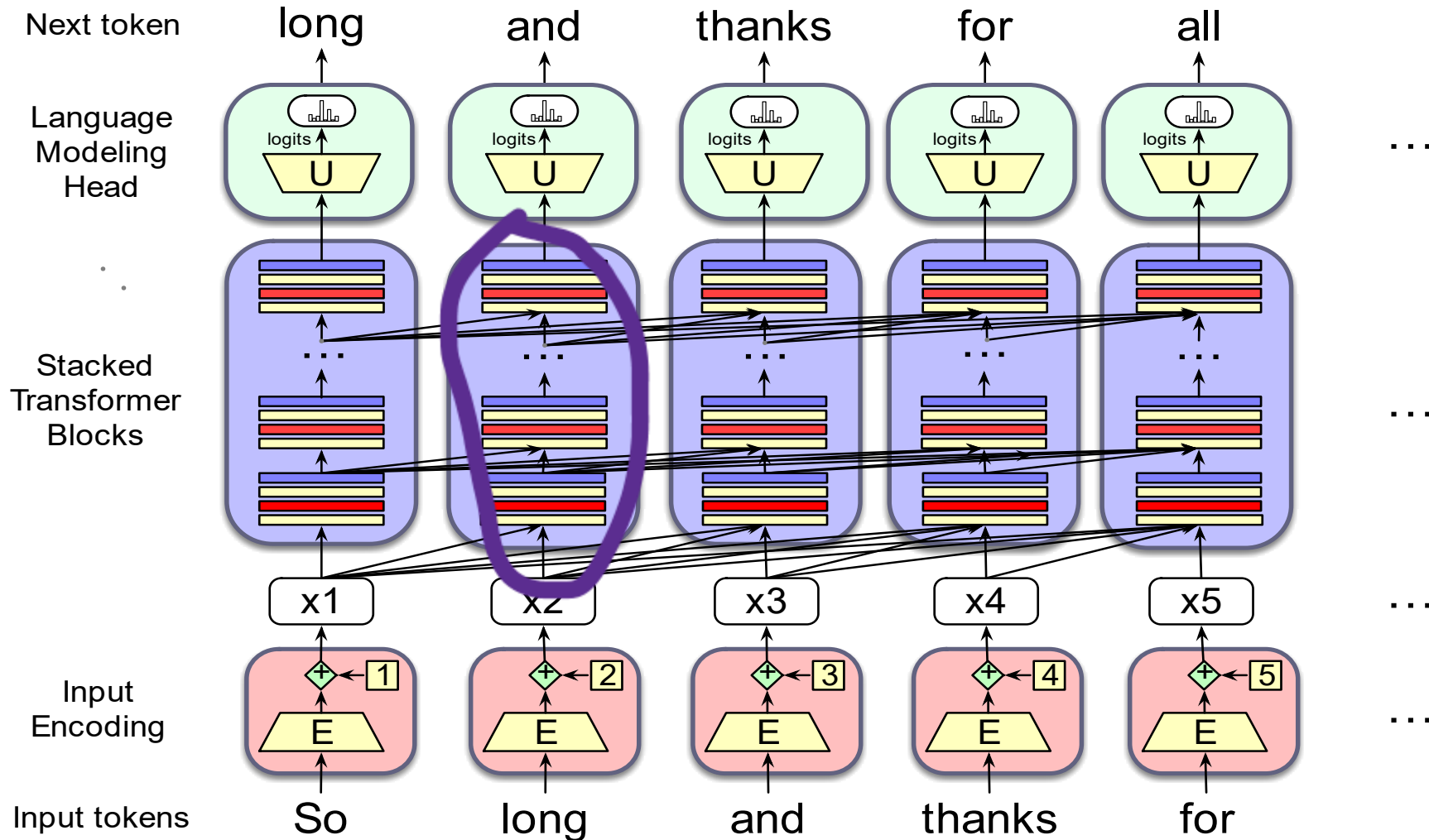
# Multi-head attention

$a_i$

$[1 \times d]$

Project down to d

$W^O \ [hd_v \times d]$

Concatenate Outputs

... $[1 \times hd_v]$

Each head attends differently to context

$[1 \times d_v]$ $[1 \times d_v]$

| Head 1 | Head 2 | ... | Head 8 |
| $W^K_1 \ W^V_1 \ W^Q_1$ | $W^K_2 \ W^V_2 \ W^Q_2$ | | $W^K_8 \ W^V_8 \ W^Q_8$ |

... $x_{i-3}$ $x_{i-2}$ $x_{i-1}$ $x_i$ $[1 \times d]$

# Summary

- Attention is a method for enriching the representation of a token by incorporating contextual information

- The result: the embedding for each word will be different in different contexts!

- Contextual embeddings: a representation of word meaning in its context.

# Transformer blocks

# Transformer language model

# In a transformer block

- In addition to the self-attention layer, there are 3 other layers: feedforward layer, residual connections, and a normalizing layer (aka. Layer norm)

- And we can stack these transformer blocks, but each block is consists of these things.

- One way of thinking about the block is called the **residual stream** (Elhage et al, 2021).

# The residual stream

- Each individual token has their own stream for processing

- A single stream starts with the original input vector, then branch out to different components (layer norm, MHA, feedforward), but we add their output back into the stream

- Initial embedding of a token $x_i$ at position $i$ is of dimensionality d.

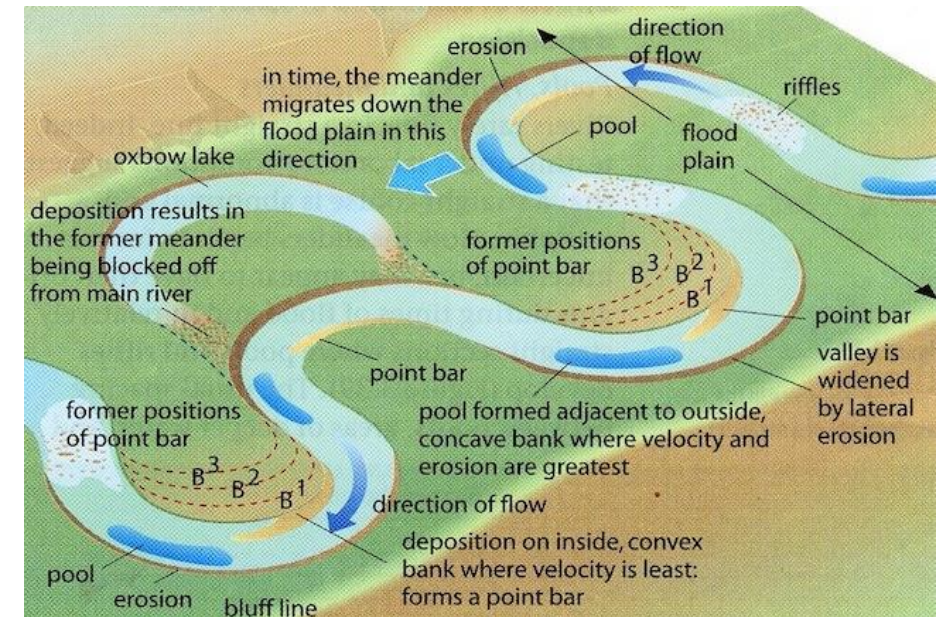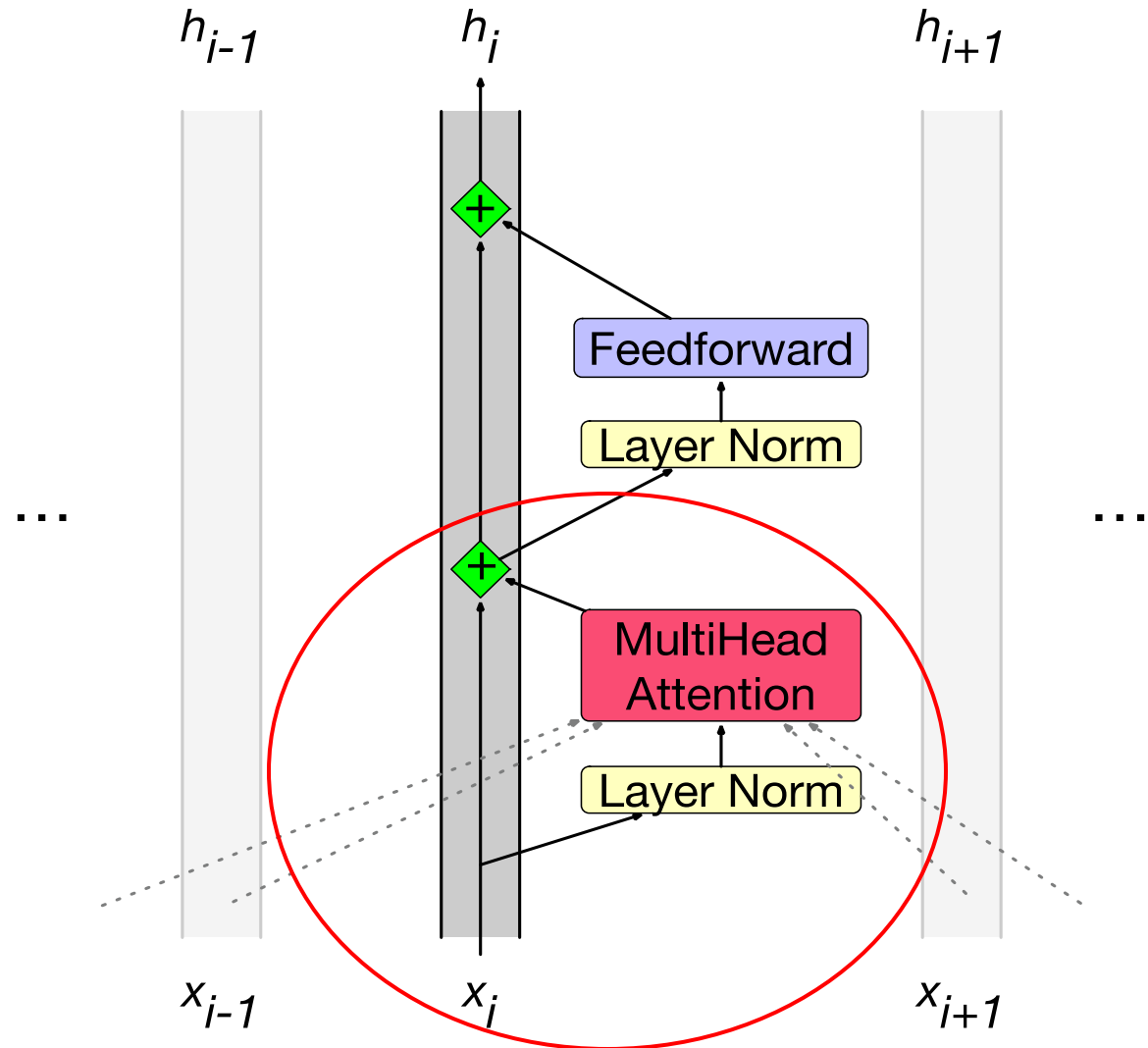- The final output of the transformer block (end of this stream) for token $i$ is $h_i$
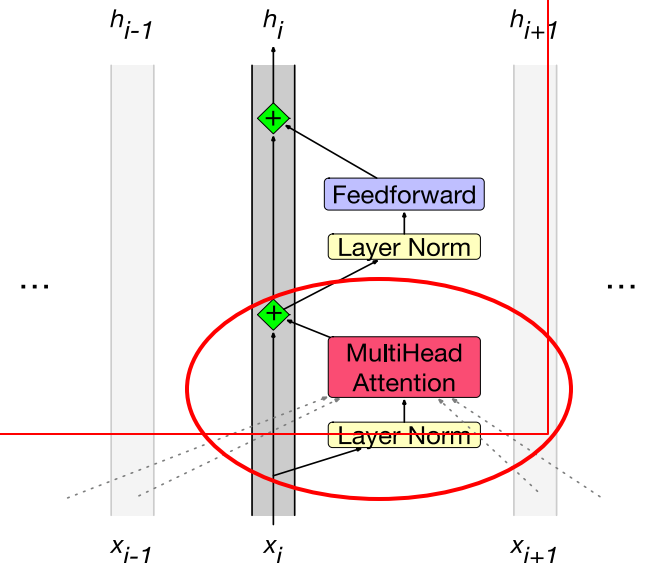


Figure from https://www.ausableriver.org/blog/why-do-streams-meander

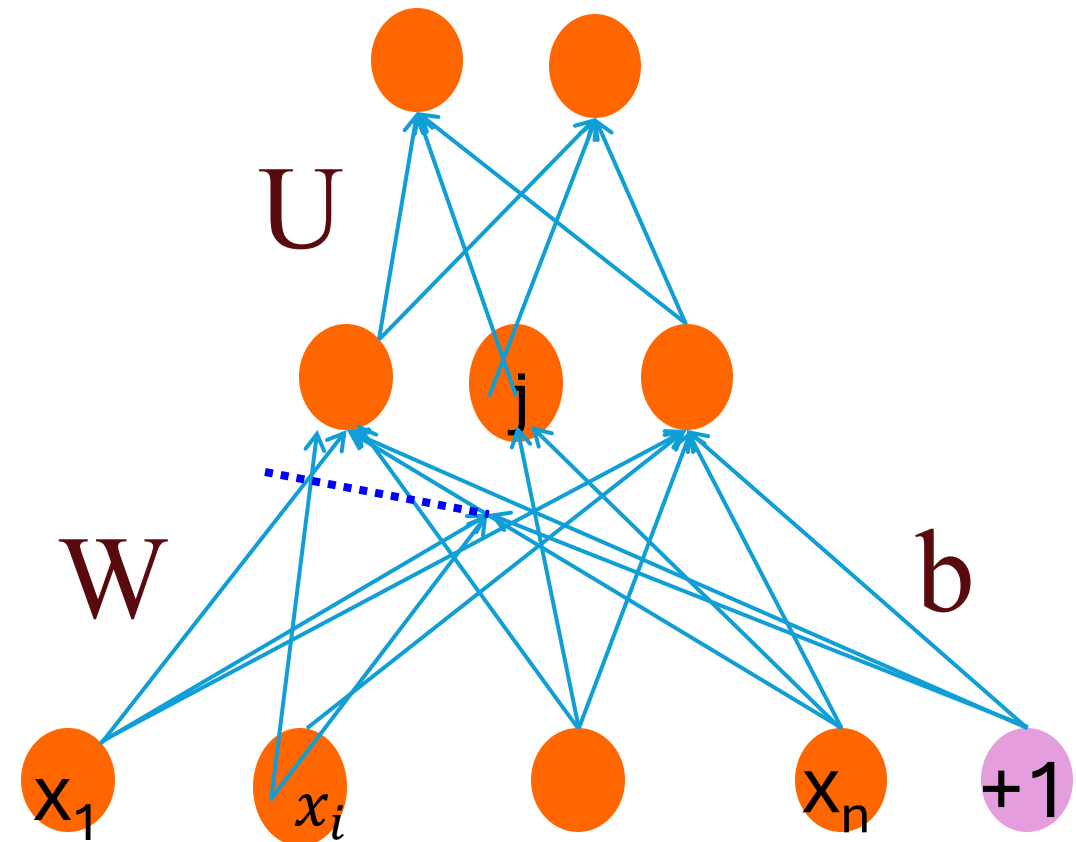# The residual stream: each token gets passed up and modified

# Minor clarification: Pre-norm vs post-norm

- Just different variations

- Post-norm: input → MHA → Residual Add → layer norm
  - OG transformer in their paper (Vaswani et al, 2017)
  - Takes raw hidden states of token, no layer norm yet

- Pre-norm: input → layer norm → MHA → residual add
  - GPT 2/3, some BERT, etc.
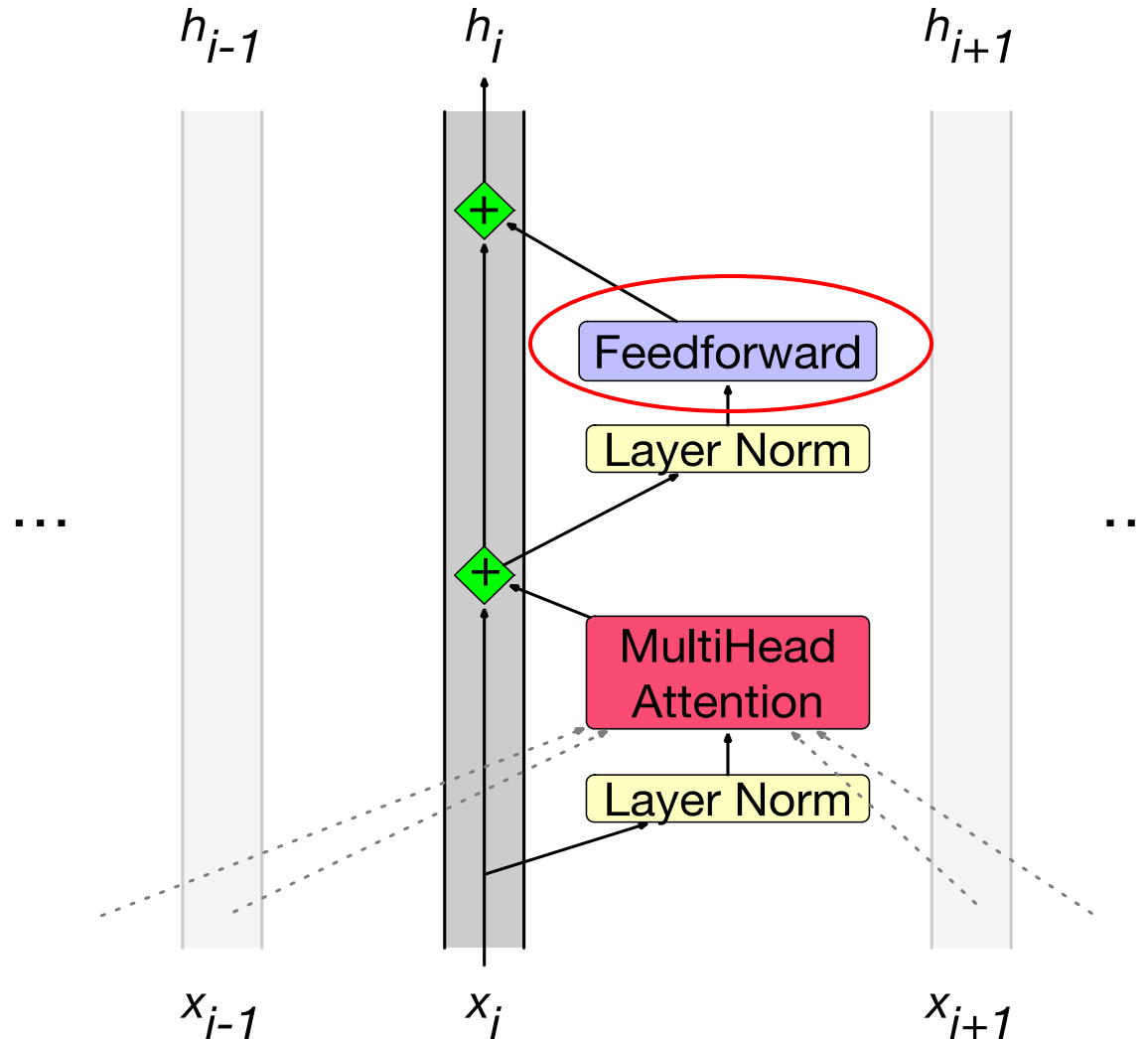  - Take the layer normalized hidden states of tokens

# The feedforward layer in residual stream/ a transformer block

- It's a fully-connected 2-layer network
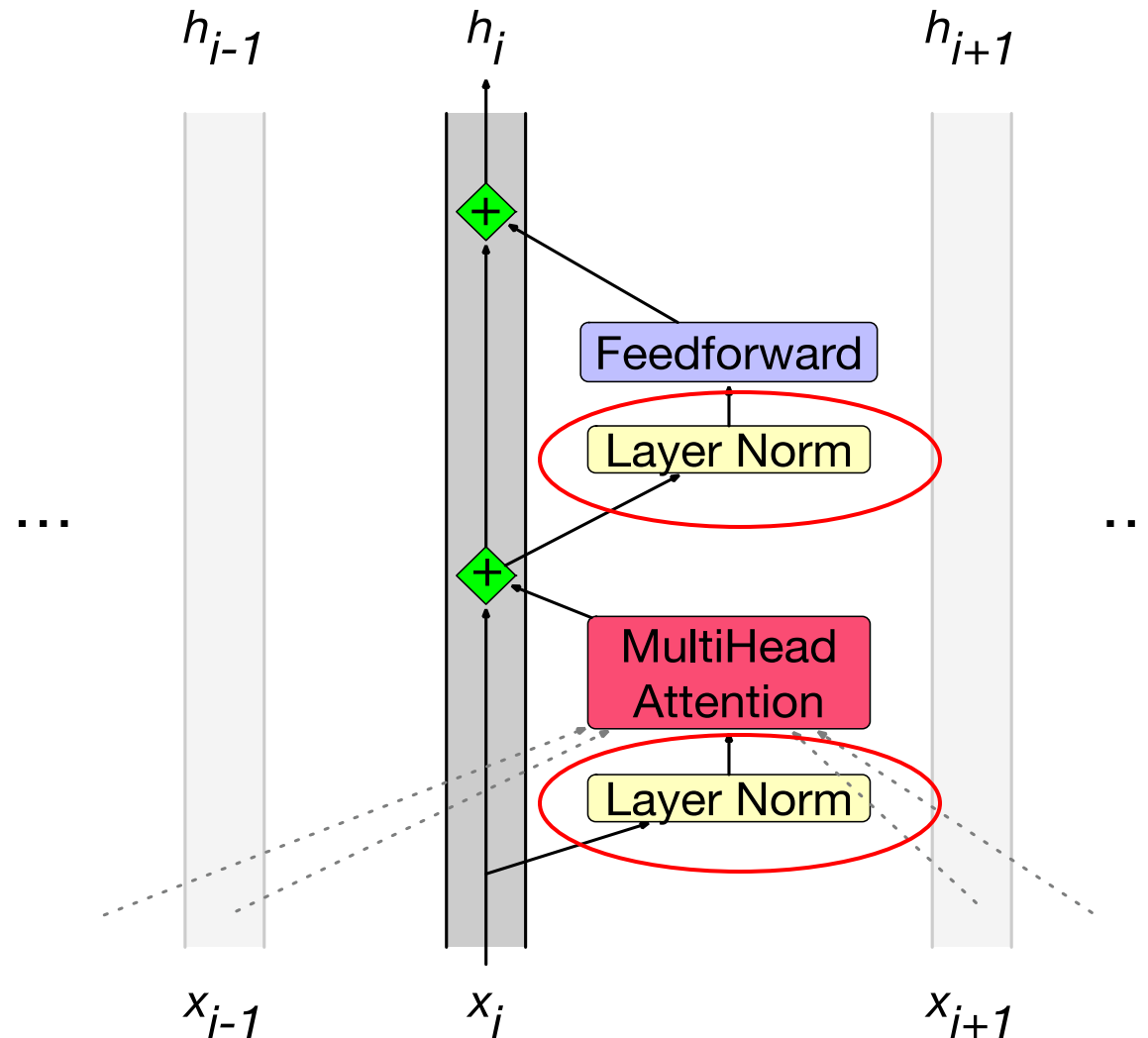  - 1 hidden layer
  - 2 weight matrices

# We'll need nonlinearities, so a feedforward layer

$$\text{FFN}(\mathbf{x}_i) = \text{ReLU}(\mathbf{x}_i \mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2$$

# Layer norm: the vector $\mathbf{x}_i$ is normalized twice

# Layer Norm

Layer norm is a variation of the z-score from statistics, applied to a single vector in a hidden layer

$$\mu \;=\; \frac{1}{d}\sum_{i=1}^{d} x_i$$

$$\sigma \;=\; \sqrt{\frac{1}{d}\sum_{i=1}^{d}(x_i - \mu)^2}$$

One of many form of normalization that can help improve training performance in deep neural network.

Input and output of layer norm are both of dimensionality $d$
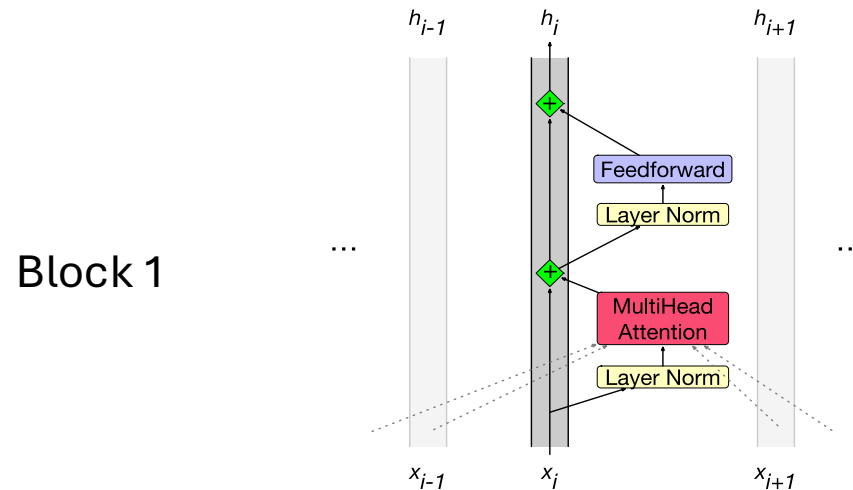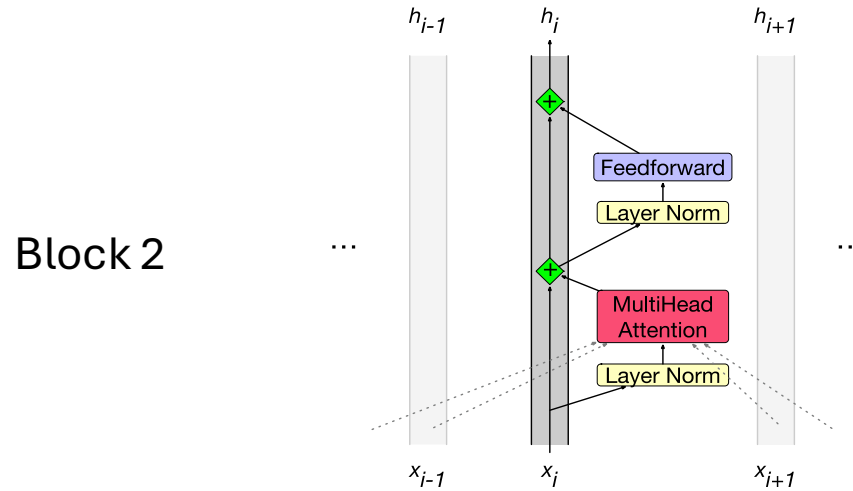
mean

$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma}$$

Standard deviation

$$\text{LayerNorm}(\mathbf{x}) = \gamma\,\frac{(\mathbf{x} - \mu)}{\sigma} + \beta$$

# Putting together a single transformer block



$$\mathbf{t}_i^1 = \text{LayerNorm}(\mathbf{x}_i)$$

$$\mathbf{t}_i^2 = \text{MultiHeadAttention}(\mathbf{t}_i^1, [\mathbf{x}_1^1, \cdots, \mathbf{x}_N^1])$$

$$\mathbf{t}_i^3 = \mathbf{t}_i^2 + \mathbf{x}_i$$

$$\mathbf{t}_i^4 = \text{LayerNorm}(\mathbf{t}_i^3)$$

$$\mathbf{t}_i^5 = \text{FFN}(\mathbf{t}_i^4)$$
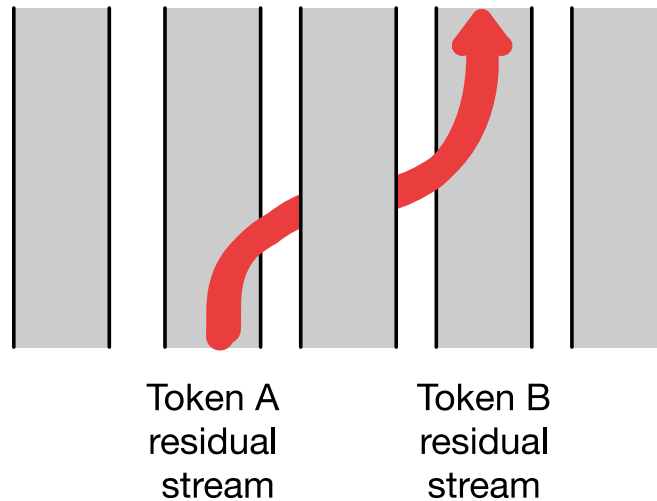
$$\mathbf{h}_i = \mathbf{t}_i^5 + \mathbf{t}_i^3$$

MHA is the only component that takes input from other tokens

# A transformer is a stack of these blocks
## so all the vectors are of the same dimensionality d

# Residual streams and attention

- Notice that all parts of the transformer block apply to 1 same token.

- Except multi-head attention component, which takes information from other tokens as well. (because we need context!)

- Elhage et al. (2021) show that we can view attention heads as literally moving information from the residual stream of a neighboring token into the current stream .



Token A
residual
stream

Token B
residual
stream

# LLMs and transformer blocks

- In large language models, there are usually stacks of these transformer blocks
  - From 12 layers: T5, GPT-3-small
  - To 96 layers: GPT-3-large
  - Even more in recent models
- Once we stack many transformer blocks, at the very end of the last transformer block, there's a single extra layer norm after the last $h_i$ (output of that block)

# Parallelizing attention

# Parallelizing attention

- The computation of a token's $a_i$ is independent of the computation of other token's

- So are all the computations in a transformer block

- That means we can easily parallelize this entire computation and take advantage of the efficient matrix multiplication routines.

# Input embeddings for N tokens

- Each row is a token embedding of $d$ dimensions.

- We have N tokens

- So a matrix $X$ with $N$ tokens is of size $[N \times d]$ (row x col)

- So a matrix with 32k tokens (in other words, input length 32k) is of size $[32k \times d]$

# Each weight matrices

- $W^Q$ → model learn how to ask questions
- $W^K$ → model learn how to match relevant information
- $W^V$ → model learn what content to pass on once we matched the keys


- These are trainable matrices. Without them we are just attending to token based on their raw embedding similarity
- You can project the X to lower-dimension space $d_k$ which is cheaper

# Let's start with a single attention head

- We have input matrix X, and we multiply X by query, key, and value matrices

$$\mathbf{Q} = \mathbf{XW}^{\mathbf{Q}}; \;\; \mathbf{K} = \mathbf{XW}^{\mathbf{K}}; \;\; \mathbf{V} = \mathbf{XW}^{\mathbf{V}}$$

# QK$^T$

- Now can do a single matrix multiplication to combine Q and K$^T$

| | | | |
|---|---|---|---|
| q1·k1 | q1·k2 | q1·k3 | q1·k4 |
| q2·k1 | q2·k2 | q2·k3 | q2·k4 |
| q3·k1 | q3·k2 | q3·k3 | q3·k4 |
| q4·k1 | q4·k2 | q4·k3 | q4·k4 |

N (left axis label)

N (bottom axis label)

# Parallelizing attention

- Scale the scores, take the softmax, and then multiply the result by V resulting in a matrix of shape *N × d*
  - An attention vector for each input token

$$\mathbf{A} = \mathrm{softmax}\left(\mathrm{mask}\left(\frac{\mathbf{QK}^{\top}}{\sqrt{d_k}}\right)\right)\mathbf{V}$$

# Masking out the future

$$A = \text{softmax}\left(\text{mask}\left(\frac{QK^\top}{\sqrt{d_k}}\right)\right)V$$

- What is this mask function?
  $QK^\top$ has a score for each query dot every key, *including those that follow the query*.

- Guessing the next word is pretty simple if you already know it!

# Masking out the future

$$A = \text{softmax}\left(\text{mask}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}}\right)\right)\mathbf{V}$$

- set $-\infty$ to cells in upper triangle
- The softmax will turn it to 0

N

| q1·k1 | $-\infty$ | $-\infty$ | $-\infty$ |
|-------|-----------|-----------|-----------|
| q2·k1 | q2·k2 | $-\infty$ | $-\infty$ |
| q3·k1 | q3·k2 | q3·k3 | $-\infty$ |
| q4·k1 | q4·k2 | q4·k3 | q4·k4 |

N

# Another point: Attention is quadratic in length

$$\mathbf{A} \;=\; \mathrm{softmax}\left(\mathrm{mask}\left(\frac{\mathbf{QK}^{\top}}{\sqrt{d_k}}\right)\right)\mathbf{V}$$
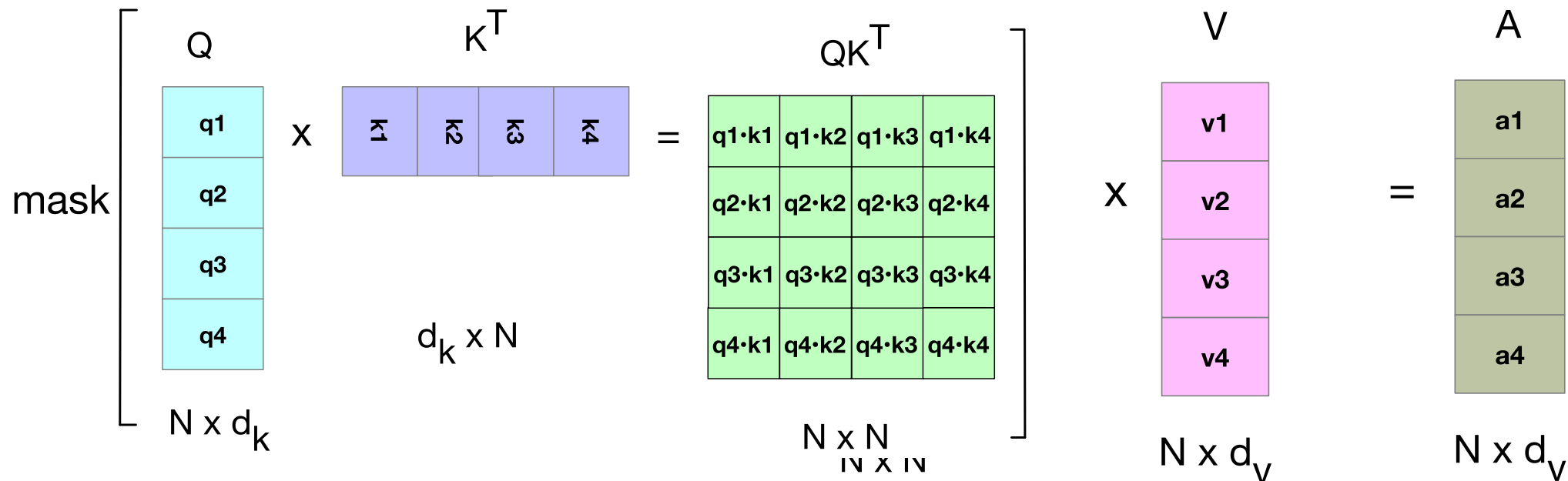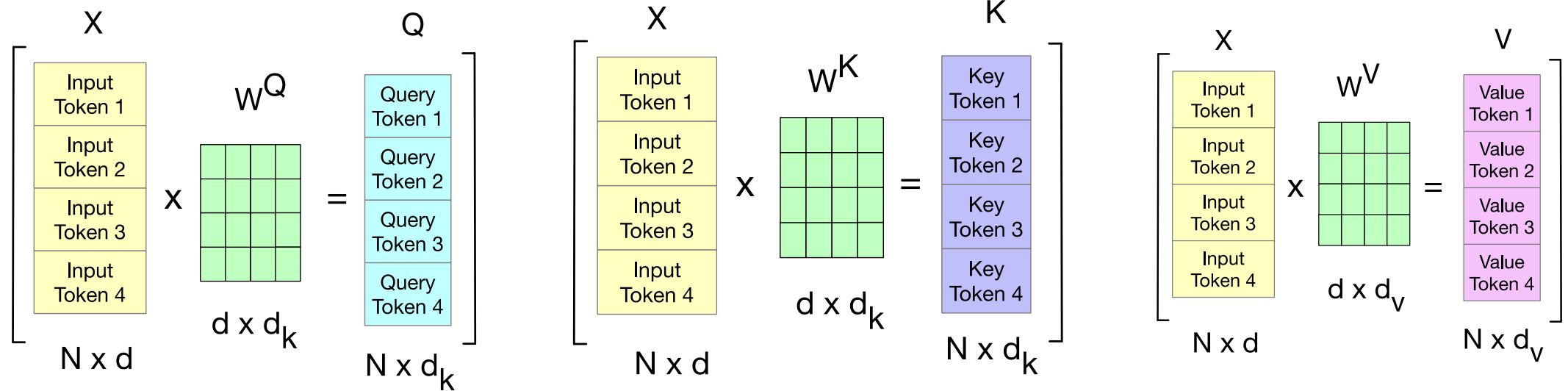
dot products
between each pair
of tokens in the
input.

N

| | | | |
|---|---|---|---|
| q1·k1 | −∞ | −∞ | −∞ |
| q2·k1 | q2·k2 | −∞ | −∞ |
| q3·k1 | q3·k2 | q3·k3 | −∞ |
| q4·k1 | q4·k2 | q4·k3 | q4·k4 |

N

# Single-head attention with N inputs

# Parallelizing Multi-head Attention

Each attention head performs attention independently, and we concatenate them at the end, and multiply by $W^o$

$$\mathbf{Q^i} = \mathbf{XW^{Qi}} ;\quad \mathbf{K^i} = \mathbf{XW^{Ki}} ;\quad \mathbf{V^i} = \mathbf{XW^{Vi}}$$

$$\mathbf{head}_i = \mathrm{SelfAttention}(\mathbf{Q^i}, \mathbf{K^i}, \mathbf{V^i}) = \mathrm{softmax}\left(\frac{\mathbf{Q^i K^{i\top}}}{\sqrt{d_k}}\right)\mathbf{V^i}$$

$$\mathrm{MultiHeadAttention}(\mathbf{X}) = (\mathbf{head}_1 \oplus \mathbf{head}_2 \ldots \oplus \mathbf{head}_h)\mathbf{W^O}$$

Attention heads can be parallelized, but not the stacked blocks, these are sequential.

# Parallelizing Multi-head Attention

$$\mathbf{O} = \text{LayerNorm}(\mathbf{X} + \text{MultiHeadAttention}(\mathbf{X}))$$

$$\mathbf{H} = \text{LayerNorm}(\mathbf{O} + \text{FFN}(\mathbf{O}))$$

- or

$$\mathbf{T}^1 = \text{MultiHeadAttention}(\mathbf{X})$$

$$\mathbf{T}^2 = \mathbf{X} + \mathbf{T}^1$$

$$\mathbf{T}^3 = \text{LayerNorm}(\mathbf{T}^2)$$

$$\mathbf{T}^4 = \text{FFN}(\mathbf{T}^3)$$

$$\mathbf{T}^5 = \mathbf{T}^4 + \mathbf{T}^3$$

$$\mathbf{H} = \text{LayerNorm}(\mathbf{T}^5)$$
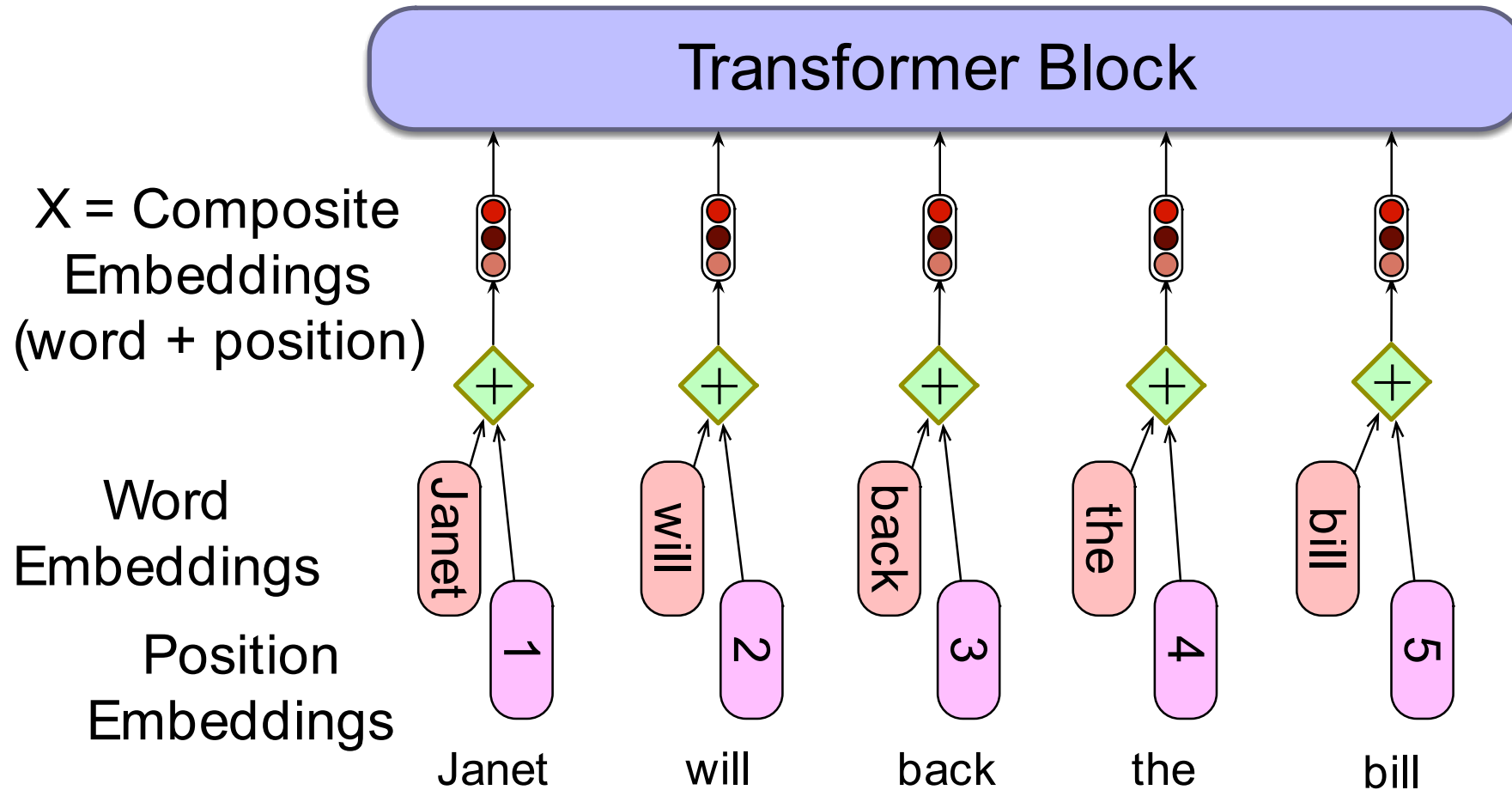
# Input embeddings

# Token and Position Embeddings

- The matrix X (of shape [*N* × *d*]) has an embedding for each word in the context.

- This embedding is created by adding two distinct embedding for each input

  - token embedding

  - positional embedding

# Position Embeddings

- There are many methods, but we'll just describe the simplest: absolute position.

- Goal: learn a position embedding matrix $E$pos of shape $[1 \times N]$.

- Start with randomly initialized embeddings

- one for each integer up to some maximum length.

- i.e., just as we have an embedding for token *fish*, we'll have an embedding for position 3 and position 17.

- As with word embeddings, these position embeddings are learned along with other parameters during training.

# Each **x** is just the sum of word and position embeddings

# Pointers

- We will continue this in the next lecture.

- If you haven't read it in your deep learning class, the original paper is worth reading! Though it is not a required reading for this class.

- Once again, the textbook is wonderful if you want to slowly go over any concepts with more details