

Neural Networks: A Theory Lab

Attention!

Hadi Daneshmand



UNIVERSITY
of VIRGINIA

Outline

2

- ▶ Introduction to LLMs and self-supervised learning
- ▶ Attention layers in LLMs
- ▶ A group activity



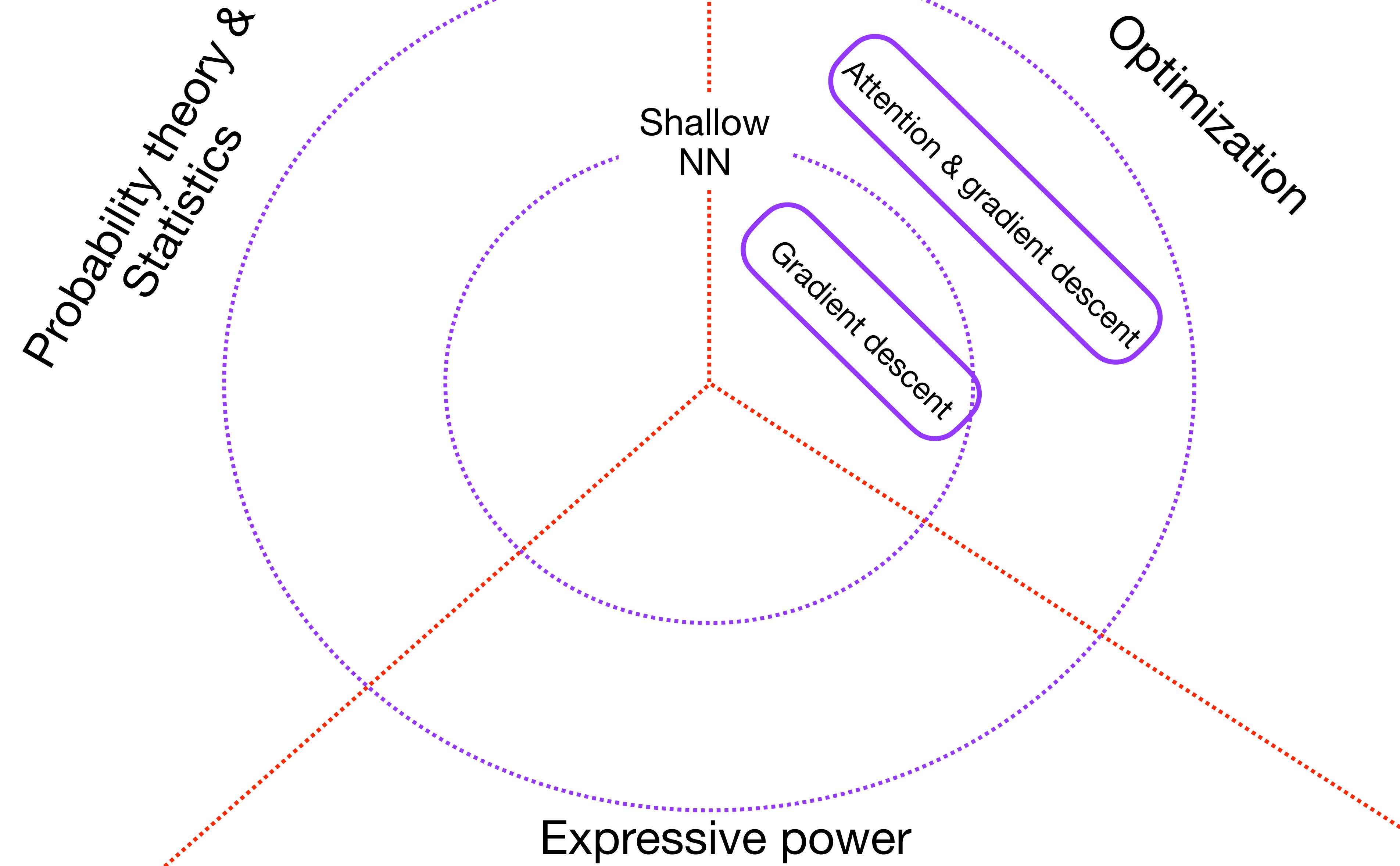
Question about accessibility of slides

3

- ▶ Are slides easy to read in .pdf?



Big picture



Large Language Models (LLMs)



Natural Language Processing before LLMs

6

Sentiment Analysis

Negative. I do not like this product 

Positive. I like this product 

Translation

English: This is a book

Dutch: Dit is een boek



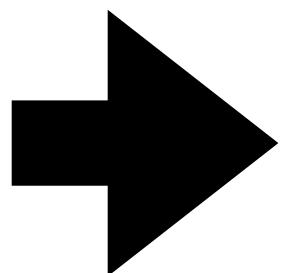
Supervised learning formulation

7

I do not like this product



x



negative



$$y \in \{0,1\}$$

- ▶ Two challenges:
 - How to represent x as a vector?
 - Labeling data is labor-intensive and time consuming



LLMs do not need human supervision

8



How can I help you today?



Self-supervised learning

9

- ▶ Self-supervised learning aims to generate labels from data itself.

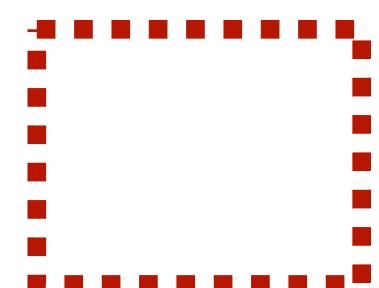
This text is

used to teach a model

- ▶ Input x



- ▶ Output y



The power of self-supervised learning

10

Sentiment Analysis

Prompt. “I do not like this product”
is positive?

Output. No

Translation

Prompt. Translate “this is a
book” to Dutch

Output: Dit is een boek



Traditional word embedding



Traditional NLP: Word Embedding

12

► How to represent a sentence/word?

- Example: I do not like this product

► One solution:

- construct a vector of the size of dictionary

$$this = [0, \dots, 0, \underbrace{1}, 0, \dots, 0]$$

• index of "this" in dictionary

► Question: What is the issue with the above binary representation?

► Solution: Does not capture semantic



Traditional word embedding

13

- ▶ w is a word with vector representations $v_w \in \mathbb{R}^d$
- ▶ All text in world can be presented as a sequence of words

$$w_1, w_2, \dots, w_i, \dots,$$

$$\triangleright p(w_i | w_j) = \begin{cases} 1 & w_i = \arg \max_w \langle v_w, v_{w_j} \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$\triangleright \max_{v_{w_i}} \sum_{j-c \leq k \leq j+c} \log p(w_{j+k} | w_j)$$

Reference: Distributed Representations of Words and Phrases and their Compositionalities: <https://arxiv.org/abs/1310.4546>



Softmax

14

- ▶ Consider model $p(w_i | w_j) = \begin{cases} 1 & w_i = \arg \max_w \langle v_w, v_{w_j} \rangle \\ 0 & \text{otherwise} \end{cases}$
- ▶ **Challenge:** The above model is not differentiable
- ▶ **Softmax:** $p(w_i | w_j) = \frac{e^{\langle v_{w_i}, v_{w_j} \rangle}}{\sum_k e^{\langle v_{w_k}, v_{w_j} \rangle}}$



Semantic of embeddings

15

► $v_{\text{madrid}} - v_{\text{spain}} + v_{\text{france}} = ?$

- $\approx v_{\text{paris}}$

► $v_{\text{czech}} + v_{\text{currency}} = ?$

- $\approx v_{\text{koruna}}$

► $v_{\text{german}} + v_{\text{airline}} = ?$

- $\approx v_{\text{Lufthansa}}$

Reference: Distributed Representations of Words and Phrases and their Compositionalities: <https://arxiv.org/abs/1310.4546>



Traditional NLP: Inference using word embeddings

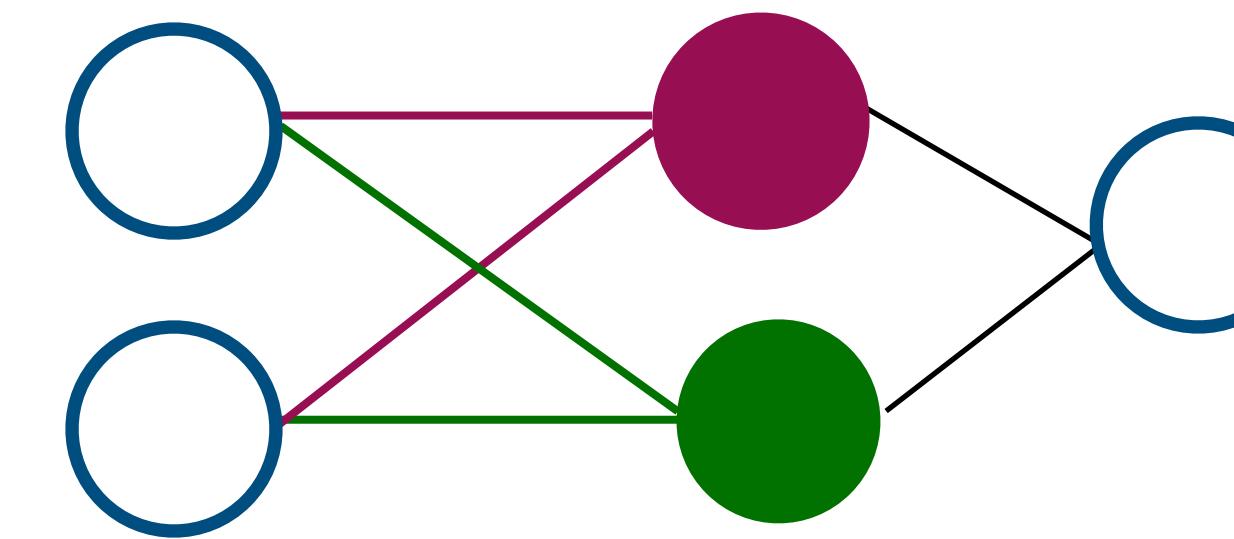
16

I do not like this product



$v_1, v_{\text{do}}, v_{\text{not}}, v_{\text{like}}, \dots$

preprocessing



negative or positive

$$y \in \{0, 1\}$$

$$\arg \min_{\theta} \mathbb{E}_{v_i, y} \|NN_{\theta}(v_1, \dots) - y\|^2$$



UNIVERSITY
of VIRGINIA

The issue with word embedding

17

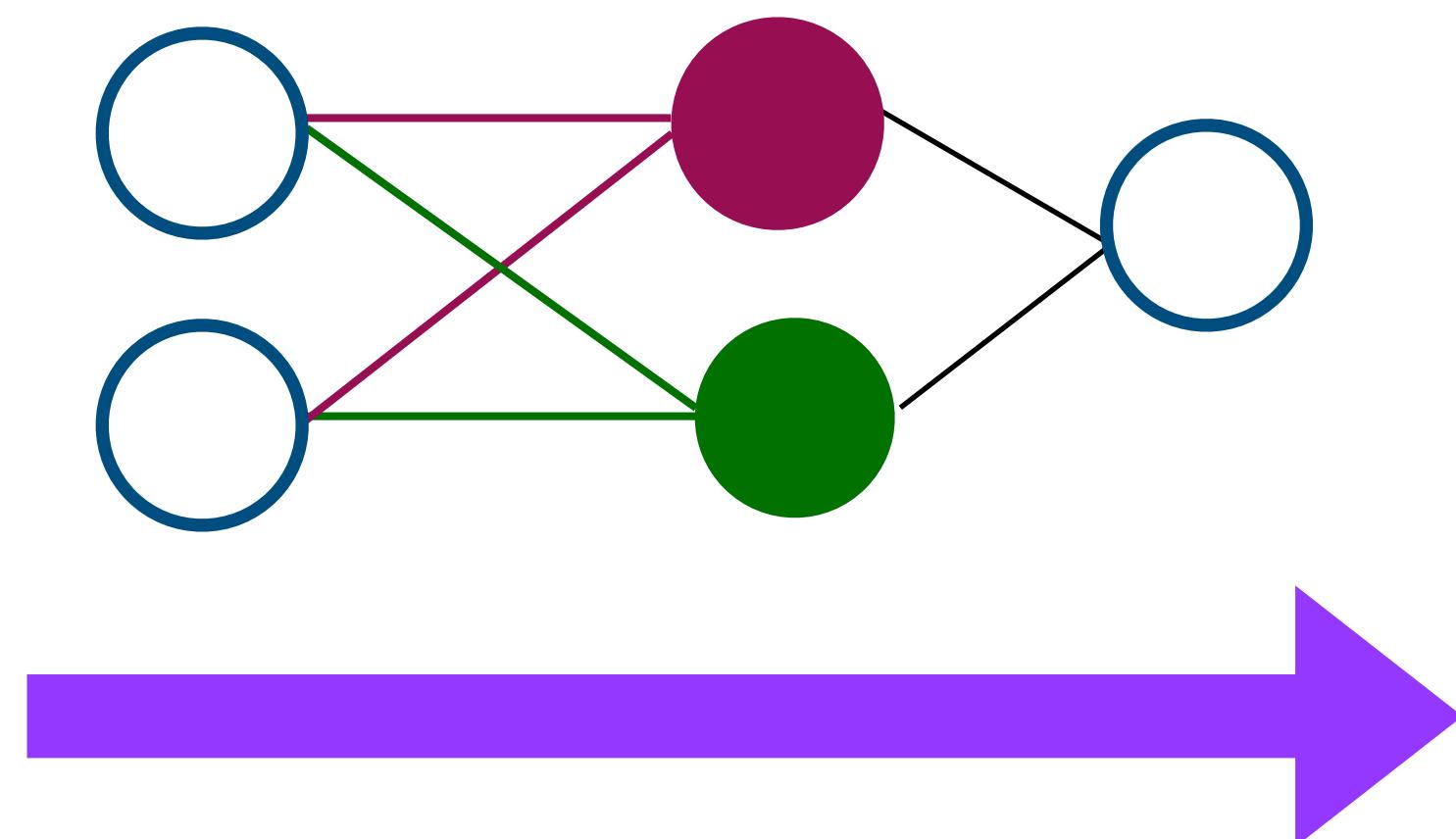
- ▶ **Question:** Word embedding is context free, why this is an issue?
- ▶ Consider word “light”
 - This bag is very (light: not heavy)
 - She turned on the (**light**: brightness)
- ▶ Word embedding for light is a fixed vector v_{light}
- ▶ **Conclusion:** while semantic is **context sensitive**, traditional word embedding is **context free**



LLMs: context-sensitive word embedding

18

She turned on the **light**
Context free
 v_{words}
 v_{light}



Transformers
Context sensitive
 v'_{words}
 v'_{light}

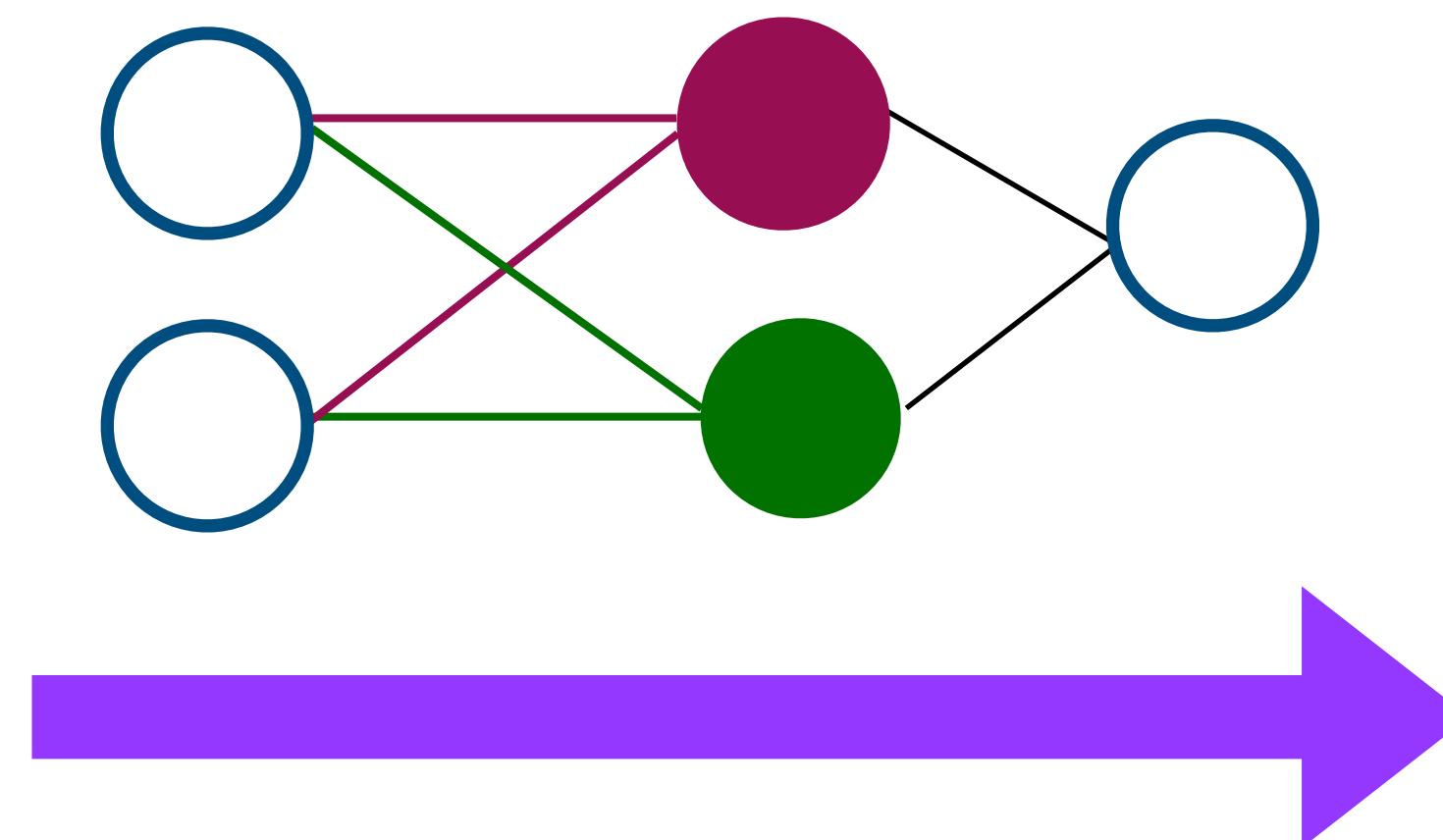


UNIVERSITY
of VIRGINIA

LLMs: context-sensitive word embedding

19

This bag is very **light**
Context free
 v_{words}
 v_{light}

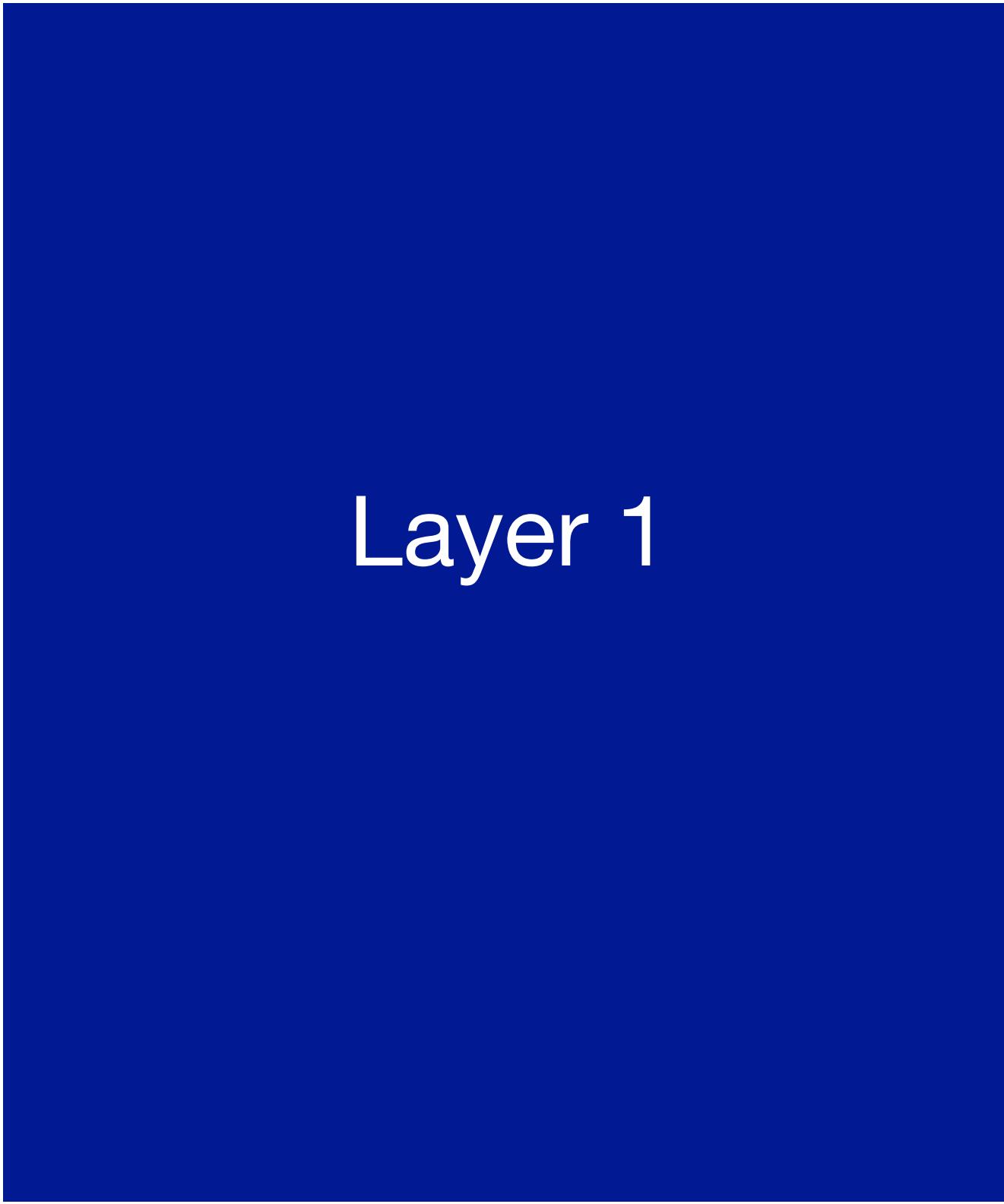


~~v_{light}~~
 v'_{words}
 v'_{light}
Context sensitive

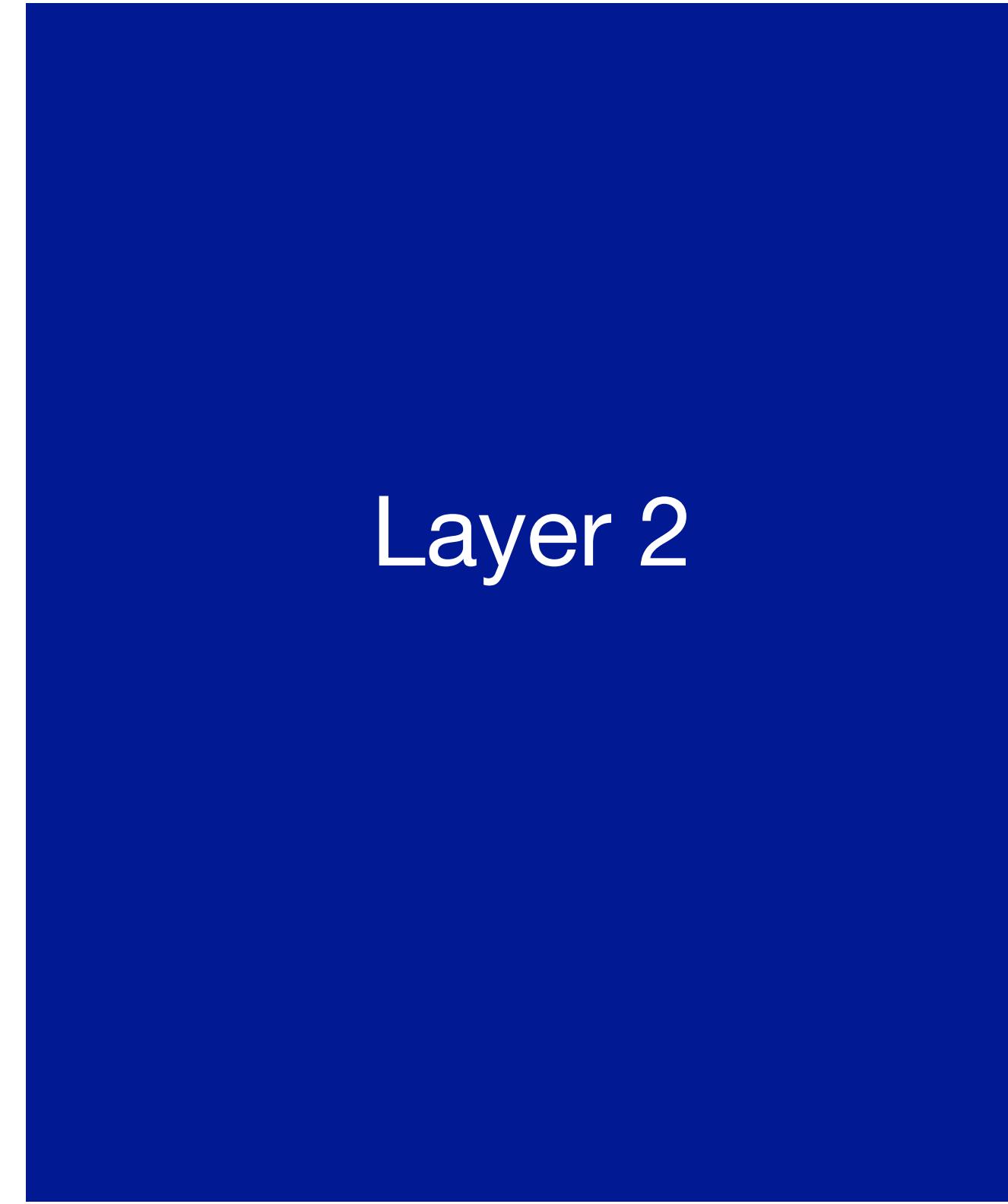
Context sensitive transformers

20

$$v_1, v_2, v_3, \dots \in \mathbb{R}^d$$



$$v'_1, v'_2, v'_3, \dots \in \mathbb{R}^d$$



$$v''_1, v''_2, v''_3, \dots \in \mathbb{R}^d$$



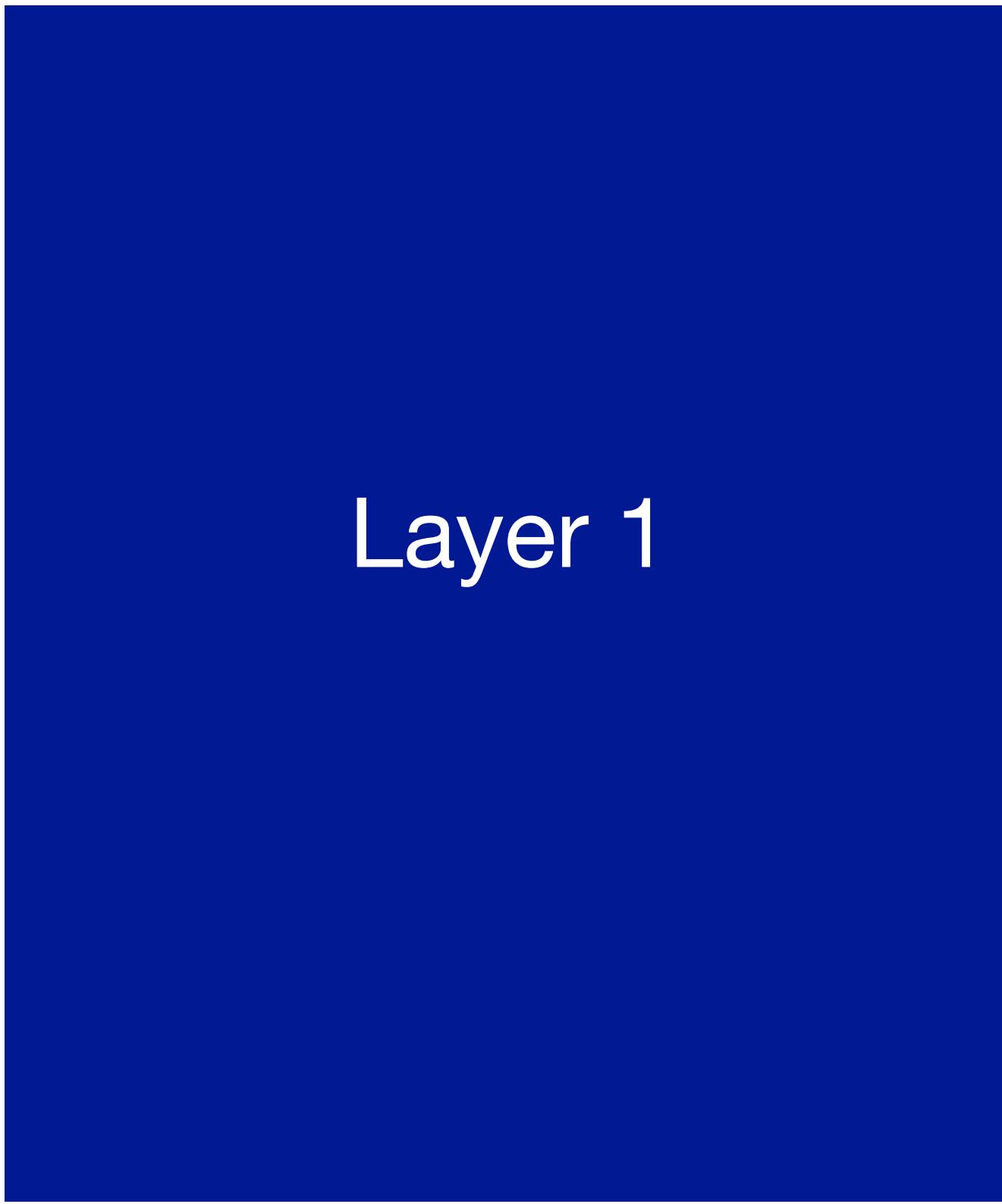
How to generate context-sensitive embeddings?



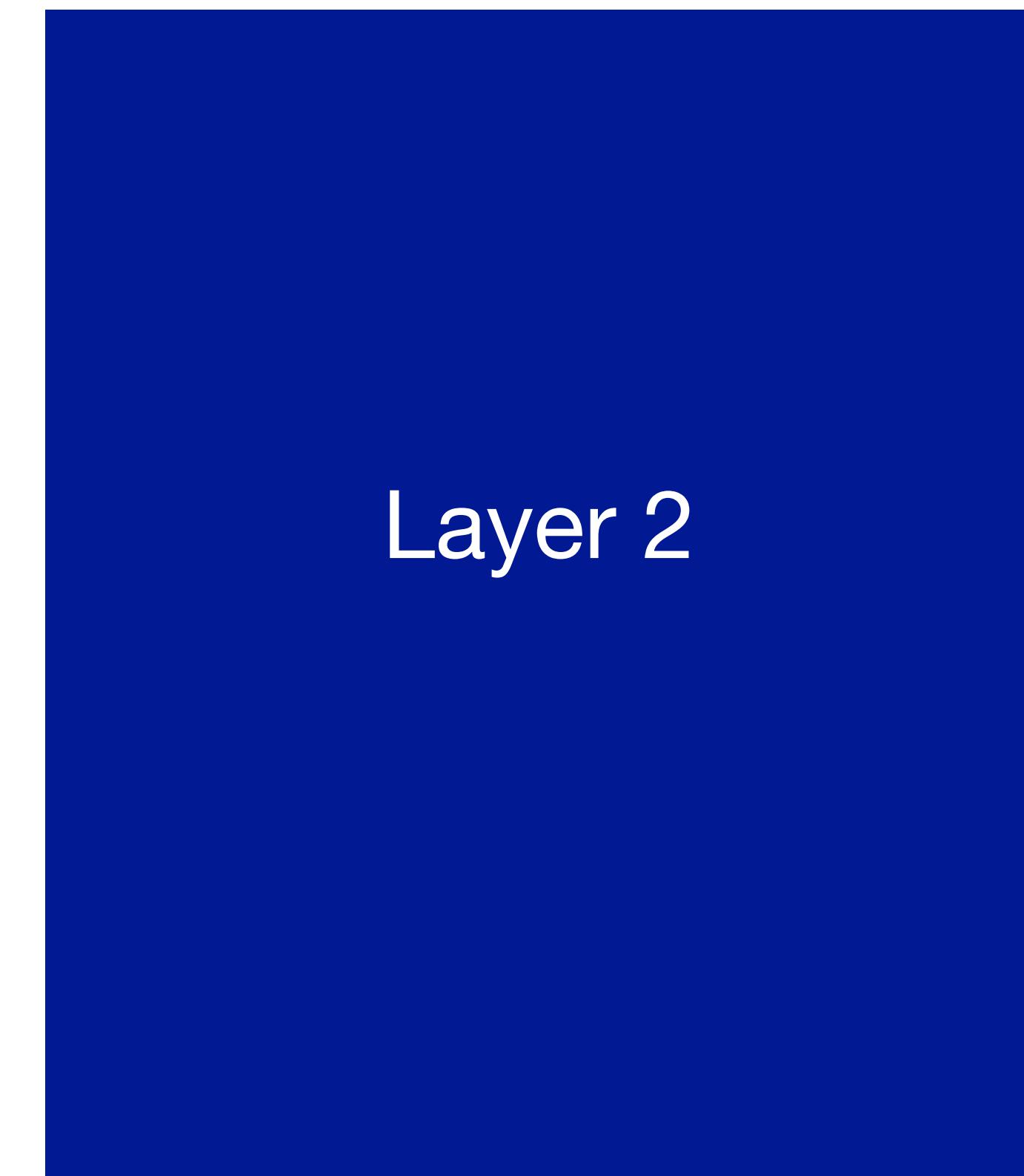
Context sensitive transformers

22

$$v_1, v_2, v_3, \dots \in \mathbb{R}^d$$



$$v'_1, v'_2, v'_3, \dots \in \mathbb{R}^d$$



$$v'_1, v'_2, v'_3, \dots \in \mathbb{R}^d$$



UNIVERSITY
of VIRGINIA

Outline

23

- ▶ Introduction to LLMs and self-supervised learning
- ▶ Attention component in LLMs
- ▶ A group activity



Development of attention over time

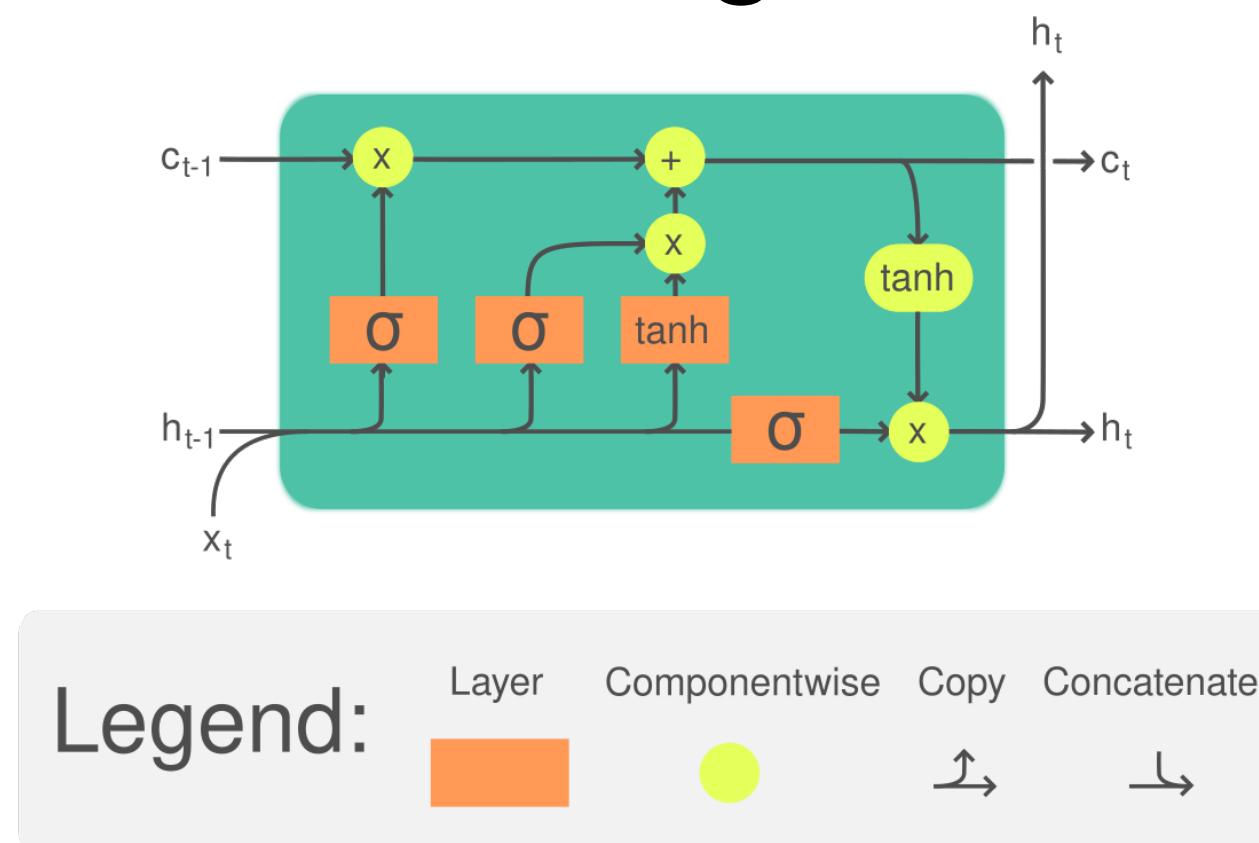
24

Model: LSTM
Long short-term memory

Goal: Talking signal vanishing

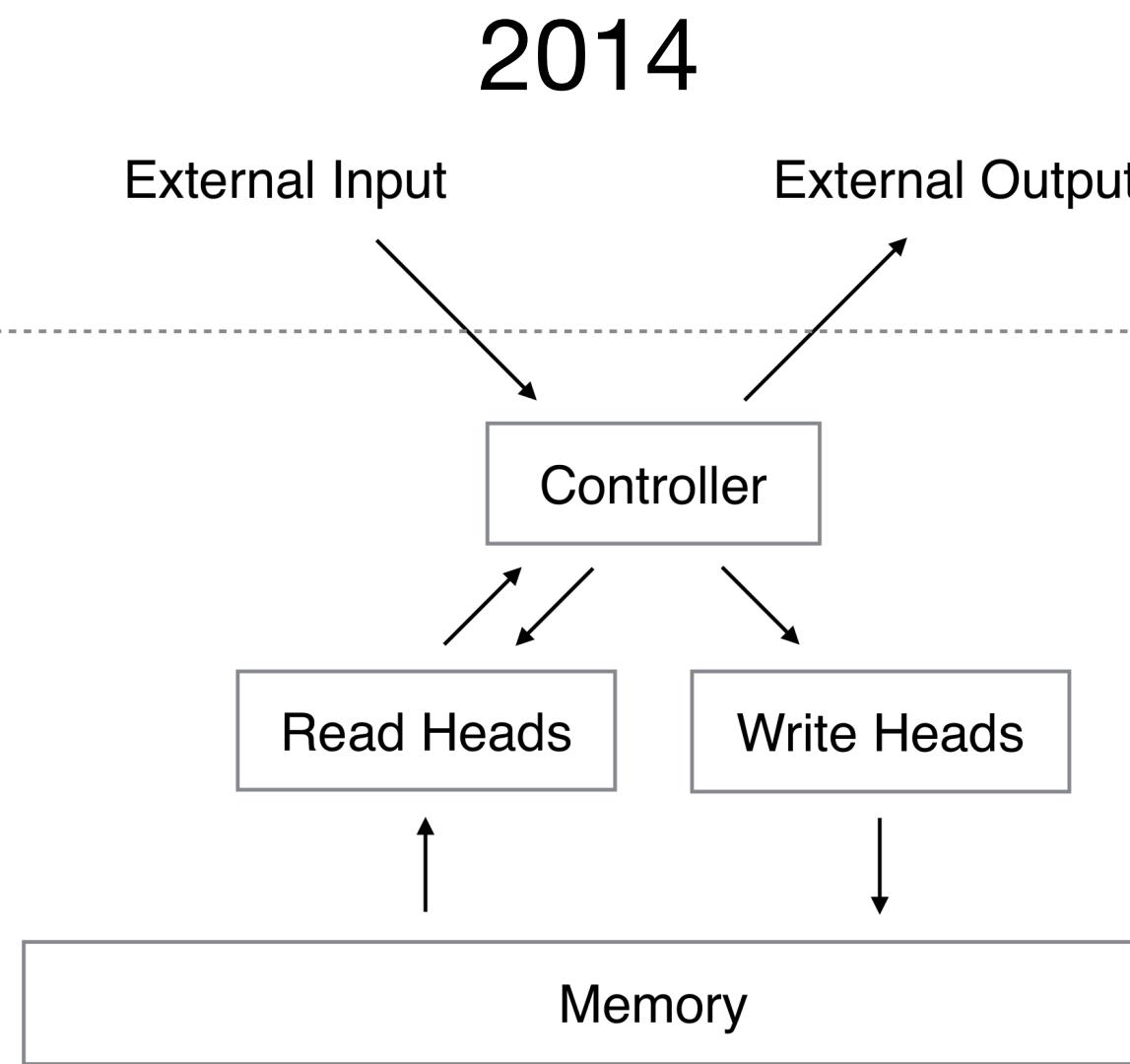
1993 +120K citations

Used in Google translate



Neural Turing Machine

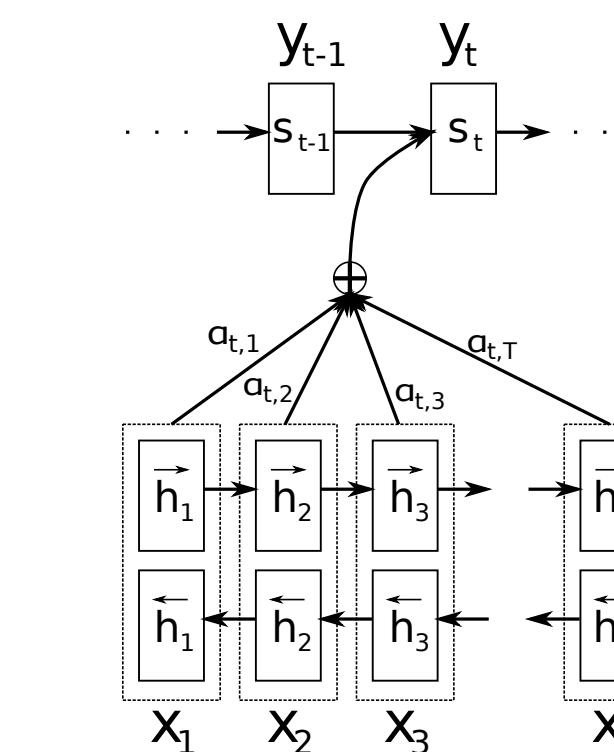
Designing memory for NNs



Neural Translation

Translation

2015
+30K citations



[Bahdanau, Cho, Bengio]

Transformers

next token pred.



UNIVERSITY
of VIRGINIA

Language model and attention mechanism

25

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

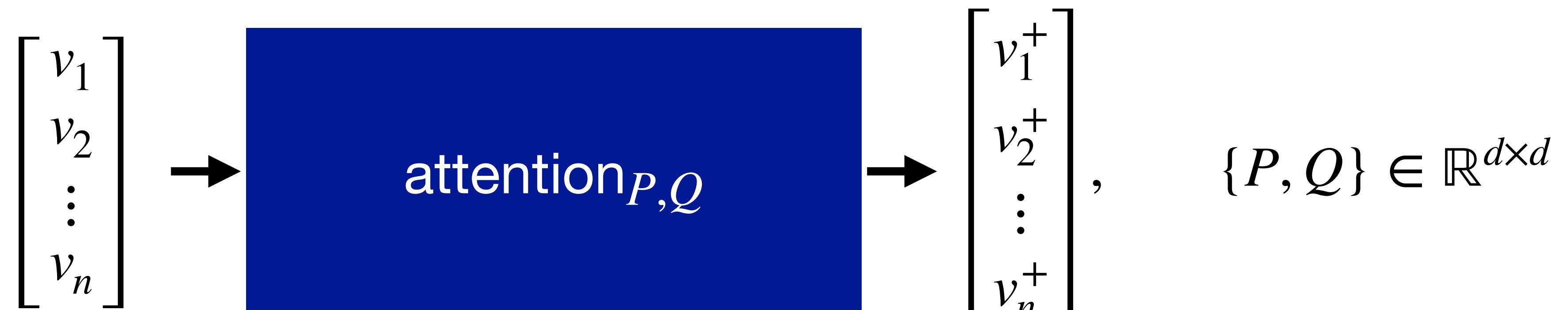


UNIVERSITY
of VIRGINIA

Attention layer

26

- ▶ The $\underbrace{v_1}_{} \underbrace{v_2}_{} \dots v_n \in \mathbb{R}^d$ is chasing a criminal on the run



$$v_i^+ = v_i + P \sum_{j=1}^n \alpha_{ij} v_j, \quad \alpha_{ij} = \frac{e^{v_j^\top Q v_i}}{\sum_k e^{v_k^\top Q v_i}}$$

[softmax]



UNIVERSITY
of VIRGINIA

Example

27

The FBI is chasing a criminal on the run .

$$v_{chasing}^+ = v_{chasing} + P(\alpha_1 v_{the} + \alpha_2 v_{FBI} + \alpha_3 v_{is})$$

$$\alpha_1 = \frac{e^{v_{the}^\top Q v_{chasing}}}{e^{v_{the}^\top Q v_{chasing}} + e^{v_{FBI}^\top Q v_{chasing}} + e^{v_{is}^\top Q v_{chasing}}} \quad [\text{softmax}]$$



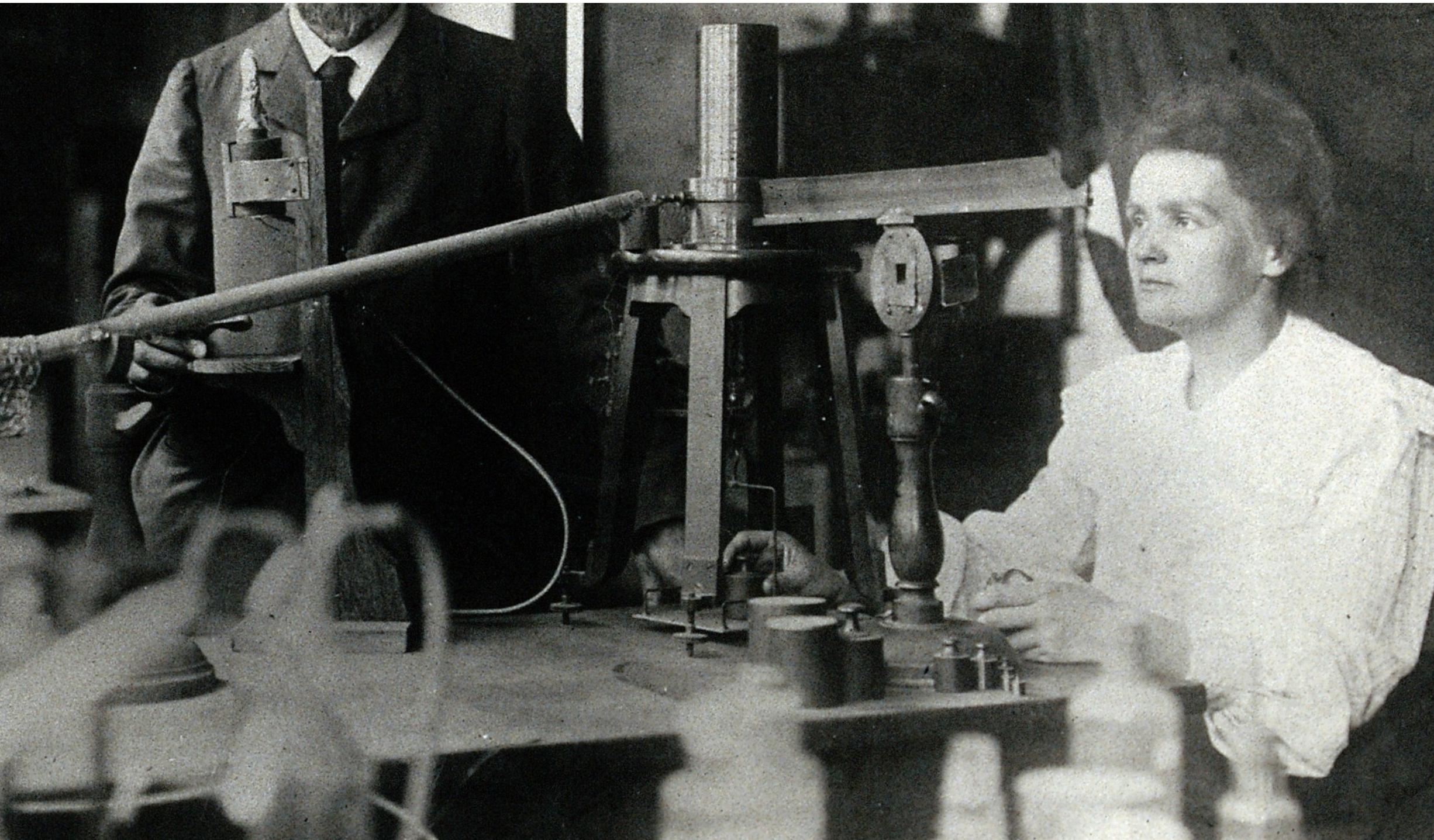
UNIVERSITY
of VIRGINIA

Outline

28

- ▶ Introduction to LLMs and self-supervised learning
- ▶ Attention component in LLMs
- ▶ A group activity





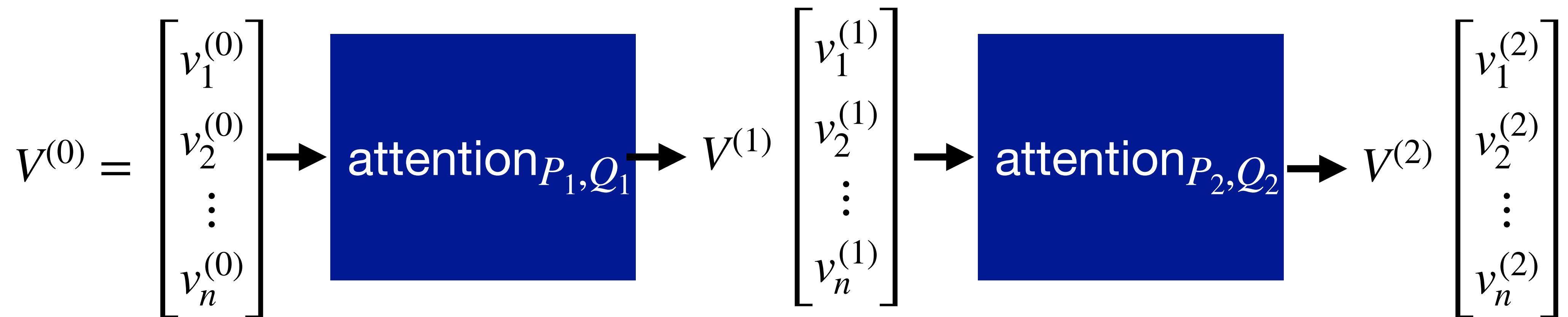
wikipedia: Marie Curie

Experiments

Get ready for hands-on group activity

Goal: word embeddings after attention layers

30



- ▶ Let $P_i = -0.01I_{n \times n}$ (identity matrix), $Q_i = GG^\top$, $G \in R^{d \times d}$, $G_{ij} \sim N(0, 1/d)$
- ▶ Task: Given $V_{ij}^{(0)} \sim N(0, 1/2\sqrt{d})$, plot $V^{(2)}$.



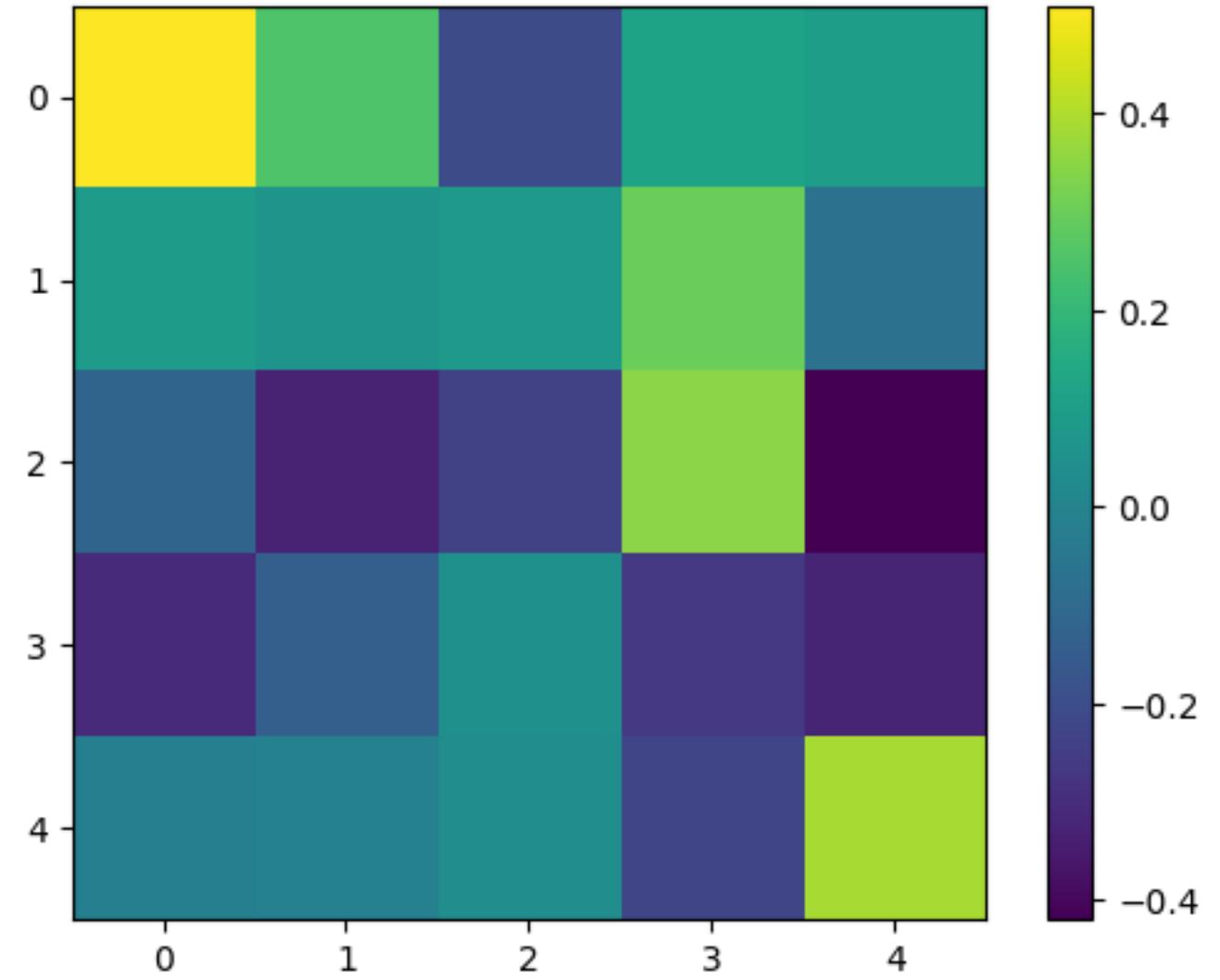
Colab

31

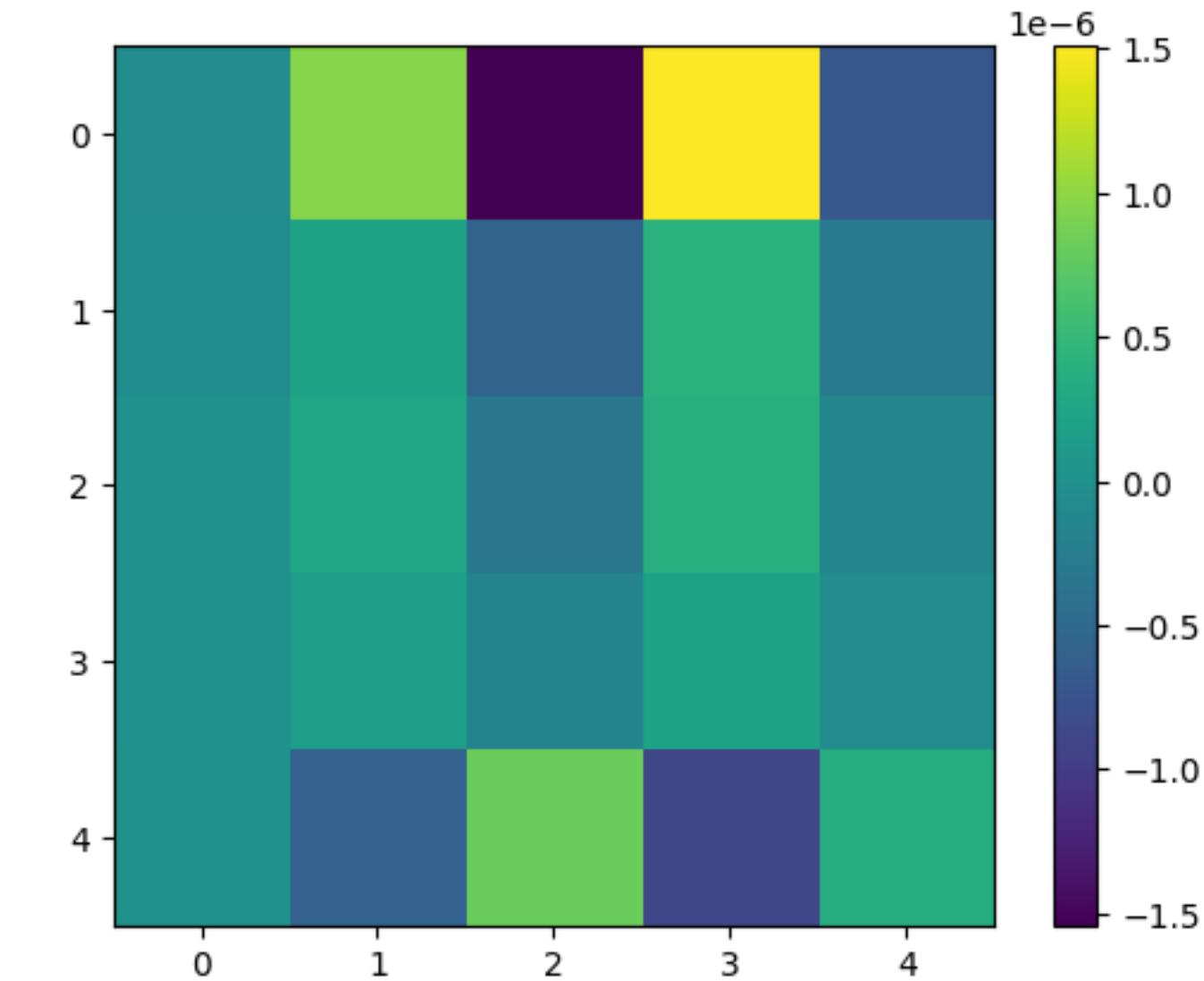
- ▶ <https://shorturl.at/2clWn>
- ▶ https://colab.research.google.com/drive/1w_O9xFkfvsgrPmQXBAZqVhHOLexriEqZ?usp=sharing
- ▶ Details are provided in the colab



An observation



$$V^{(0)} = \begin{bmatrix} v_1^{(0)} \\ v_2^{(0)} \\ \vdots \\ v_n^{(0)} \end{bmatrix}$$



$$V^{(2)} = \begin{bmatrix} v_1^{(2)} \\ v_2^{(2)} \\ \vdots \\ v_n^{(2)} \end{bmatrix}$$

Reference: Geshkovski, Borjan, et al. "The emergence of clusters in self-attention dynamics." *Advances in Neural Information Processing Systems* 36 (2024).



The next lecture

33

- ▶ Why word embeddings become low-rank? (Explaining the observation)
- ▶ Why attention layers are linked to gradient descent?



Thank you very much!

34



UNIVERSITY
of VIRGINIA

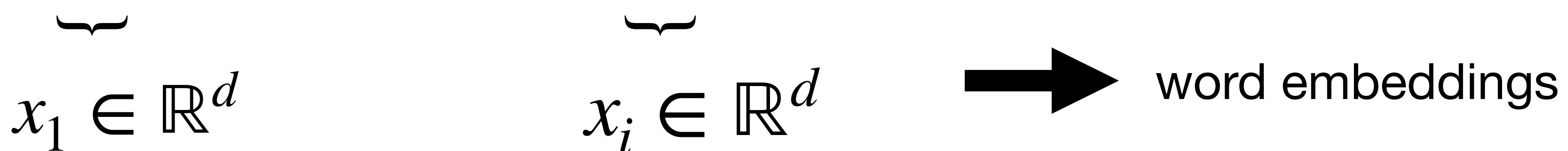
Next word predication and supervised learning

35

- ▶ This text is used to teach a model

- ▶ Questions:

- What are response variables y (outputs)?
 - All words or characters. Would you prefer characters or words?
- What are input features x ?
 - This text is used to teach a model



The power of self-supervised learning

36

Sentiment Analysis

Prompt. “I do not like this product”
is positive?

Output. No

Translation

Prompt. Translate “this is a
book” to Dutch

Output: Dit is een boek



The skill of "in-context learning"

37

- ▶ **In-context learning:** Language models can learn from in-context data

Example: swapping characters

Prompt: ; ; ?

GPT3:

- ▶ Reference: Language Models are Few-Shot Learners, [Tom B. Brown et al.](#)