

Self-Attention, MLP, and Interpretability

Haw-Shiuan Chang

Logistics

- TA Nguyen Tran permanently moved his TA office hour
 - from Fri 3pm-4pm
 - to Fri 4pm-5pm (1 hour later)

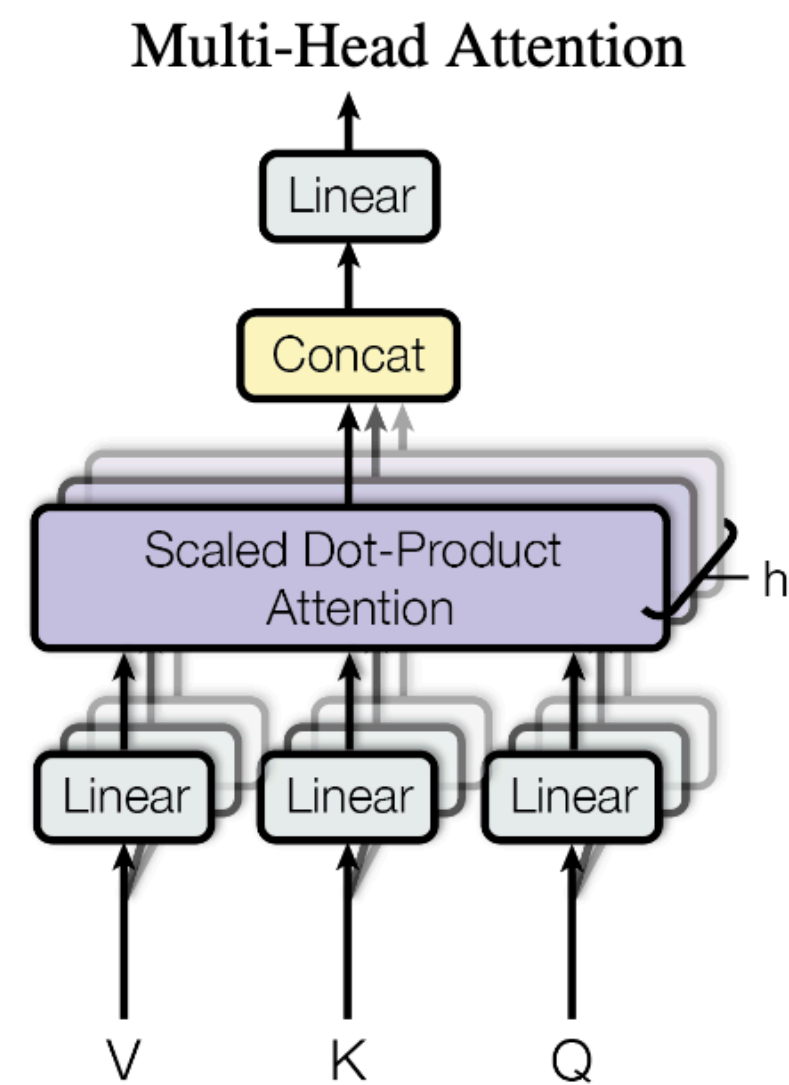
Deadlines

- **<https://people.cs.umass.edu/~hschang/cs685/schedule.html>**
- **3/3:** Quiz 2 due
 - I might extend the deadline on Piazza (if I cannot finish teaching cross-attention this Wednesday)
- **3/7:** Project proposals due
 - If you cannot reach all of your group members this week, please report the situation on Piazza privately
 - If you have some project ideas, you can go to TA office hours to seek some feedback.
- **3/14:** HW 1 due
 - It's about BERT, but you should start the annotation early

**Please stop me if you don't
understand!**

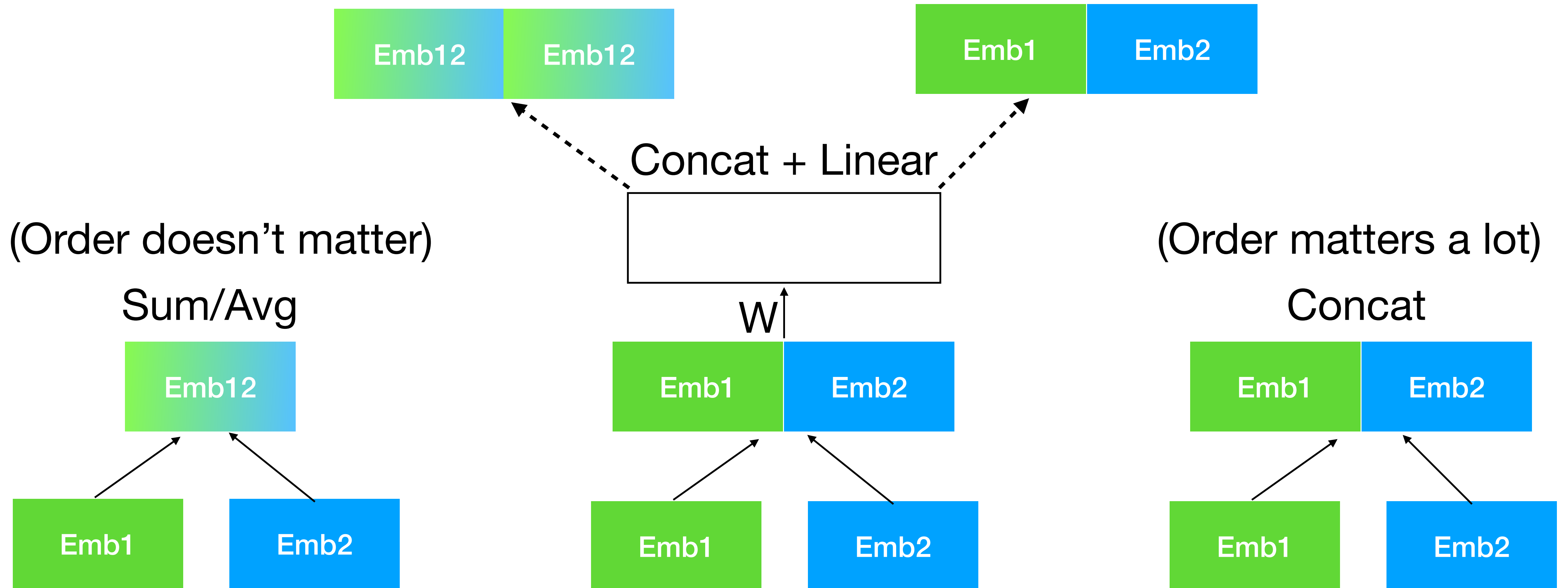
This Lecture will be more advanced, I will explain more slowly

Last Year Note Review



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

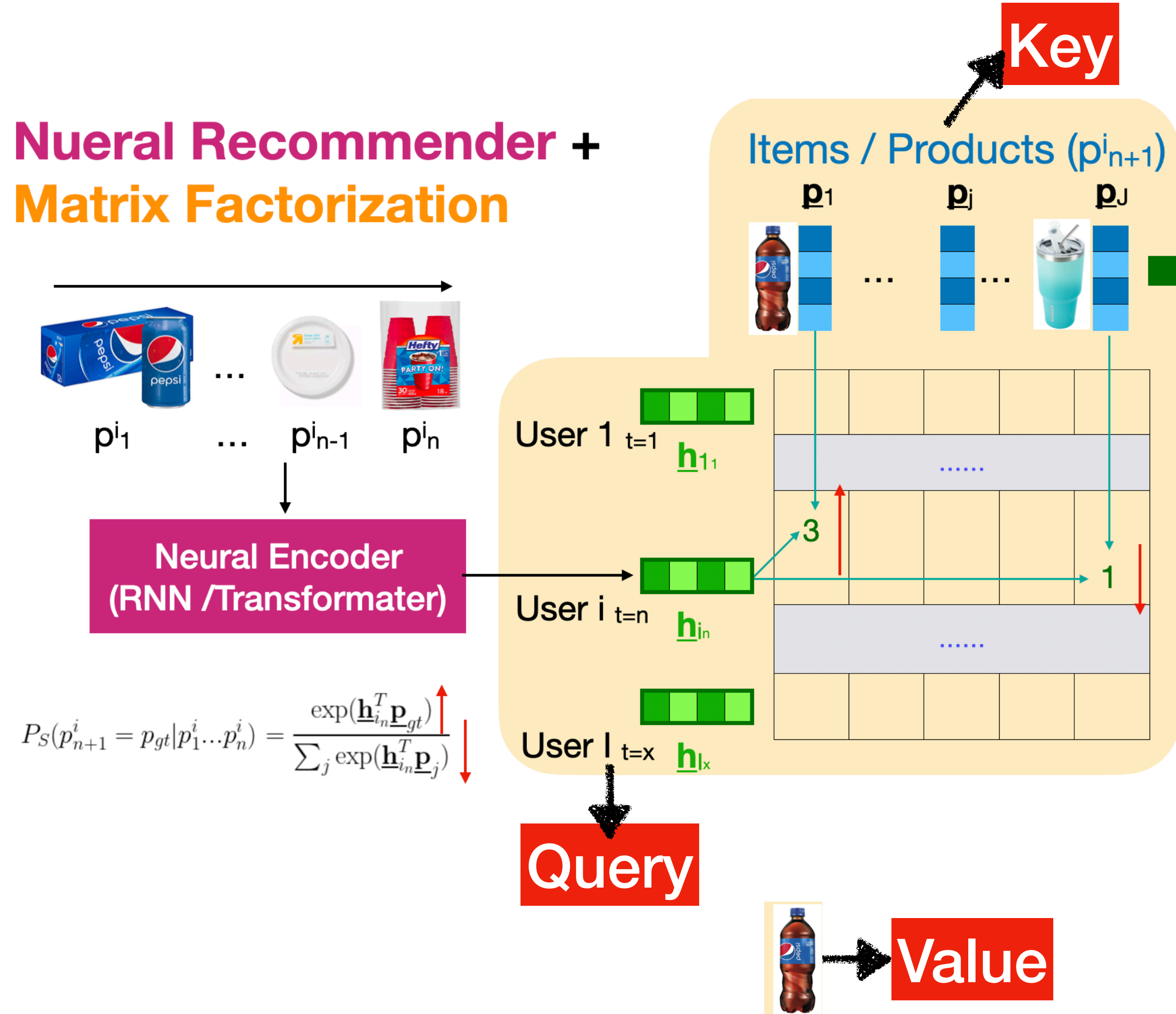
Merging Two Embeddings



Other less common merging includes max, min, product, minus, ...

A Metaphor

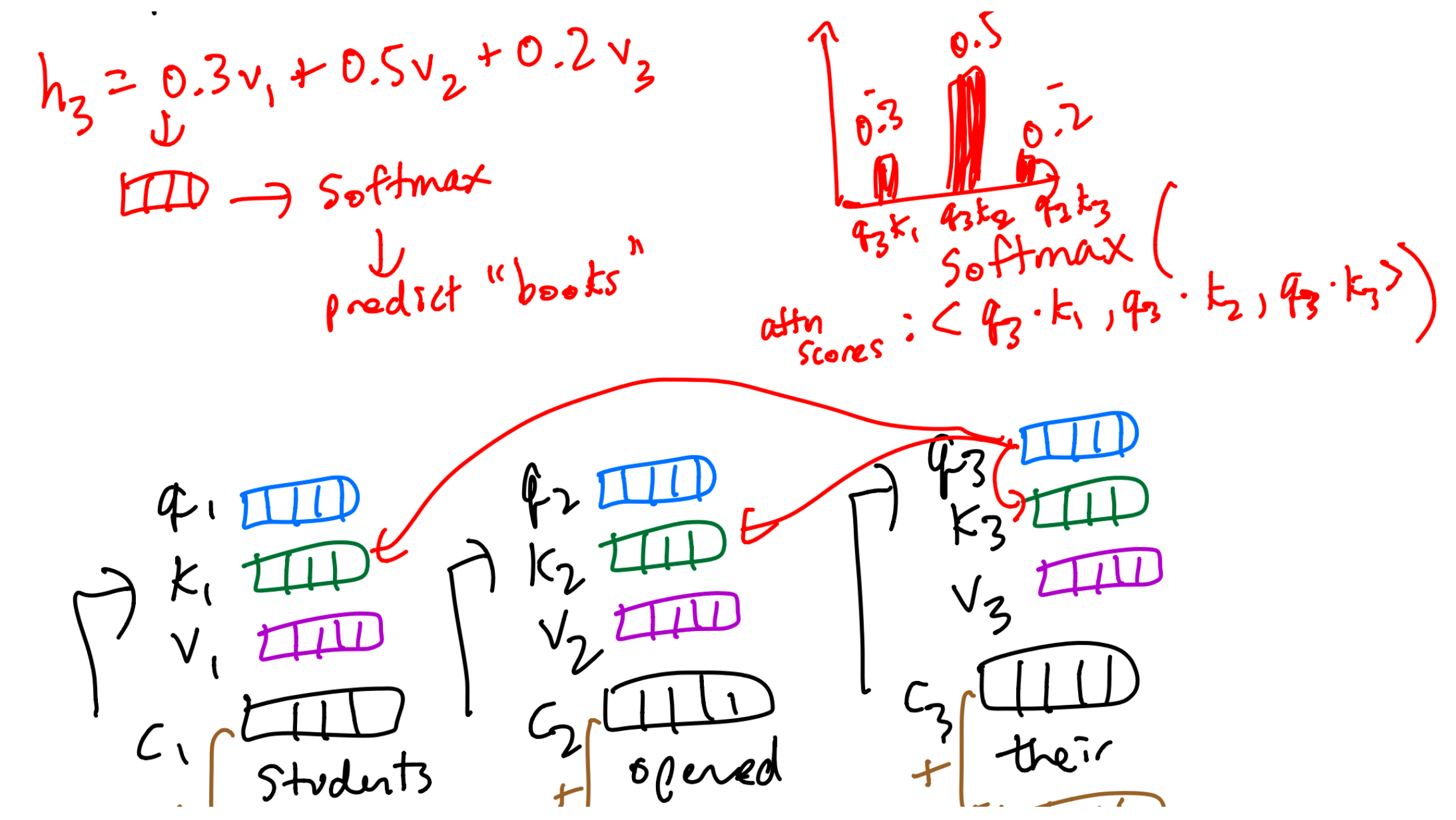
Neural Recommender + Matrix Factorization



$$P_S(p_{n+1}^i = p_{gt}^i | p_1^i \dots p_n^i) = \frac{\exp(\mathbf{h}_{i,n}^T \mathbf{p}_{gt}^i)}{\sum_j \exp(\mathbf{h}_{i,n}^T \mathbf{p}_j^i)}$$

Why do queries and keys need to be different?

1. Matrix factorization's need
2. The diagonal value cannot be low
 1. $q_1^T q_1 > q_1^T q_i$



Note: This is a cross-domain association (not common knowledge among NLP researchers)

Self-Attention Illustrative Example

Why do we need multiple heads?

How do they learn these?

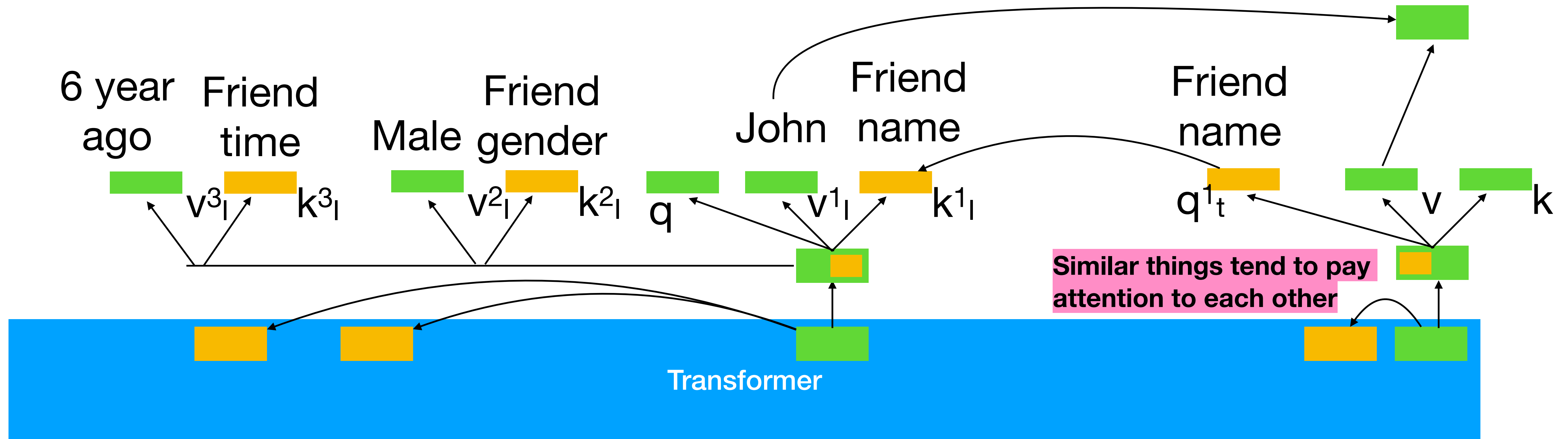
Why are keys and values different?

Retrieving only partial information

Gradient descent

Action != "property"

John



... Please call the friend of your main character John ...

... Mary's friend, _____

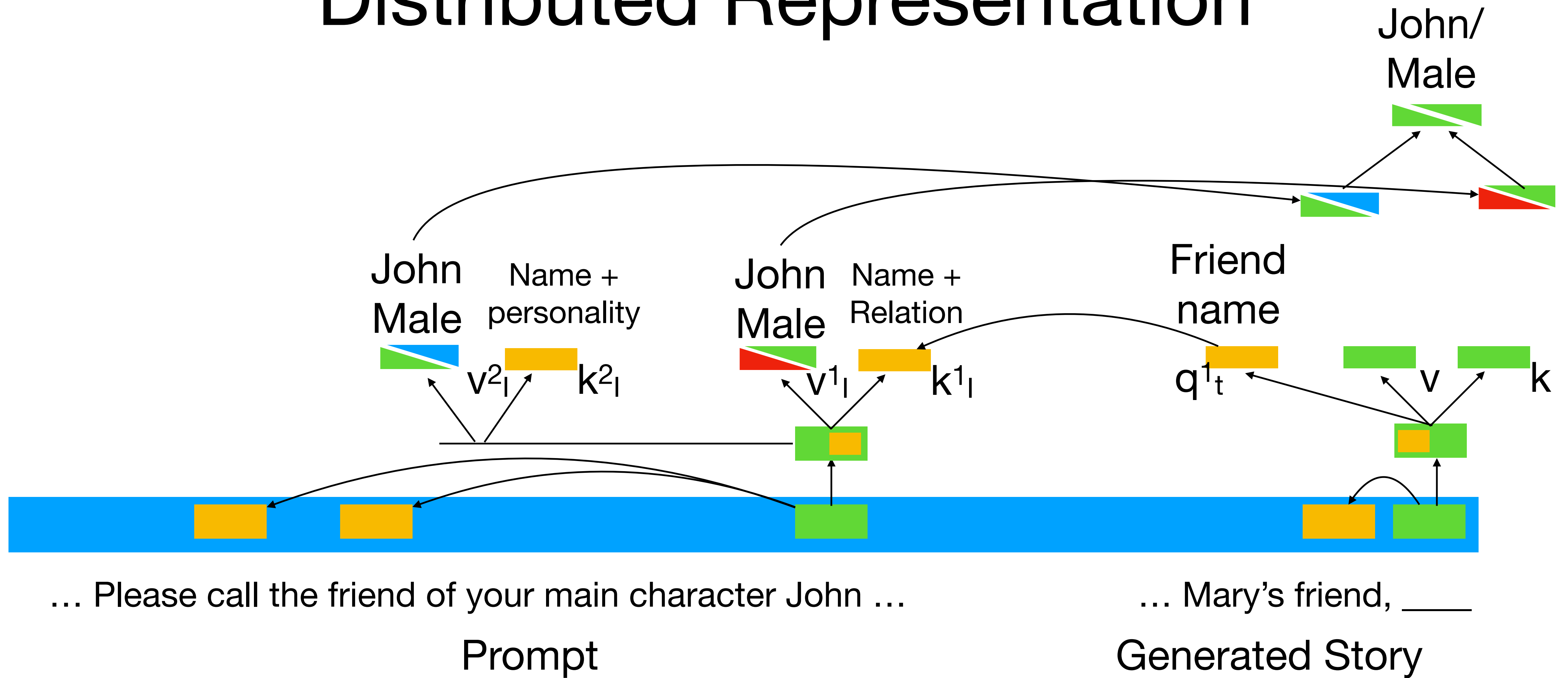
your main character met her friend 6 years ago

Prompt

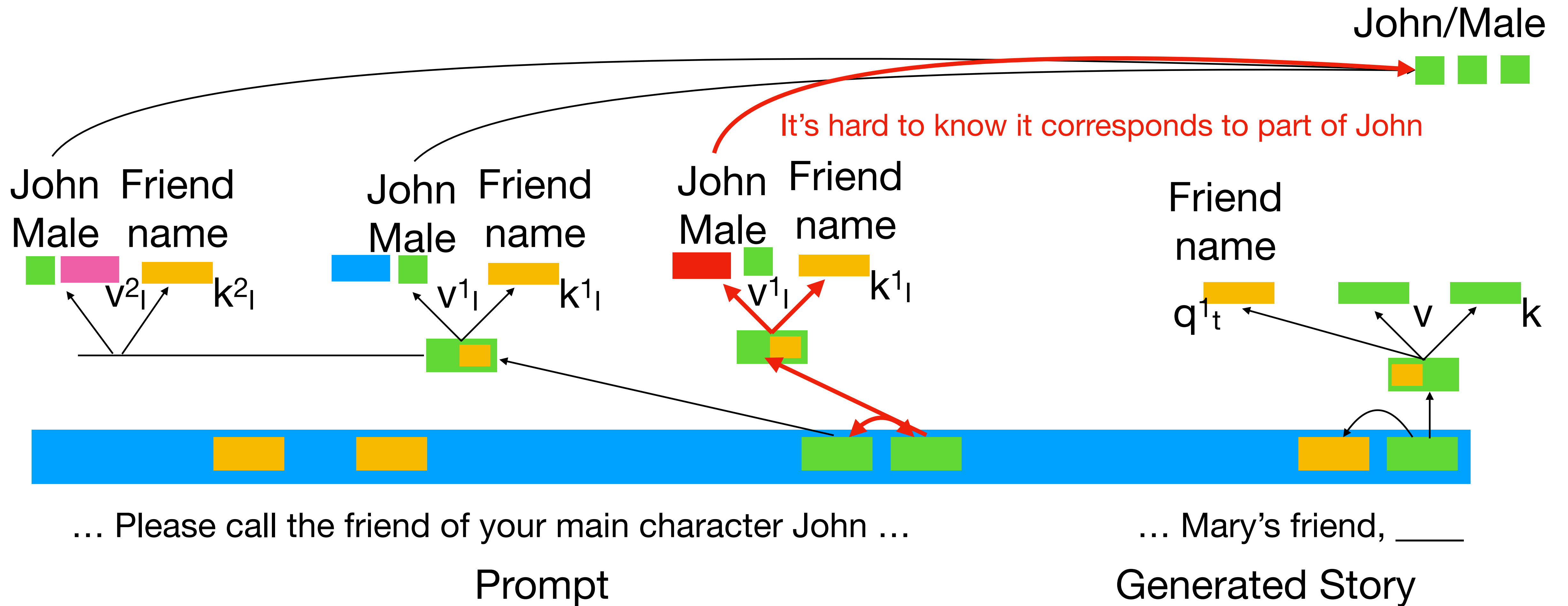
Generated Story

However, most heads are not very interpretable in practice.

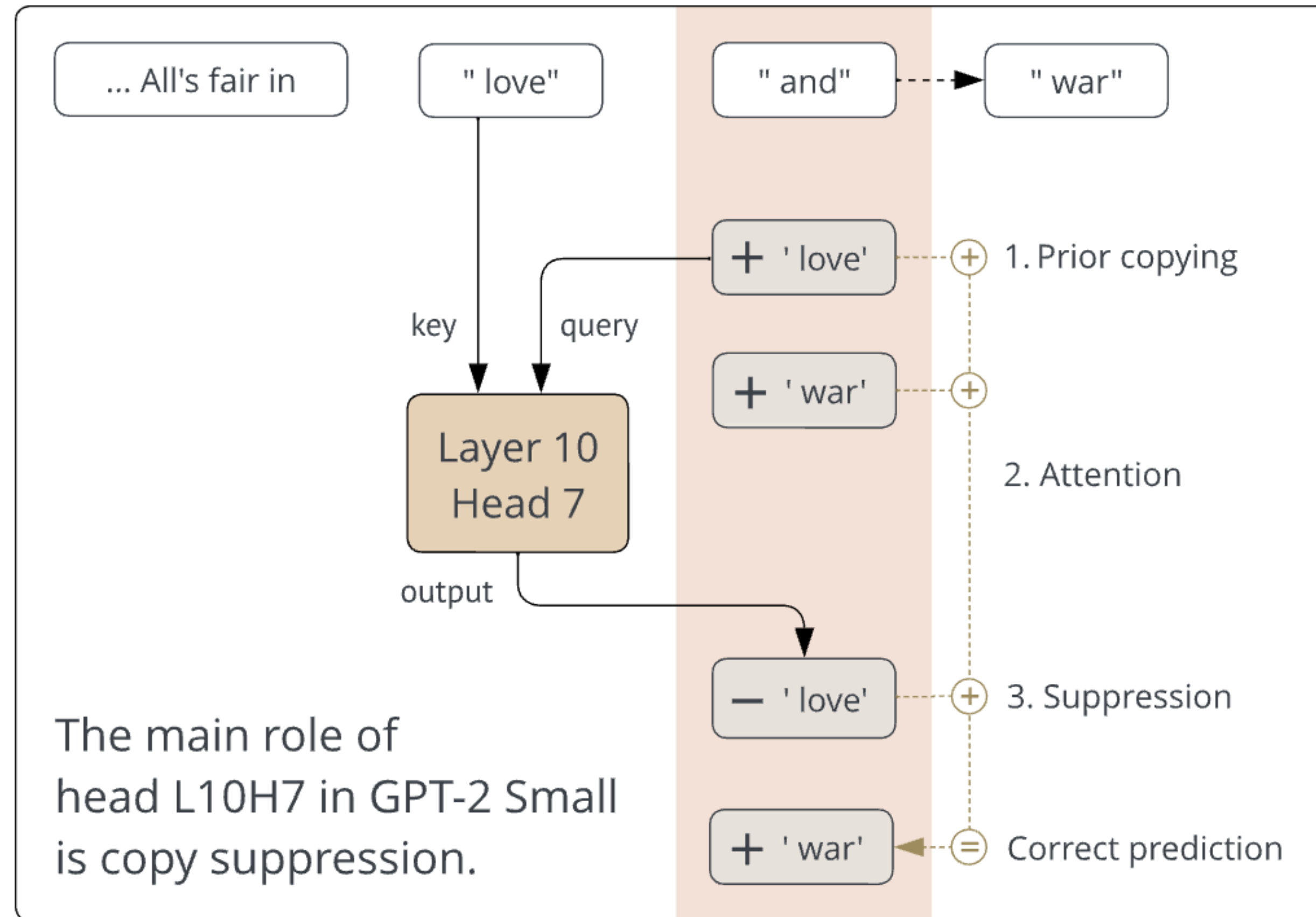
Distributed Representation



Distributed Representation



Interpretable Head Exists but Rare



COPY SUPPRESSION: COMPREHENSIVELY UNDERSTANDING AN ATTENTION HEAD
<https://arxiv.org/pdf/2310.04625>

A RNN Language Model

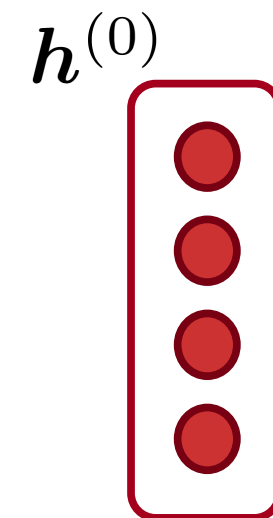
Testing time: Autoregressive LM

Inference limitation for self attention

hidden states

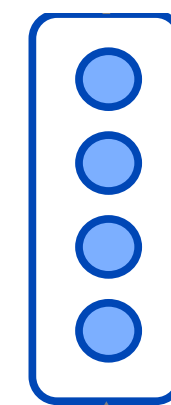
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!



word embeddings

c_1, c_2, c_3, c_4



the

c_1

A RNN Language Model

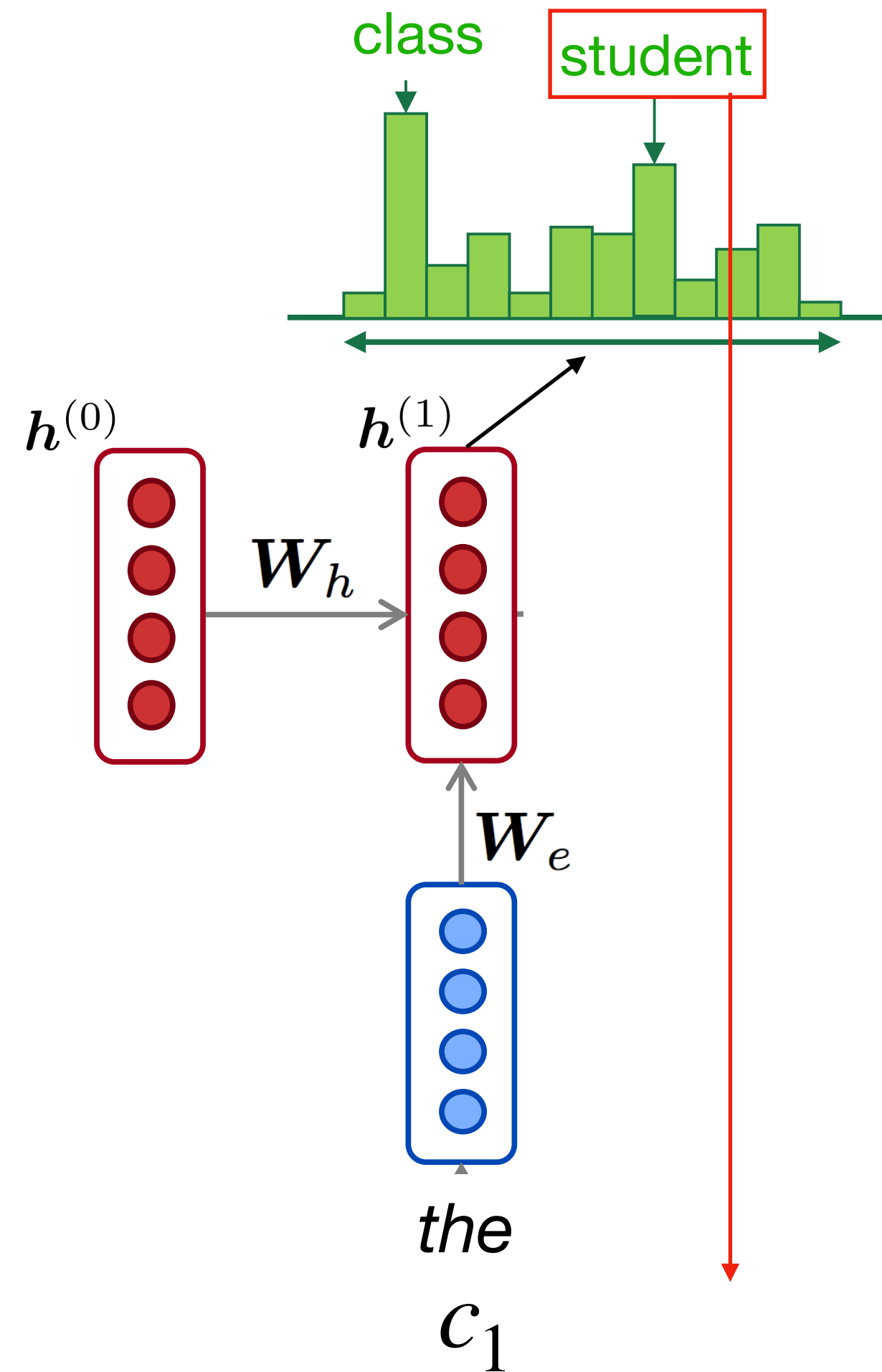
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!

word embeddings

c_1, c_2, c_3, c_4



A RNN Language Model

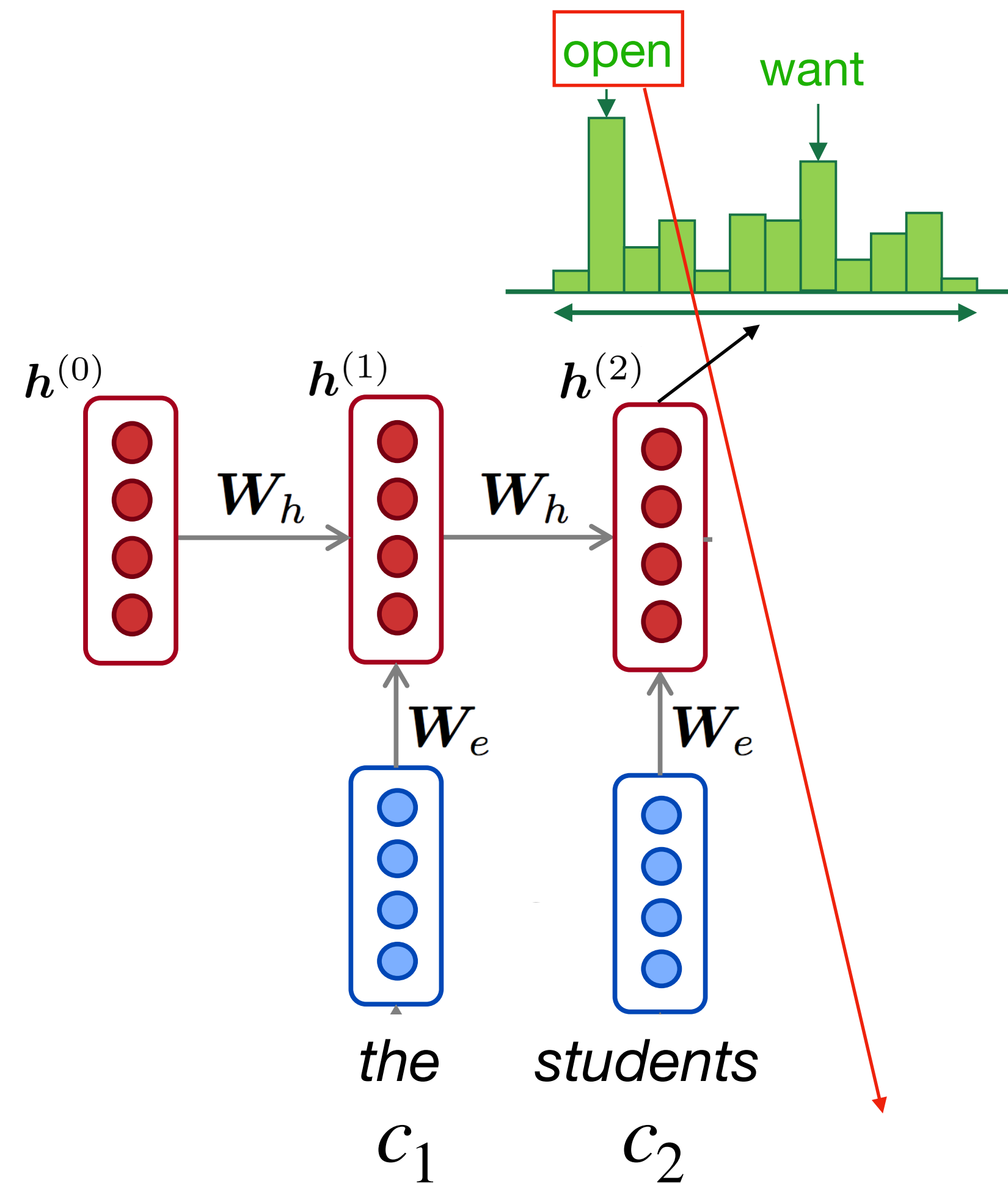
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!

word embeddings

c_1, c_2, c_3, c_4



A RNN Language Model

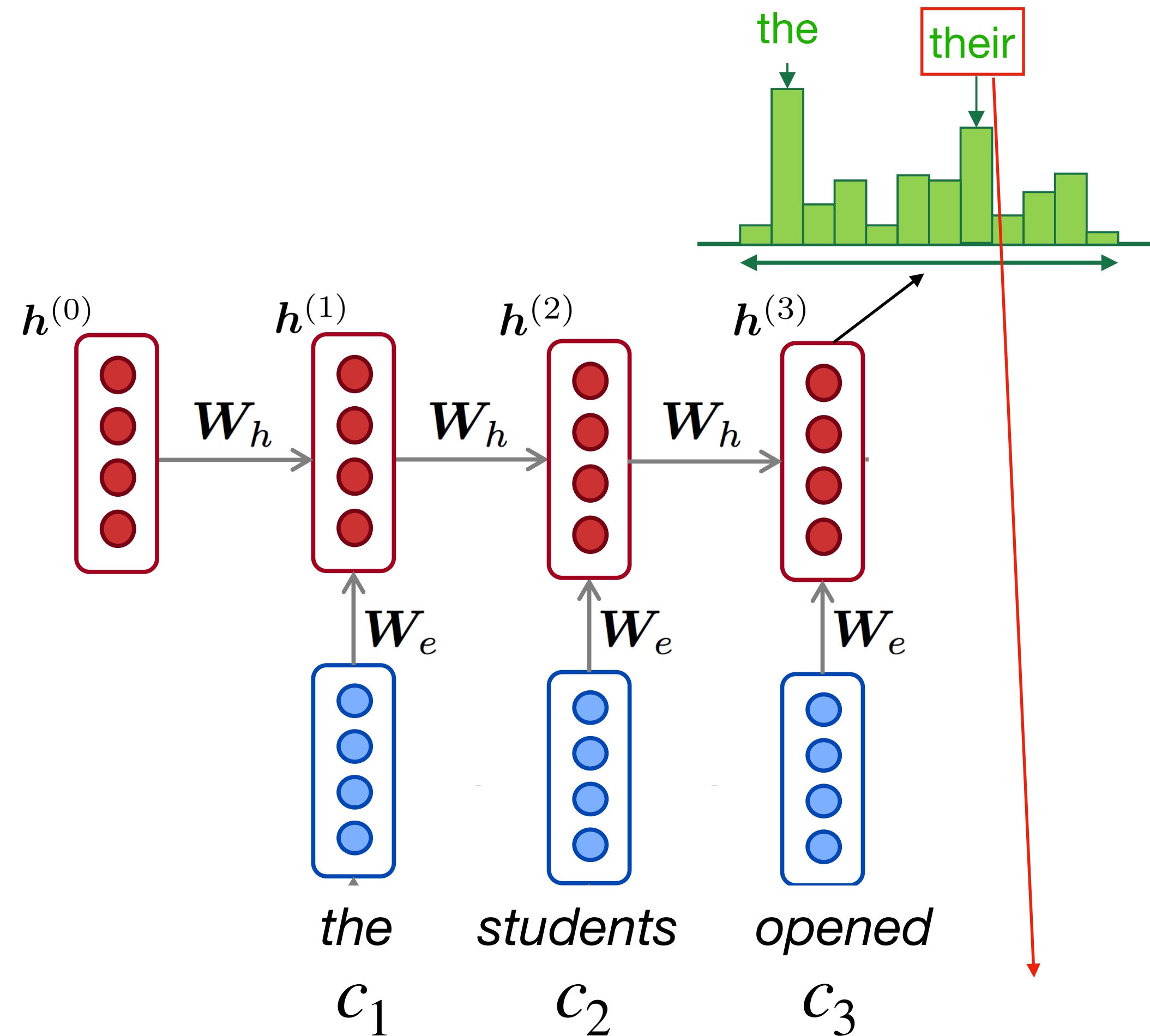
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!

word embeddings

c_1, c_2, c_3, c_4



A RNN Language Model

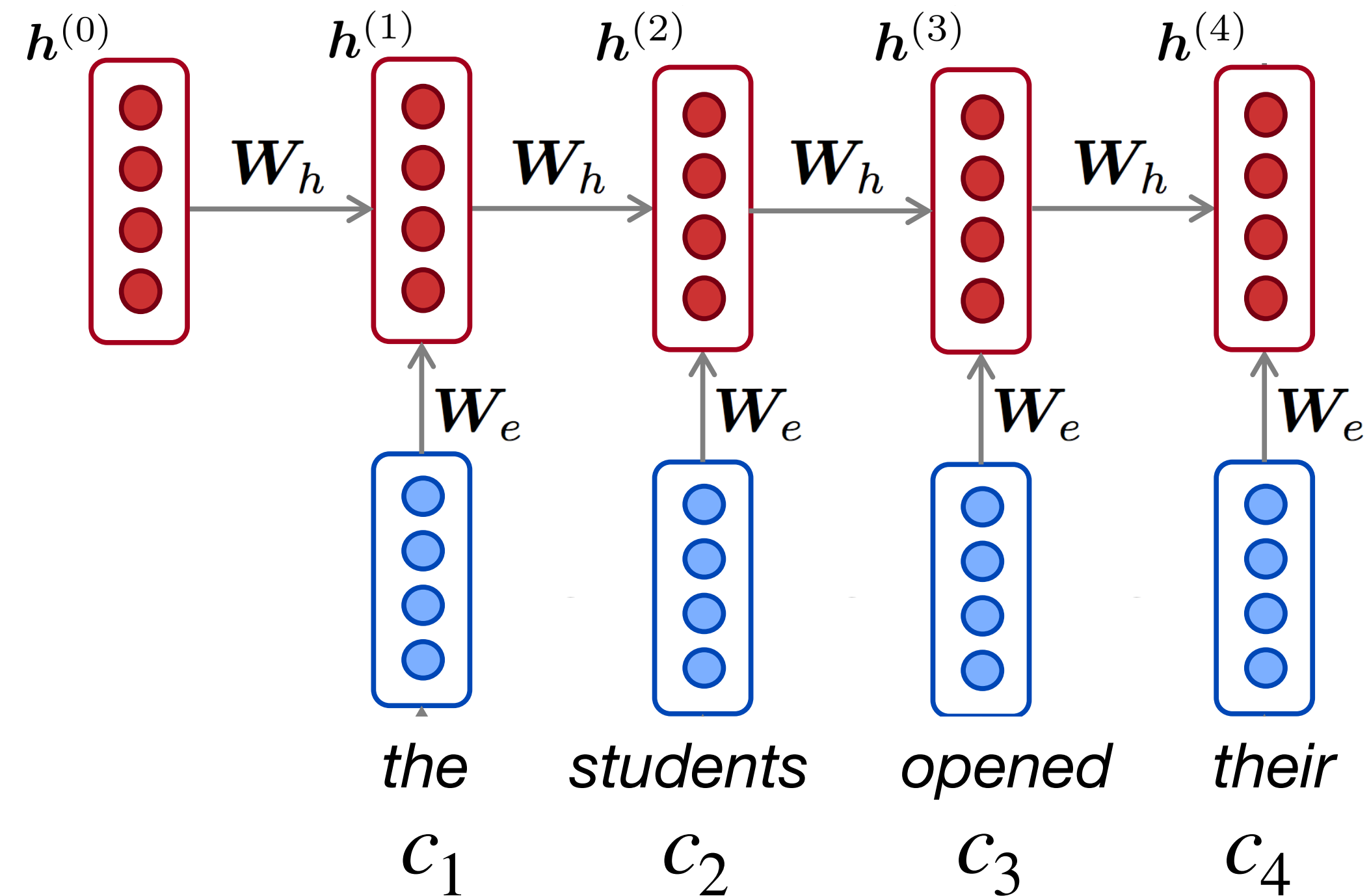
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!

word embeddings

c_1, c_2, c_3, c_4



A RNN Language Model

output distribution

$$\hat{y} = \text{softmax}(W_2 h^{(t)})$$

hidden states

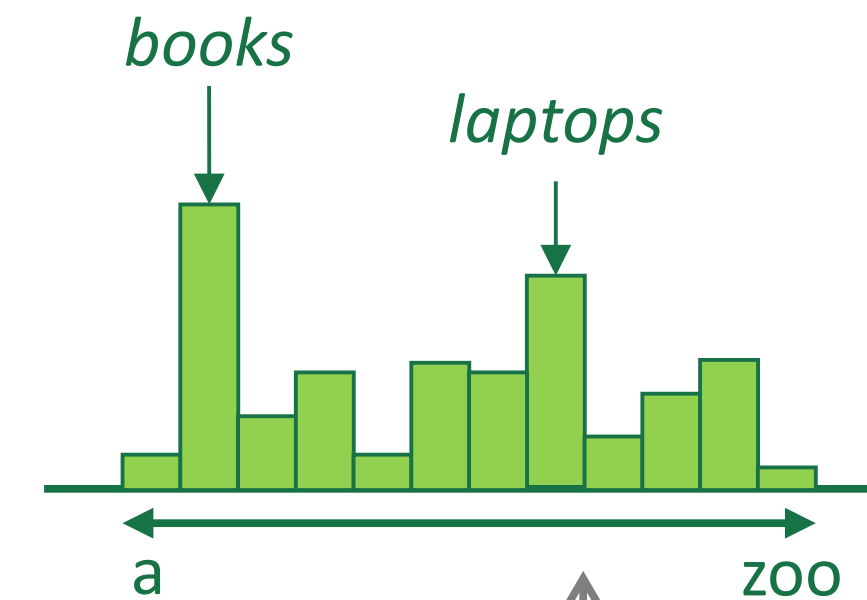
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$ is initial hidden state!

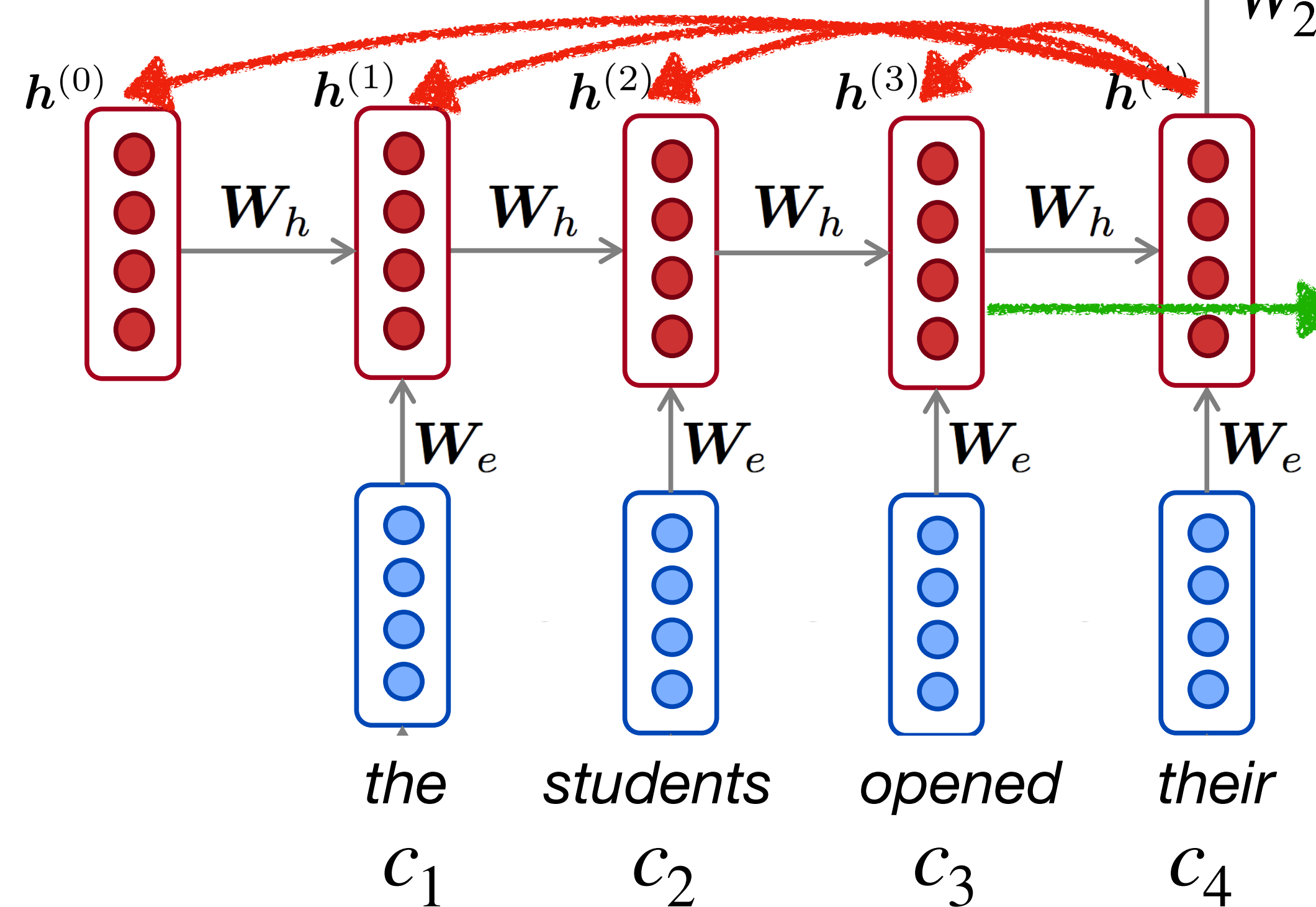
word embeddings

$$c_1, c_2, c_3, c_4$$

$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



Self-attention needs to store all KV cache



RNN only needs to store this for decoding

Multi-Head Latent Attention (MLA)

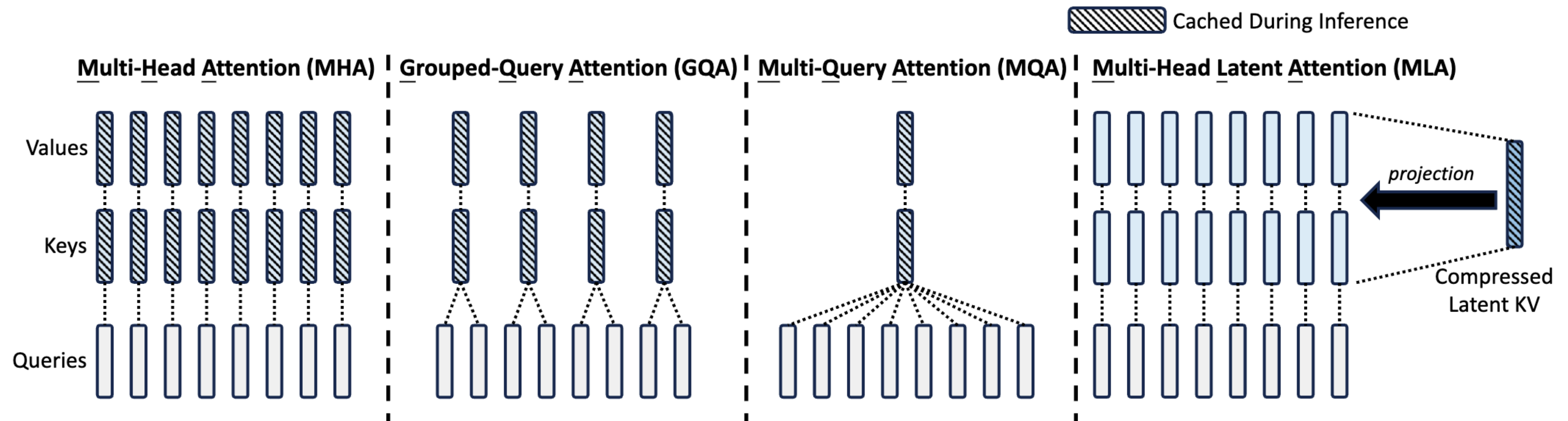
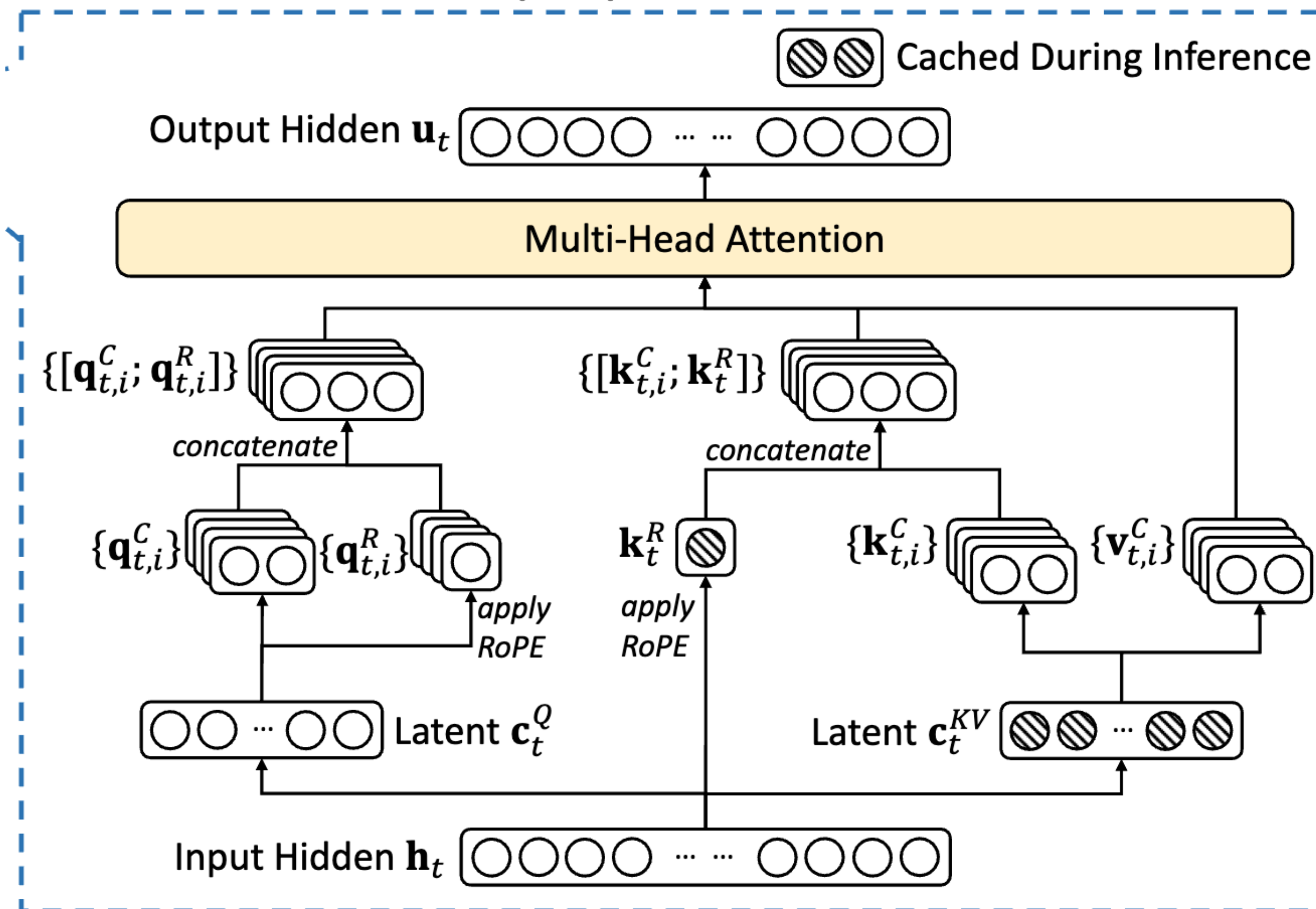


Figure 3 | Simplified illustration of Multi-Head Attention (MHA), Grouped-Query Attention (GQA), Multi-Query Attention (MQA), and Multi-head Latent Attention (MLA). Through jointly compressing the keys and values into a latent vector, MLA significantly reduces the KV cache during inference.

Deepseek V2 (<https://arxiv.org/pdf/2405.04434>)

Multi-Head Latent Attention (MLA)

Multi-Head Latent Attention (MLA)



Why not store h?

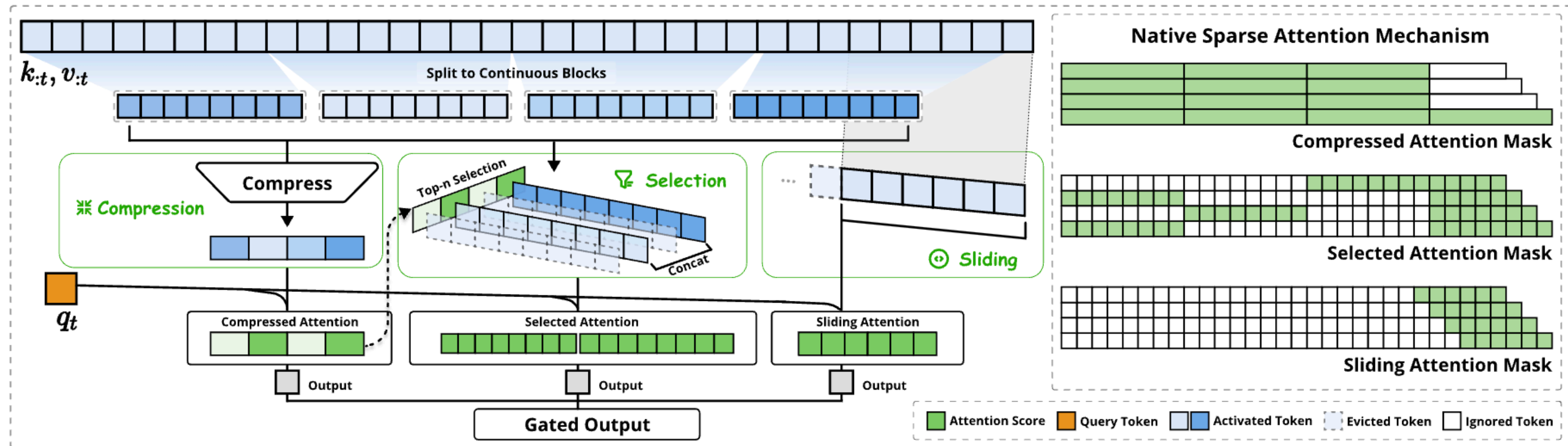
$$q^T k = (W_q h)^T (W_k h) = h^T W_q^T W_k h$$

Attention Mechanism	KV Cache per Token (# Element)	Capability
Multi-Head Attention (MHA)	$2n_h d_h l$	Strong
Grouped-Query Attention (GQA)	$2n_g d_h l$	Moderate
Multi-Query Attention (MQA)	$2d_h l$	Weak
MLA (Ours)	$(d_c + d_h^R) l \approx \frac{9}{2} d_h l$	Stronger

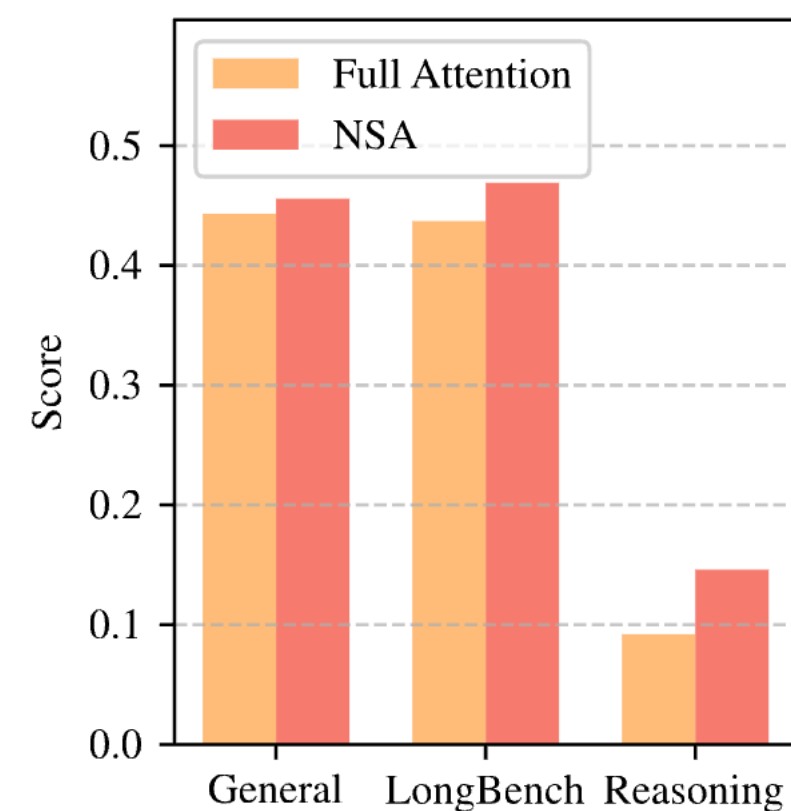
Table 1 | Comparison of the KV cache per token among different attention mechanisms. n_h denotes the number of attention heads, d_h denotes the dimension per attention head, l denotes the number of layers, n_g denotes the number of groups in GQA, and d_c and d_h^R denote the KV compression dimension and the per-head dimension of the decoupled queries and key in MLA, respectively. The amount of KV cache is measured by the number of elements, regardless of the storage precision. For DeepSeek-V2, d_c is set to $4d_h$ and d_h^R is set to $\frac{d_h}{2}$. So, its KV cache is equal to GQA with only 2.25 groups, but its performance is stronger than MHA.

<https://github.com/Zefan-Cai/Awesome-LLM-KV-Cache>

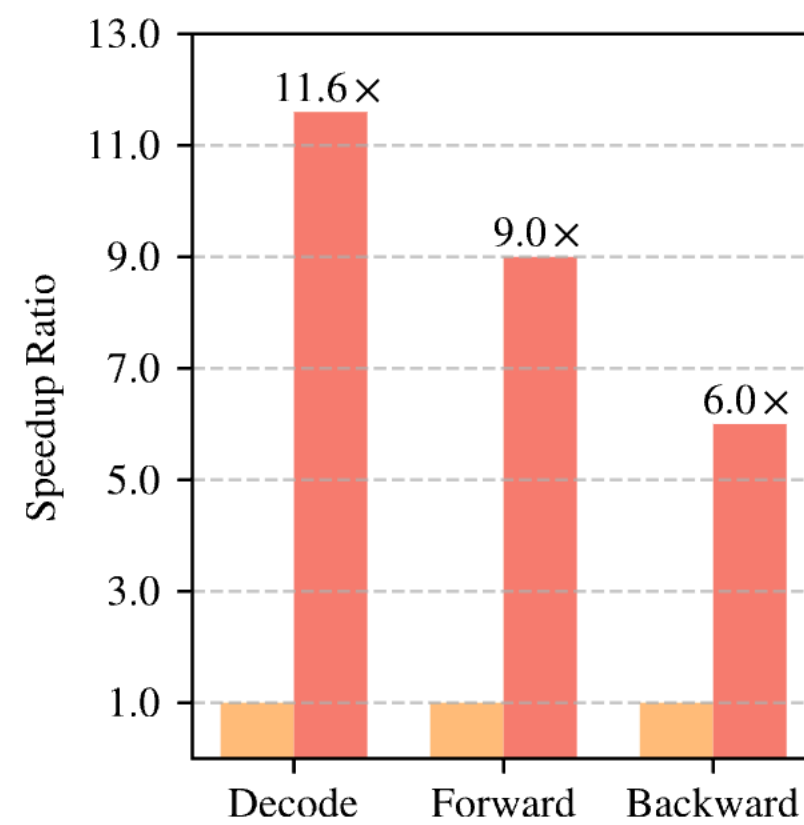
Hierarchical Attention



Performance on Benchmarks



Speed on Stages



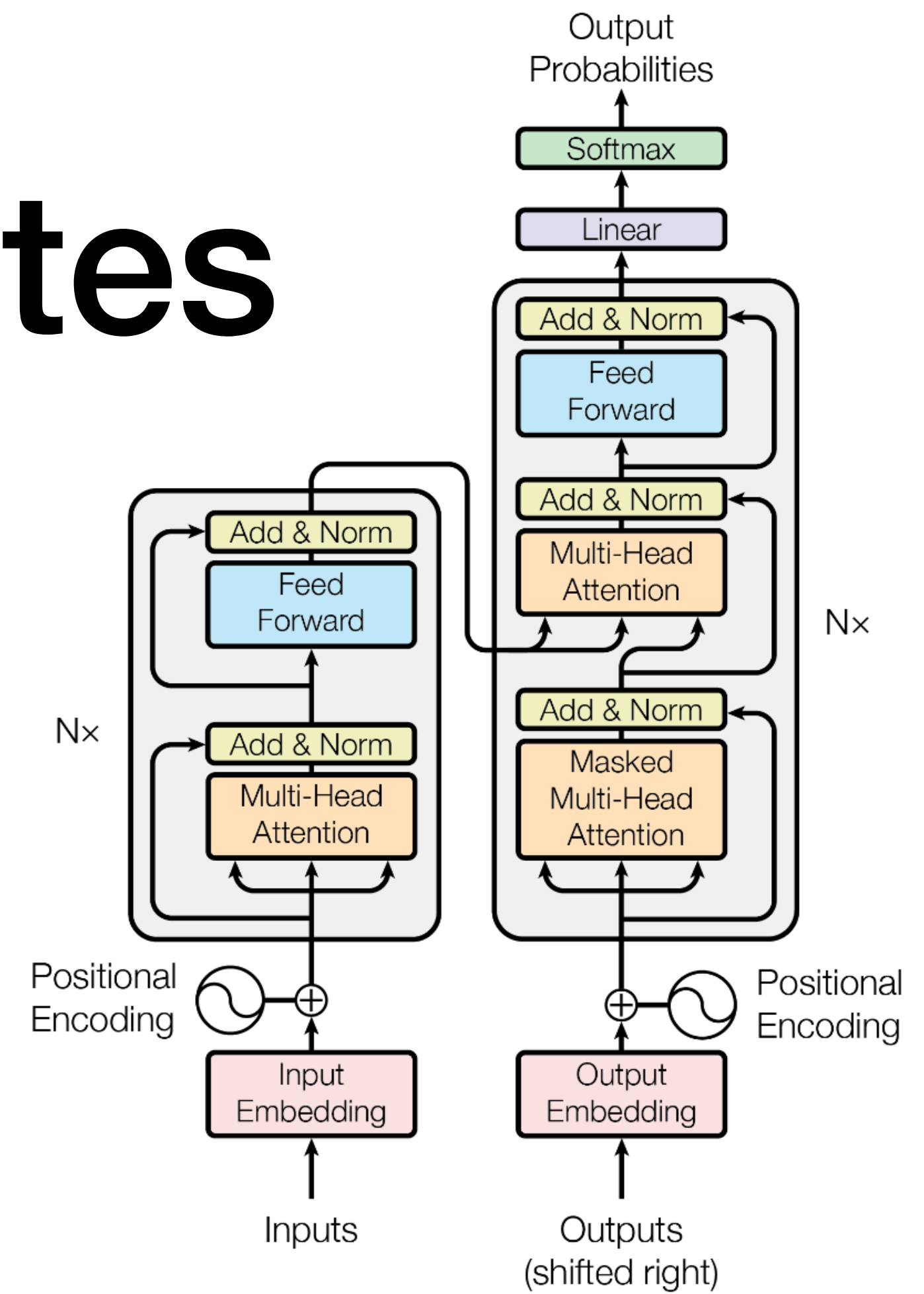
Similar to RAG

Native Sparse Attention: Hardware-Aligned and Natively Trainable Sparse Attention (<https://arxiv.org/pdf/2502.11089>)

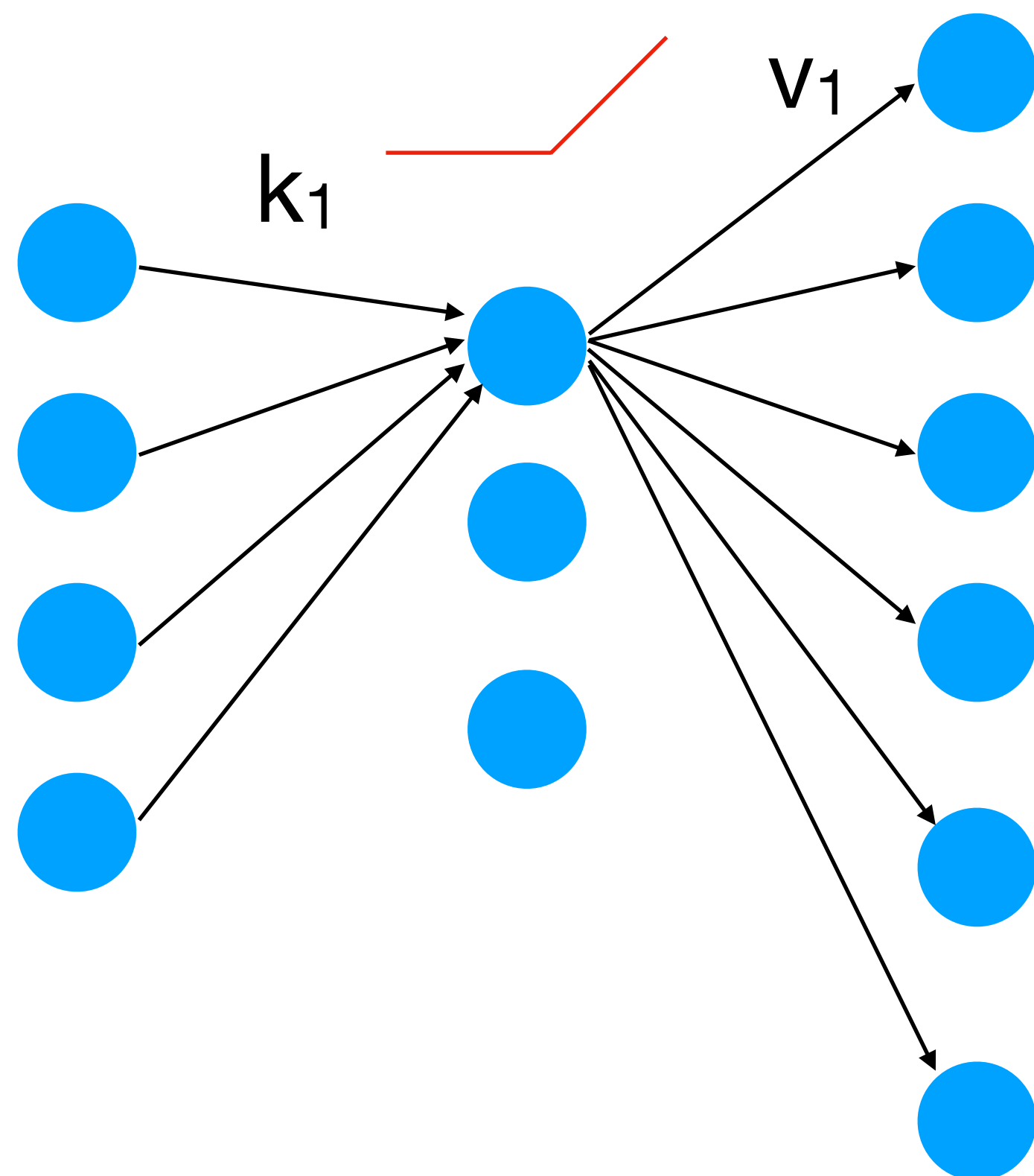
(From deepseek AI)

Last Year Notes

$$\text{Layernorm}(x) = \gamma \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$



MLP / Feed forward NN



Transformer Feed-Forward Layers Are Key-Value Memories (<https://arxiv.org/pdf/2012.14913>)

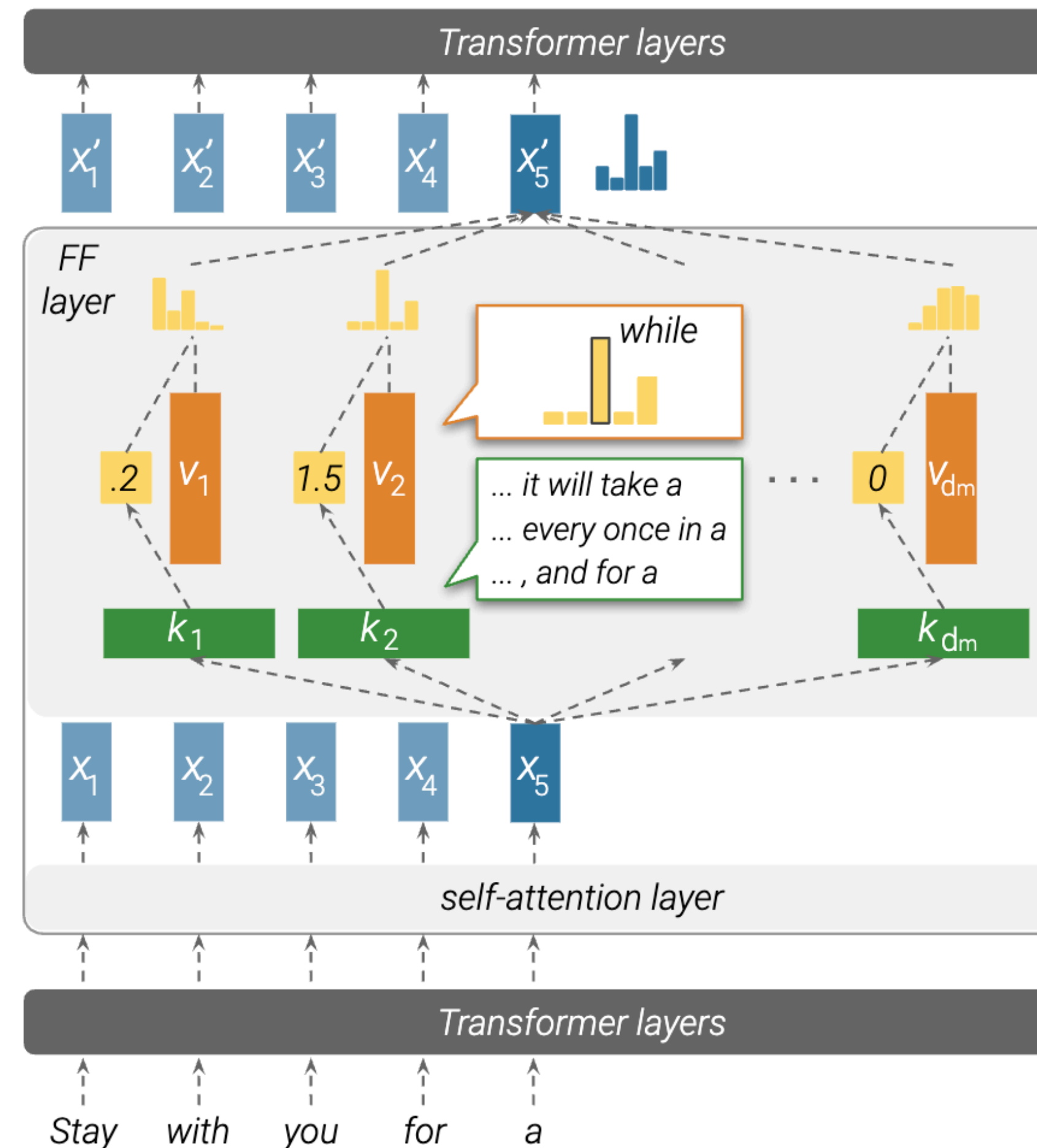


Figure 1: An illustration of how a feed-forward layer emulates a key-value memory. Input vectors (here, x_5) are multiplied by *keys* to produce *memory coefficients* (e.g., the memory coefficient for v_1 is 0.2), which then weigh distributions over the output vocabulary, stored in the *values*. The feed-forward layer's output is thus the weighted sum of its values.

Non-negativity \rightarrow Meaningful Dimensions



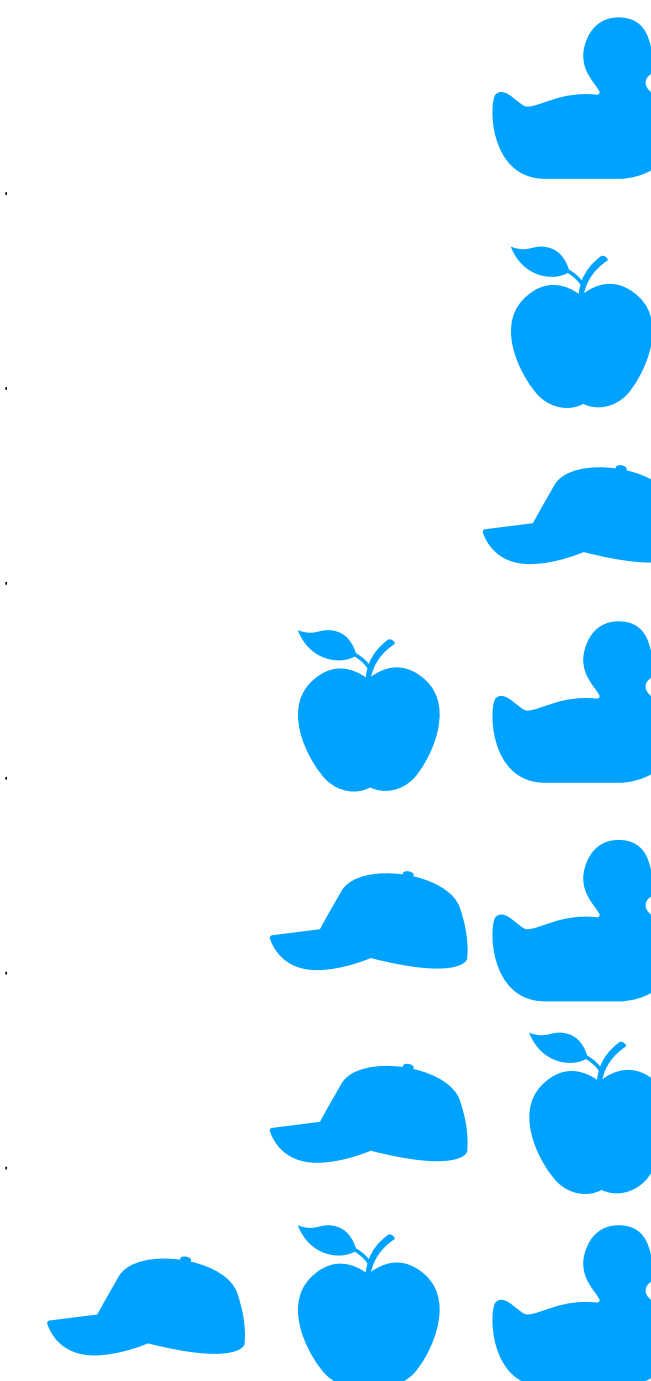
Hidden state

0	0	1
0	1	0
1	0	0
0	1	1
1	0	1
1	1	0
1	1	1

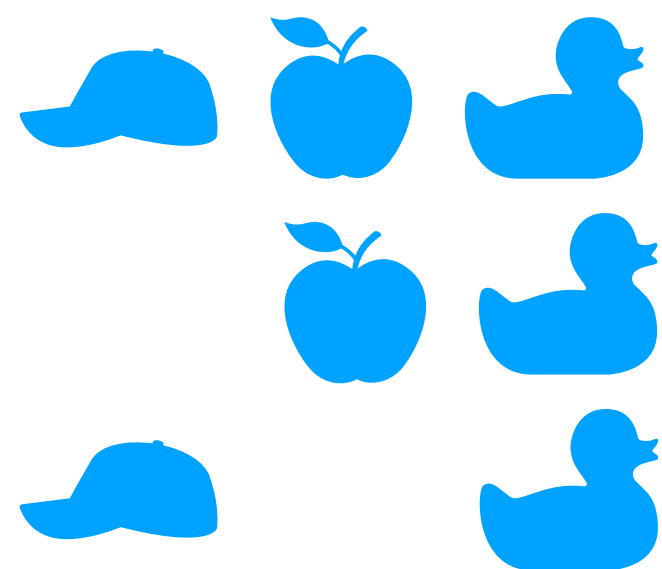
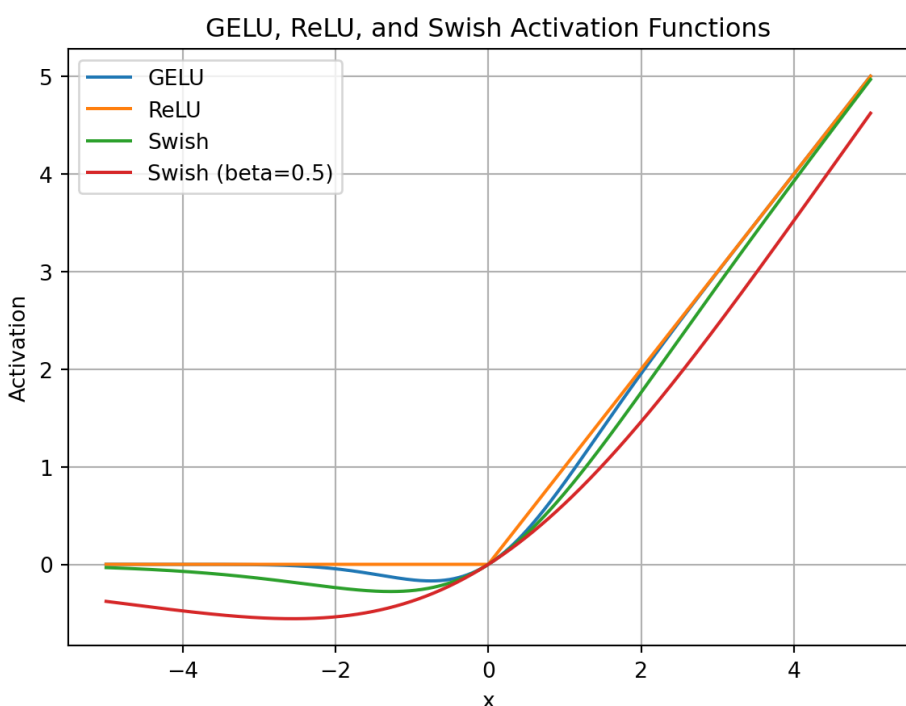
v

	1	-1	0	0	0	0
	0	0	2	-2	0	0
	0	0	0	0	3	-3

0	0	0	0	3	-3
0	0	2	-2	0	0
1	-1	0	0	0	0
0	0	2	-2	3	-3
1	-1	0	0	3	-3
1	-1	2	-2	0	0
1	-1	2	-2	3	-3



Distributed Representation



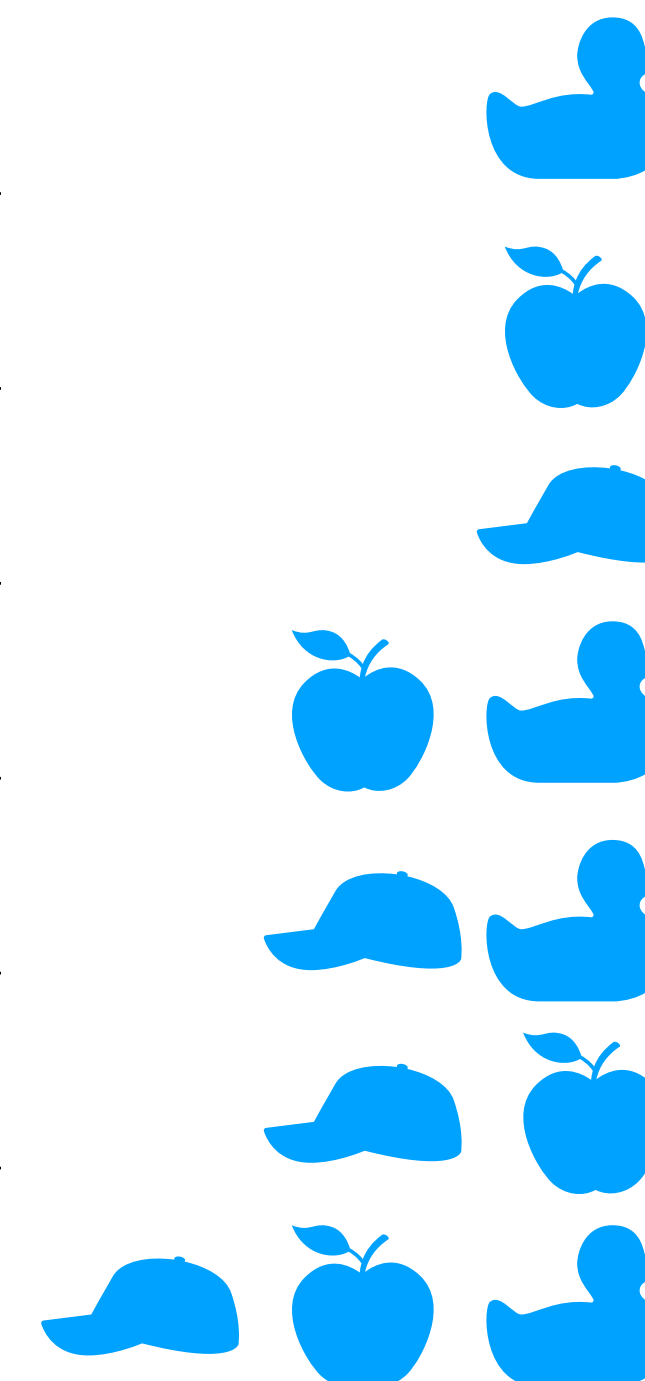
1	-1	-2	-2	3	-3
0	0	2	-2	3	-3
1	-1	0	0	3	-3

A good example showing that the mechanism interpretability highly depends on the small architecture change

Linear probes extract attributes from hidden states

-1	1	1
1	0	-1
1	-1	0
0	1	0
0	0	1
2	-1	-1
1	0	0

0	0	0	0	3	-3
0	0	2	-2	0	0
1	-1	0	0	0	0
0	0	2	-2	3	-3
1	-1	0	0	3	-3
1	-1	2	-2	0	0
1	-1	2	-2	3	-3

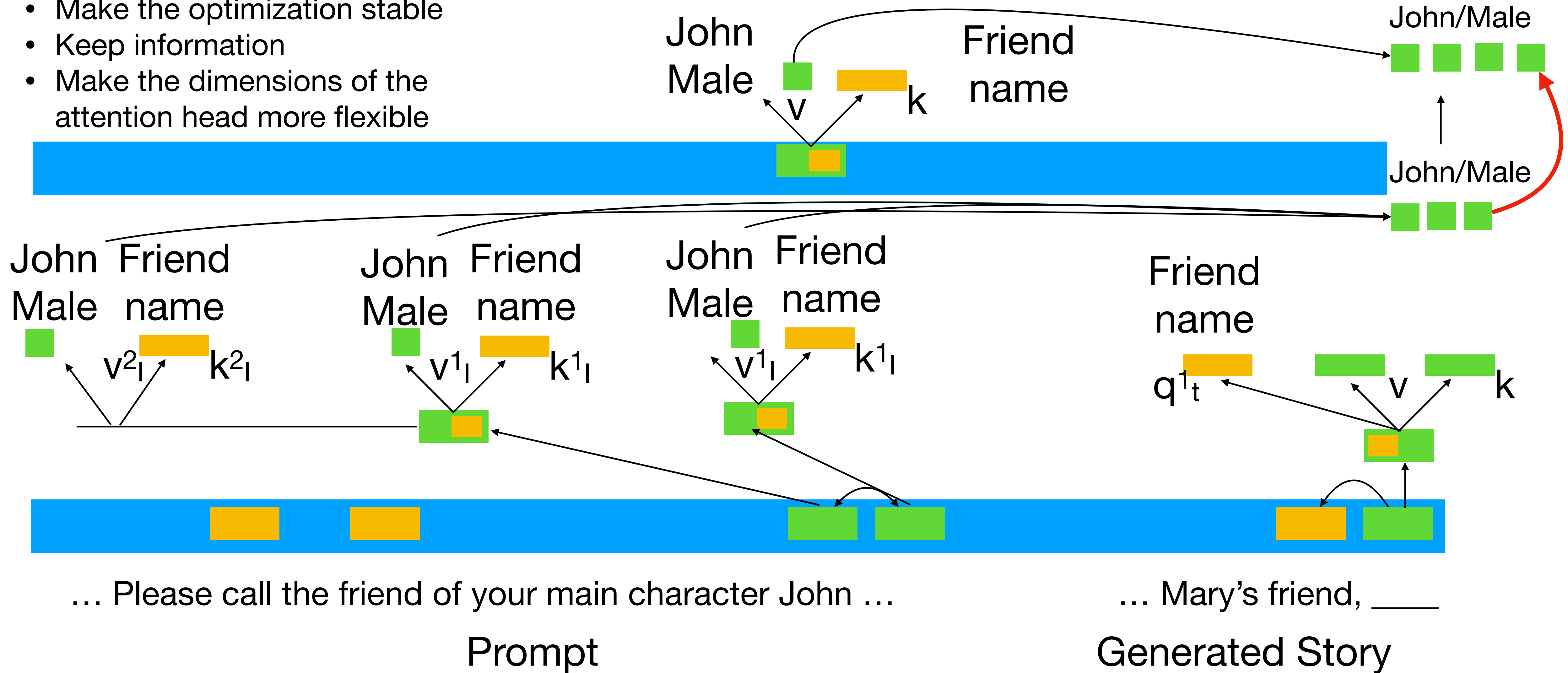


Transformer LM

Haw-Shiuan Chang

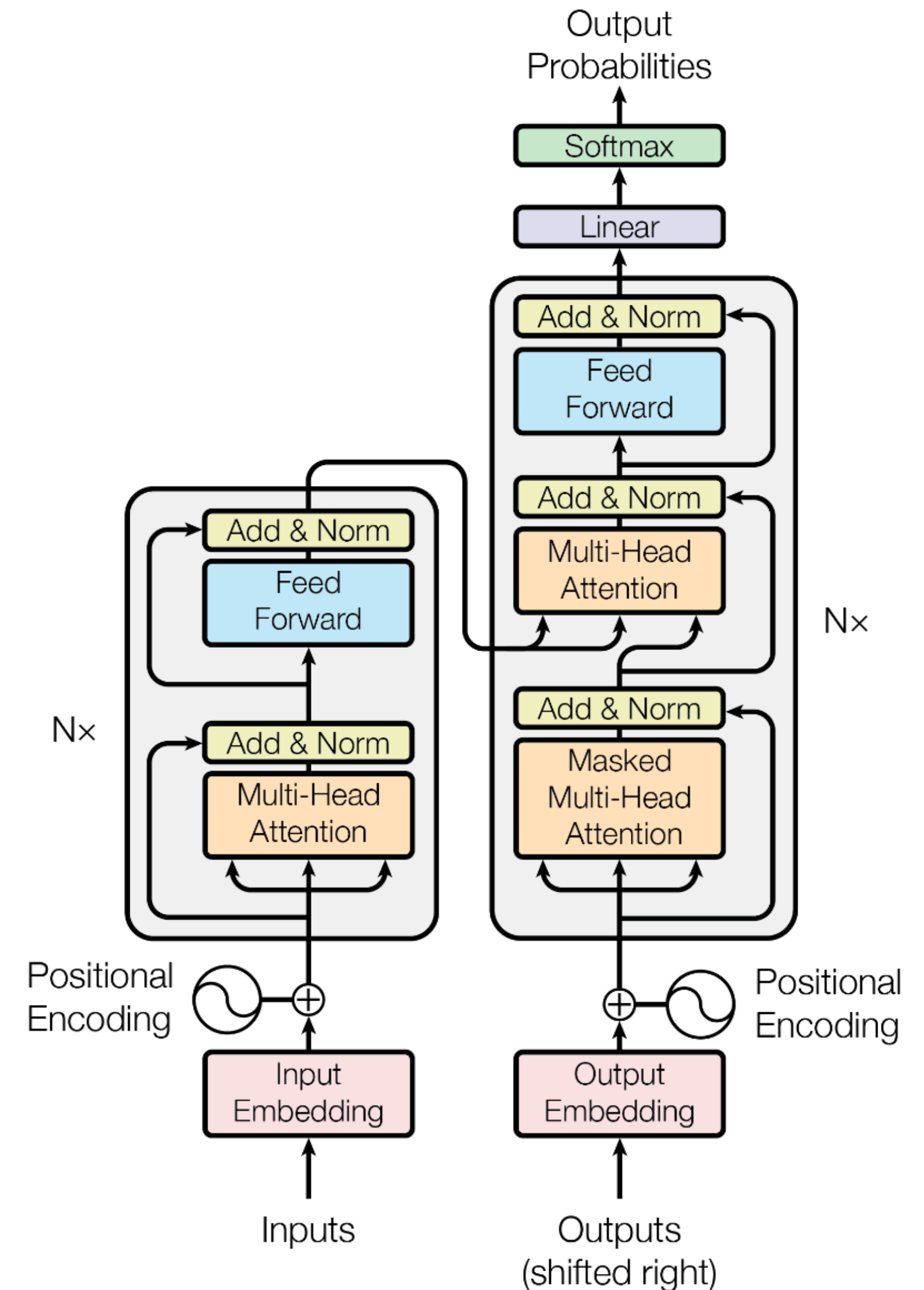
Residual Link

- Make the optimization stable
- Keep information
- Make the dimensions of the attention head more flexible

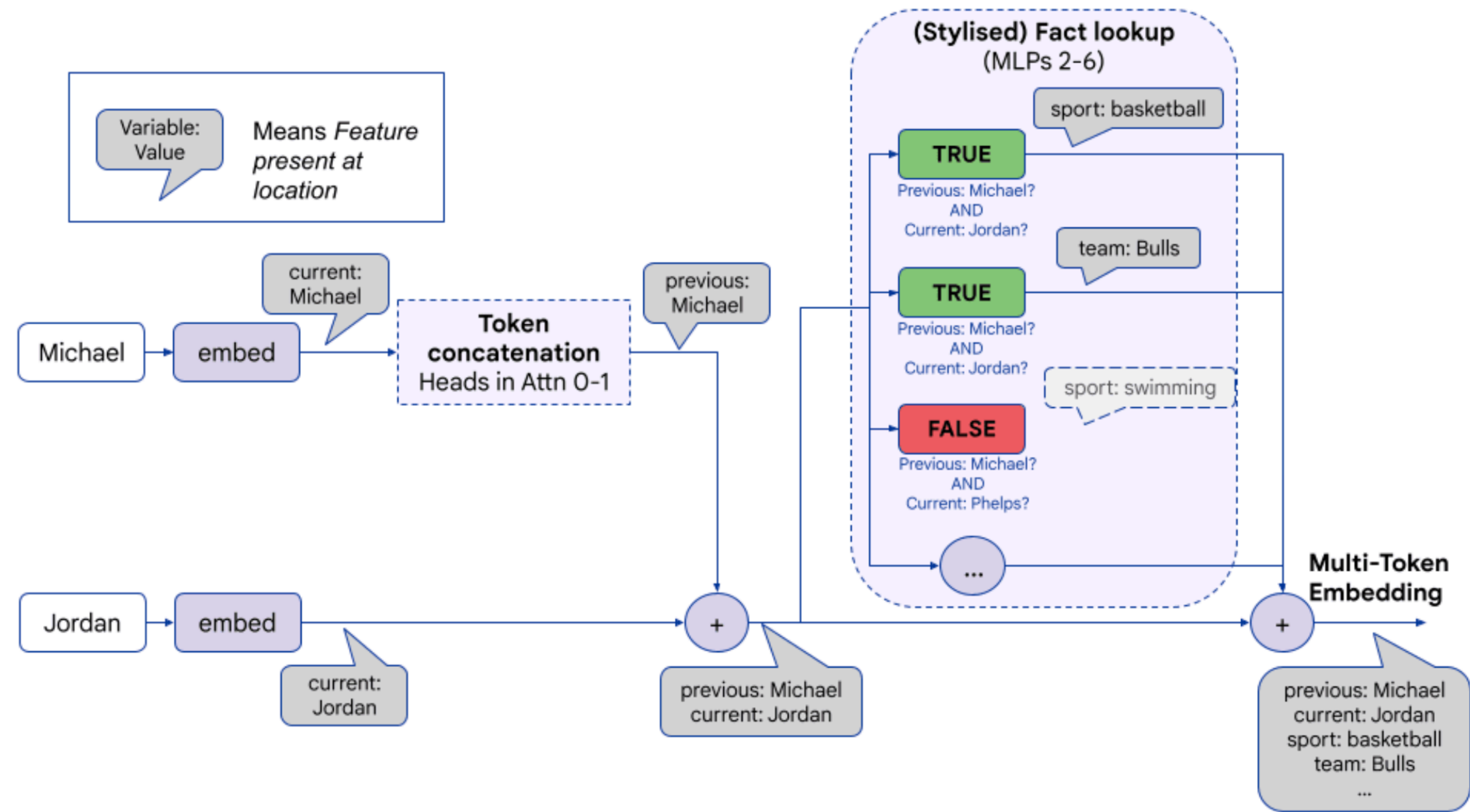
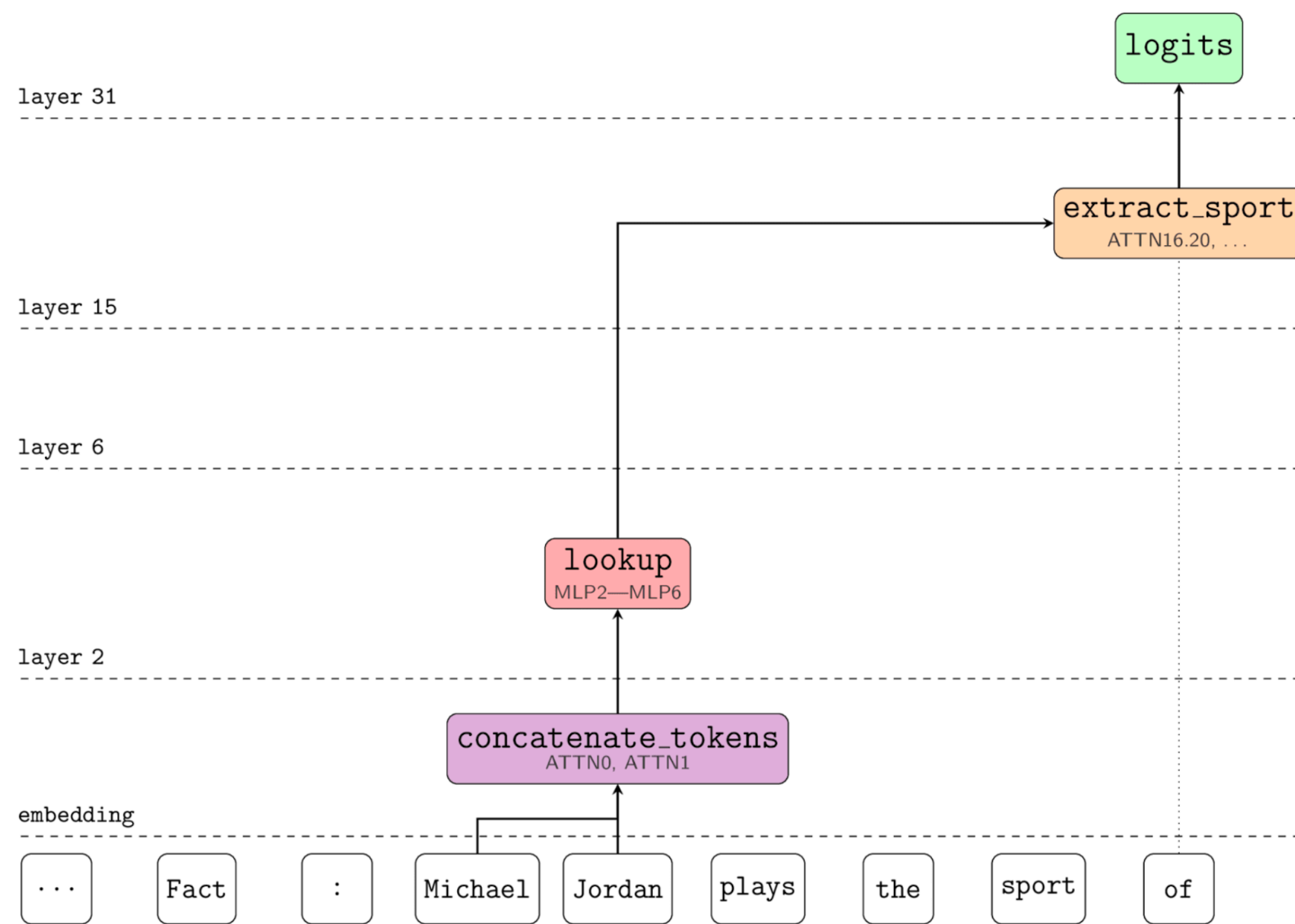


Transformer Architecture

- Self-attention needs $\sim 1/3$ parameters (GPT-3)
 - Most context processing happens here
 - You can replace this with RNN
- MLP needs $\sim 2/3$ parameters
 - Most memorization happens here
- Optionally sharing word embeddings
- The architectures are designed for GPU



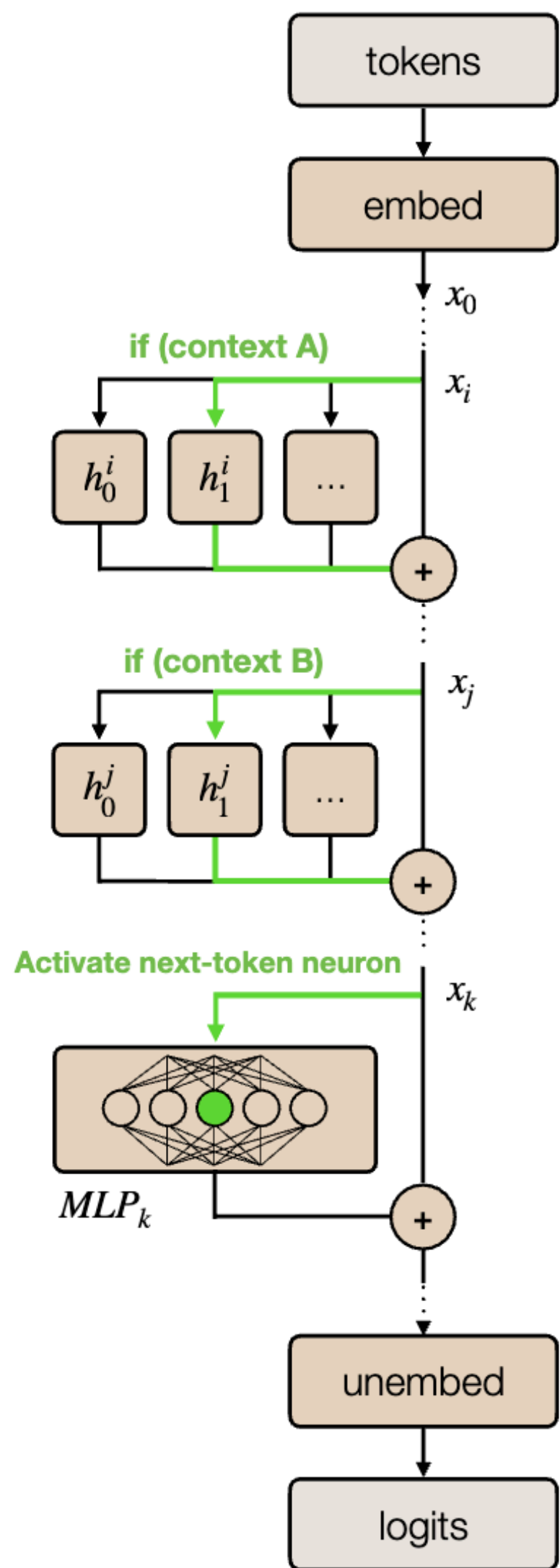
Example for Attention & MLP



Top deep-learning scientists such as Ilya could probably see these after reading the Transformer paper

<https://www.lesswrong.com/posts/iGuwZTHWb6DFY3sKB/fact-finding-attempting-to-reverse-engineer-factual-recall>

Another Perspective

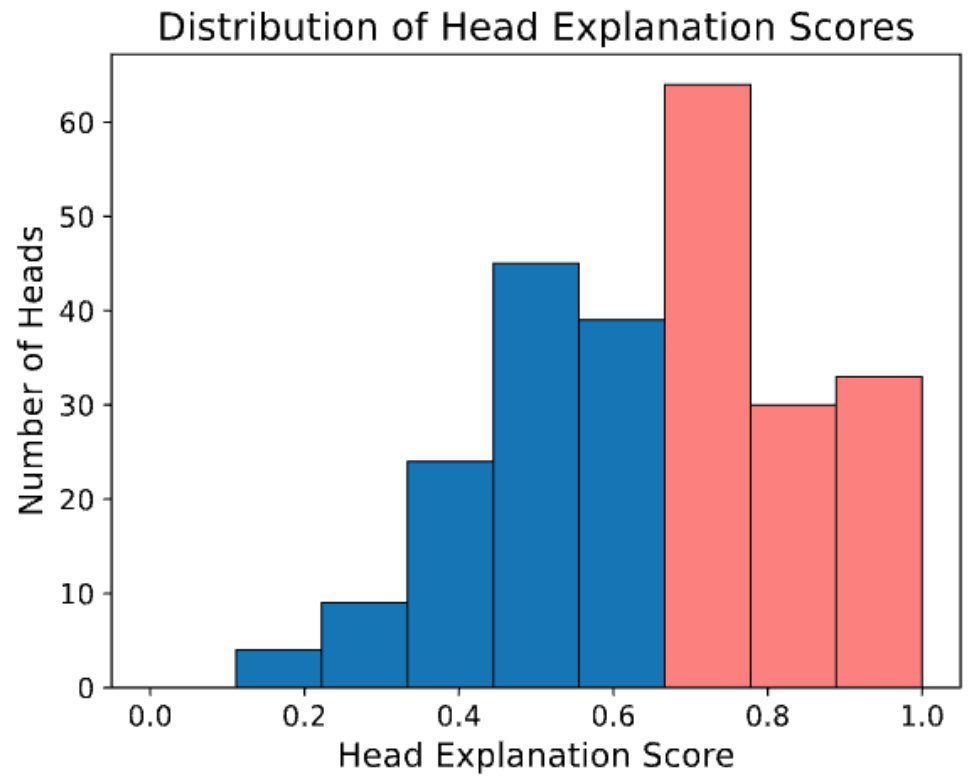
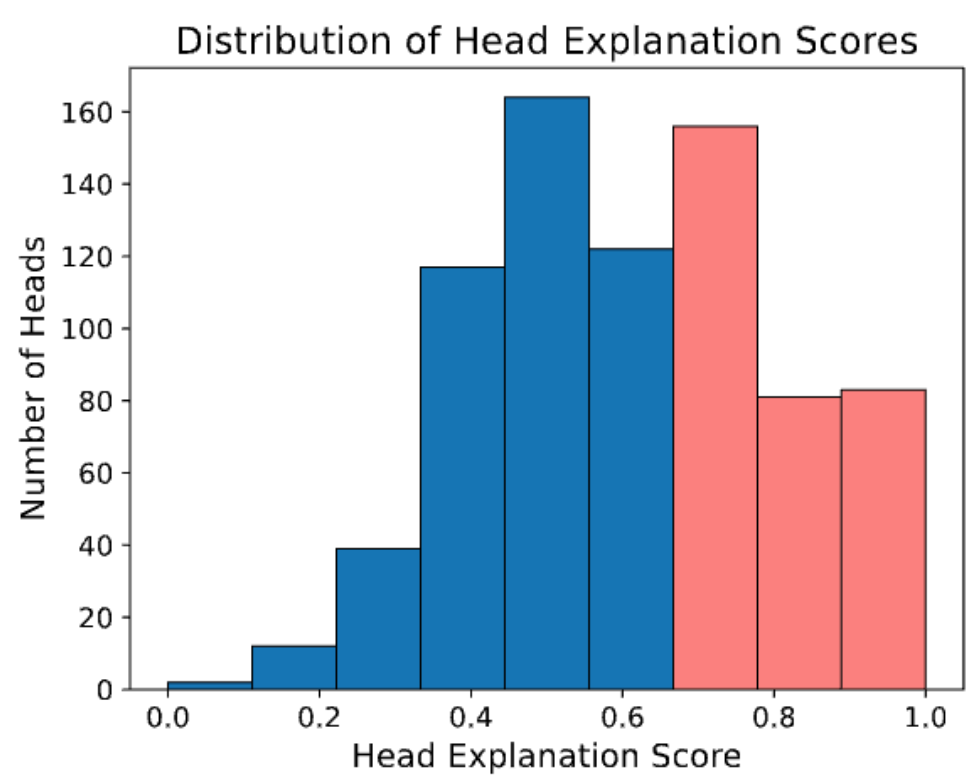
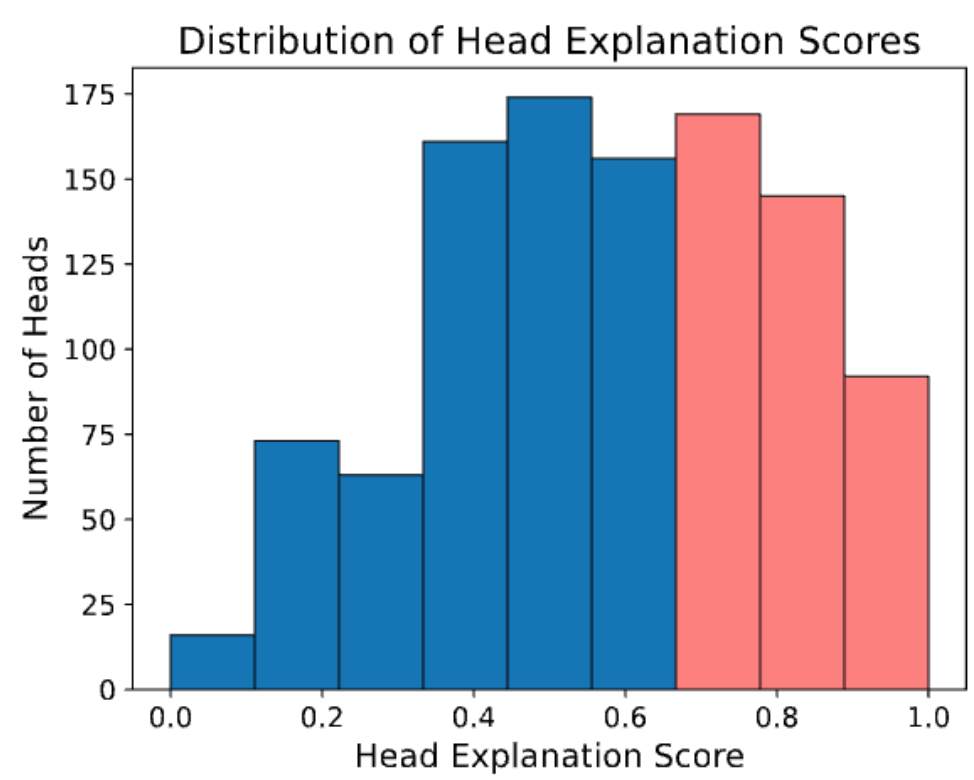
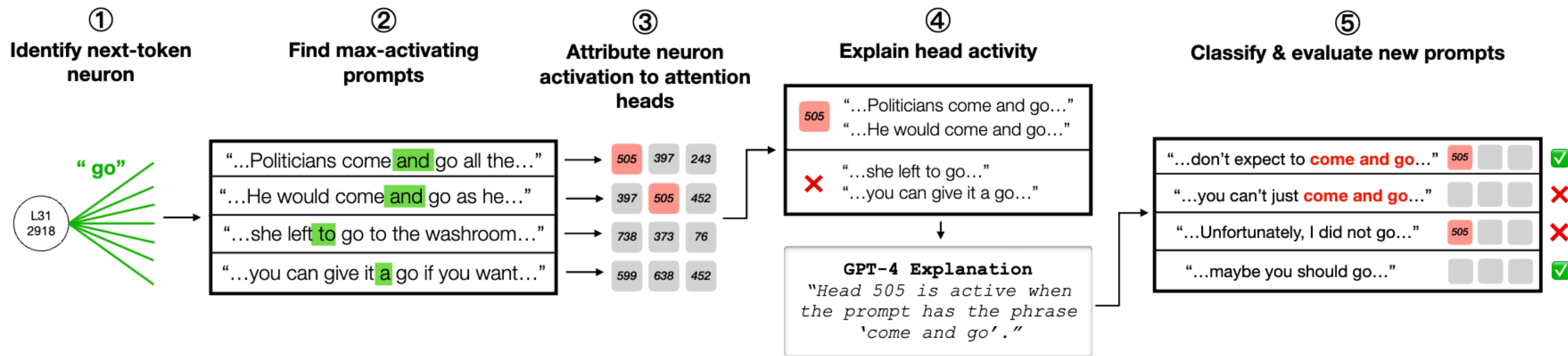


A word can be used differently in multiple contexts.

Different attention heads can pick up on different contexts, and have an output that activates a MLP neuron in later layers.

When the neuron activates, it outputs the embedding of that token into the residual stream, increasing the probability of that token.

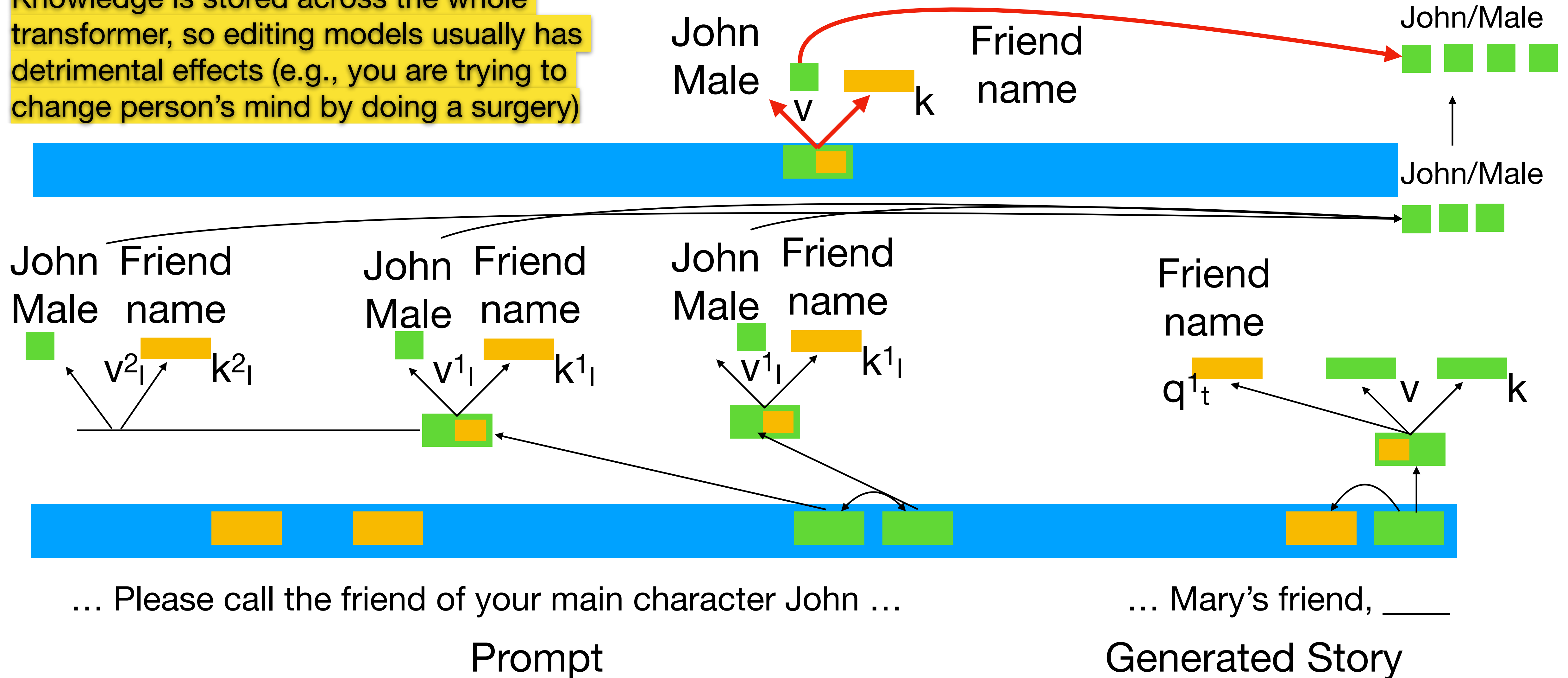
The neuron facilitates the model's ability to predict that token over different context or phrases.



Interpreting Context Look-ups in Transformers: Investigating Attention-MLP Interactions (<https://arxiv.org/pdf/2402.15055v1>)

Distributed Representation

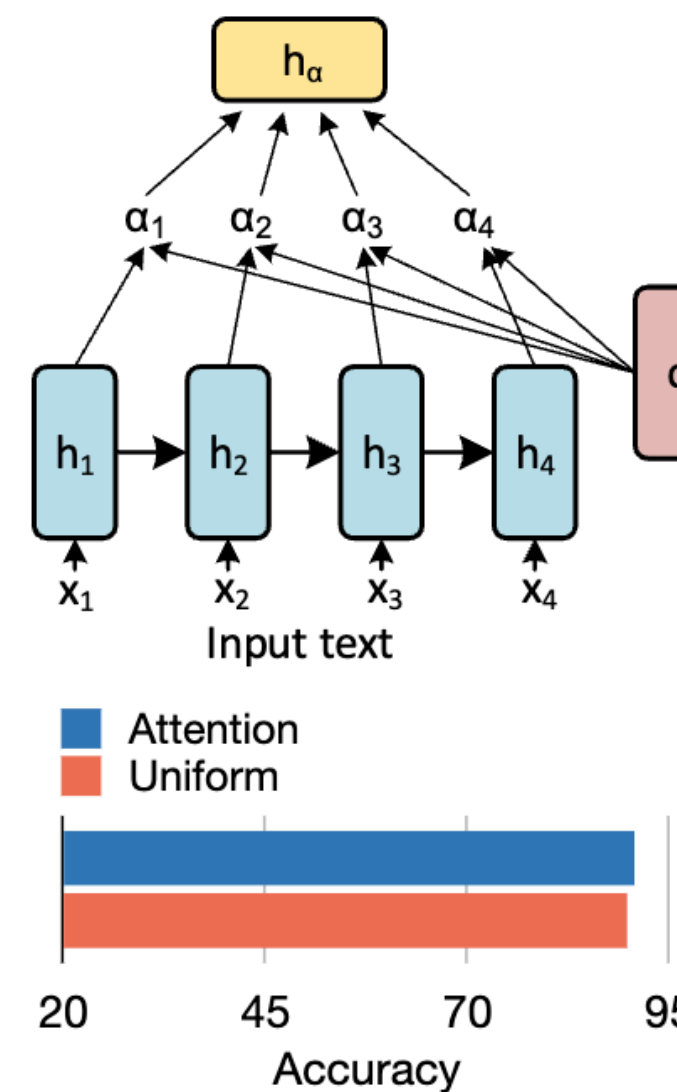
Knowledge is stored across the whole transformer, so editing models usually has detrimental effects (e.g., you are trying to change person's mind by doing a surgery)



Can Attentions Tell you Token “Importance”?

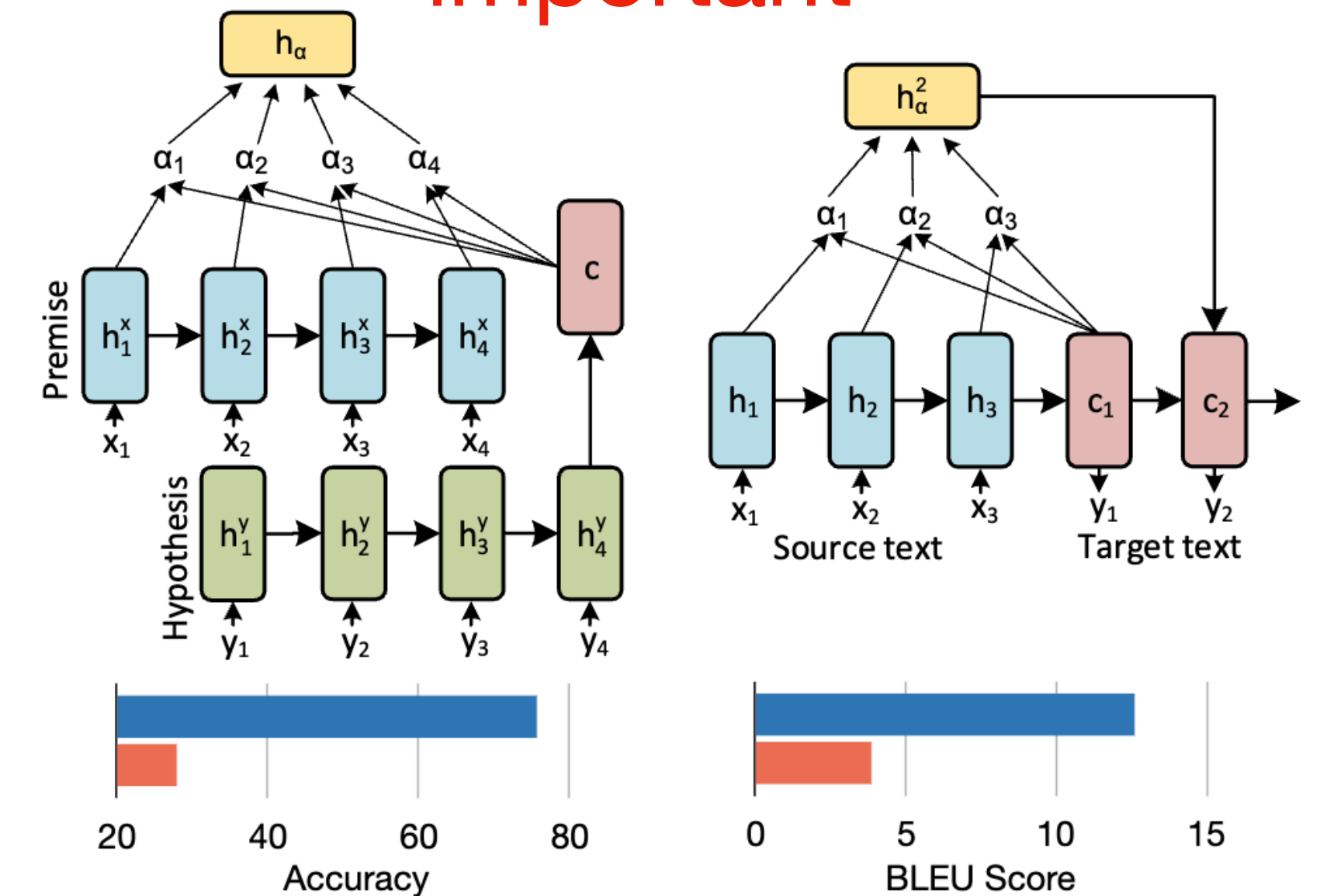
- Multi-head multi-layer self-attention
 - Unlikely
 - But could be used to compress the KV cache
- Multi-head self-attention
 - Probably not
- Single self-attention
 - Task Dependent

Attention is not important



(a) Text Classification

Attention is important



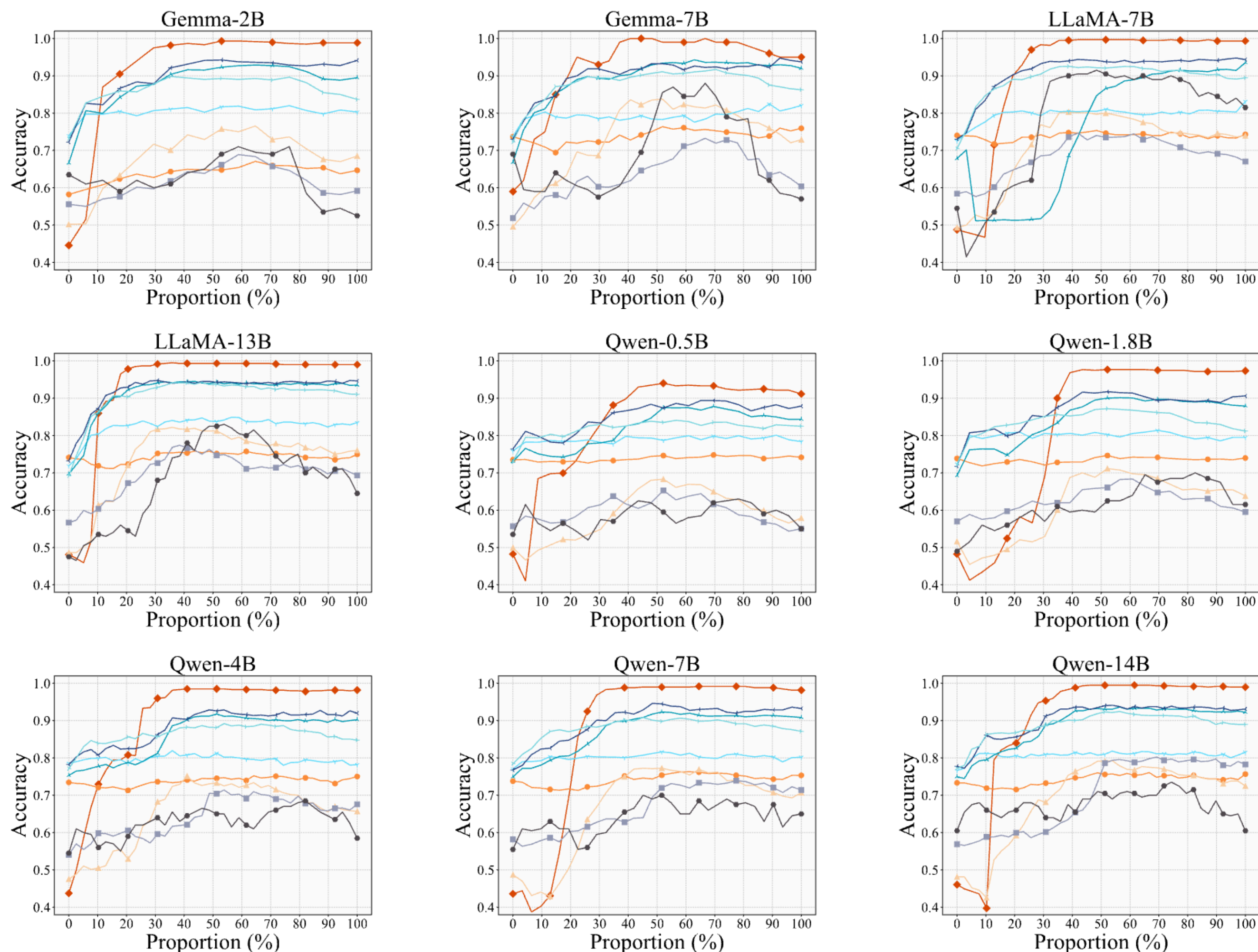
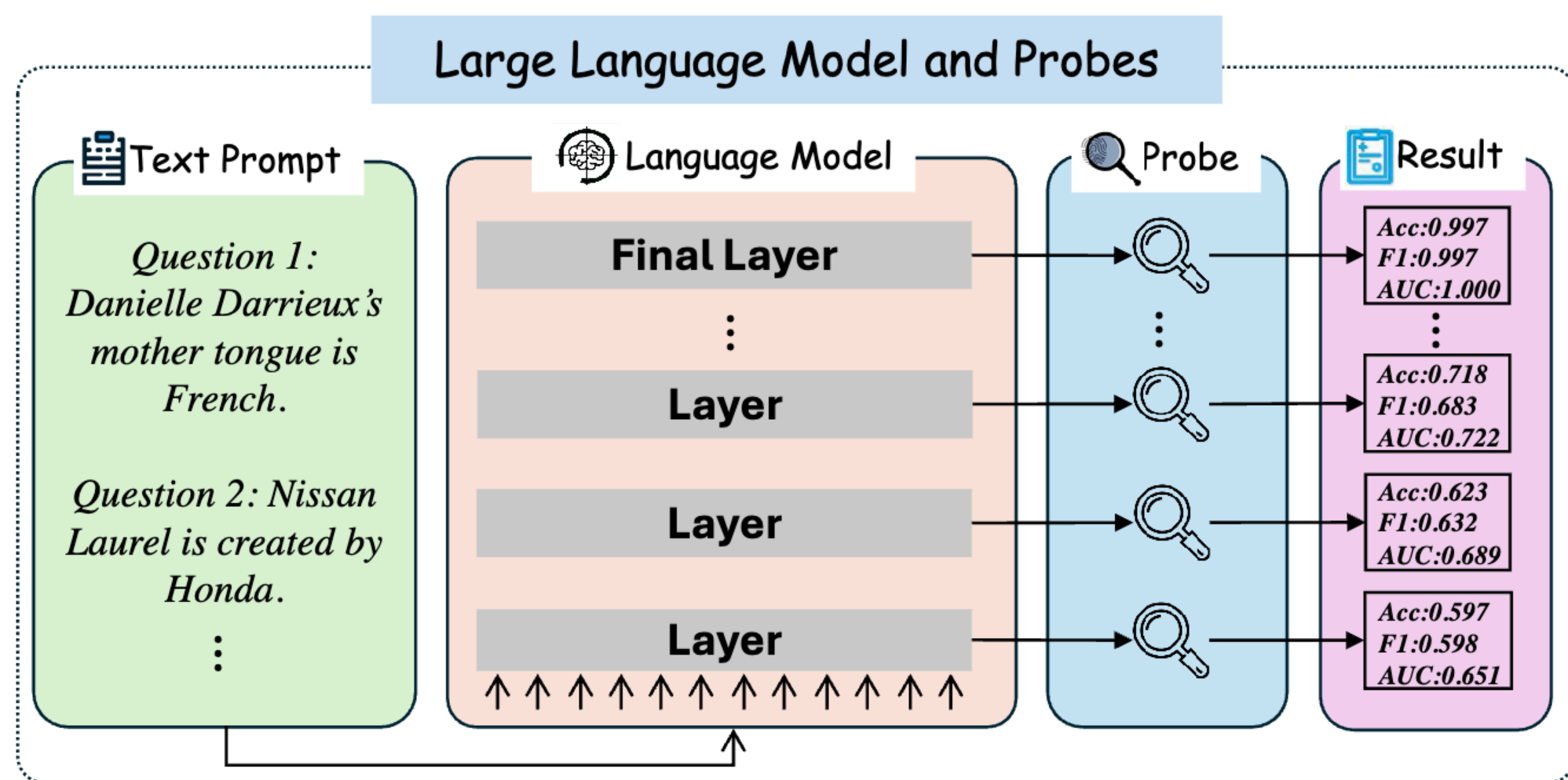
(b) Natural Language Inference

(c) Neural Machine Translation

Figure 1: Comparison of performance with and without neural attention on text classification (IMDB), Natural Language Inference tasks (SNLI) and Neural Machine Translation (News Commentary). Here, α and c denote attention weights and context vector respectively. The results show that attention does not substantially effect performance on text classification. However, the same does not hold for other tasks.

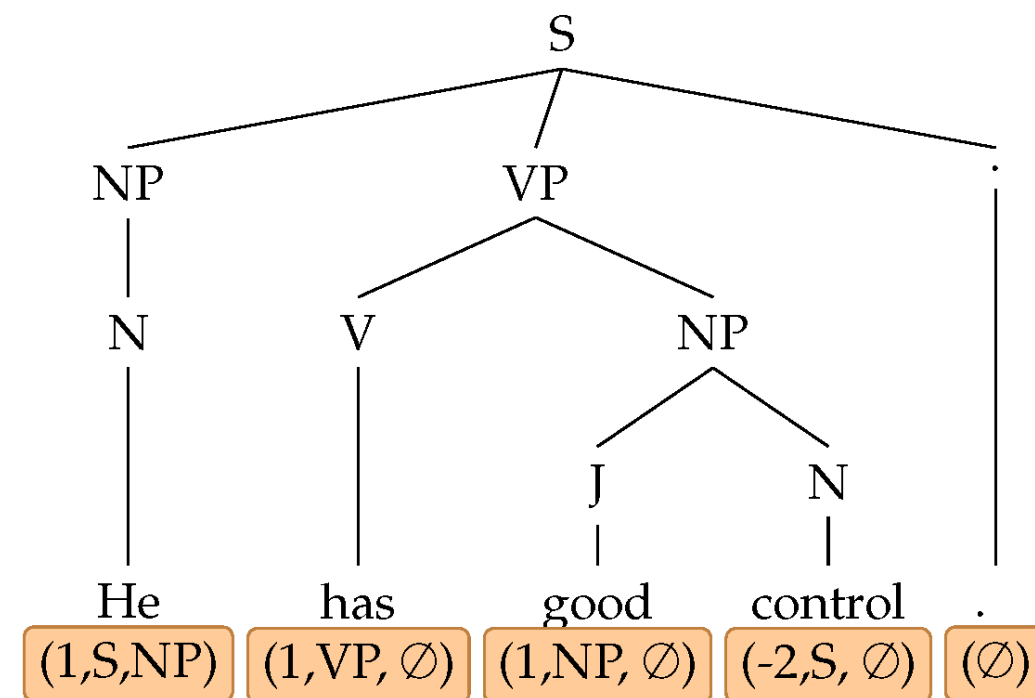
Where are the Facts Stored?

- Facts “tend to” be stored in earlier layers

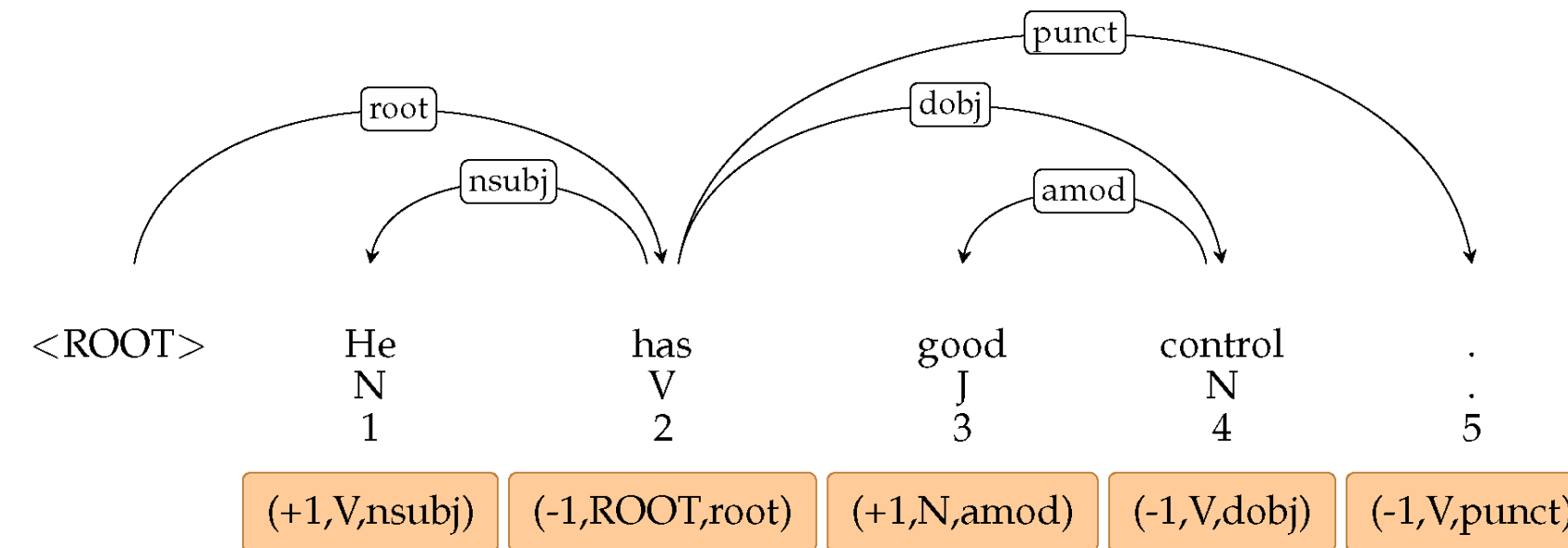


Exploring Concept Depth: How Large Language Models Acquire Knowledge and Concepts at Different Layers? (<https://arxiv.org/pdf/2404.07066>)

Classic NLP Tasks



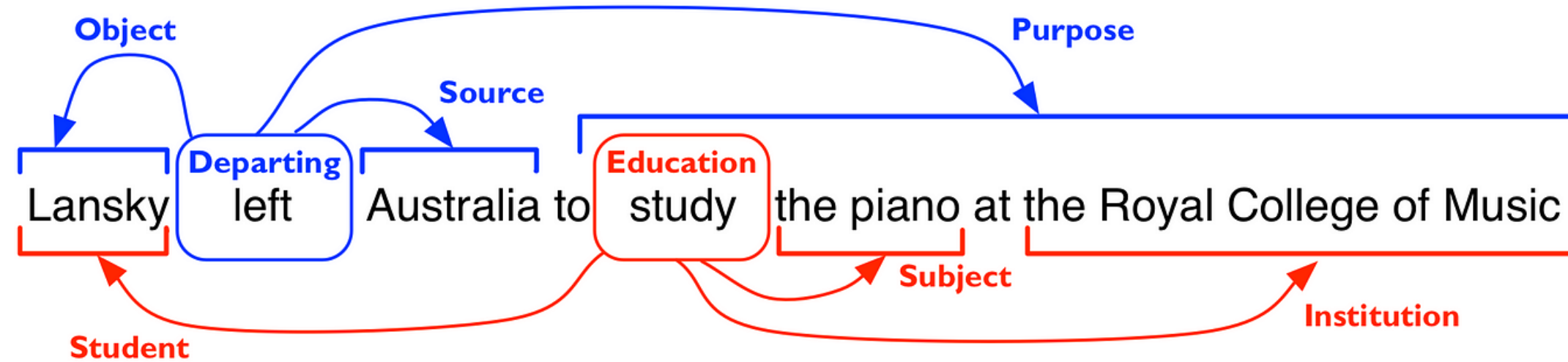
(a) A constituency tree



(b) A dependency tree

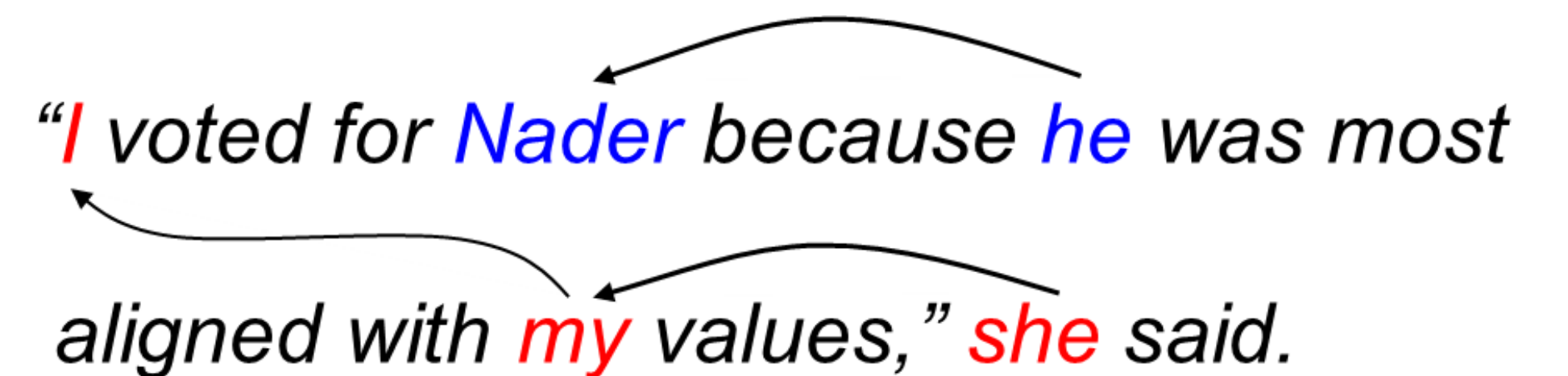
<https://www.mdpi.com/2504-3900/21/1/49>

SRL



<https://medium.com/thedeephub/deciphering-sentences-a-glimpse-into-semantic-role-labeling-with-deep-learning-6b7809bfdcbf>

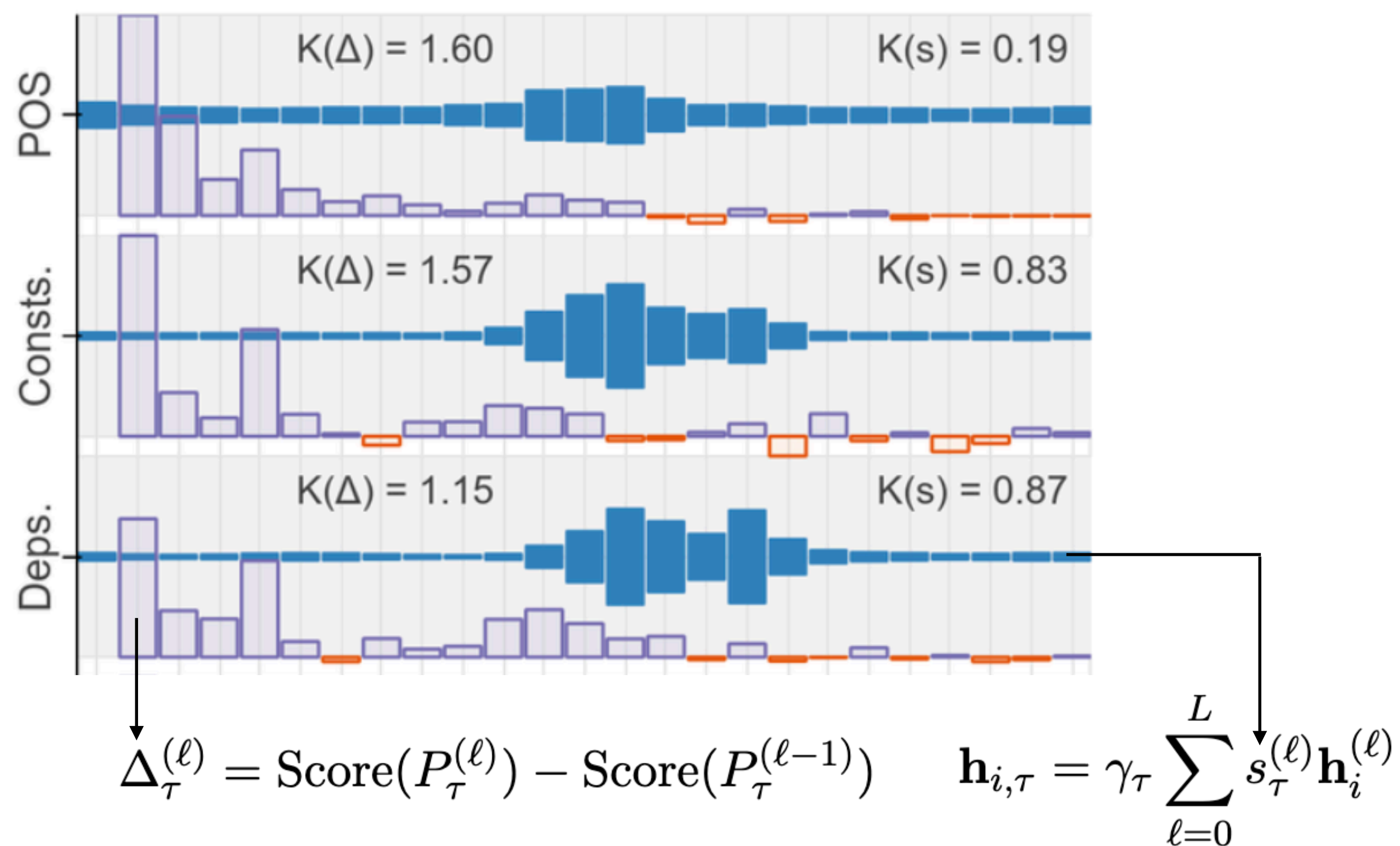
Coref



<https://nlp.stanford.edu/projects/coref.shtml>

What does each Layer Do?

- Higher layers “tend to” handle more semantic information



	F1 Scores		Expected layer & center-of-gravity										
	$\ell=0$	$\ell=24$	0	2	4	6	8	10	12	14	16		
POS	88.5	96.7	3.39	11.68									
Consts.	73.6	87.0	3.79	13.06									
Deps.	85.6	95.5	5.69	13.75									
Entities	90.6	96.1	4.64	13.16									
SRL	81.3	91.4	6.54	13.63									
Coref.	80.5	91.9	9.47	15.80									
SPR	77.7	83.7	9.93	12.72									
Relations	60.7	84.2	9.40	12.83									

BERT Rediscovered the Classical NLP Pipeline (<https://arxiv.org/abs/1905.05950>)

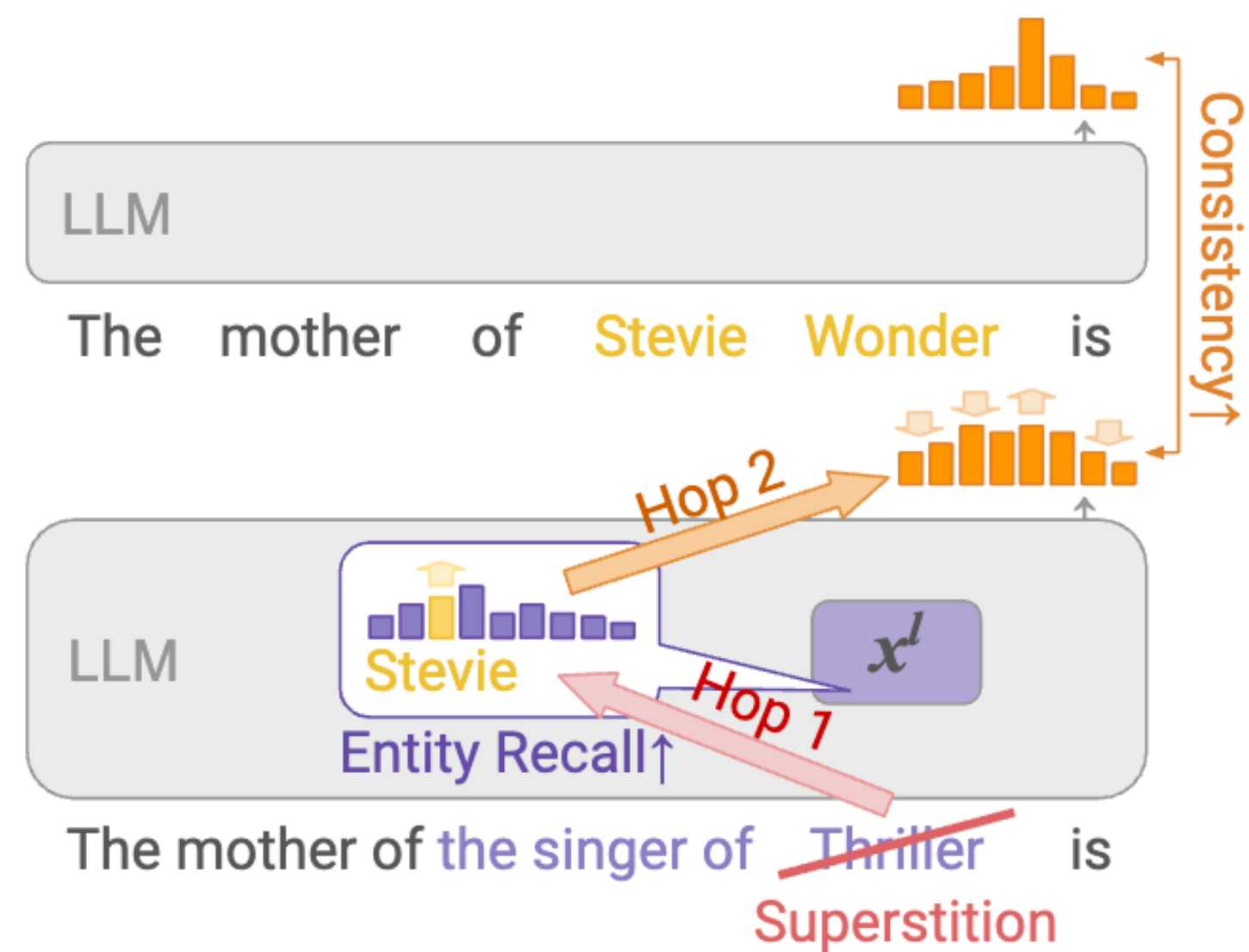


Figure 1: We investigate the latent multi-hop reasoning of LLMs. For the first hop, we change the input prompt to refer to the bridge entity (Stevie Wonder) and check how often it increases the model’s internal recall of the bridge entity. For the second hop, we check if increasing this recall causes the model output to be more consistent with respect to what it knows about the bridge entity’s attribute (mother of Stevie Wonder).

Do Large Language Models Latently Perform Multi-Hop Reasoning? (<https://arxiv.org/pdf/2402.16837>)

Input query: In the year **Scarlett Johansson** was born, the Summer Olympics were hosted in **the country of**

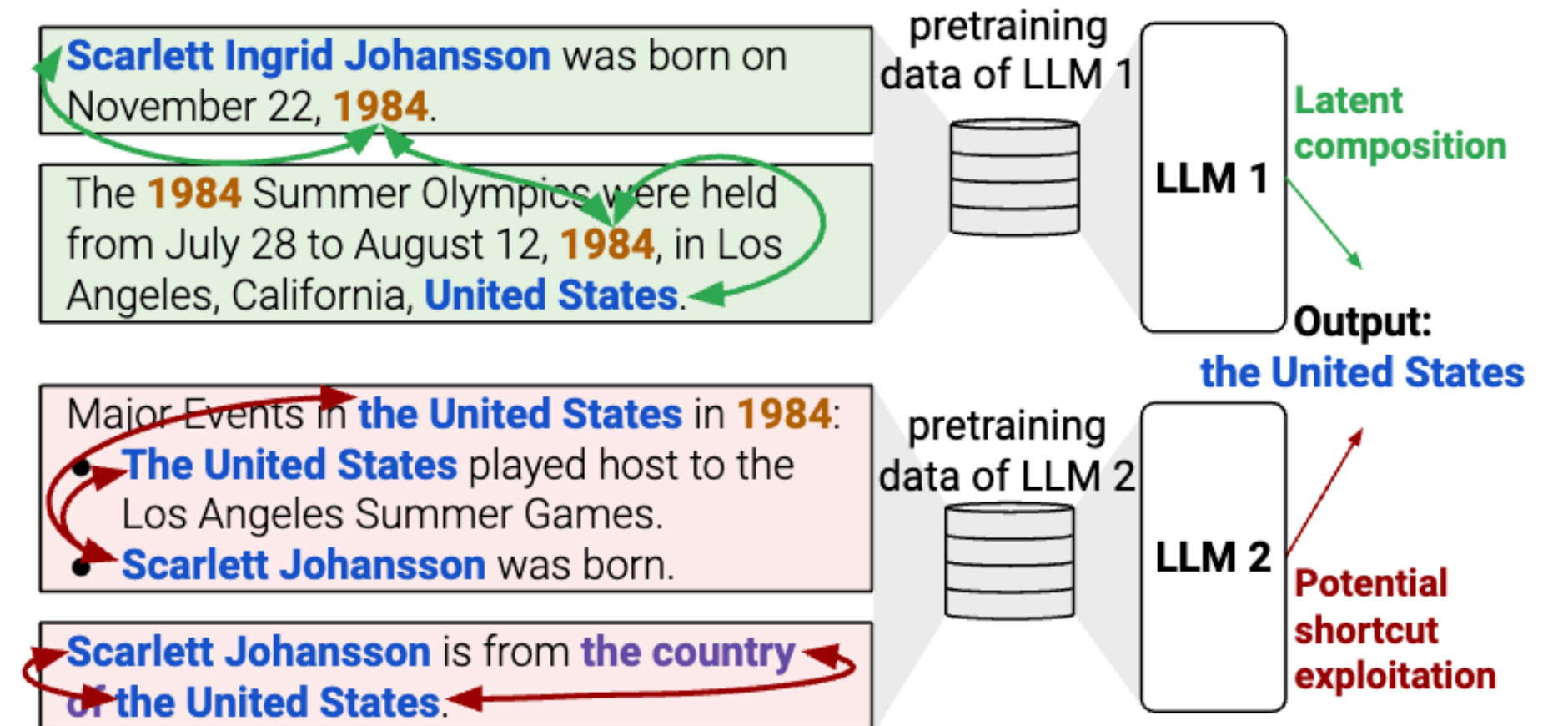
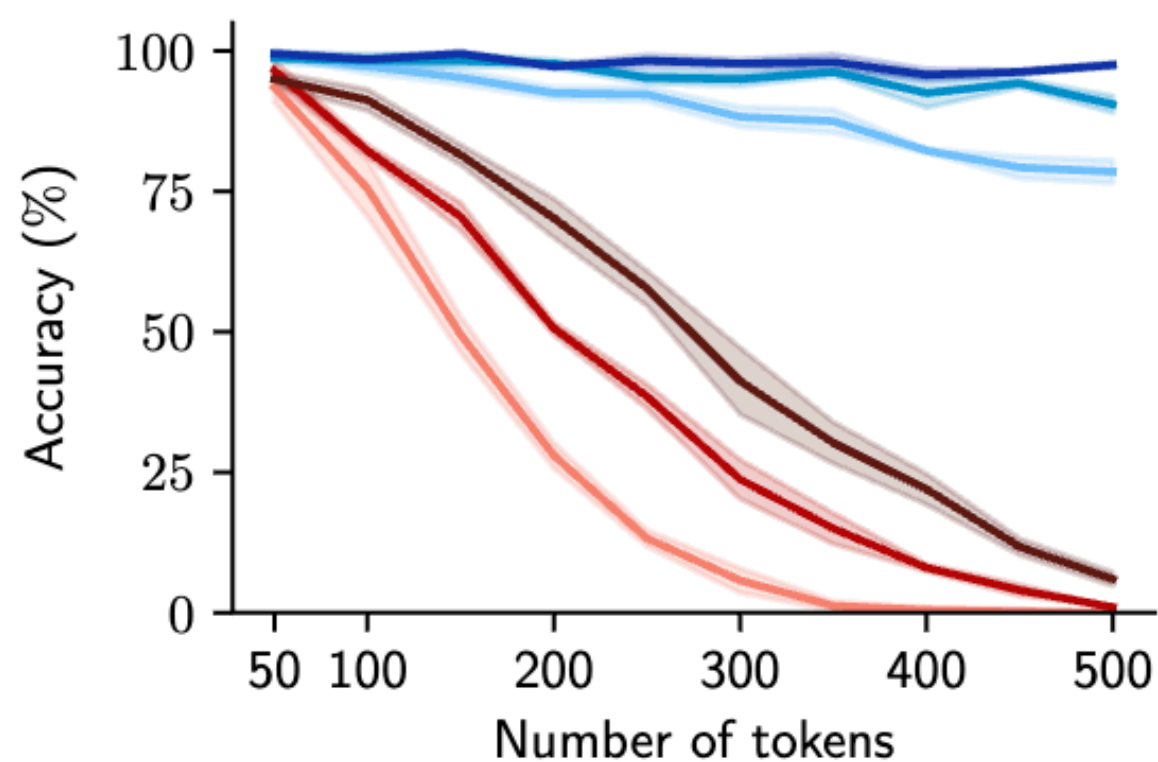


Figure 1: Evaluation of latent multi-hop reasoning should exclude cases where LLMs can bypass the process of latently composing the single-hop facts by exploiting shortcuts. LLMs can develop shortcuts when they frequently encounter the head entity (“*Scarlett Johansson*”) or the relation pattern in the query (“*the country of*”) with the answer entity (“*United States*”). We propose desiderata for shortcut-free evaluation of latent multi-hop reasoning ability to minimize the chance of shortcuts.

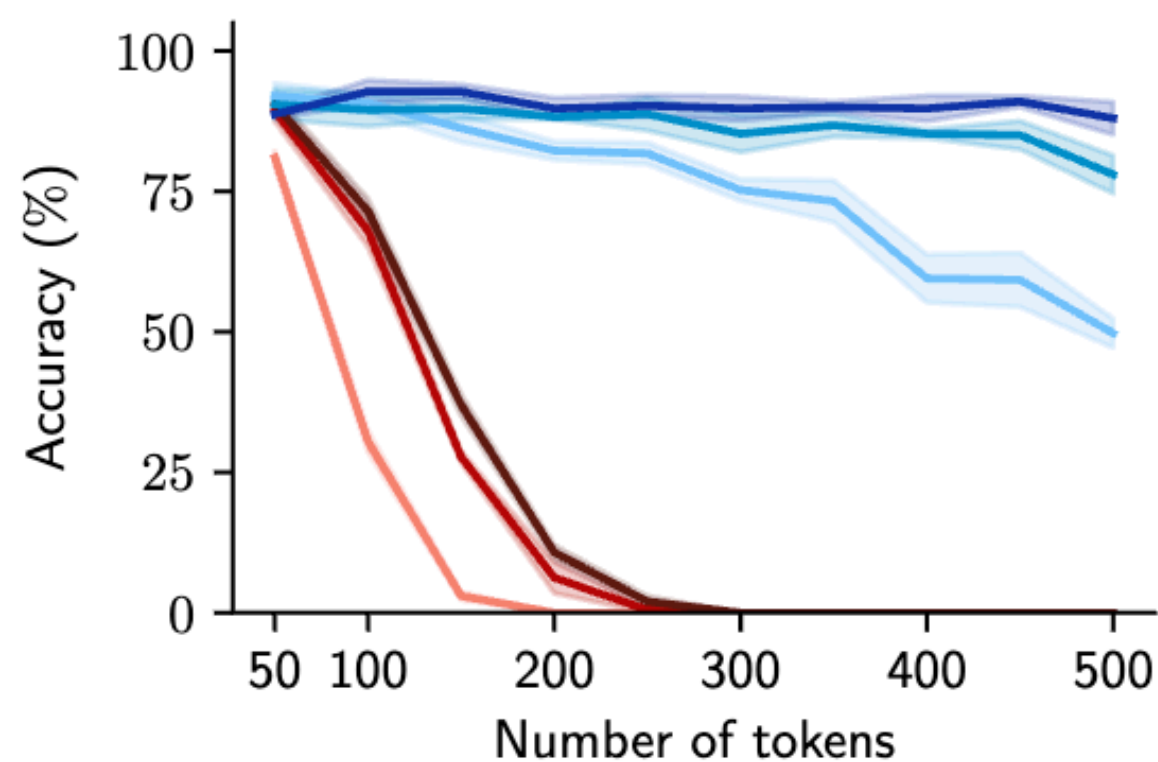
Do Large Language Models Perform Latent Multi-Hop Reasoning without Exploiting Shortcuts? (<https://arxiv.org/pdf/2411.16679>)

Why is Attention Effective?



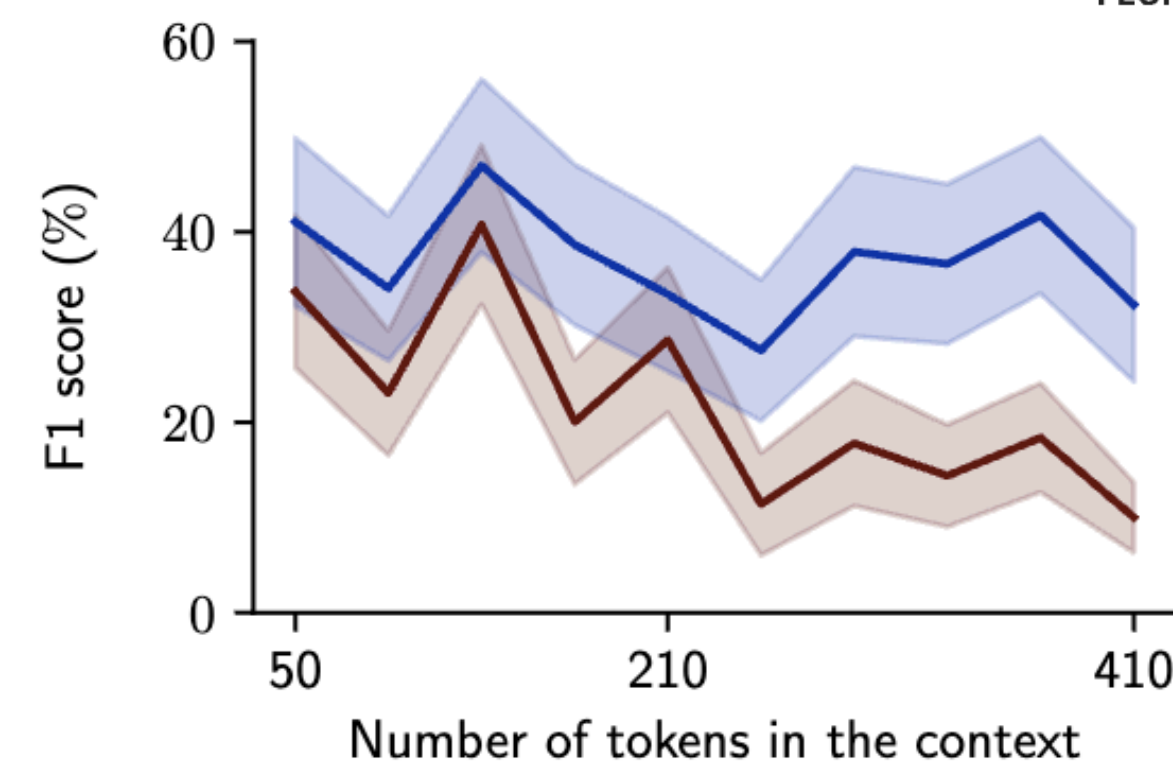
Pythia: 410M 1.4B 2.8B
Mamba: 360M 1.4B 2.8B

(a) Copy: natural language strings



Pythia: 410M 1.4B 2.8B
Mamba: 360M 1.4B 2.8B

(b) Copy: shuffled strings



Pythia: 2.8B
Mamba: 2.8B

(c) Question answering (SQUAD)

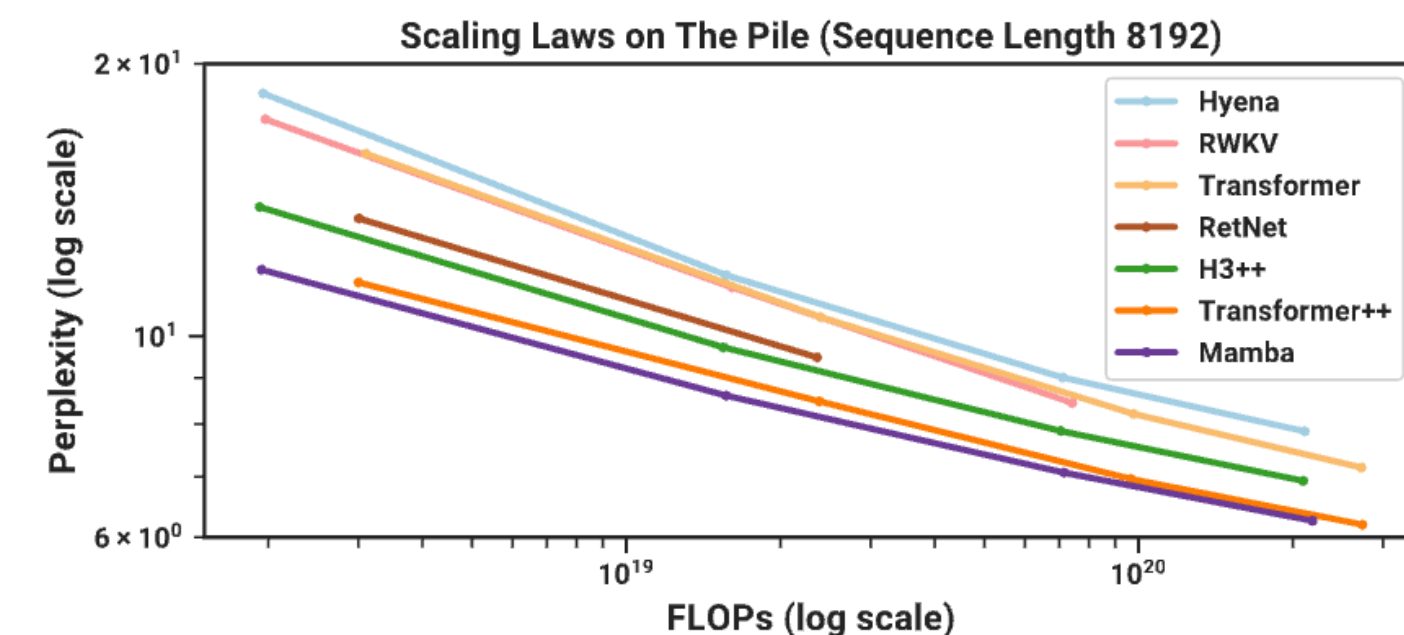
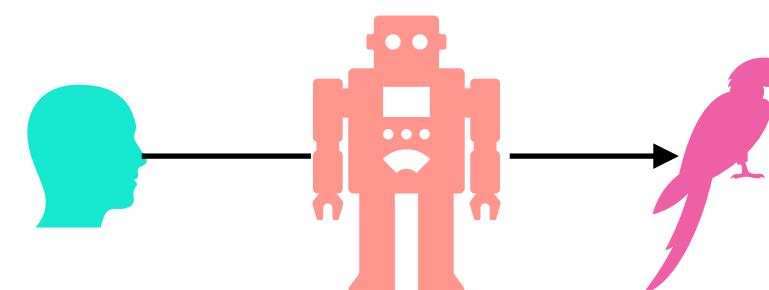
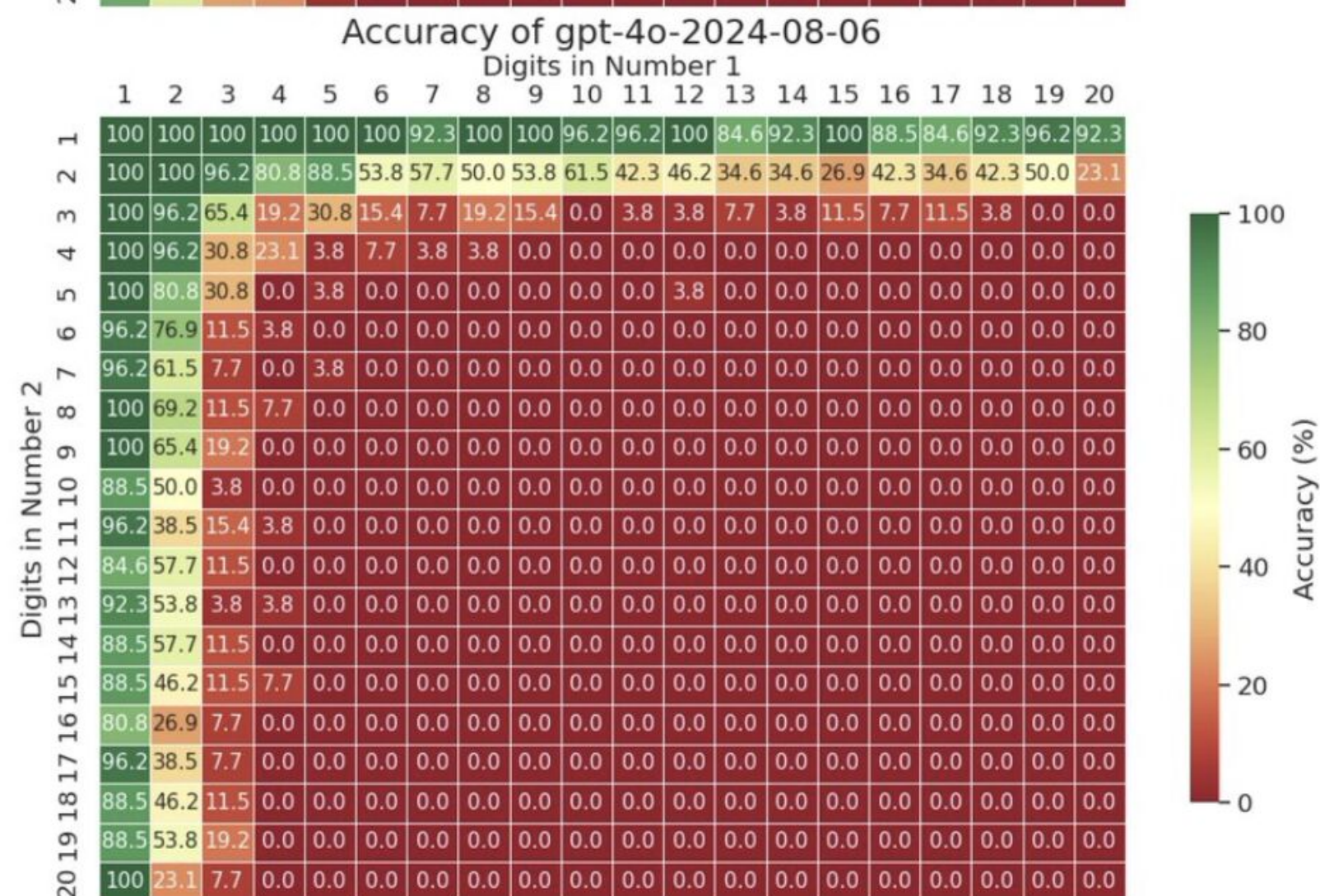
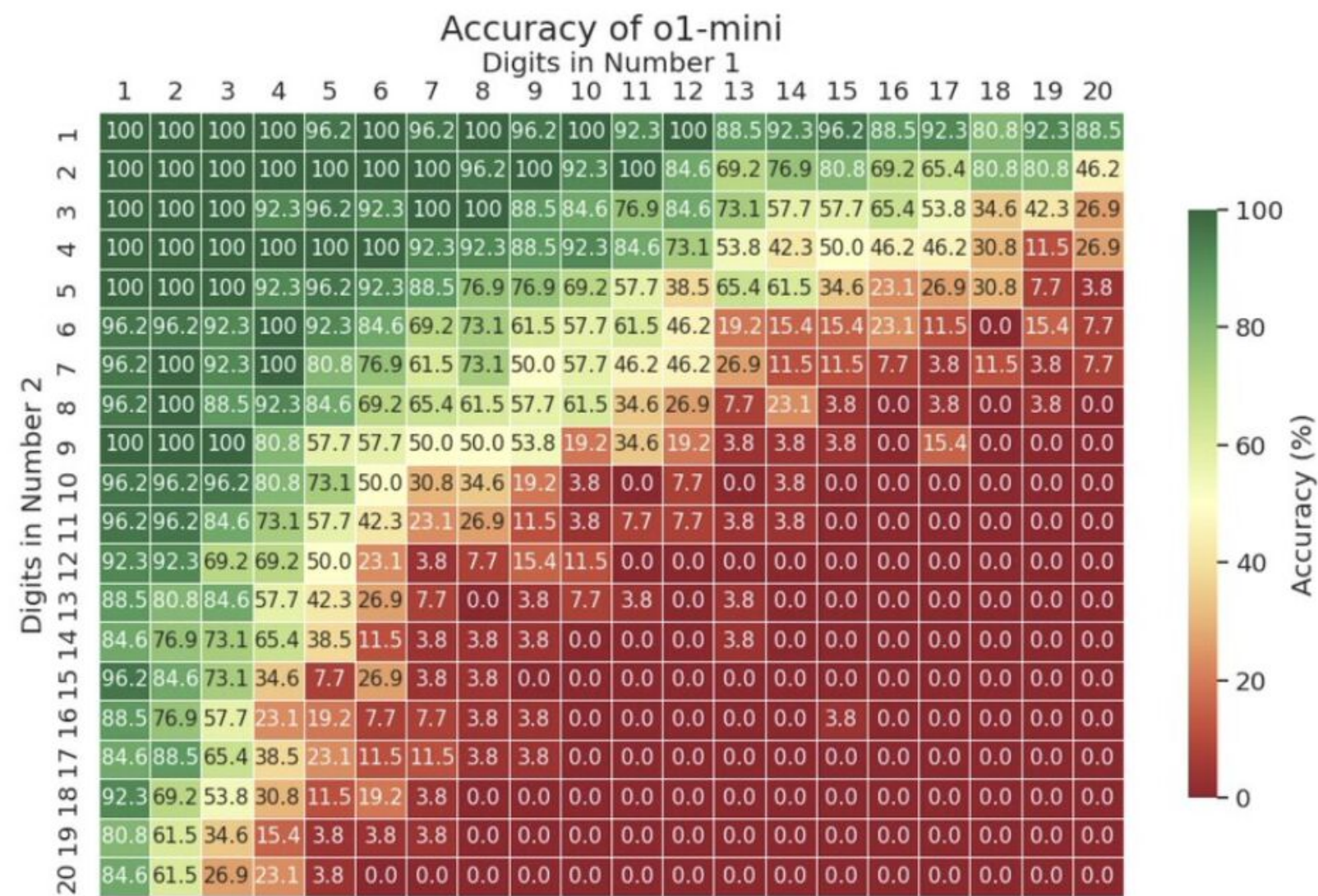


Figure 7. (a) Copy: natural language strings. We compare pretrained models on their ability to copy natural language strings sampled from C4 of varying lengths and report string-level accuracy. The transformer models substantially outperform the GSSMs. **(b) Copy: shuffled strings.** To test whether it mattered that the strings were in natural language, we randomly shuffle the word order of the strings from the previous experiment. We find that this degrades performance, especially for the Mamba models. **(c) Question answering (SQUAD).** We compare Pythia and Mamba on a standard question answering dataset where we bin the dataset based on the length of the context paragraph. We find that Mamba performance decays more quickly with the length of the context.

Repeat After Me: Transformers are Better than State Space Models at Copying (<https://arxiv.org/abs/2402.01032>)

Arithmetic Computation Limitation

- Transformer architecture design
- MLP retrieves the relevant memory from training
- Self-attention retrieves and merges the memories
MLP retrieves
- Transformers are not designed to handle logic



<https://x.com/yuntiandeng/status/1889704768135905332>

```
function multiply (x[1..p], y[1..q]):
  // multiply x for each y[i]
  for i = q to 1
    carry = 0
    for j = p to 1
      t = x[j] * y[i]
      t += carry
      carry = t // 10
      digits[j] = t mod 10
      summands[i] = digits

  // add partial results (computation not shown)
  product = \sum_{i=1}^q summands[q+1-i] \cdot 10^{i-1}
  return product
```

$A(\mathbf{x})$

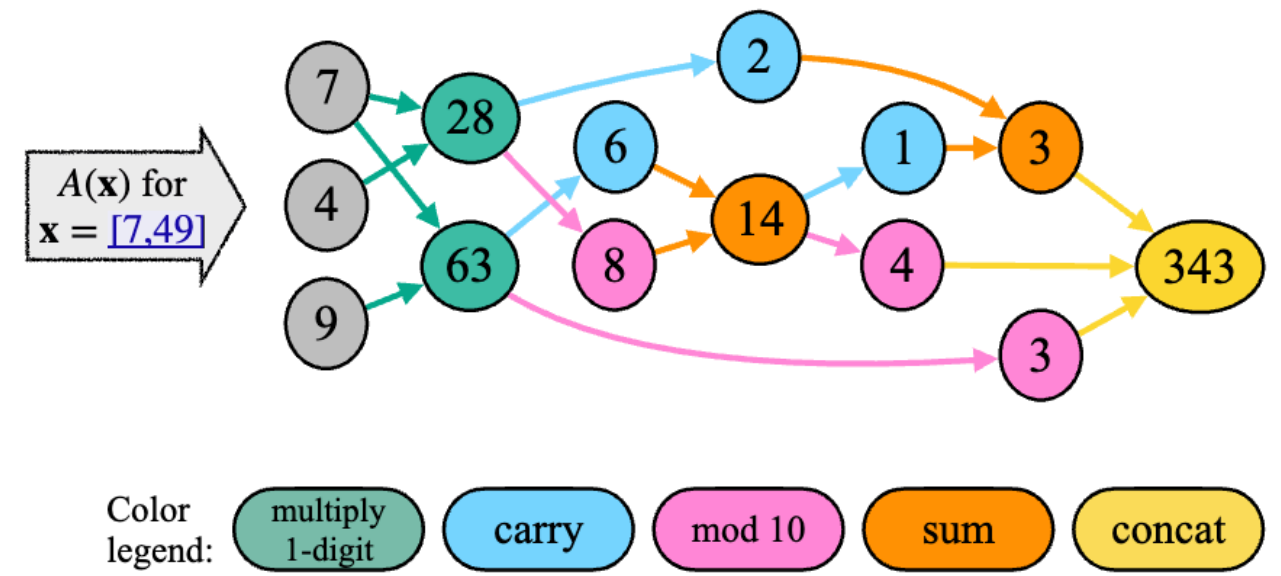
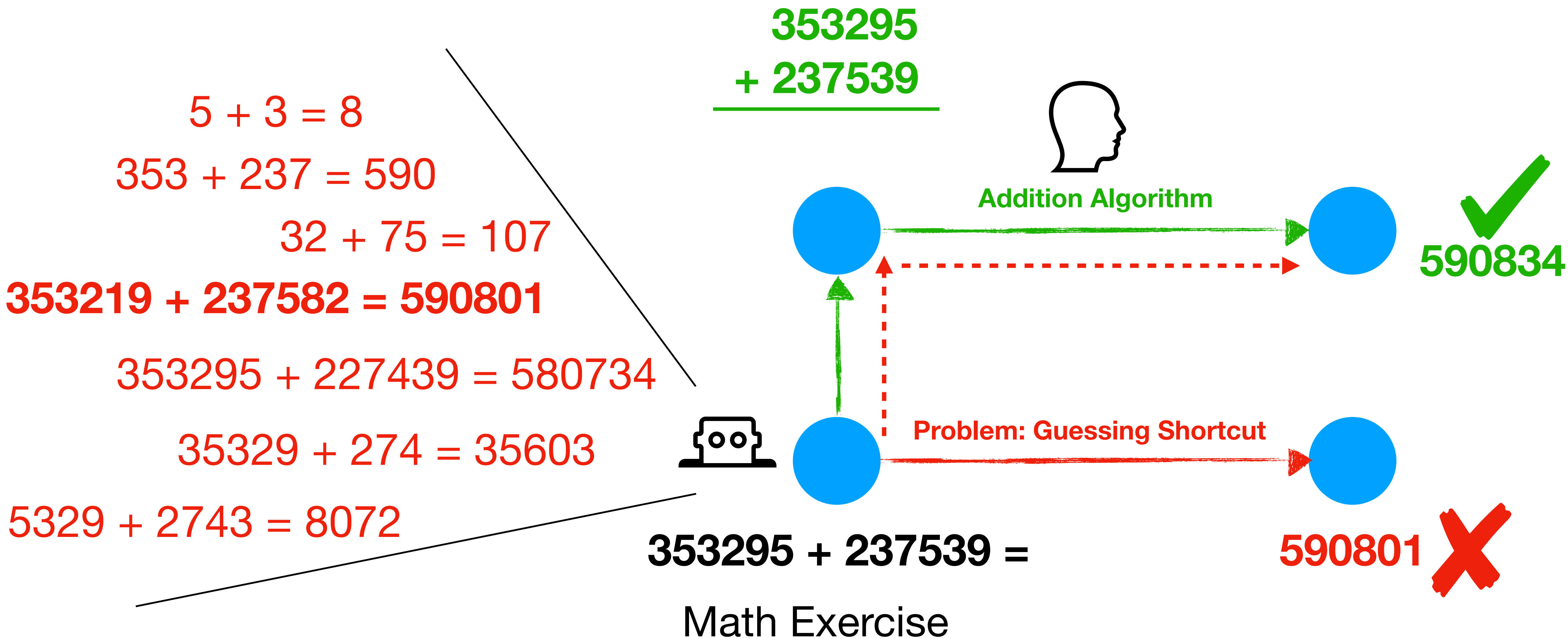


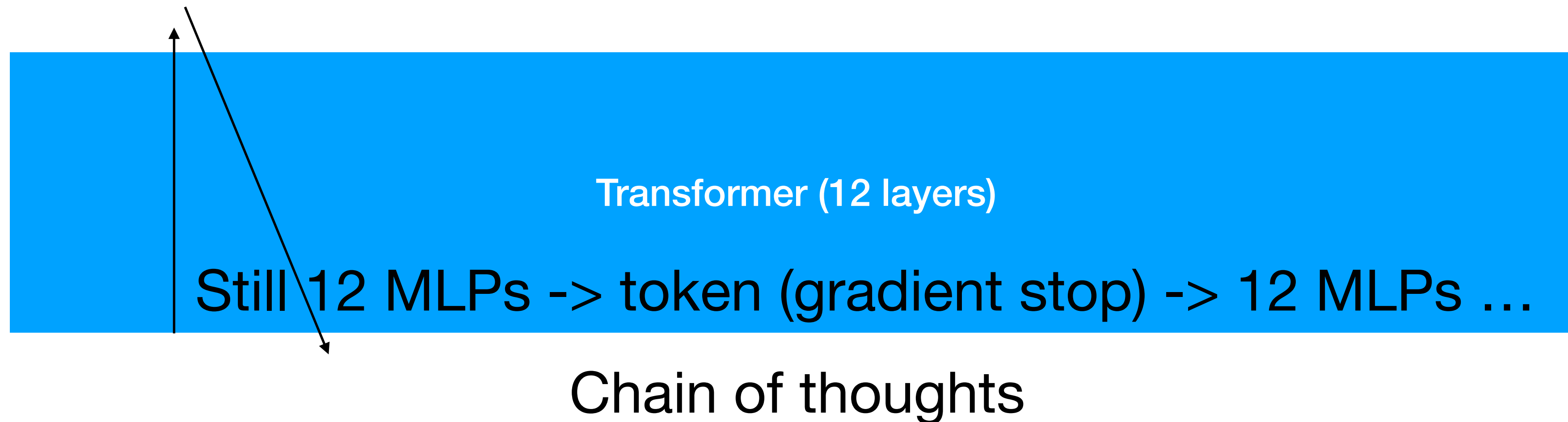
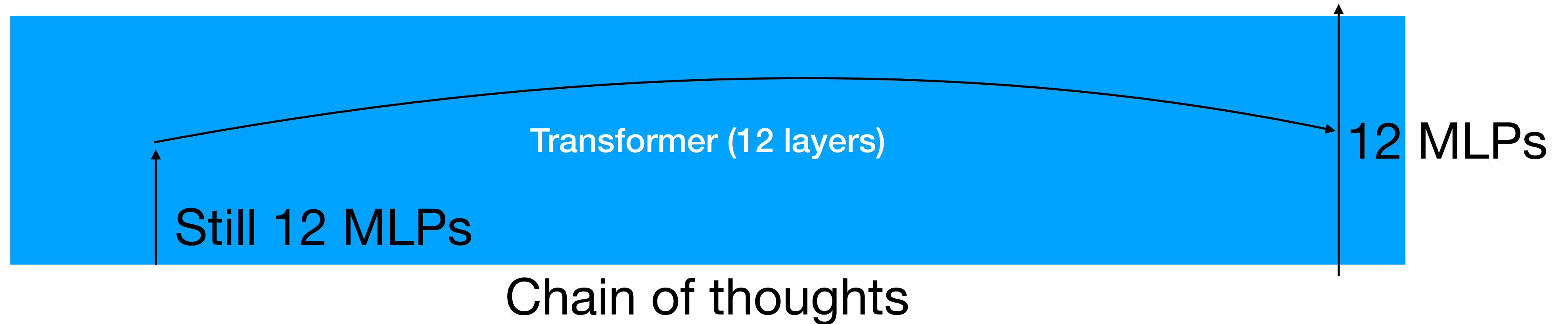
Figure 1: Transformation of an algorithm A to its computational graph $G_{A(\mathbf{x})}$. The depicted example is of long-form multiplication algorithm A , for inputs $\mathbf{x} = [7, 49]$ (i.e. computing 7×49).

Faith and Fate: Limits of Transformers on Compositionality (https://proceedings.neurips.cc/paper_files/paper/2023/file/deb3c28192f979302c157cb653c15e90-Paper-Conference.pdf)

Arithmetic Computation Limitation



Limited Reasoning Chain Length



Midterm Example Questions

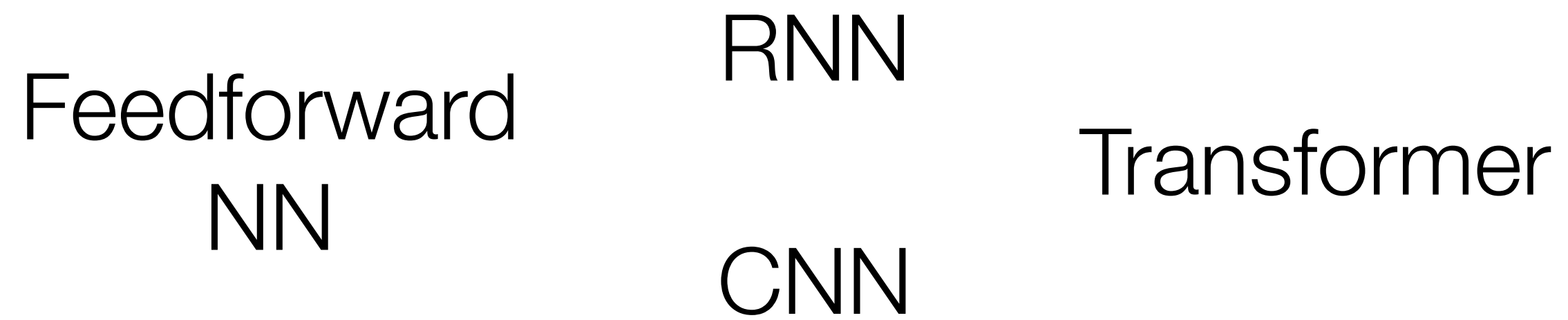
Q1: RNN represents one sequence using one embedding, but Transformer also represents one sequence using one embedding. Why does Transformer mitigate the embedding bottleneck problem?

Q2: Which model is more expensive to train? RNN or Self-attention?

Self-Attention vs RNN vs CNN

- CNN
 - Pros
 - Super Fast
 - Cons
 - Bad long distance
- RNN
 - Pros
 - Long reasoning chain length
 - Low inference memory
 - Fast Inference
 - Good at positional Information
 - Cons
 - Copy difficulty
 - Simple operation or slow parallel training
- Self-attention
 - Pros
 - Long Distance Copy
 - Parallel Training
 - Good at handling a set
 - Cons
 - Limited reasoning chain length
 - High inference memory
 - Relatively Slow Inference
 - $O(n^2)$

Sequence Processing



Computation Cost

More



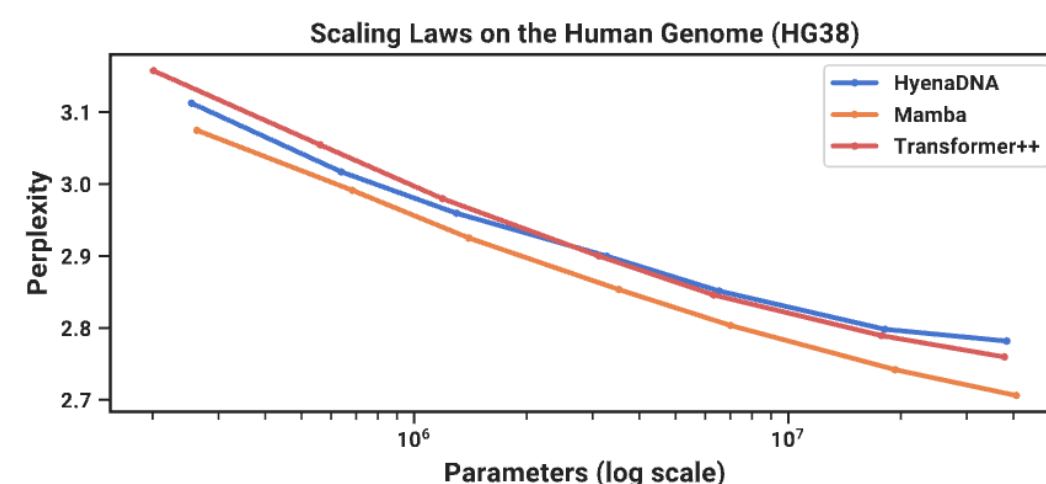
Long Dependency

More



Position Information

Less



Mamba: Linear-Time Sequence Modeling with Selective State Spaces (<https://arxiv.org/pdf/2312.00752>)