

MACHINE LEARNING FOR SIGNAL PROCESSING

24-3-2025

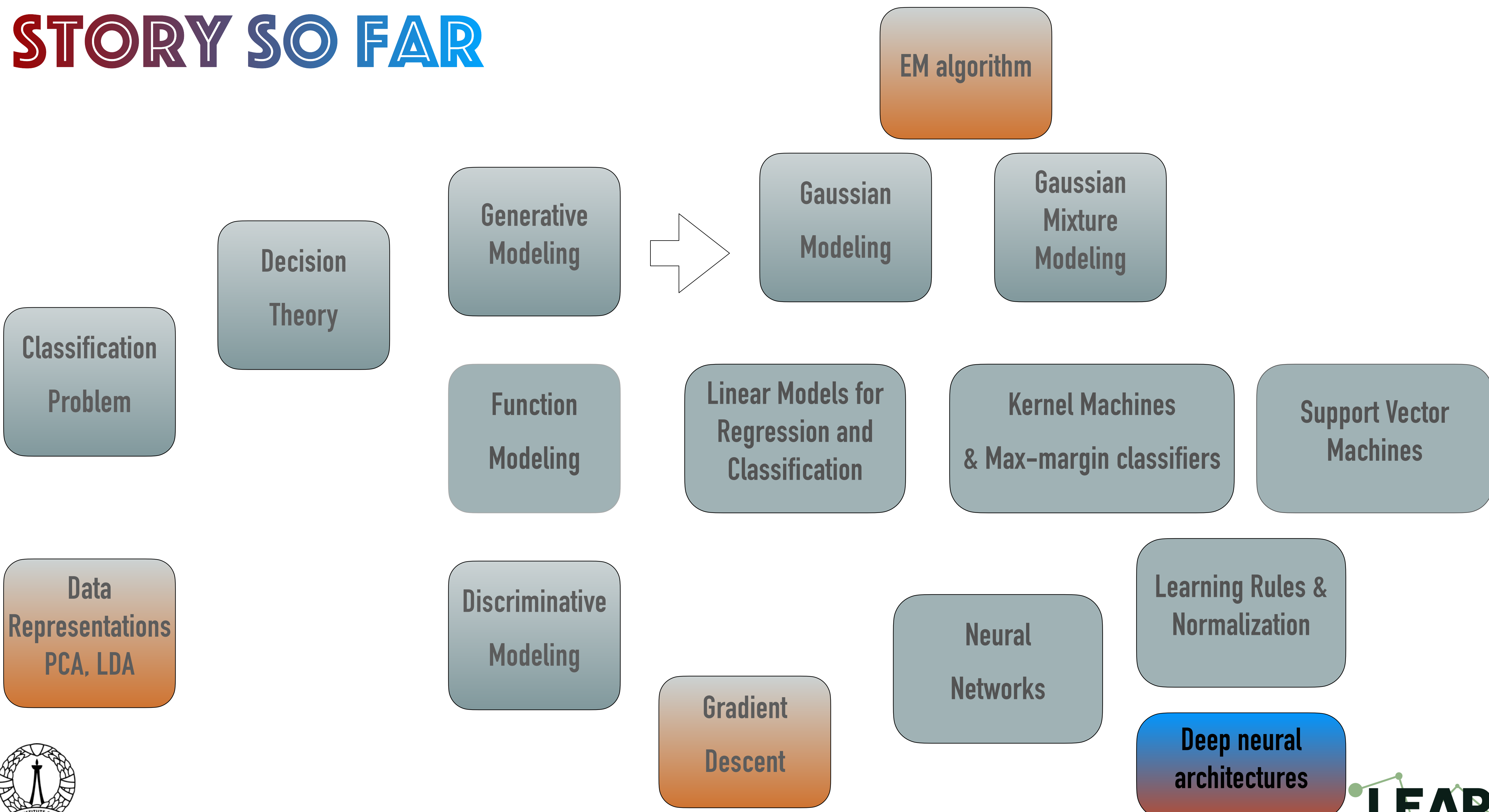
Sriram Ganapathy
LEAP lab, Electrical Engineering, Indian Institute of Science,
sriramg@iisc.ac.in

.....
Viveka Salinamakki, Varada R.
LEAP lab, Electrical Engineering, Indian Institute of Science
.....

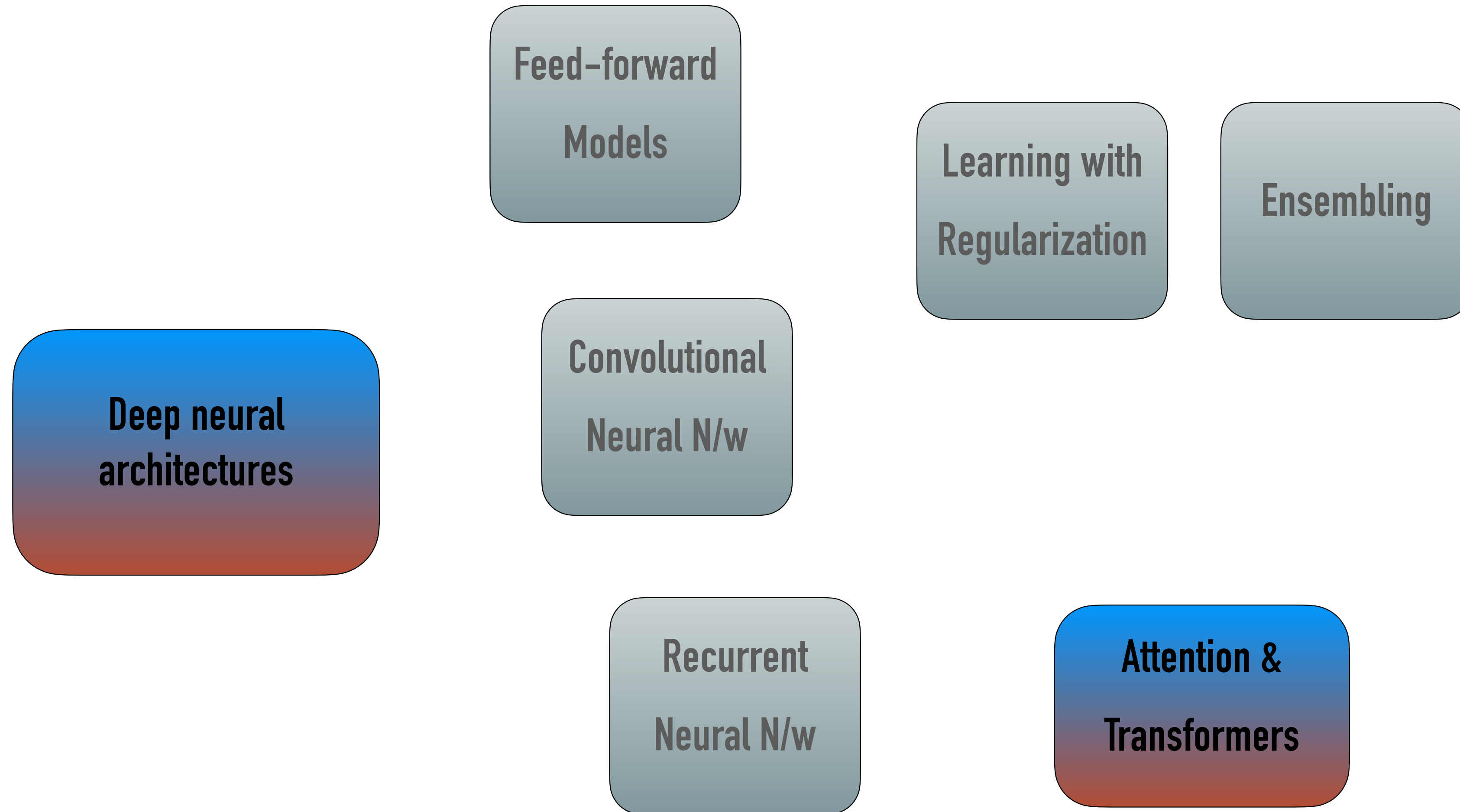
<http://leap.ee.iisc.ac.in/sriram/teaching/MLSP25/>



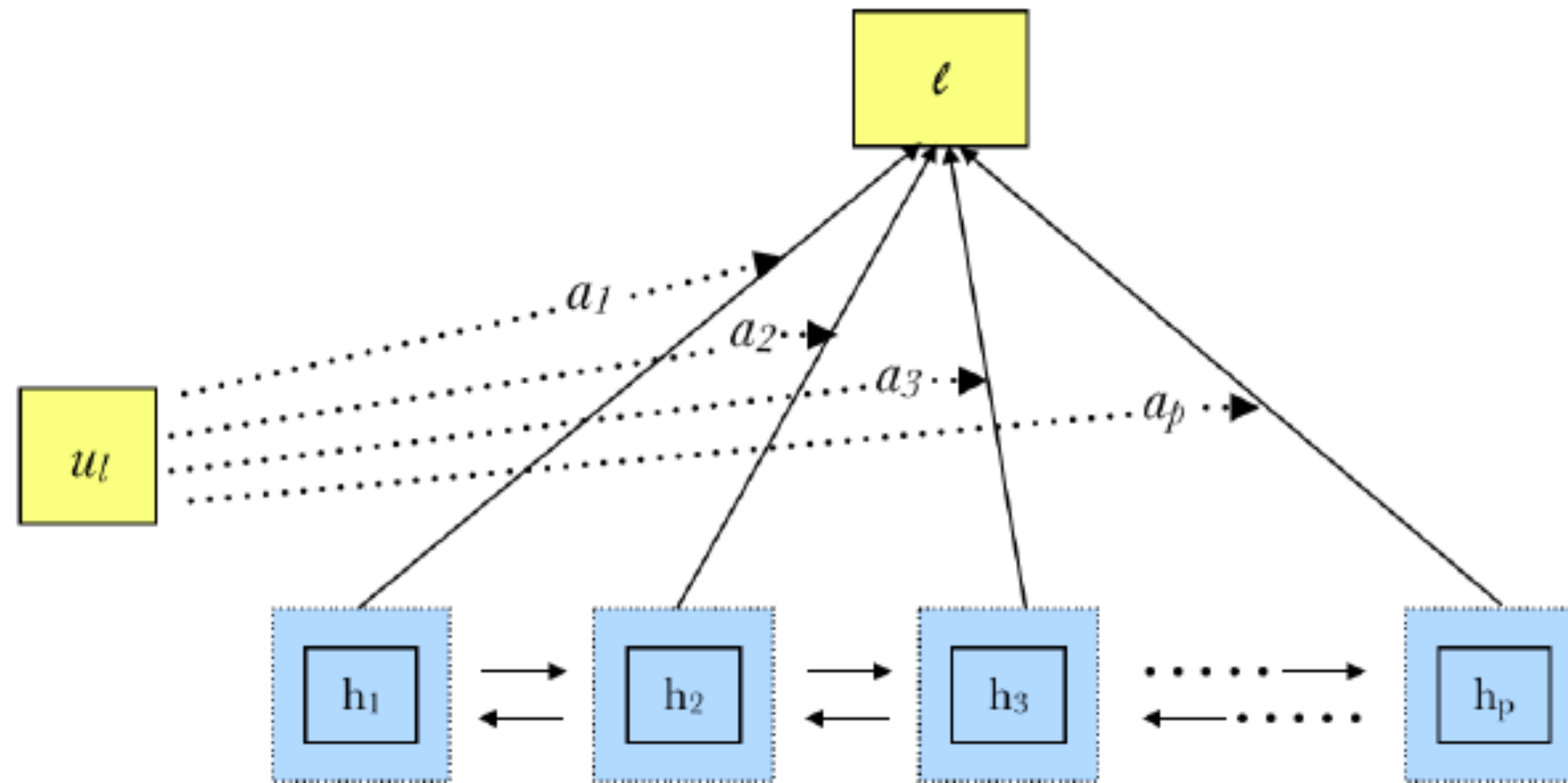
STORY SO FAR



STORY SO FAR



Attention in LSTM Networks



$$\mathbf{u}_t = \tanh(\mathbf{W}_l \mathbf{h}_t + \mathbf{b}_l)$$

$$a_t = \frac{\exp(\mathbf{u}_t^T \mathbf{u}_l)}{\sum_t \exp(\mathbf{u}_t^T \mathbf{u}_l)}$$

$$\mathbf{l} = \sum_t a_t \mathbf{h}_t$$

- ❖ Attention allows a mechanism to add relevance
- ❖ Certain regions of the audio have more importance than the rest for the task at hand.

Encoder-Decoder Attention

Encoder
hidden
state

Je

hidden
state #1

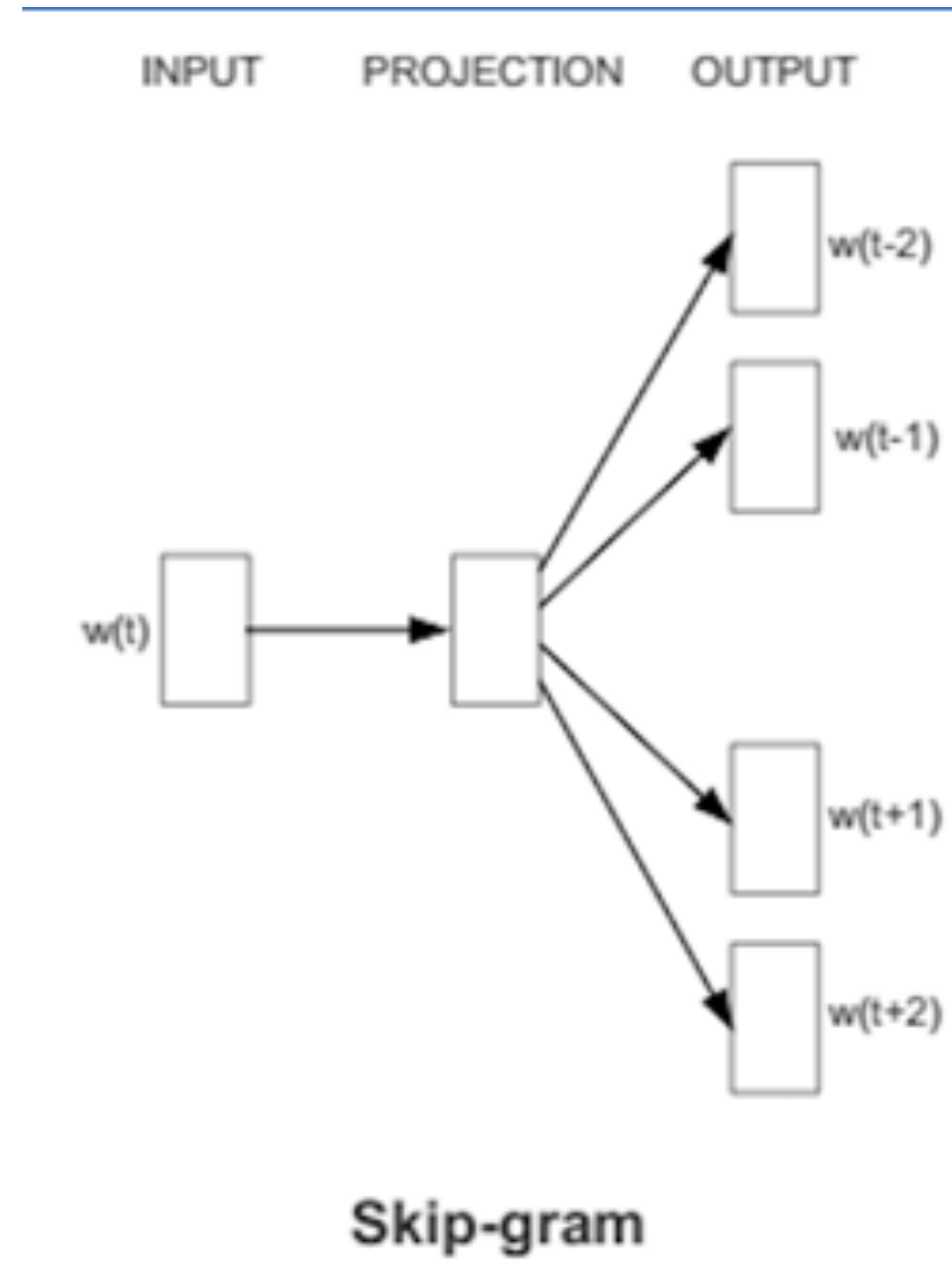
suis

hidden
state #2

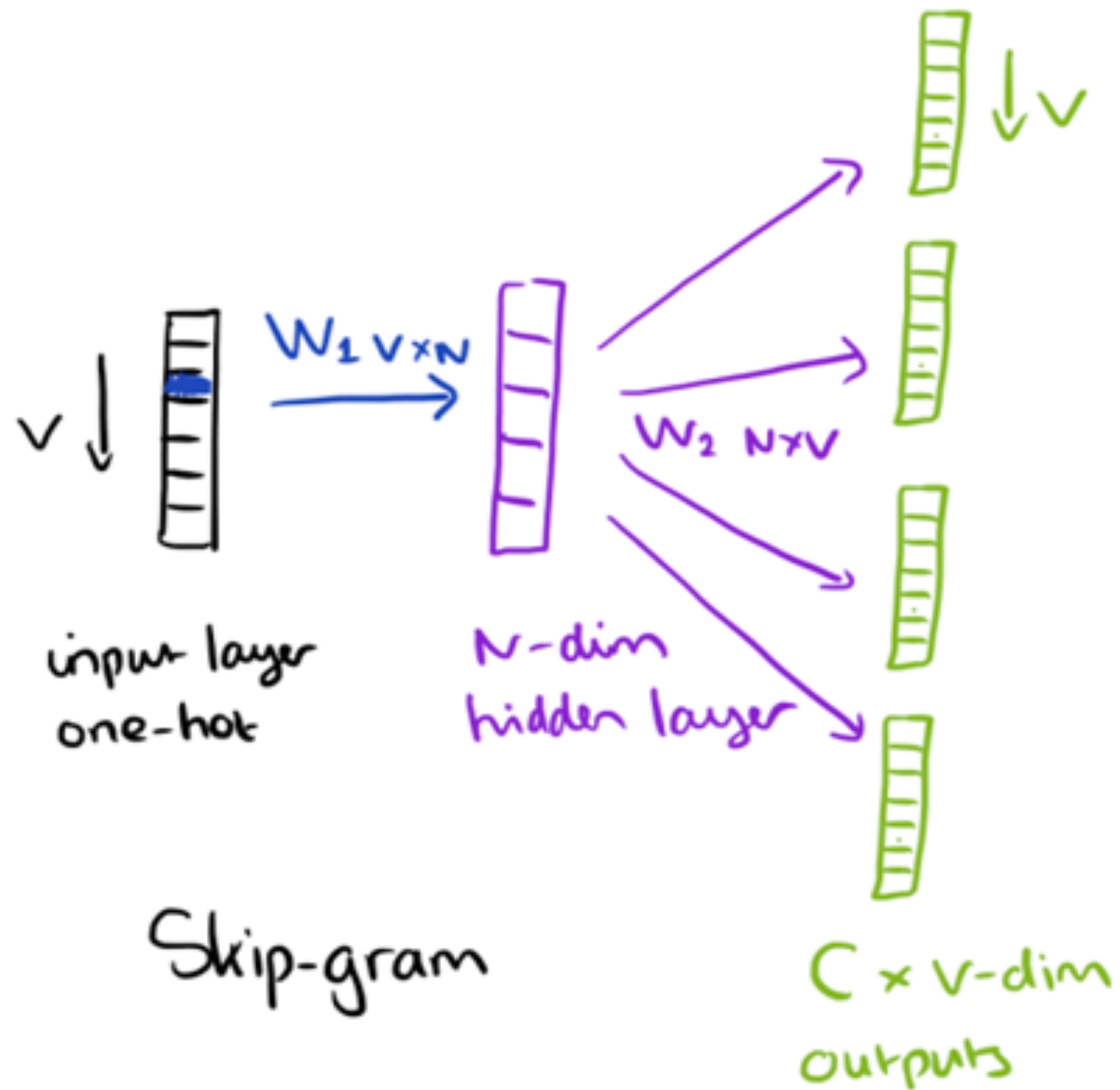
étudiant

hidden
state #3

word2vec models as text representations



word2vec representations

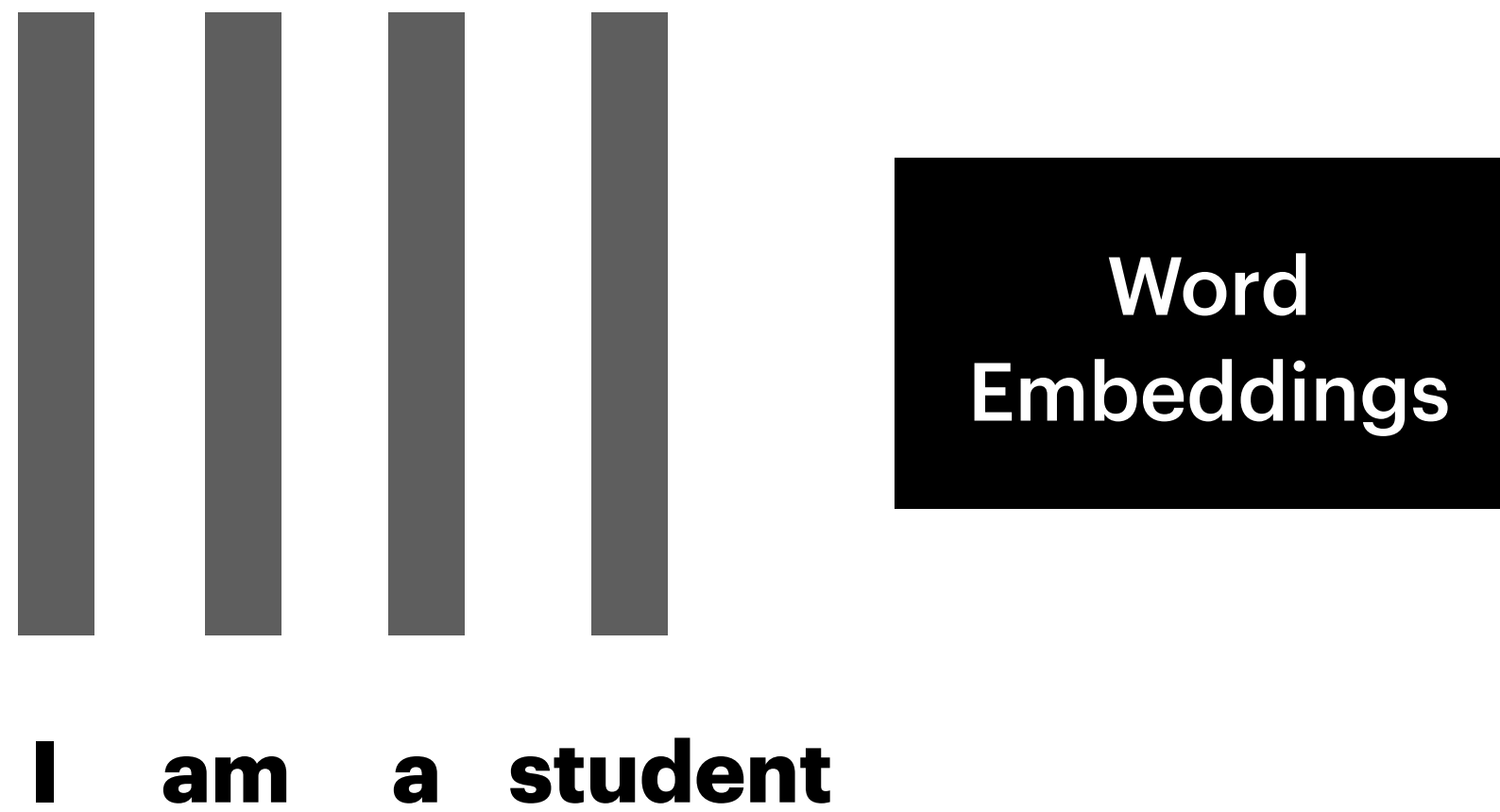


$$\begin{matrix} \text{input} \\ 1 \times v \end{matrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{matrix} W_1 \\ v \times N \end{matrix} \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & j & k & l \end{bmatrix} = \begin{matrix} \text{hidden} \\ 1 \times N \end{matrix} \begin{bmatrix} e & f & g & h \end{bmatrix}$$

W_2

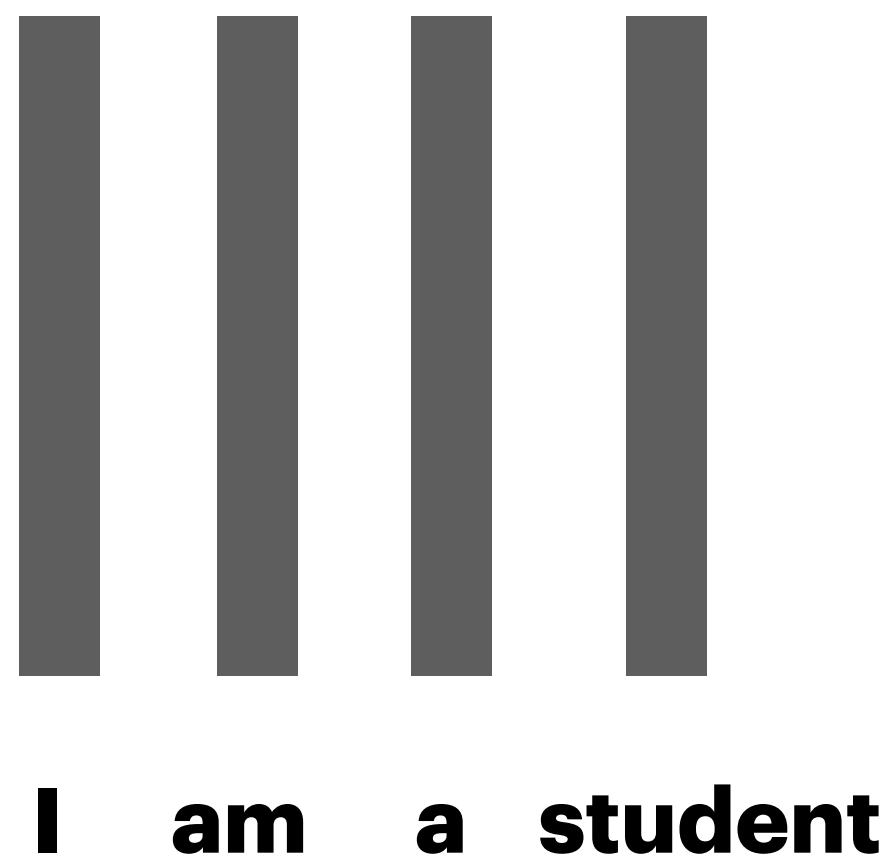
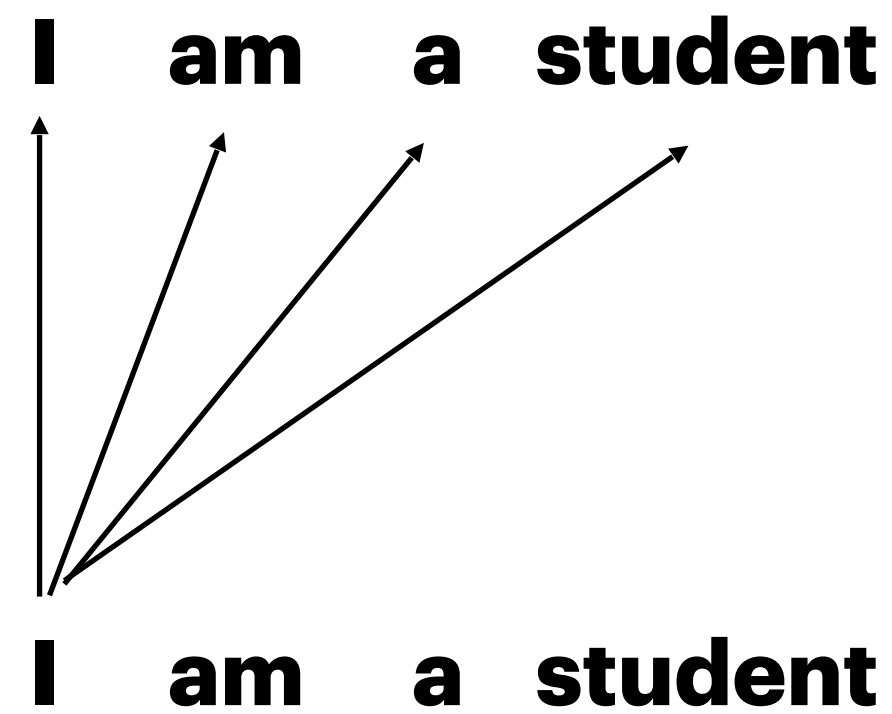
Transformers

- Embedding context in sequence inputs



Transformers - self-attention

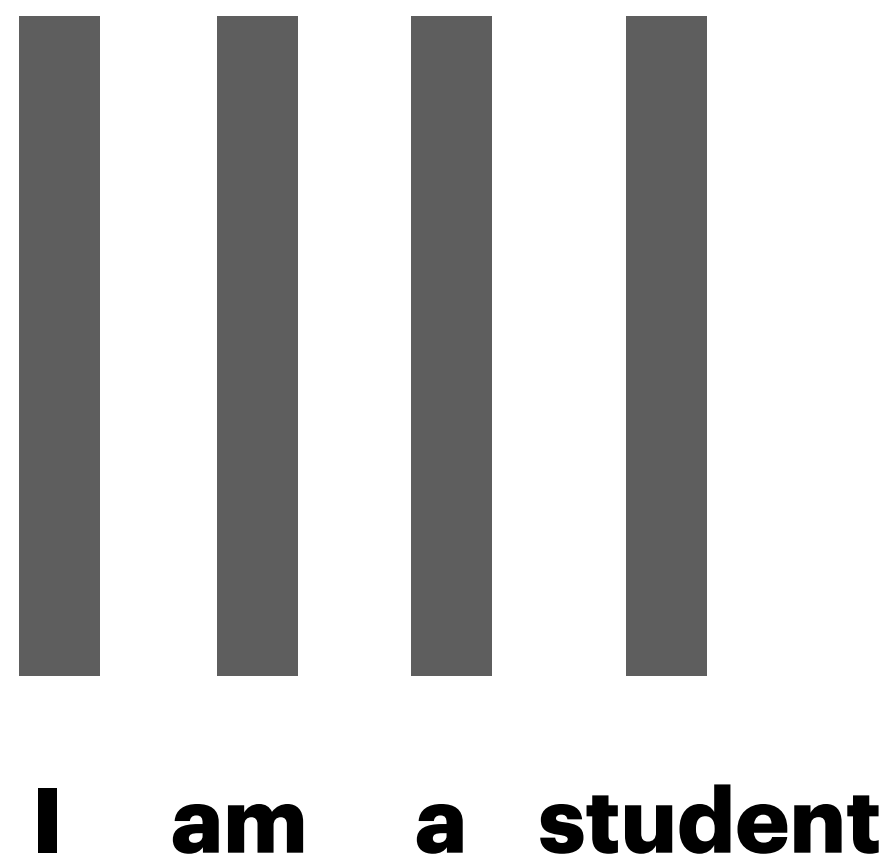
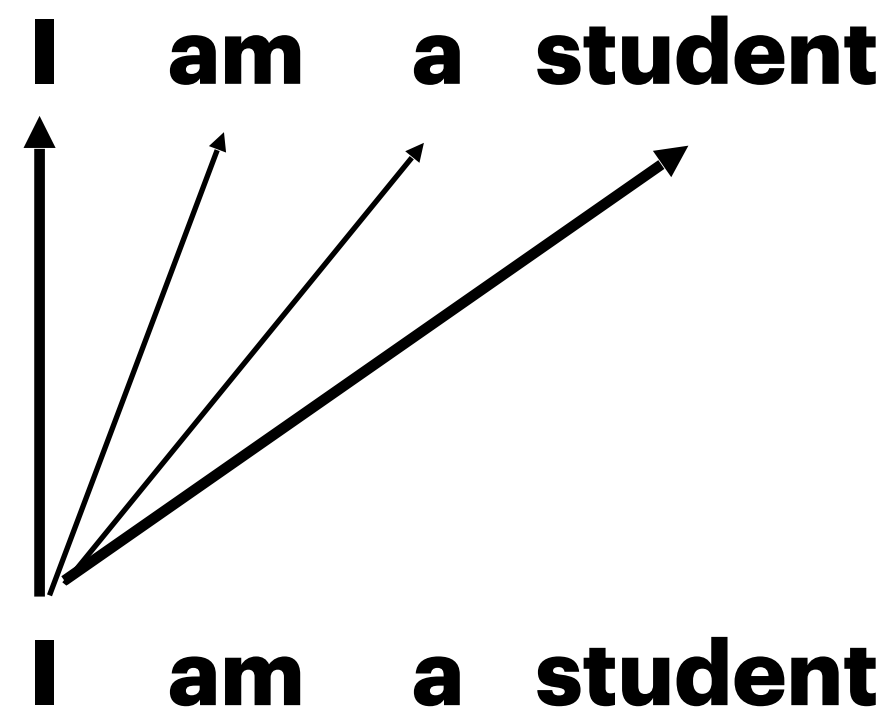
- Embedding context in sequence inputs
 - * Let us take an example



Word
Embeddings

Transformers - self-attention

- Embedding context in sequence inputs
 - * Let us take an example
 - * Using word embeddings as the input representation



$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}; \mathbf{x}_t \in \mathcal{R}^D$$

Word
Embeddings

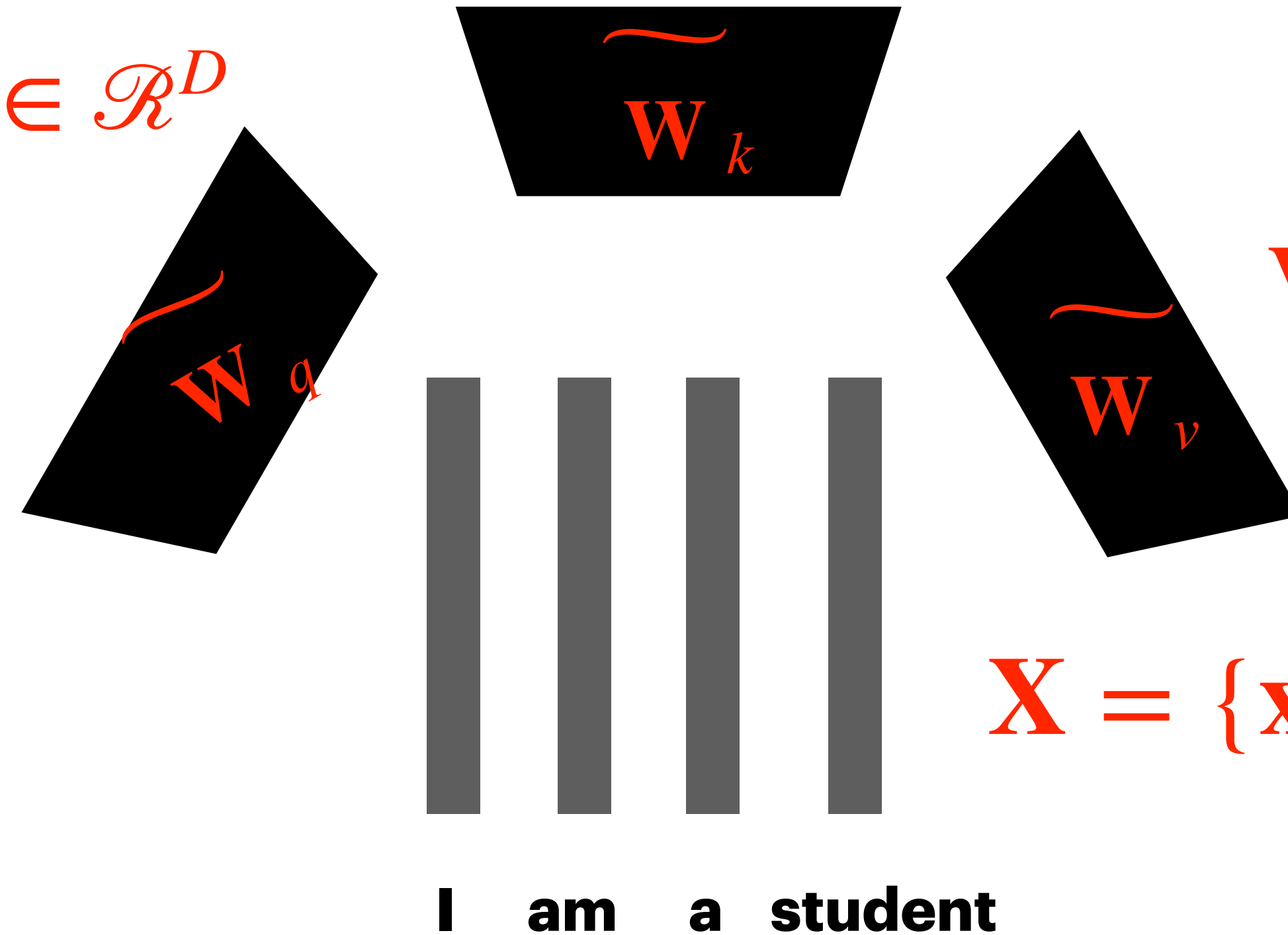
Transformers - self-attention

$$\mathbf{K} = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_T\}; \mathbf{k}_t \in \mathcal{R}^D$$

$$\mathbf{Q} = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T\}; \mathbf{q}_t \in \mathcal{R}^D$$

$$\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T\}; \mathbf{v}_t \in \mathcal{R}^D$$

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}; \mathbf{x}_t \in \mathcal{R}^D$$



Word
Embeddings

Transformers - self-attention

I am a student

$$\mathbf{K} = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_T\}; \mathbf{k}_t \in \mathcal{R}^D$$

$$\mathbf{Q} = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T\}; \mathbf{q}_t \in \mathcal{R}^D$$

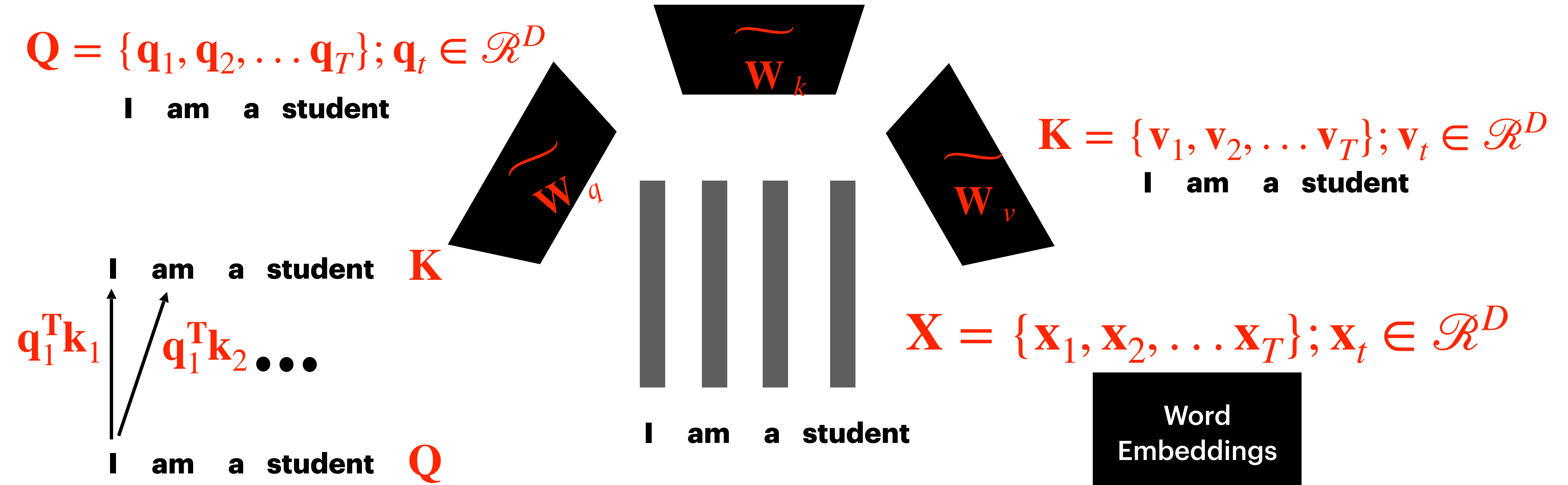
I am a student

$$\mathbf{K} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T\}; \mathbf{v}_t \in \mathcal{R}^D$$

I am a student

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}; \mathbf{x}_t \in \mathcal{R}^D$$

Word
Embeddings



Transformers - self-attention

I am a student

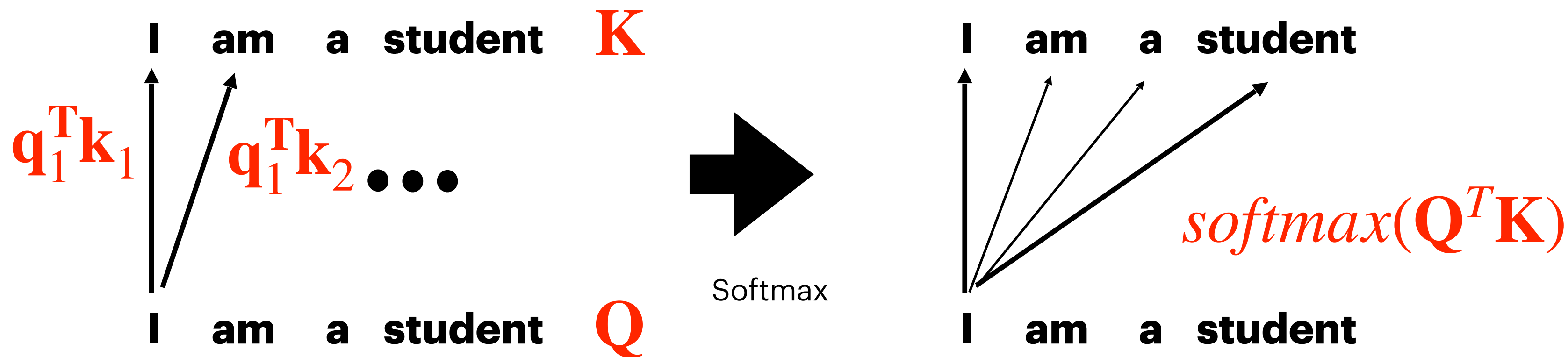
$$\mathbf{K} = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_T\}; \mathbf{k}_t \in \mathcal{R}^D$$

I am a student

$$\mathbf{Q} = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T\}; \mathbf{q}_t \in \mathcal{R}^D$$

$$\mathbf{K} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T\}; \mathbf{v}_t \in \mathcal{R}^D$$

I am a student



Transformers - self-attention

I am a student

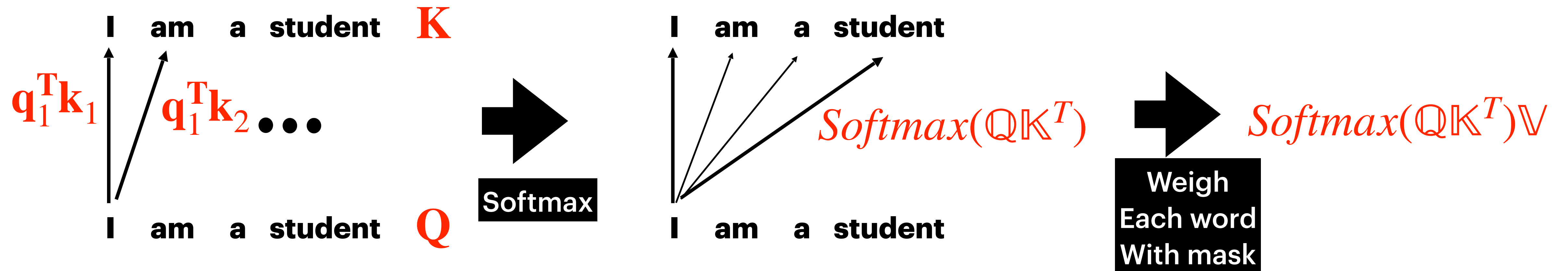
$$\mathbf{K} = \{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_T\}; \mathbf{k}_t \in \mathcal{R}^D$$

I am a student

$$\mathbf{Q} = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T\}; \mathbf{q}_t \in \mathcal{R}^D$$

$$\mathbf{K} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T\}; \mathbf{v}_t \in \mathcal{R}^D$$

I am a student



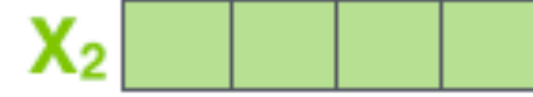
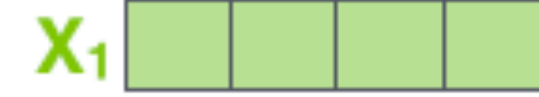
Transformers - self-attention

Input

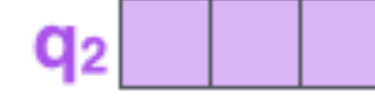
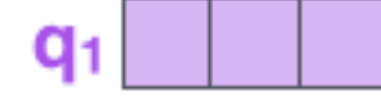
Thinking

Machines

Embedding

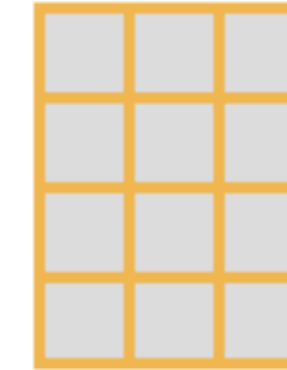
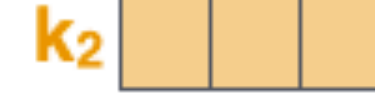
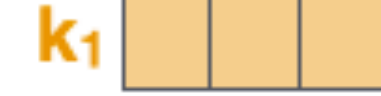


Queries



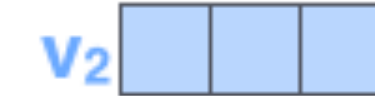
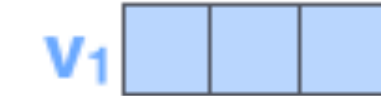
W^Q

Keys



W^K

Values



W^V

Transformers - self-attention

Input

Embedding

Queries

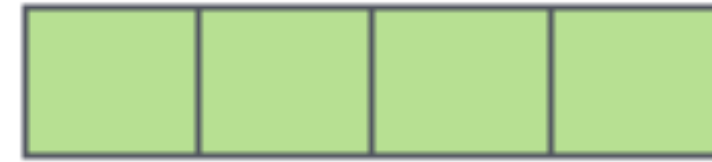
Keys

Values

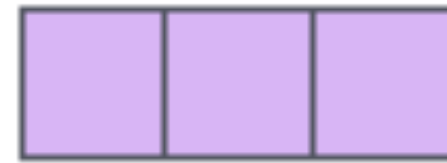
Score

Thinking

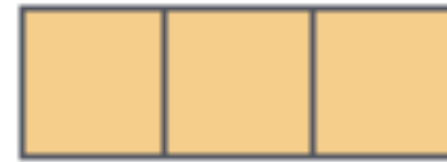
x_1



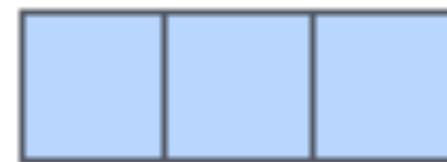
q_1



k_1



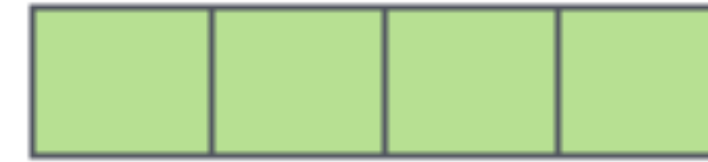
v_1



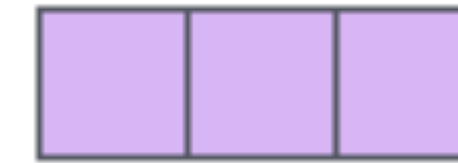
$$q_1 \cdot k_1 = 112$$

Machines

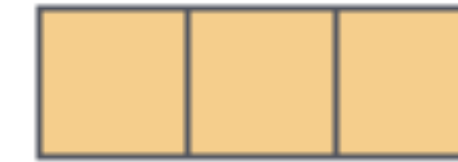
x_2



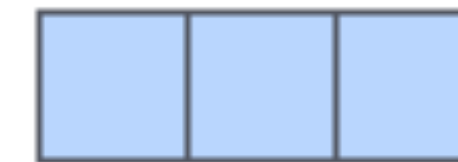
q_2



k_2

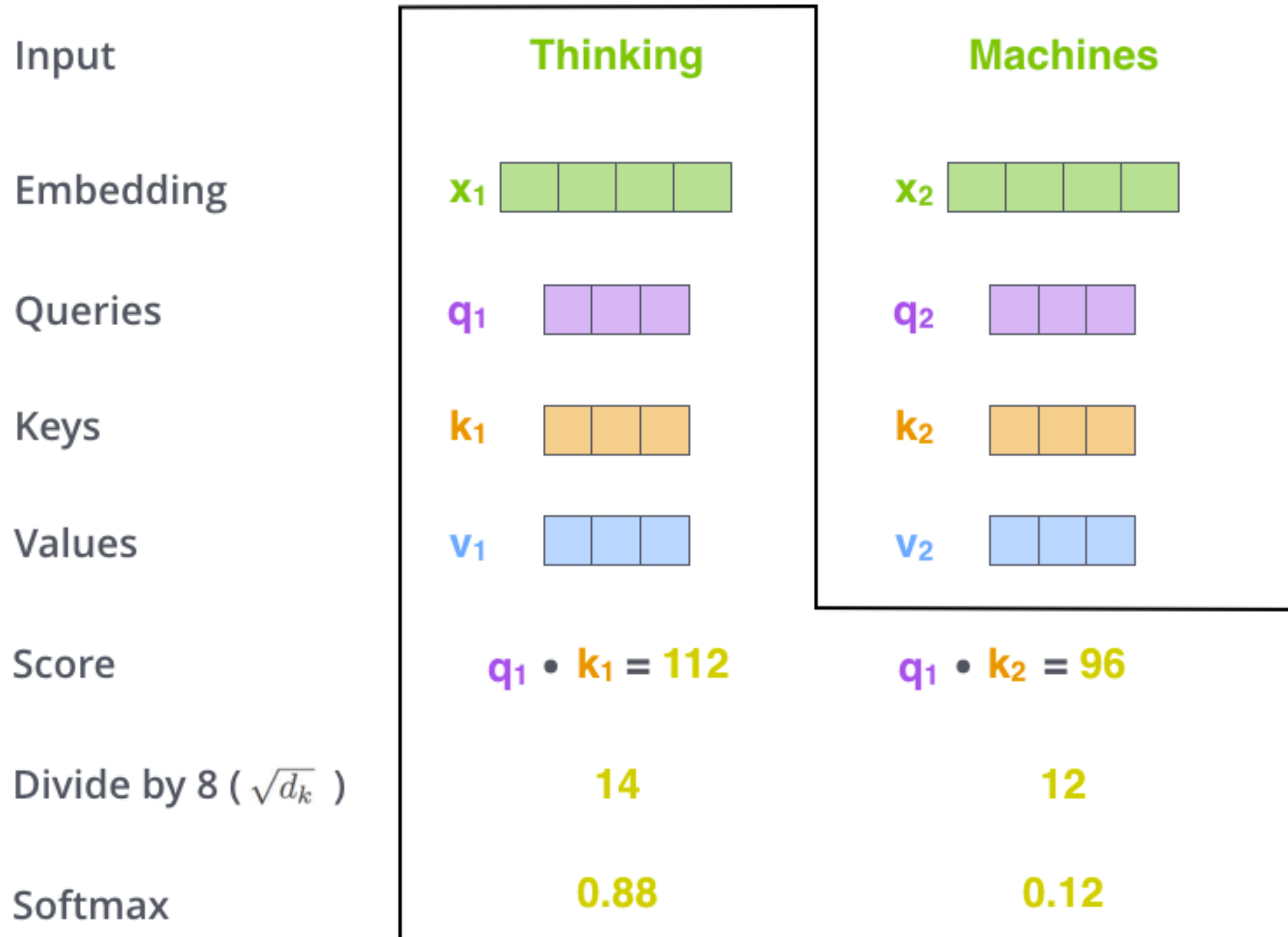


v_2

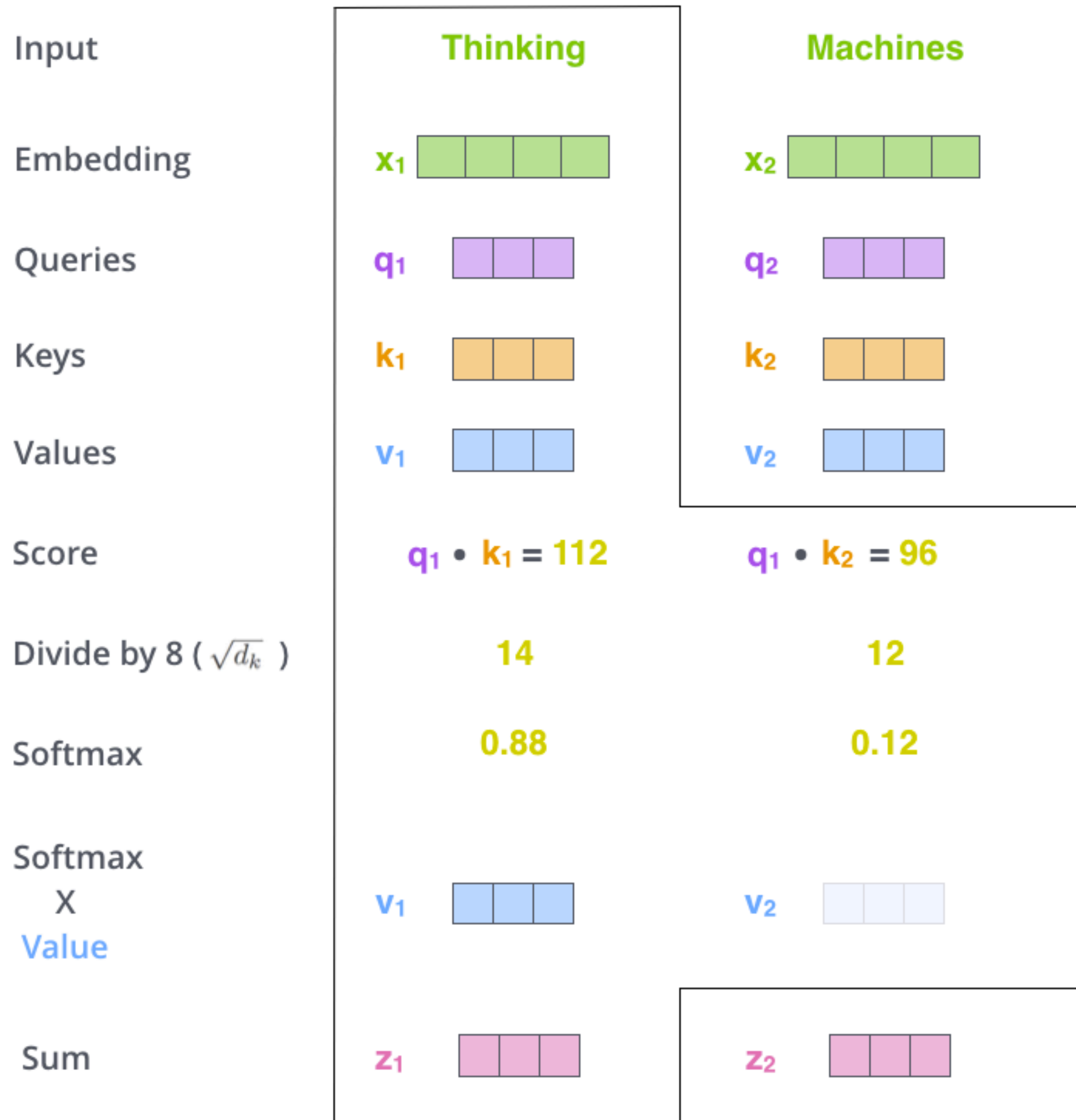


$$q_1 \cdot k_2 = 96$$

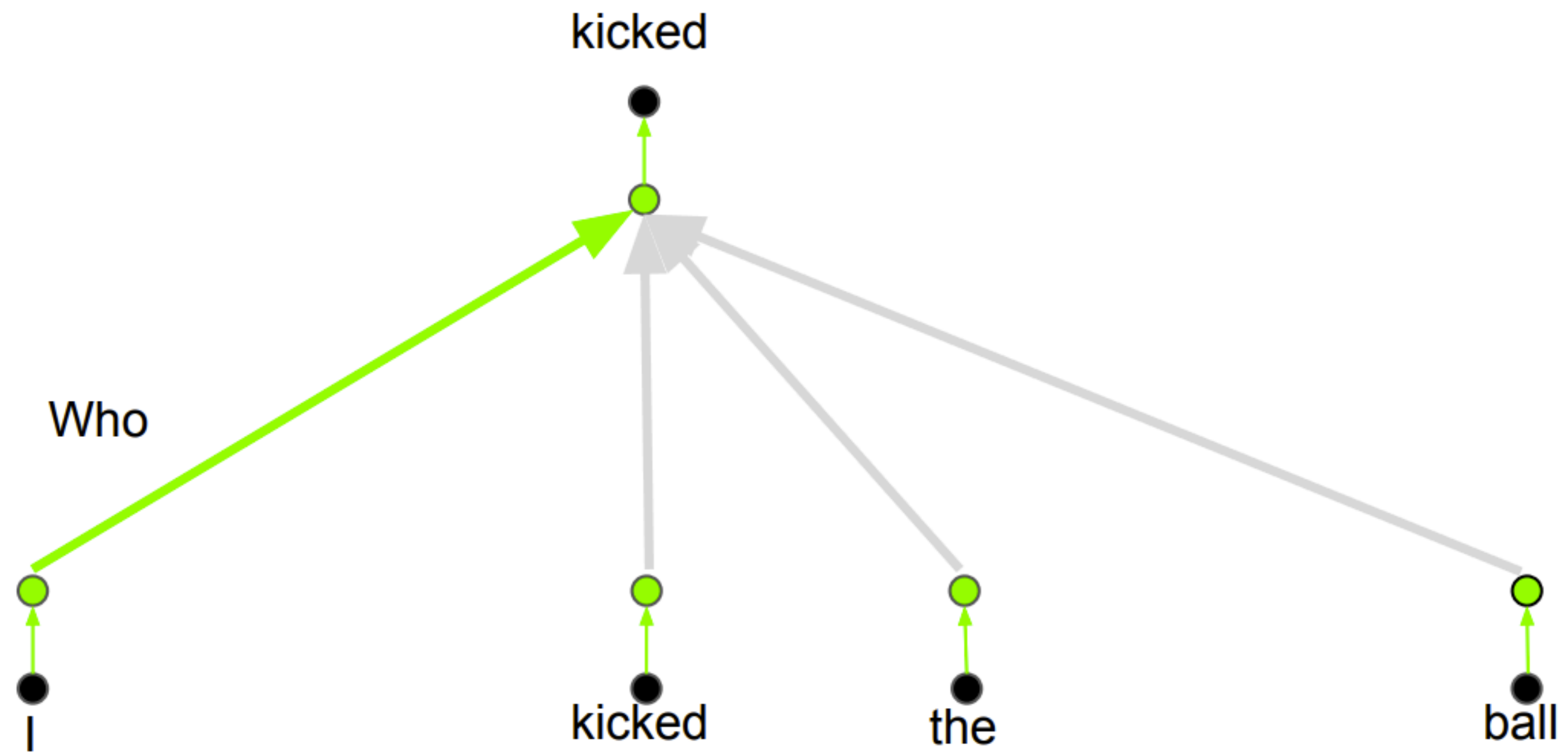
Transformers - self-attention



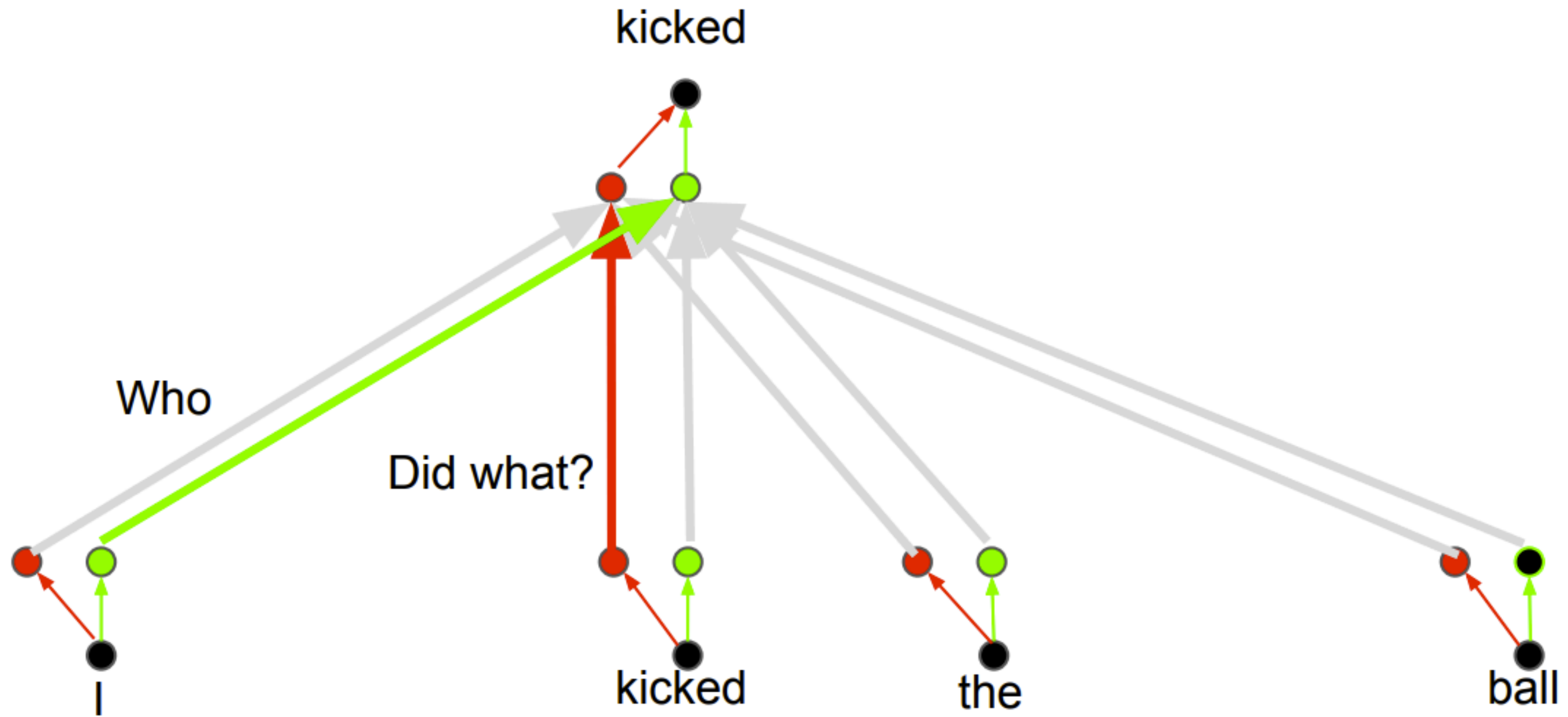
Transformers - self-attention



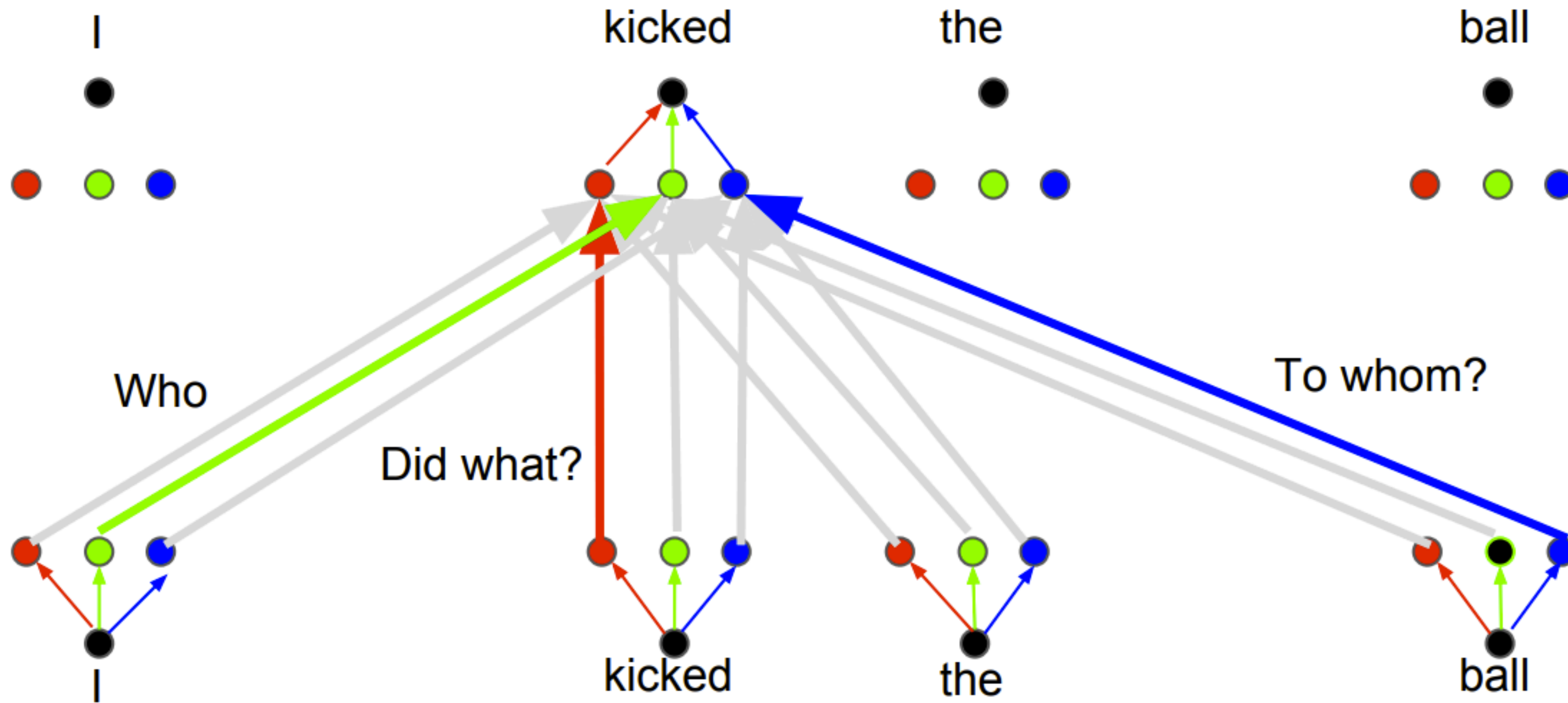
Multi-head



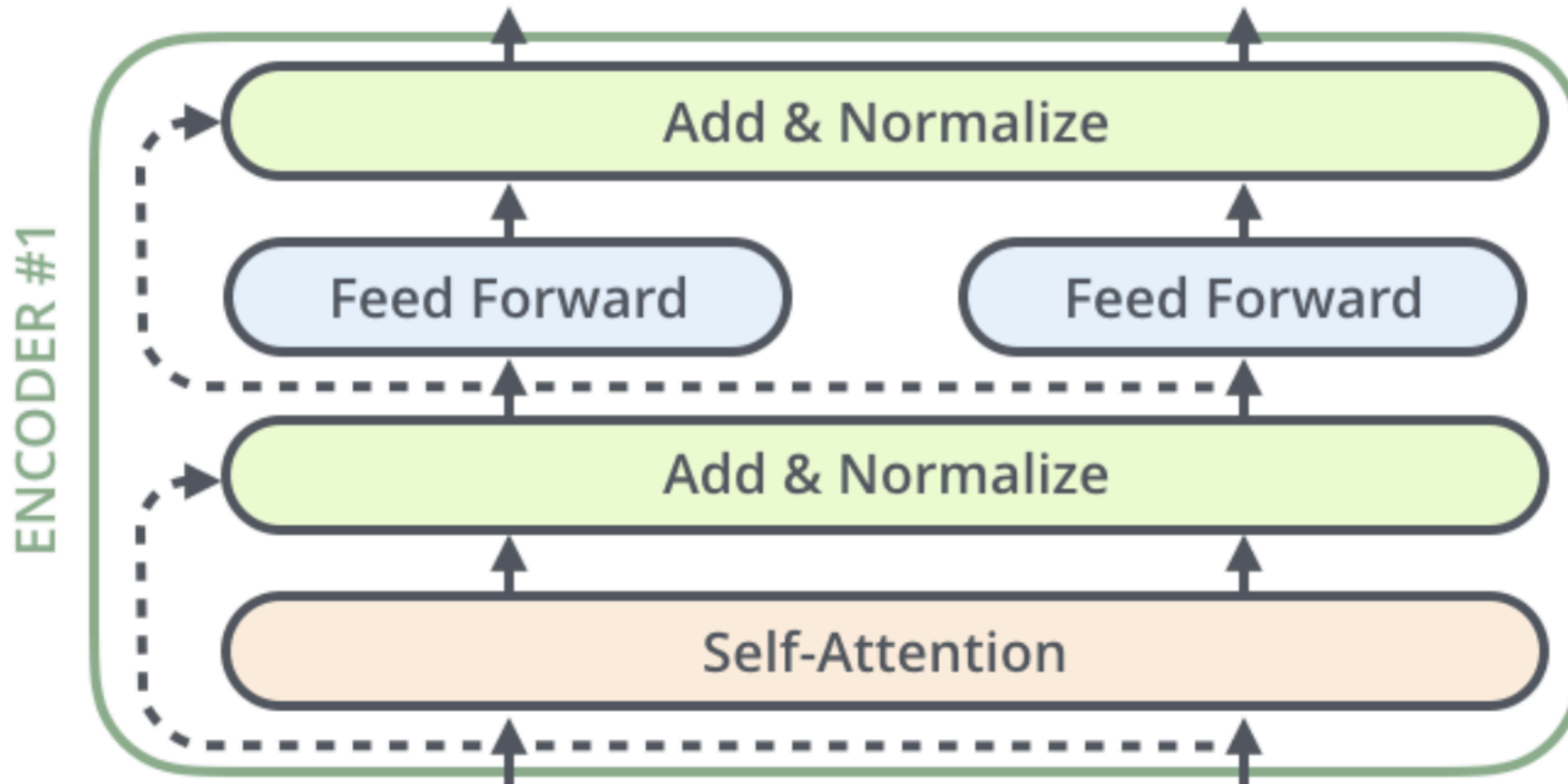
Multi-head



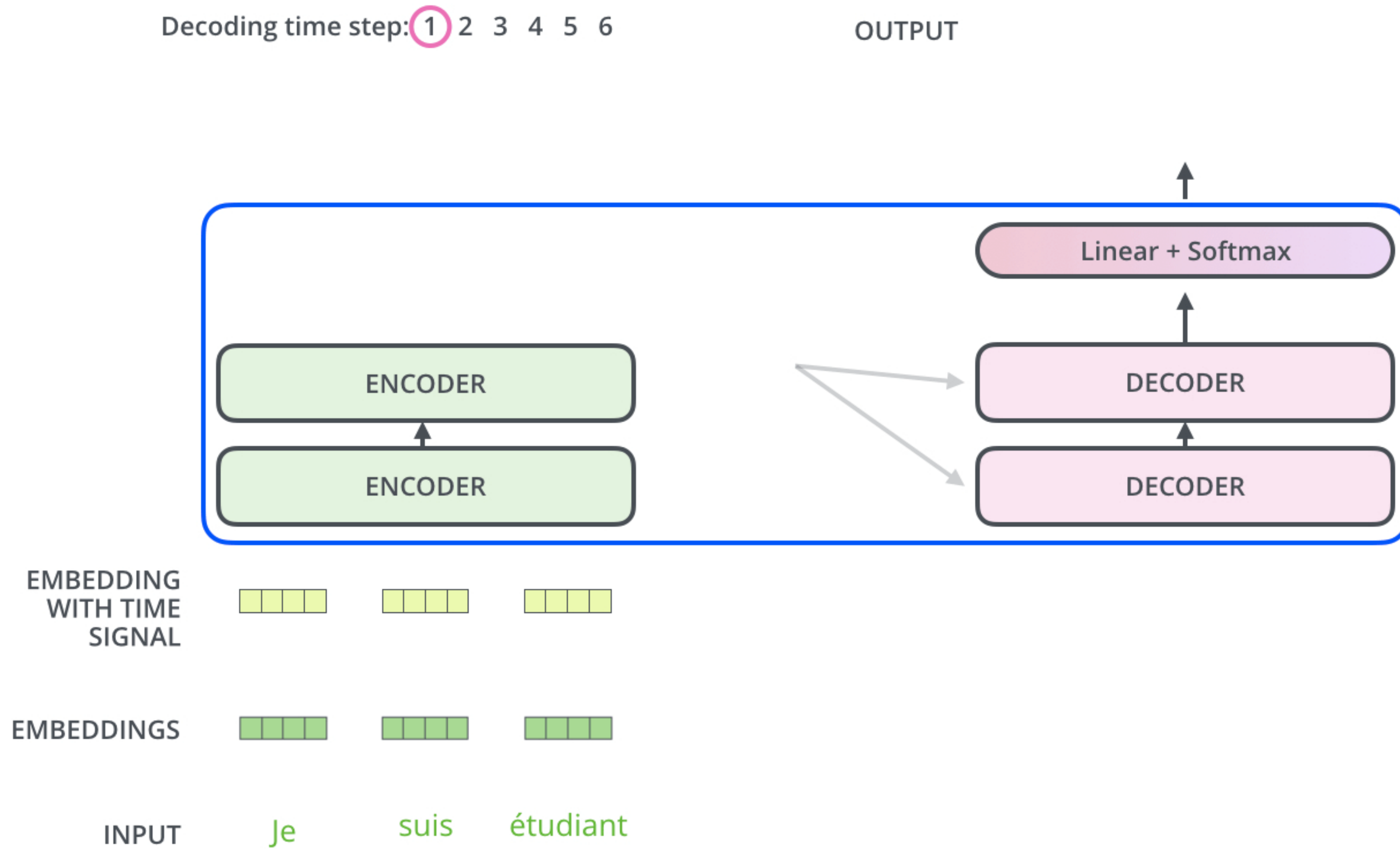
Multi-head



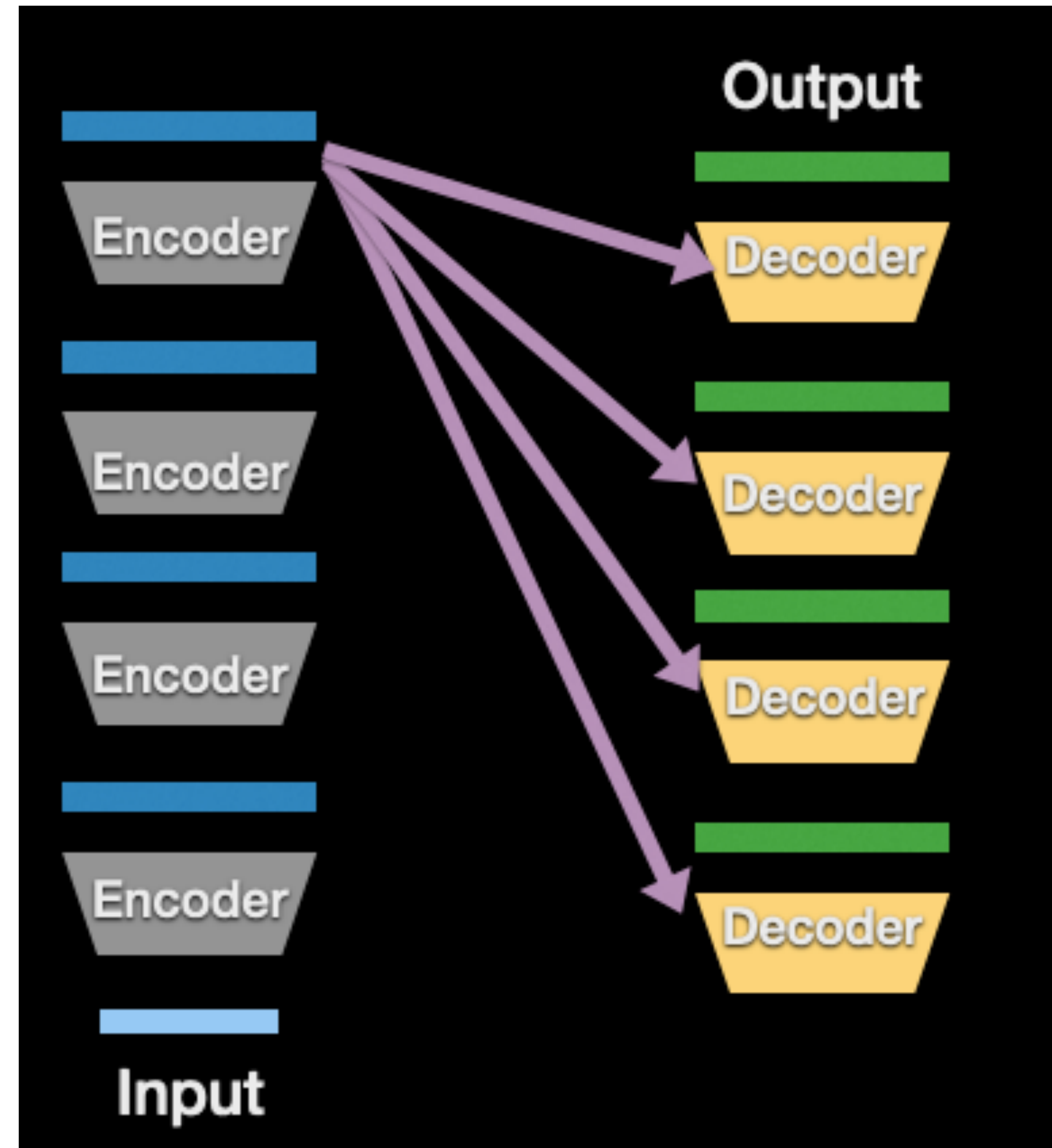
Transformer encoder



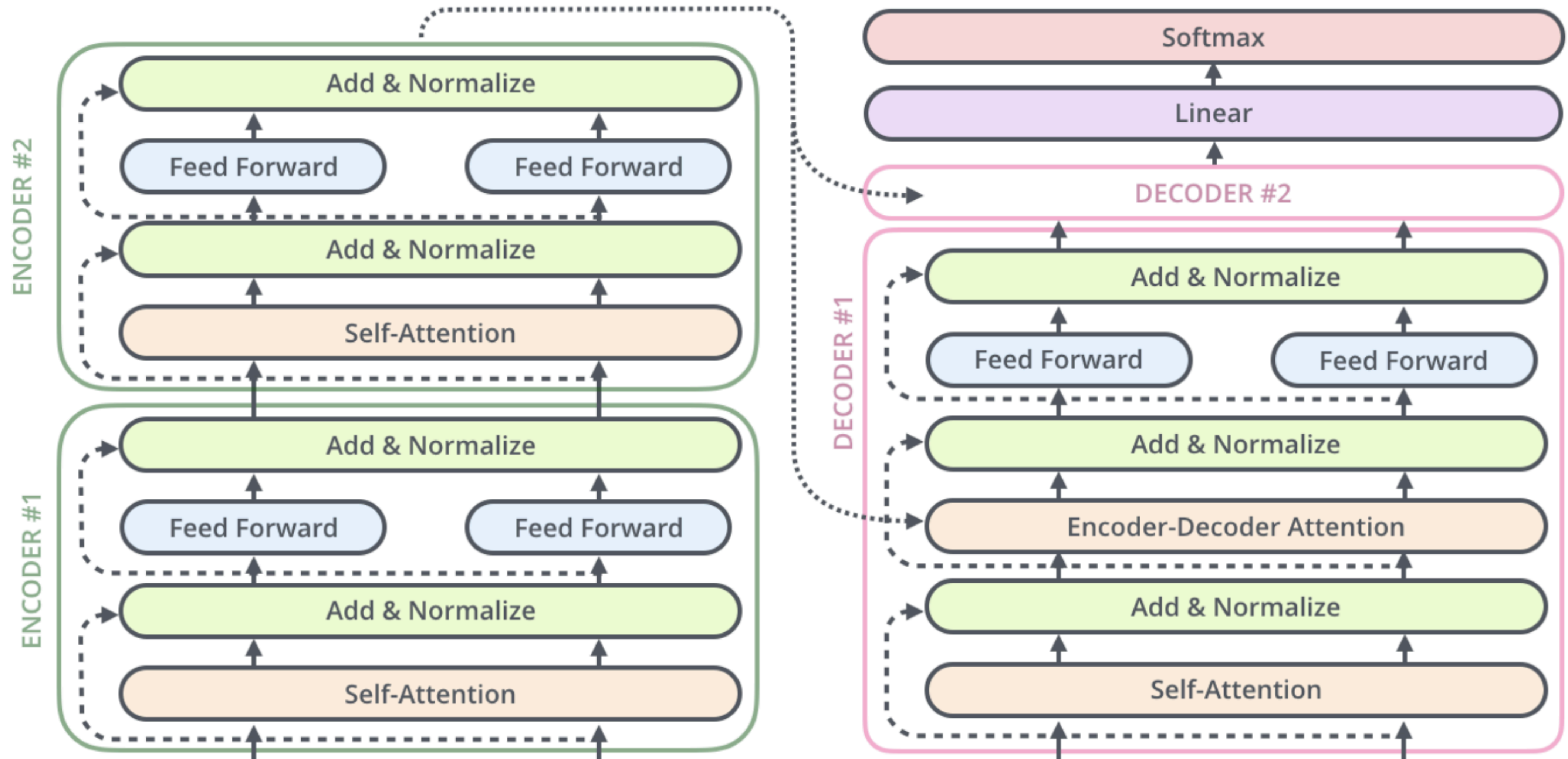
Transformer



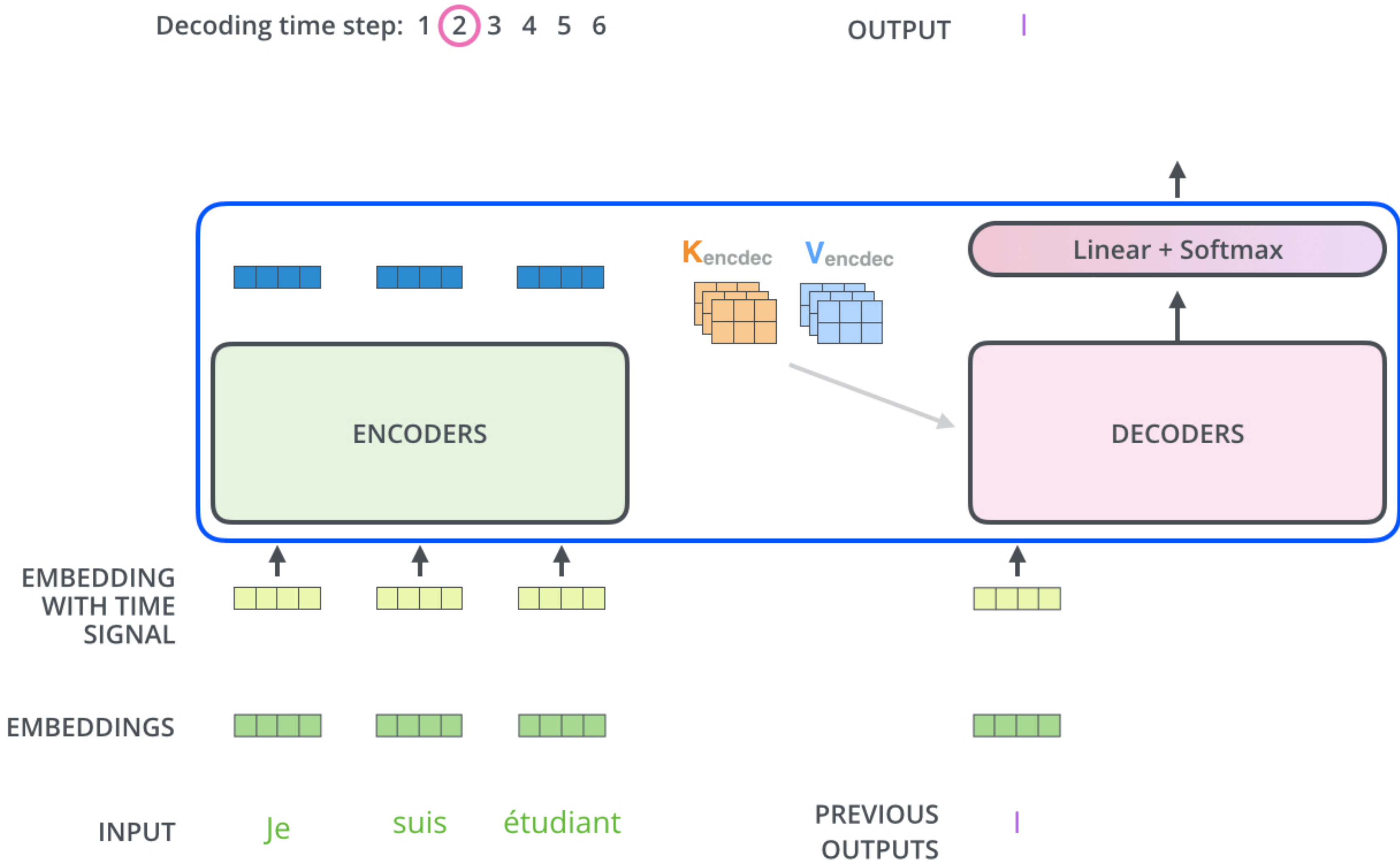
Encoder- Decoder Attention



Transformer decoder



Transformer Example

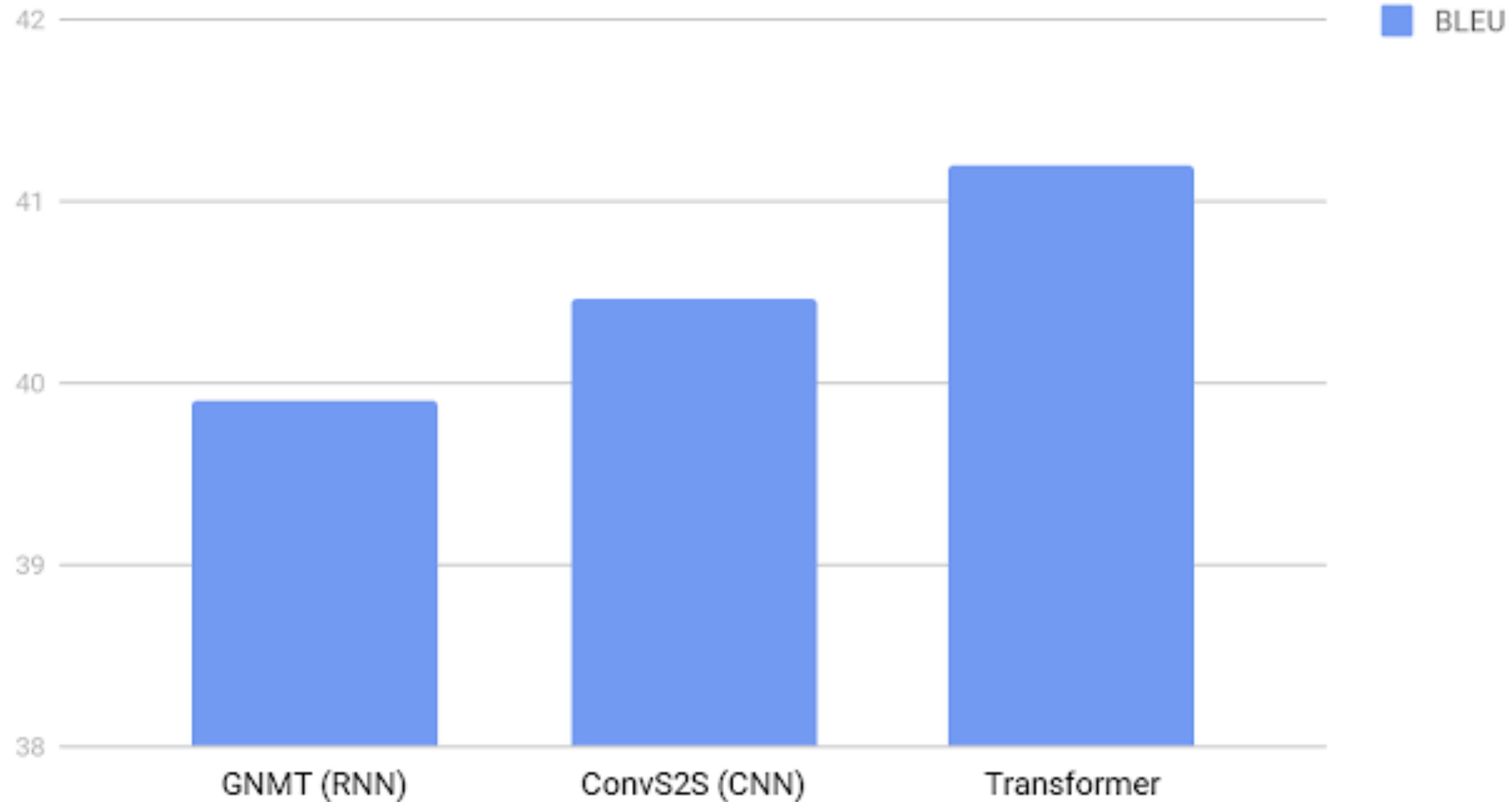


Pics taken from : <https://jalammar.github.io/illustrated-transformer/>

Neural Machine Translation Example

Neural Machine Translation Example

English French Translation Quality

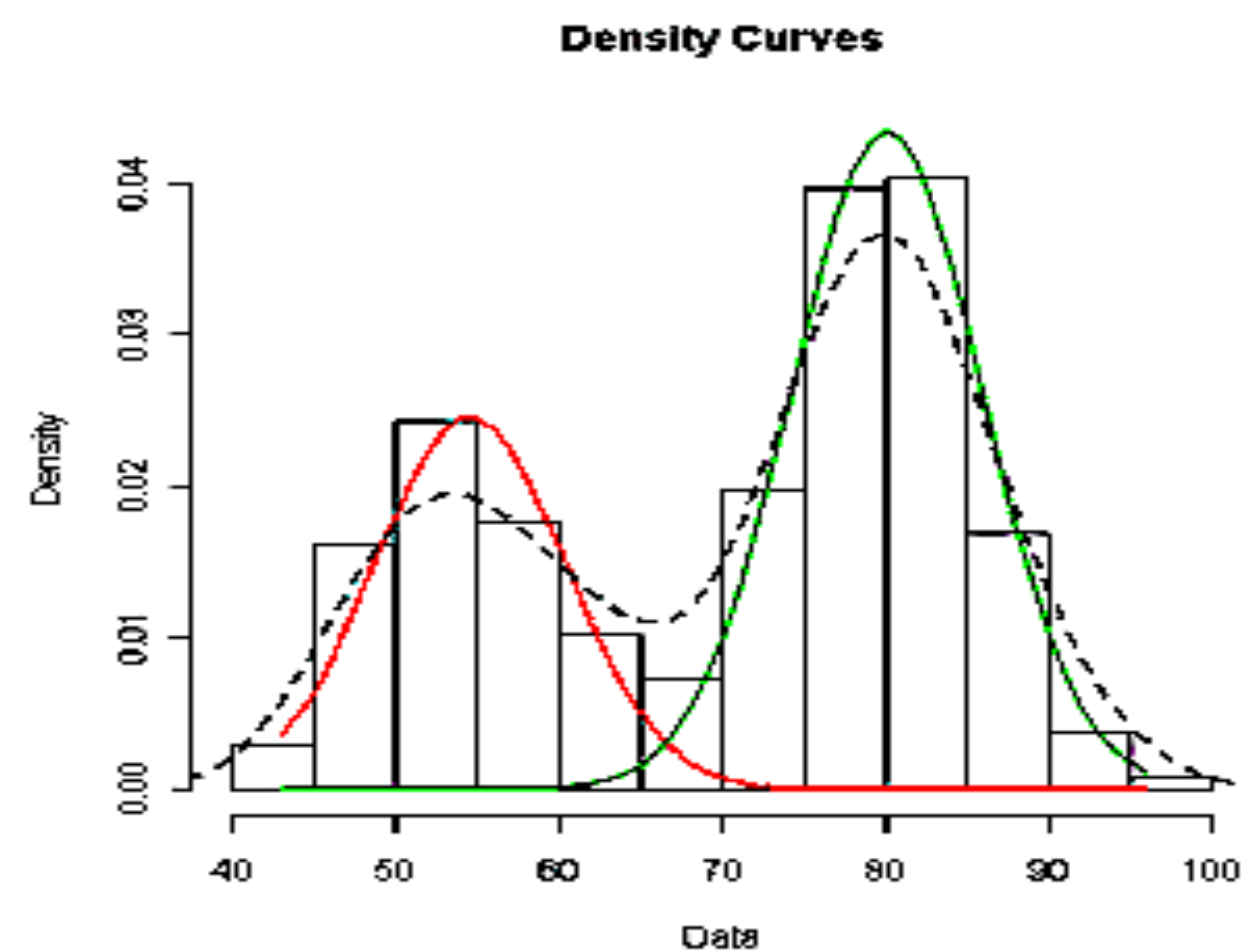


Unsupervised Learning

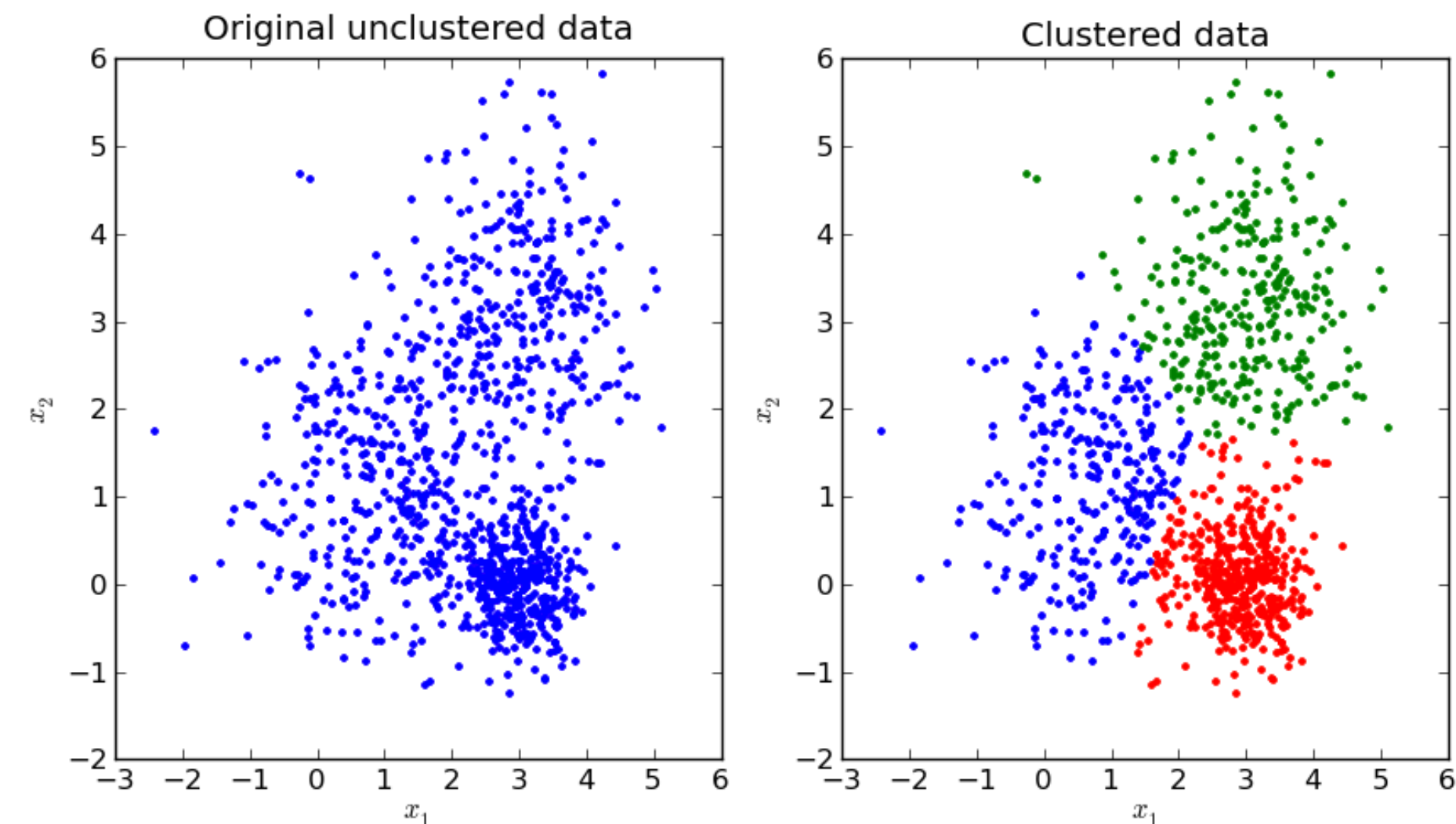
Unsupervised Learning

- Developing models that do not need labels
- May model the generation of data.
- May allow generation of new data samples
- Broad strategies for unsupervised learning

Learning the distribution of the data

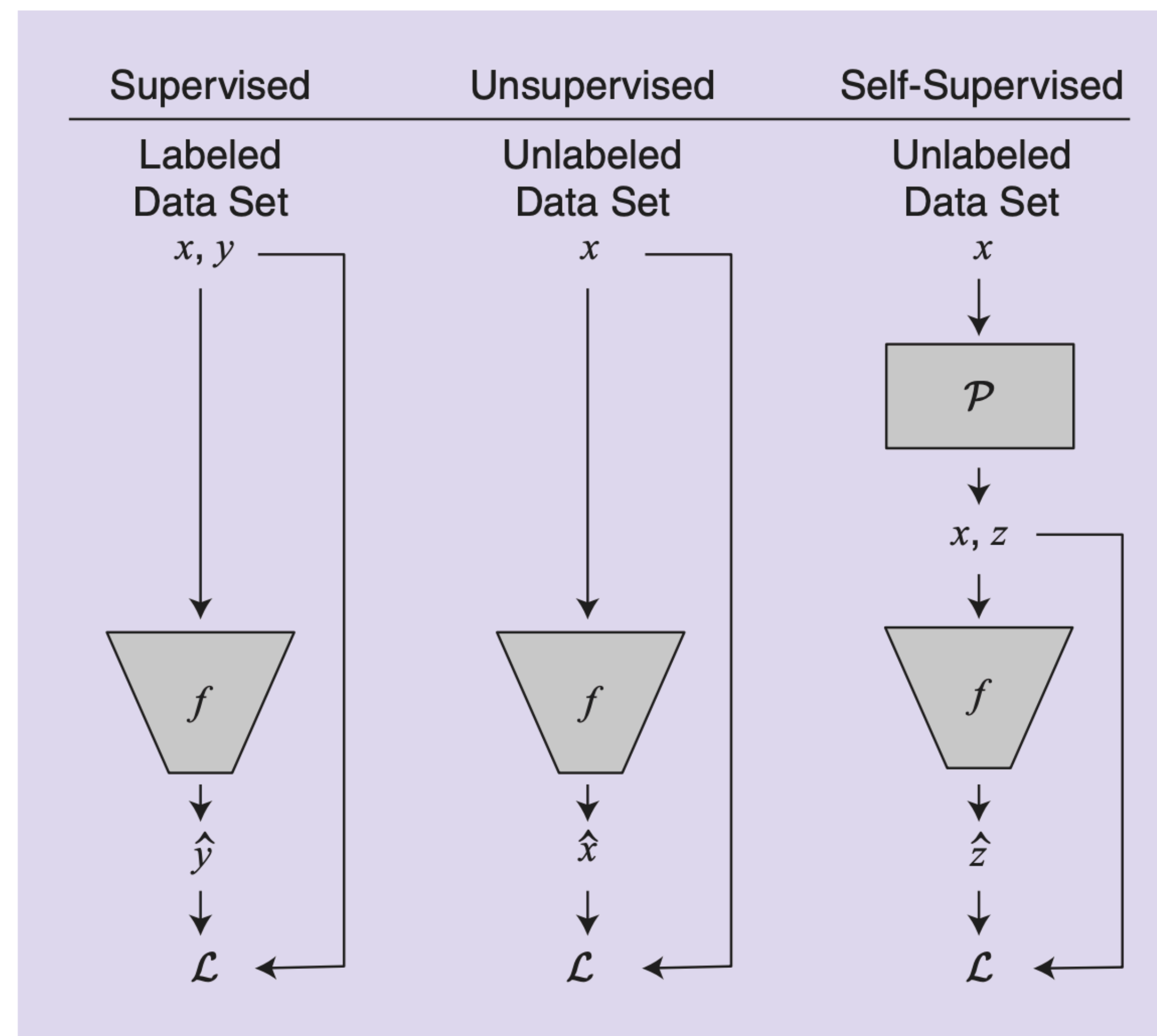


Detecting clusters in the data



Self supervision

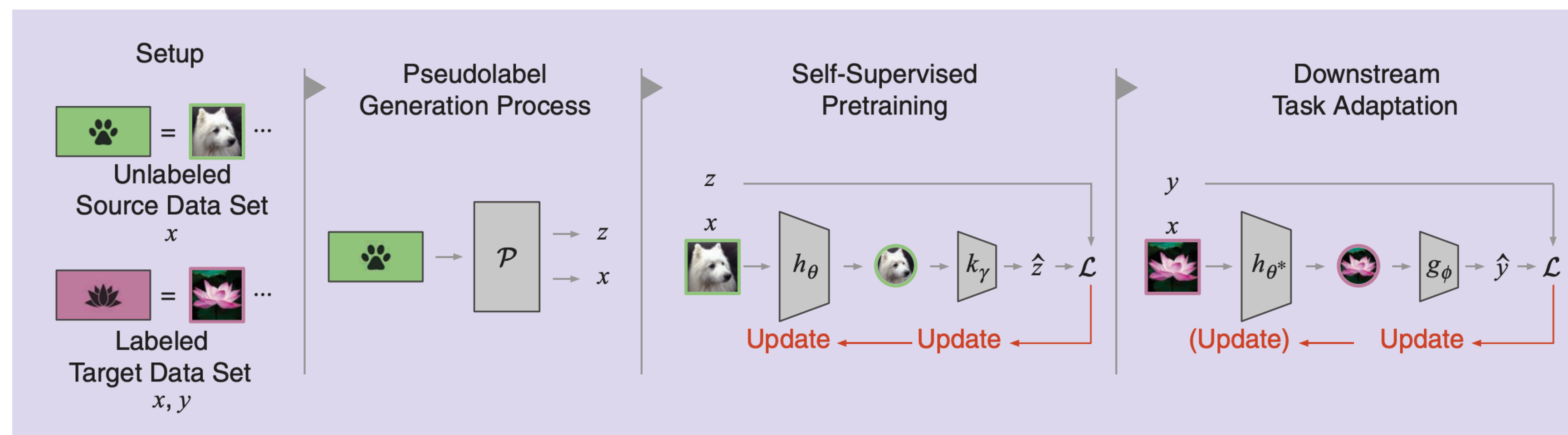
- ◆ Different from supervised and unsupervised learning
- * Does not perform distribution learning or reconstruction
- * Uses a pretext task
- * Performing contrastive or predictive learning
- ◆ Using large volumes of unsupervised data



Ericsson, Linus, et al. "Self-supervised representation learning: Introduction, advances, and challenges." *IEEE Signal Processing Magazine* 39.3 (2022): 42-62.

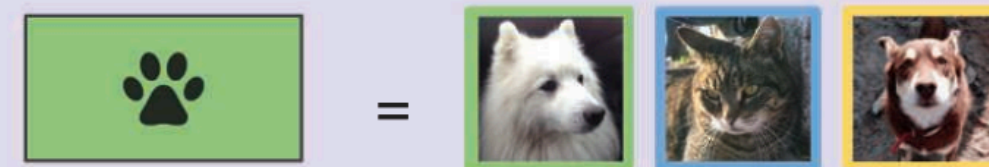
Self supervision - principle

- ◆ Two levels of modeling with unsupervised data
 - ❖ Generating a pseudo-label
 - ❖ Learning the upstream model
- ◆ Downstream task performs fine-tuning of the SSL model.

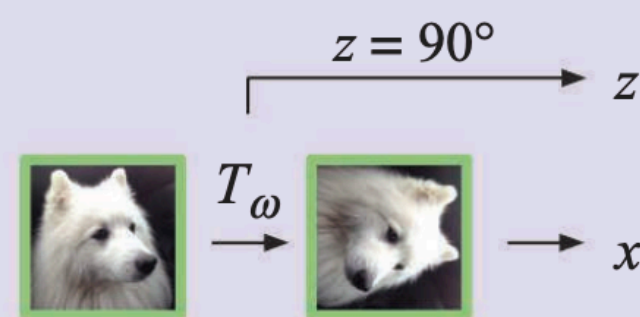


Self supervision - pre-text task

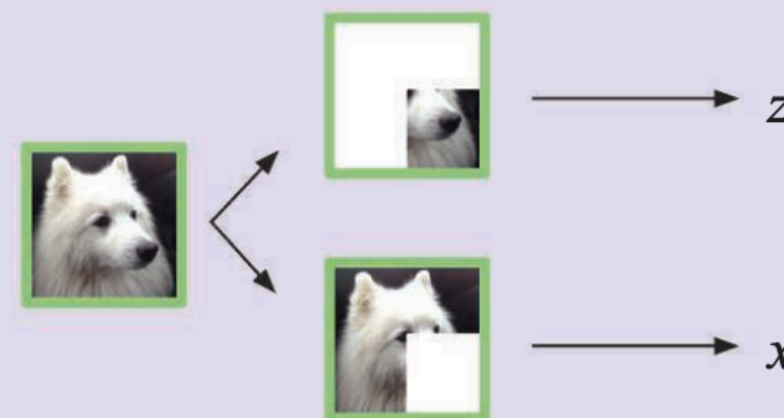
Pseudolabel Generation Processes



Transformation Prediction



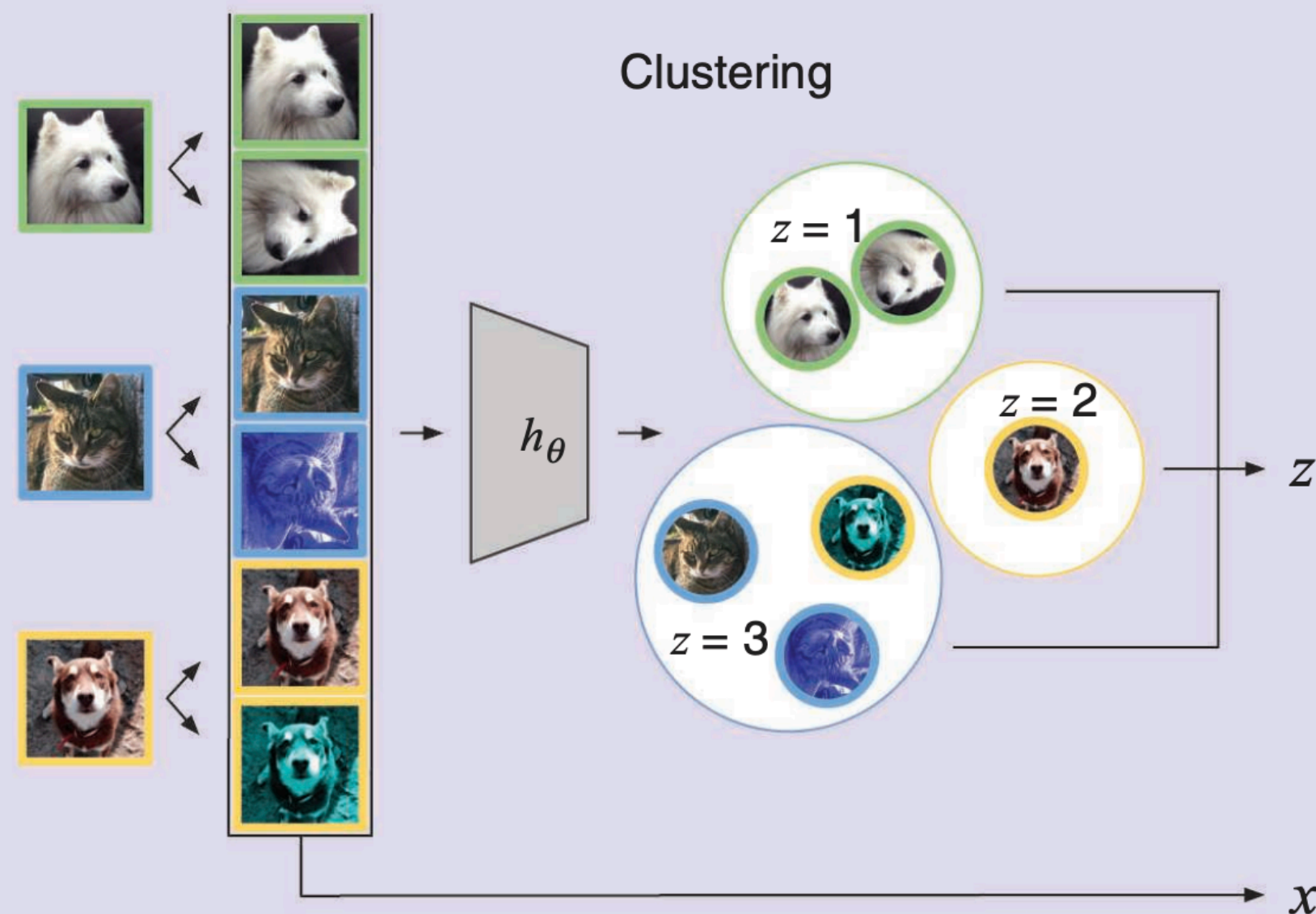
Masked Prediction



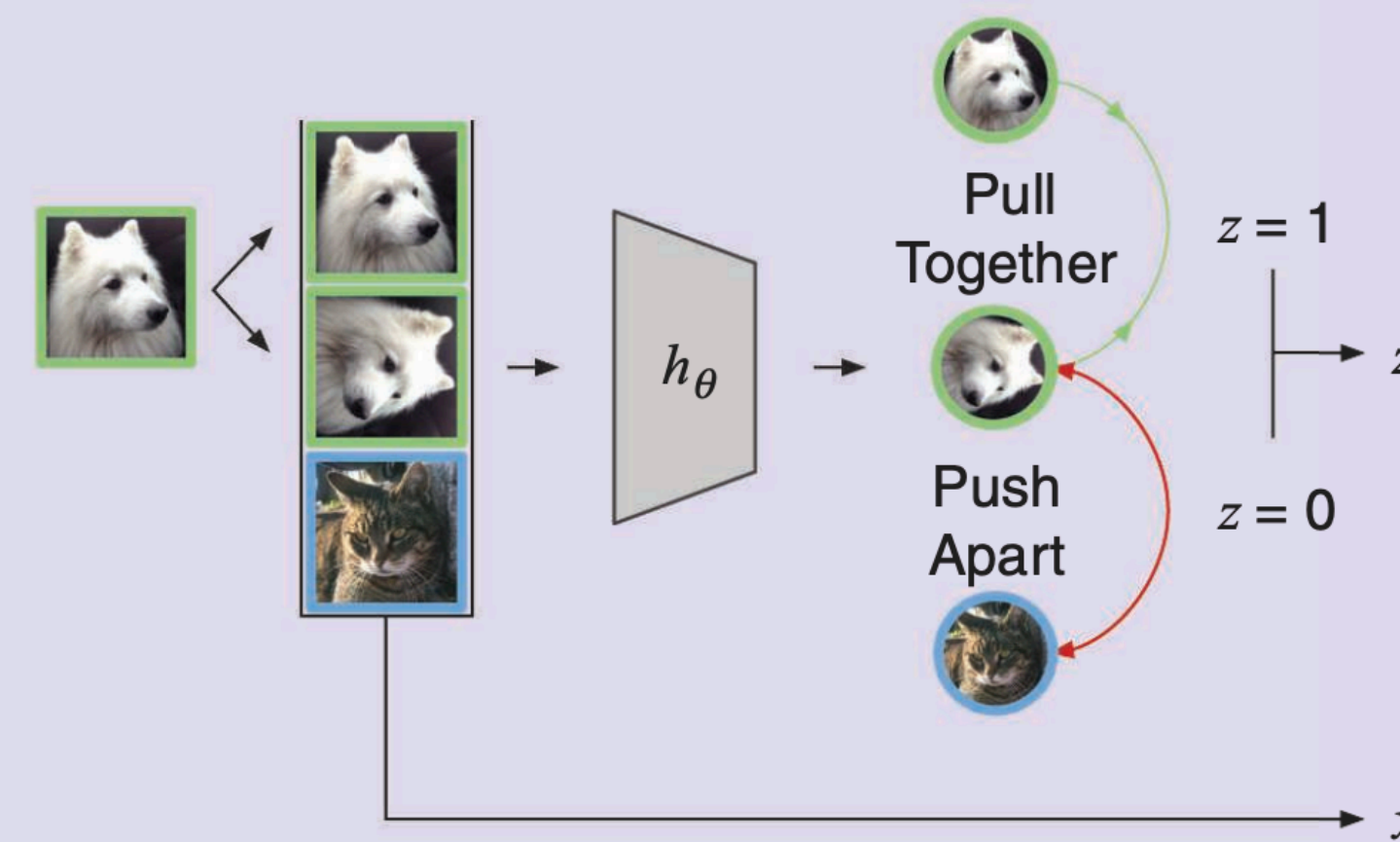
Instance Discrimination



Clustering

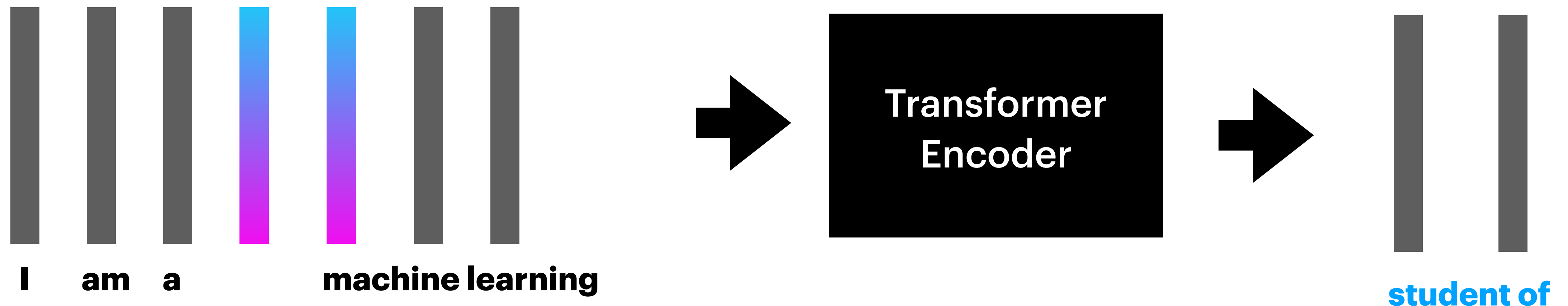


Contrastive Instance Discrimination



Self-supervision as a task

- Masking out portions of the input data
 - * Pass the rest of the embeddings (with zeros or random entries at the masked locations) to the transformer encoder
 - * Have the model predict the word tokens in the masked portions - **Masked Language Modelling (MLM)**



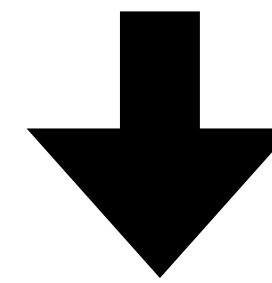
Large language models (LLMs)

- Extending the task of self-supervision
- Mine lots of text data
 - * Crawled from the web, as well as, from other resources.
- Design the **model with large capacity** (Millions → Billions of parameters)
- Pre-train the model
 - * With MLM and similar style of losses
 - * High resource of computations.
- Final trained model can be **fine-tuned for supervised tasks**
 - * Load the parameters as initialization and perform supervised learning.

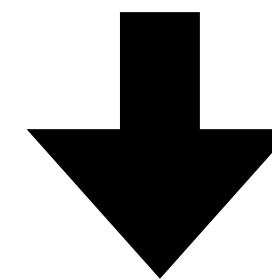
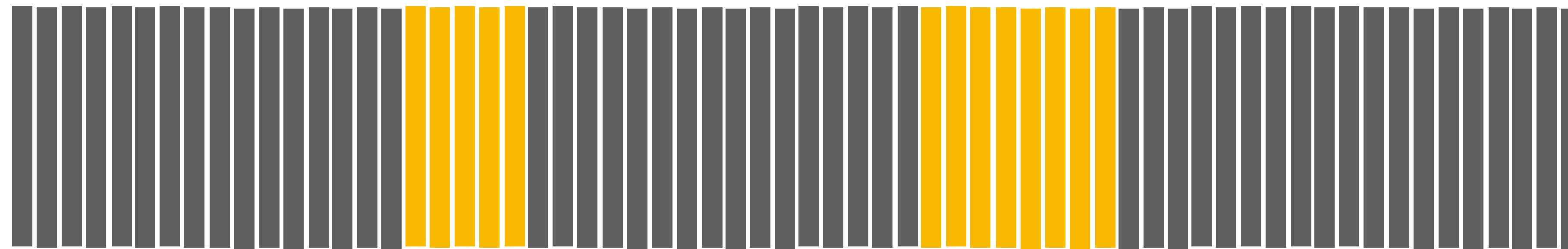
Large language models (LLMs)

- Self-supervised learning
 - * Has shown emergent abilities to generalise to wide variety of downstream tasks.
 - ✓ Tasks that the model was not trained on
 - ✓ Not seen in smaller models
 - * Enables to build reasoning capabilities in the model.
 - * **Applicable for several domains** - text, speech and images.

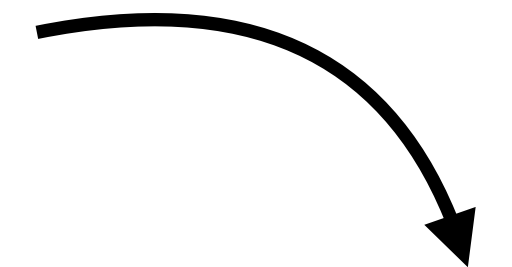
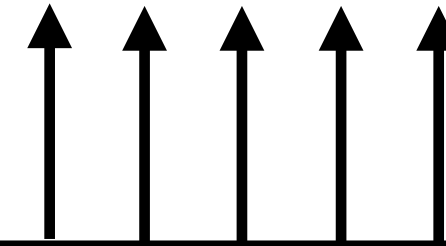
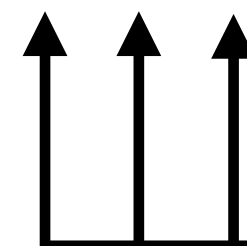
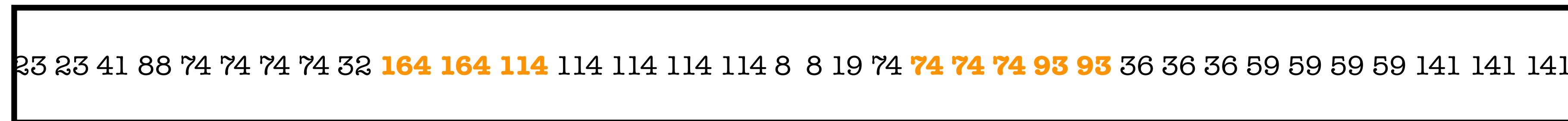
Self-supervision in audio - wav2vec



Sequence of spectral vectors



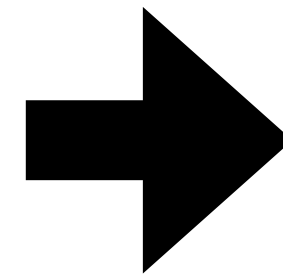
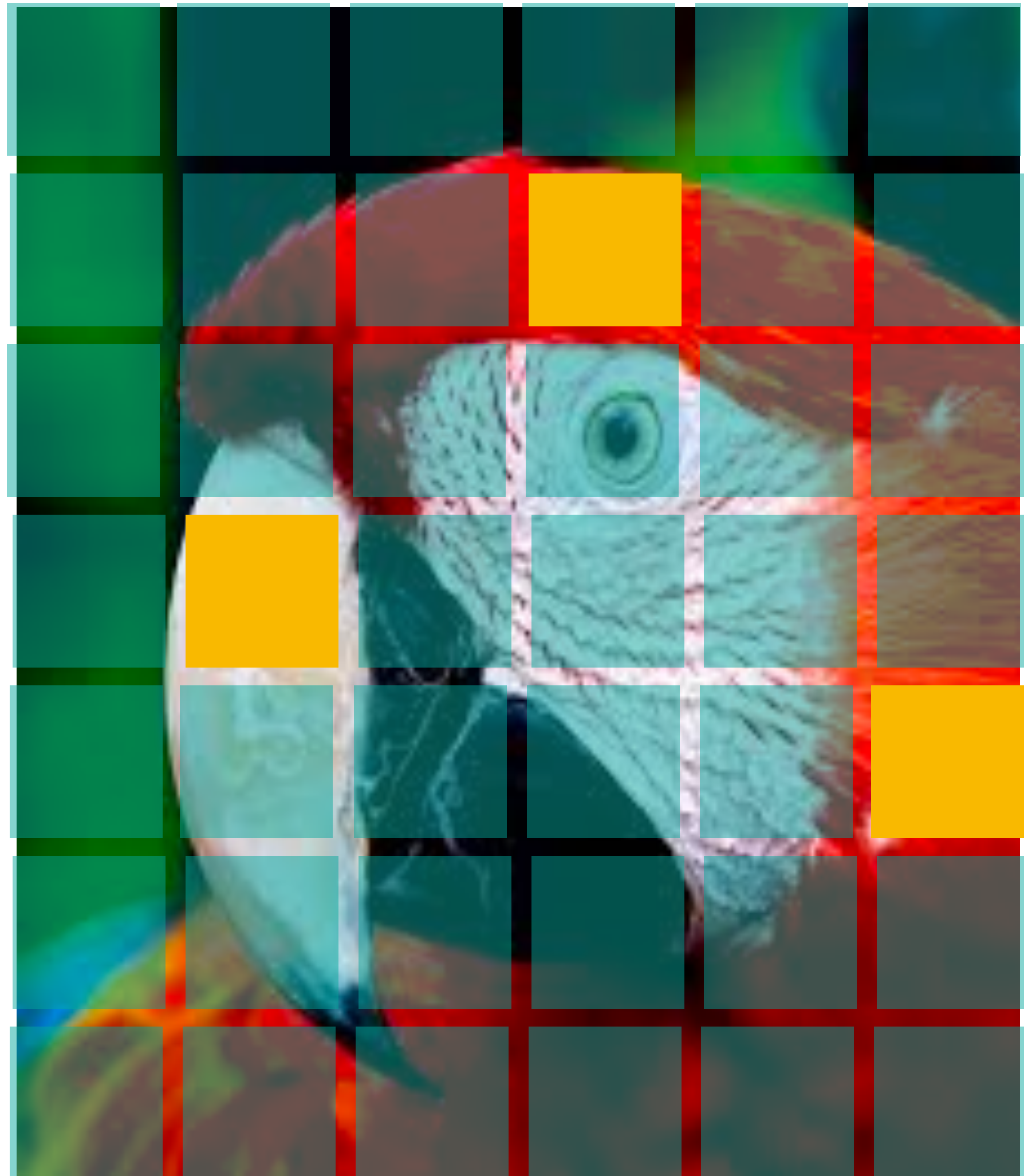
Quantization



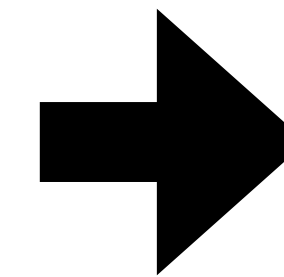
**Transformer
Encoder**



Self-supervision in images - Vision Transformer



Transformer
Encoder



Predict
Masked
Regions

LLM - Examples

- Generative Pre-trained Transformers (GPT) series

	Architecture	Data Used	Model Size
GPT-1	Transformer (12 layer, decoder only model)	Book Corpus (4.5GB)	117M
GPT-2	GPT-1 with additional normalisation layers	Web Text (40GB)	1.5B
GPT-3/3.5	GPT-2 with more layers Adding Fine-tuning tasks and human feedback	Large Web Crawl (570B)	175B
GPT-4/4o	Details Undisclosed [Trained with Text + Images]	—	—

Future works (some already underway)

- Multi-modal

- * Incorporating learning across modalities

- ✓ Create a domain specific encoder/decoder and learning the joint language model.

- Combining some labeled data with the self-supervised data to further improve the models.

- ✓ Current models like GPT use human feedback.

- Understanding the risks and vulnerabilities of these models.

THANK YOU

Sriram Ganapathy and TA team
LEAP lab, C328, EE, IISc
sriramg@iisc.ac.in

