

Transformers

Word and image embeddings with global attention.

Karel Zimmermann

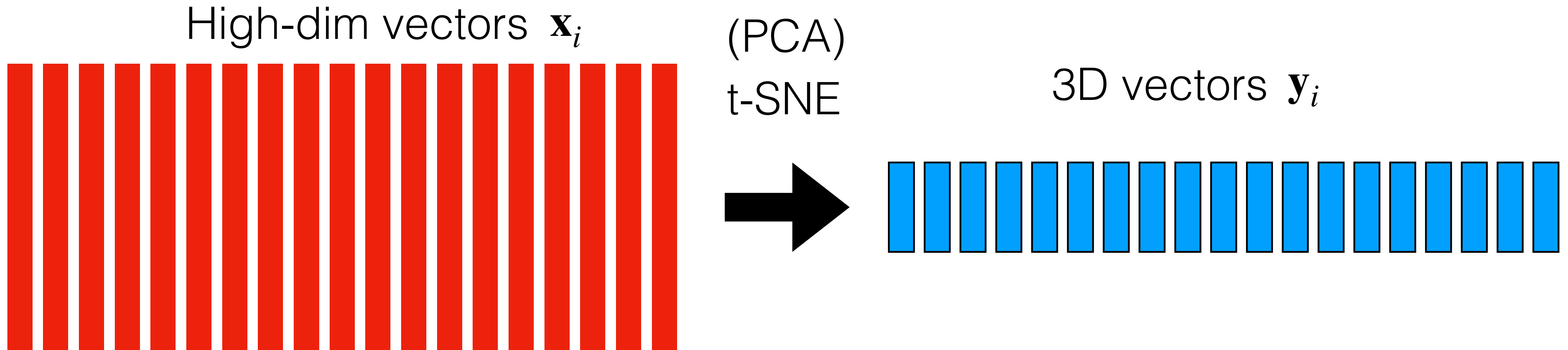
Czech Technical University in Prague

Faculty of Electrical Engineering, Department of Cybernetics



Pre-requidity: Visualizing high-dimensional data

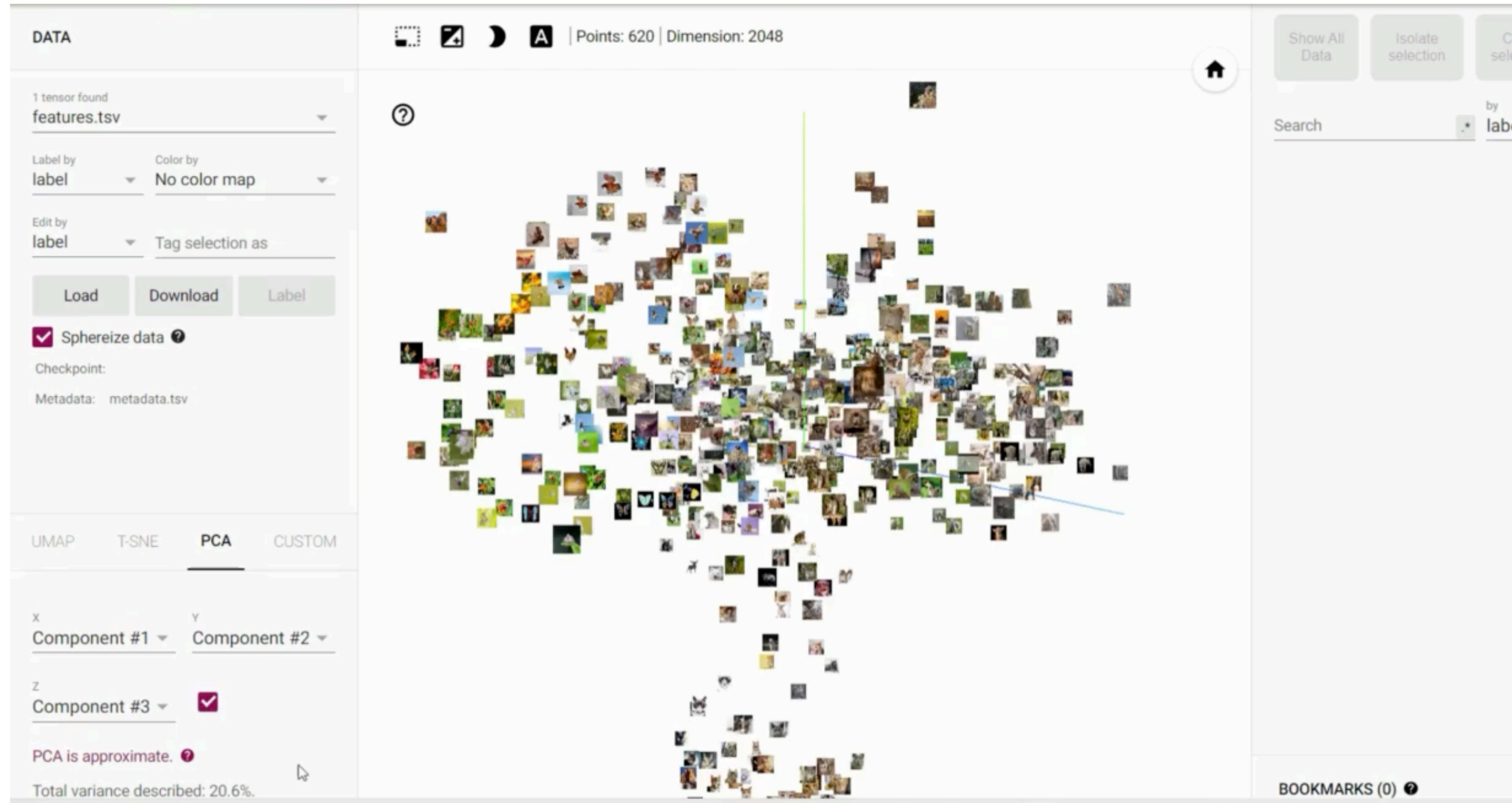
Visualizing high-dimensional embedding in 2D/3D world.

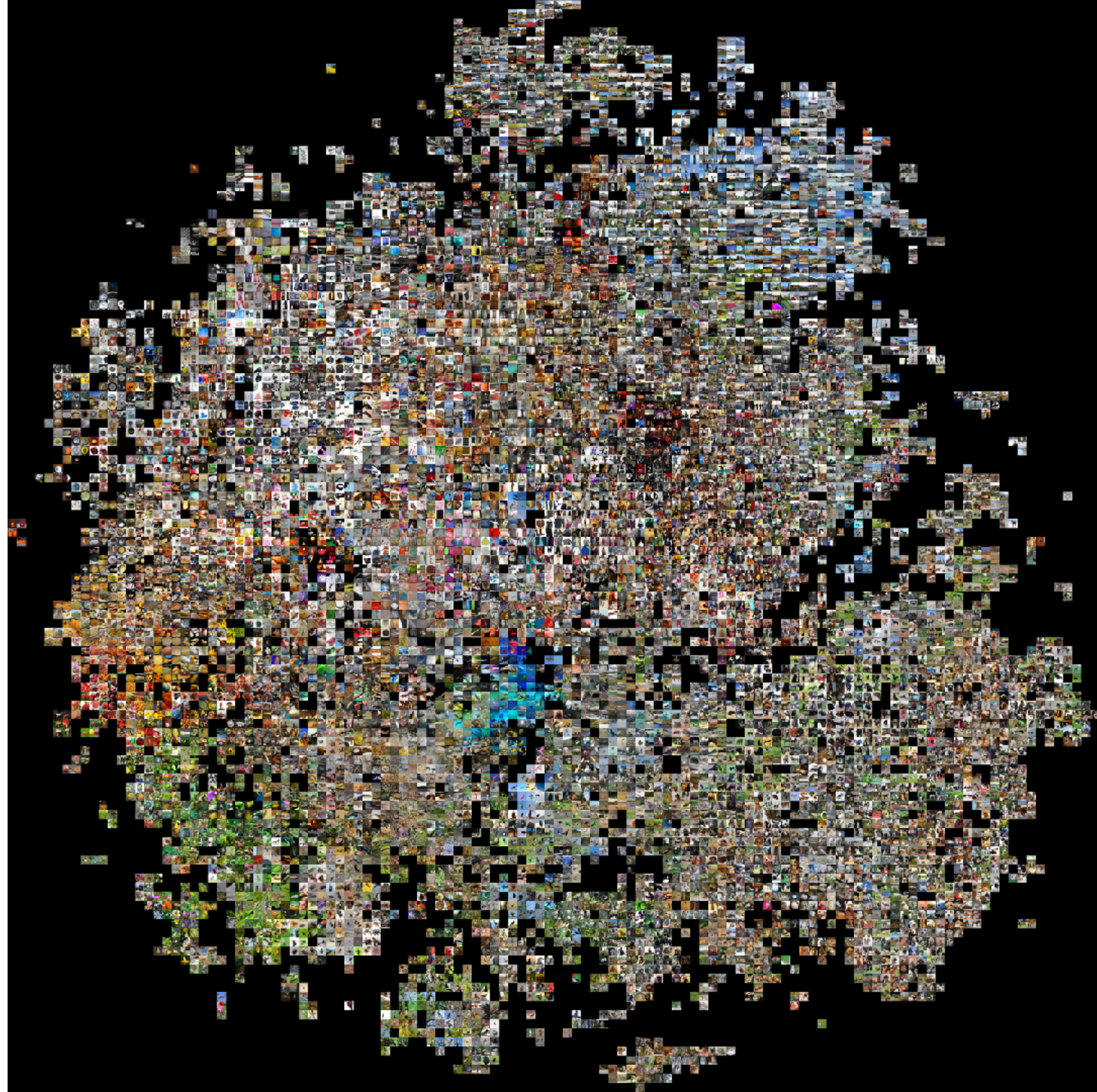


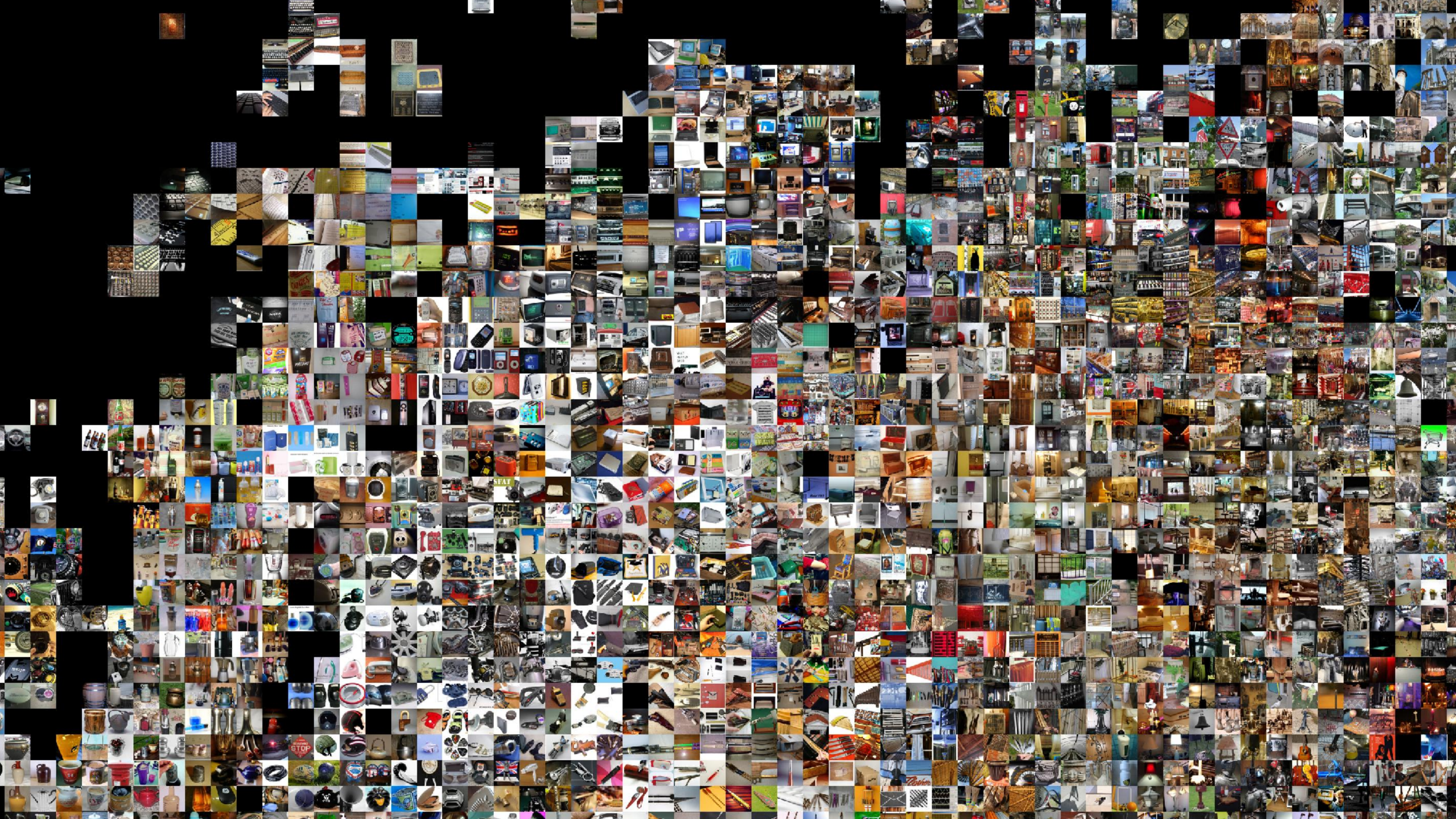
1. Randomly initialize \mathbf{y}_i by normal zero-mean noise $\mathcal{N}(0, 0.001)$
2. Compute pair-wise probabilities in \mathbf{x}_i :
$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$
3. Compute pair-wise probabilities in \mathbf{y}_i :
$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$
4. Optimize \mathbf{y}_i to get similar distribution
$$\text{KL}(P \| Q) = \sum_{i \neq j} p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right)$$

Gaussian
distribution
t-distr. with
heavier-tails
(crowding)

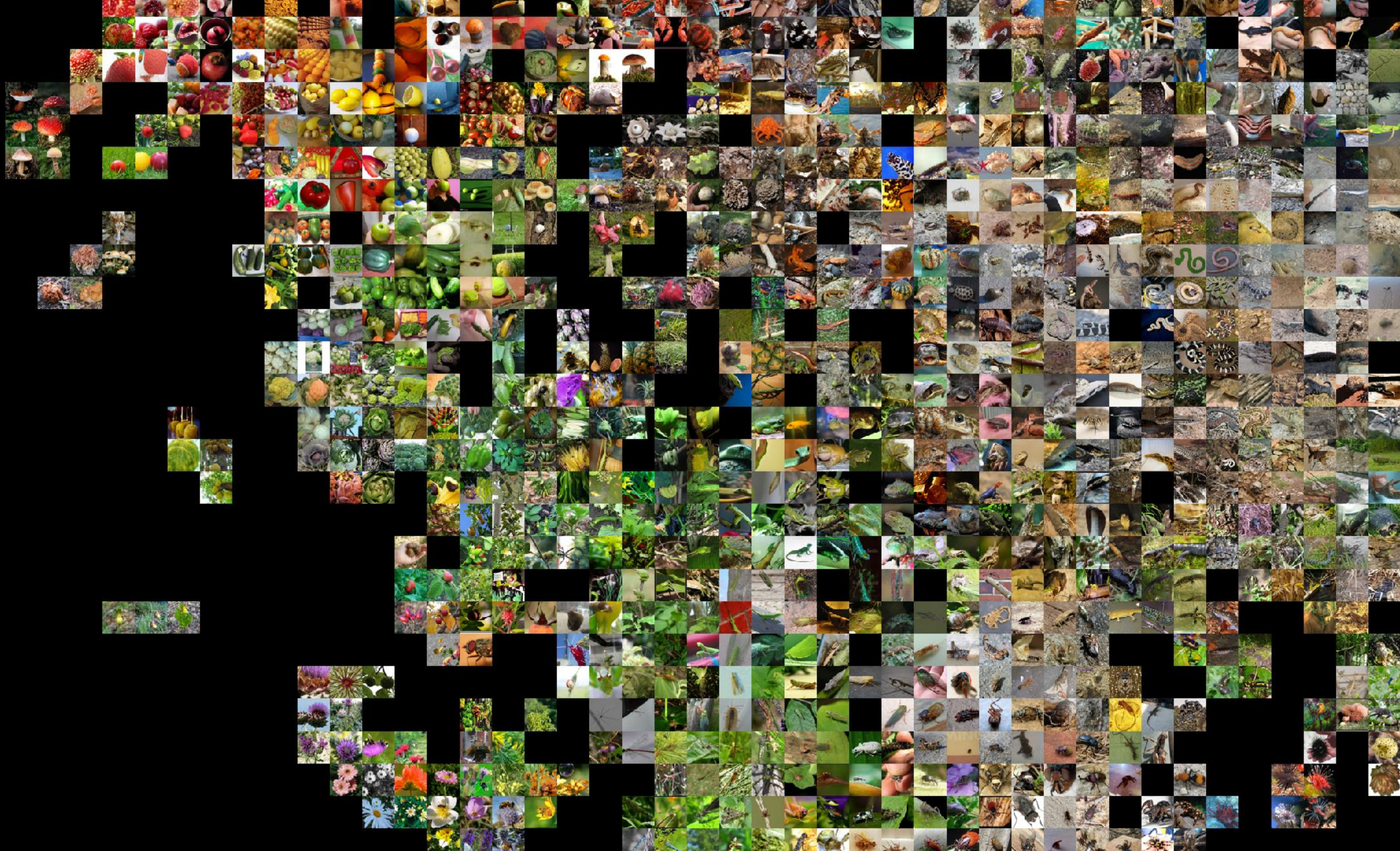
Visualizing high-dimensional embedding in 2D/3D world.



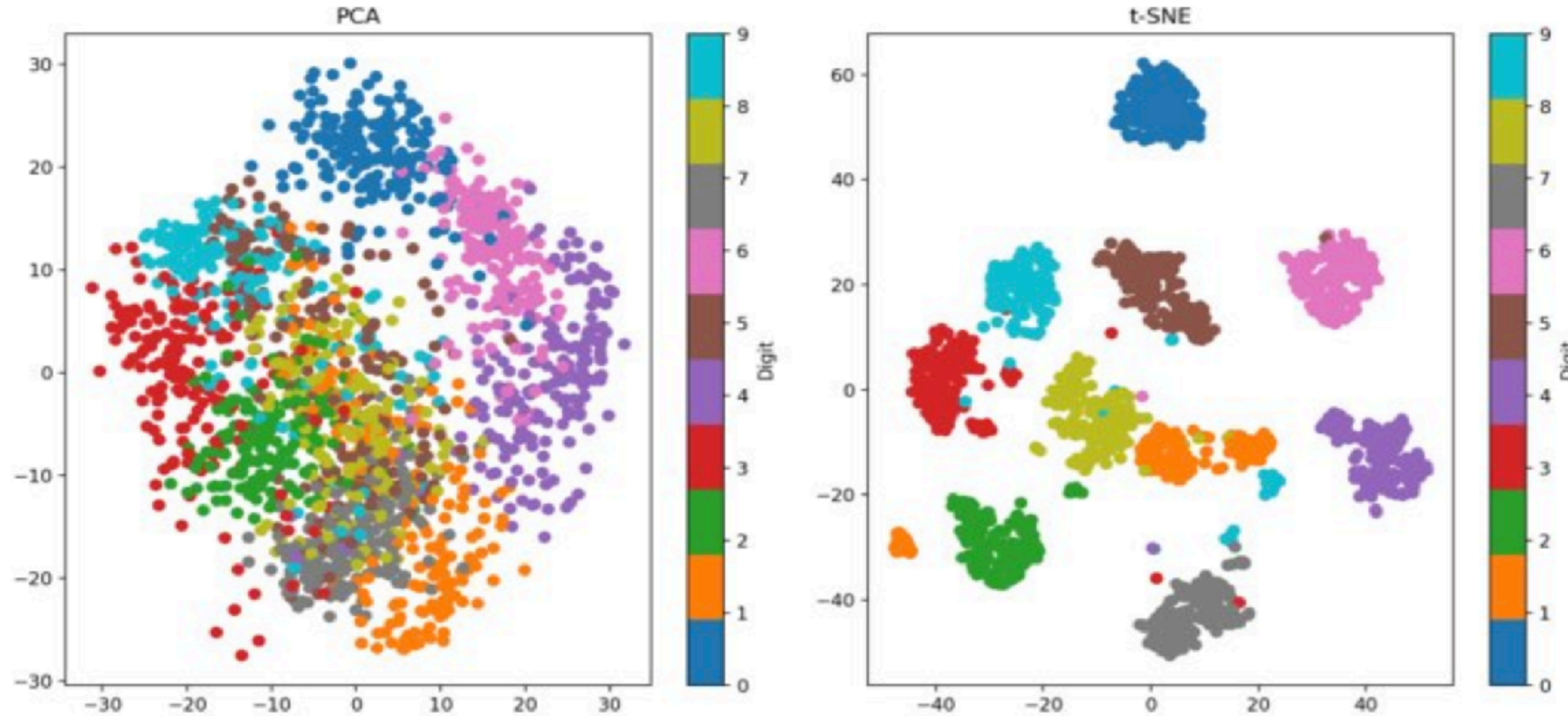








Visualizing high-dimensional embedding in 2D/3D world.



- t-SNE (t-distributed Stochastic Neighbor Embedding)
 - Captures **non-linear** relationships in data
 - Separate clusters based on their high-dimensional proximity
 - Outcome is stochastic and depends of perplexity σ
- PCA
 - Captures **linear** relationships in data
 - Deterministic and useful for preprocessing

Transformers in language (NLP)

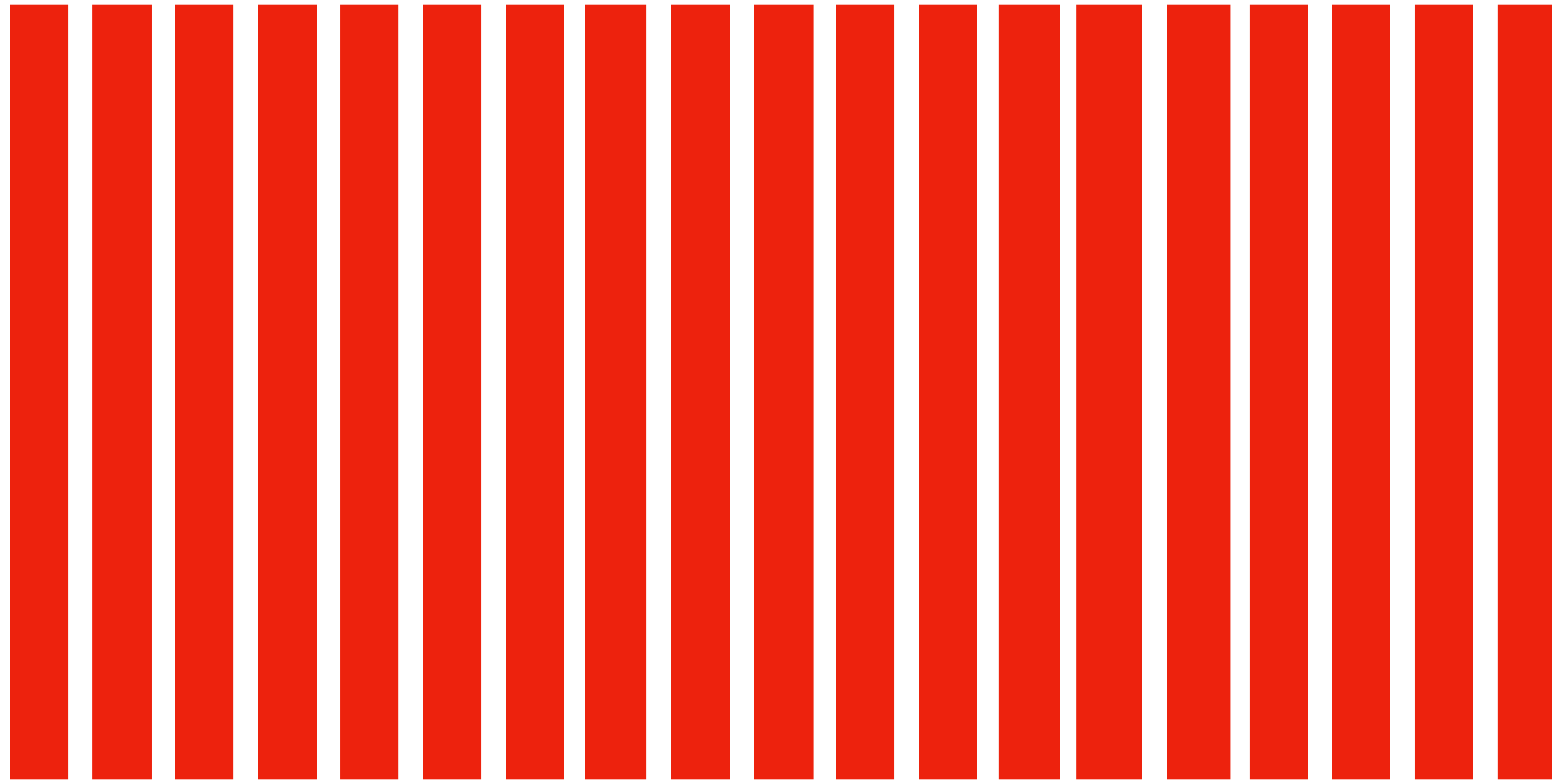
Word2vec represents words as low-dimensional continuous vectors
[Mikolov NIPS 2013]

Word2vec represents words as low-dimensional continuous vectors

N-word vocabulary (one-hot enc.)

Karel is the best teacher in the whole world

N-dim inputs

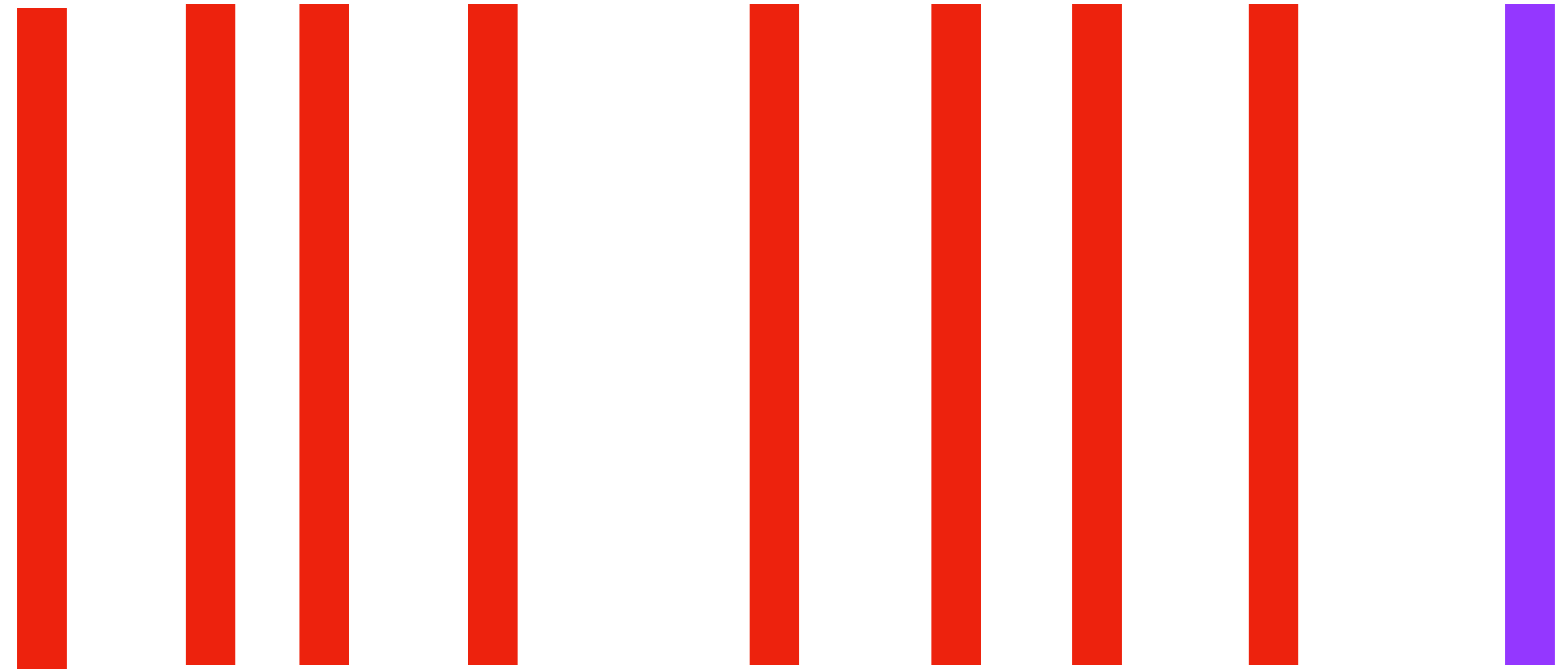
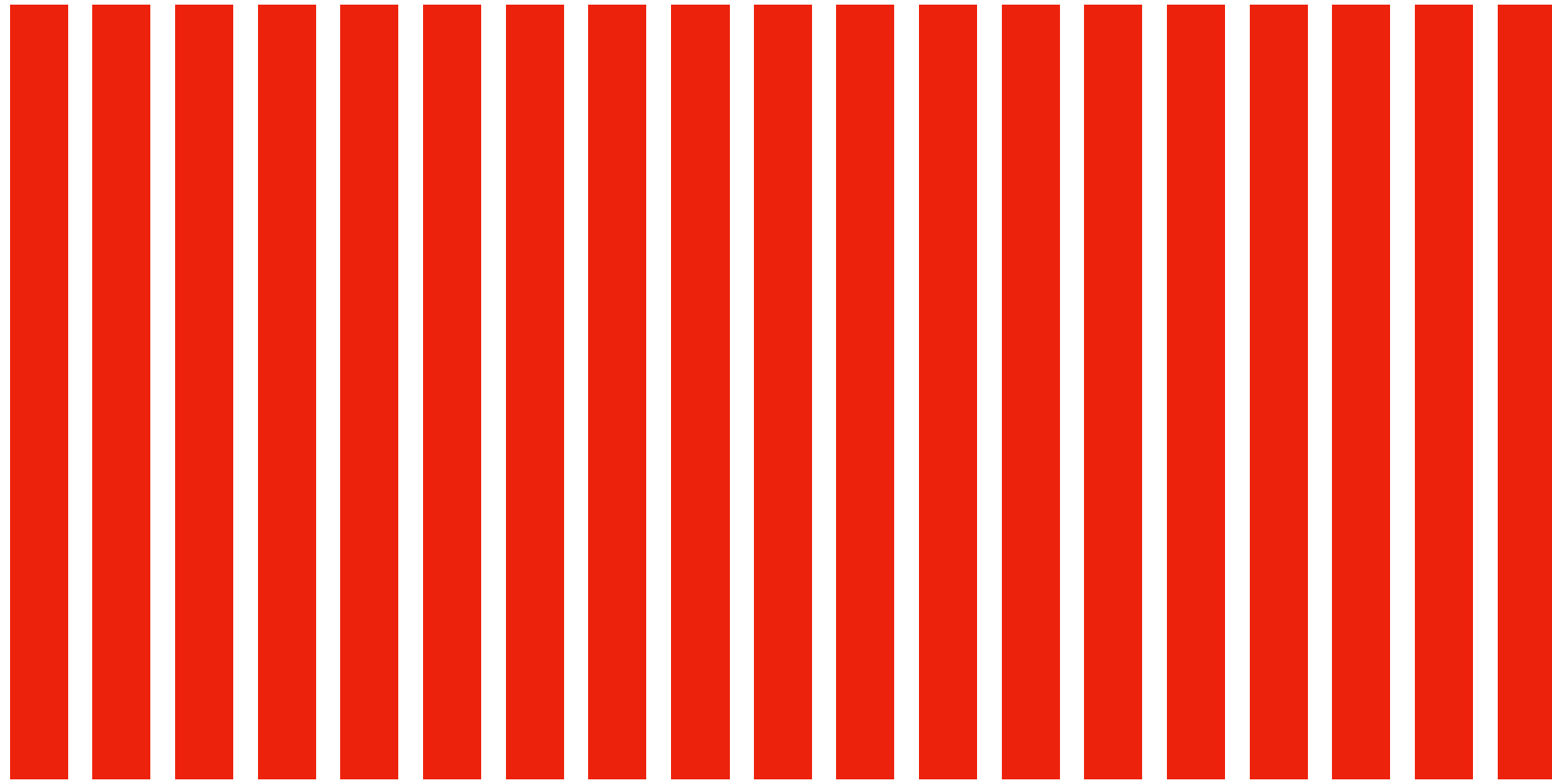


Word2vec represents words as low-dimensional continuous vectors

N-word vocabulary (one-hot enc.)

Karel is the best teacher in the whole world

N-dim inputs



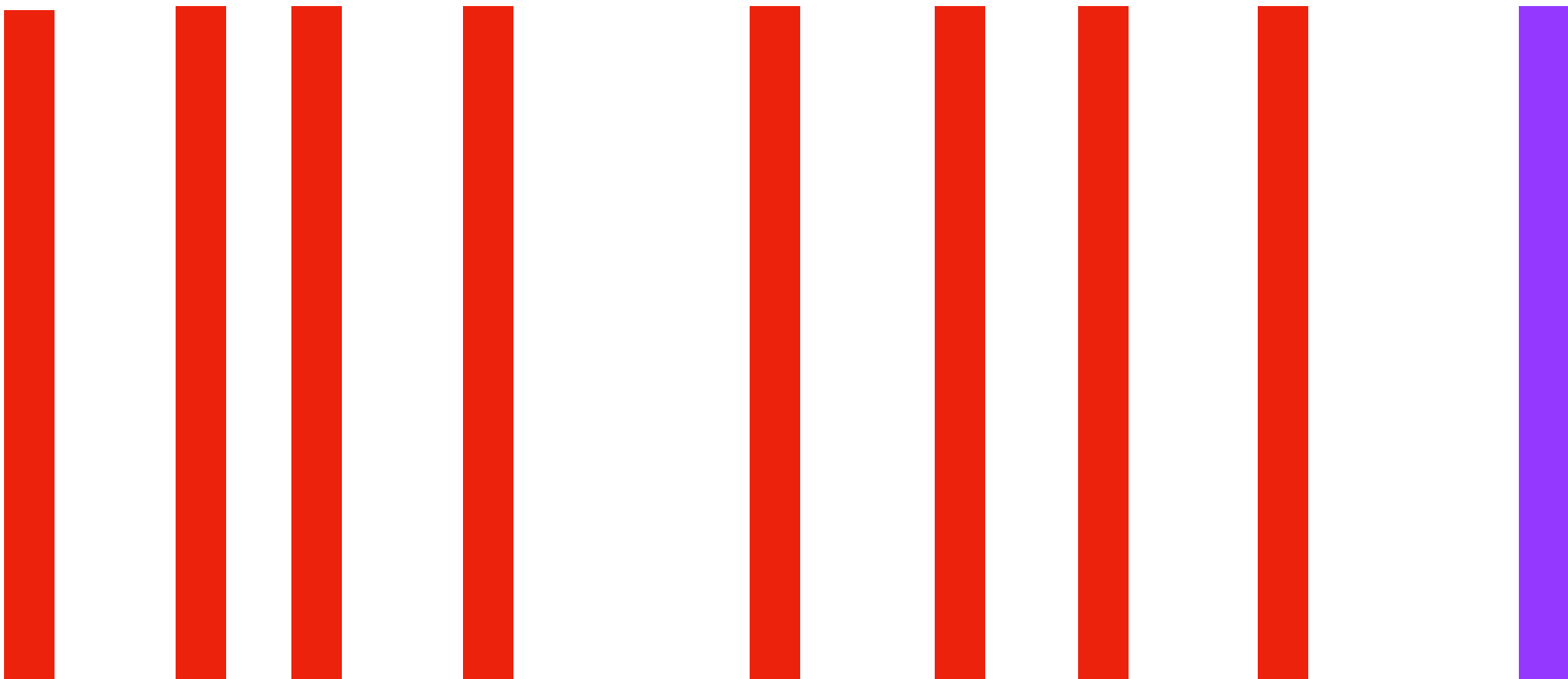
Word2vec represents words as low-dimensional continuous vectors

N-word vocabulary (one-hot enc.)

N-dim inputs

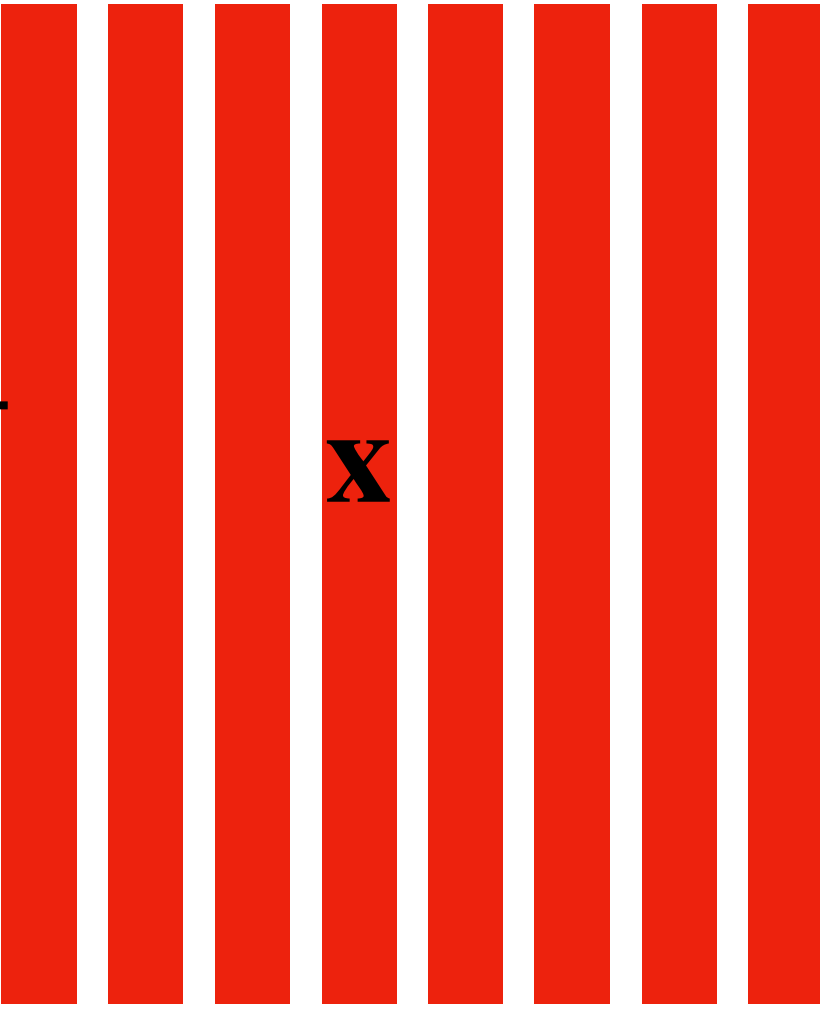


Karel is the best teacher in the whole world



S input words

N-dim inputs



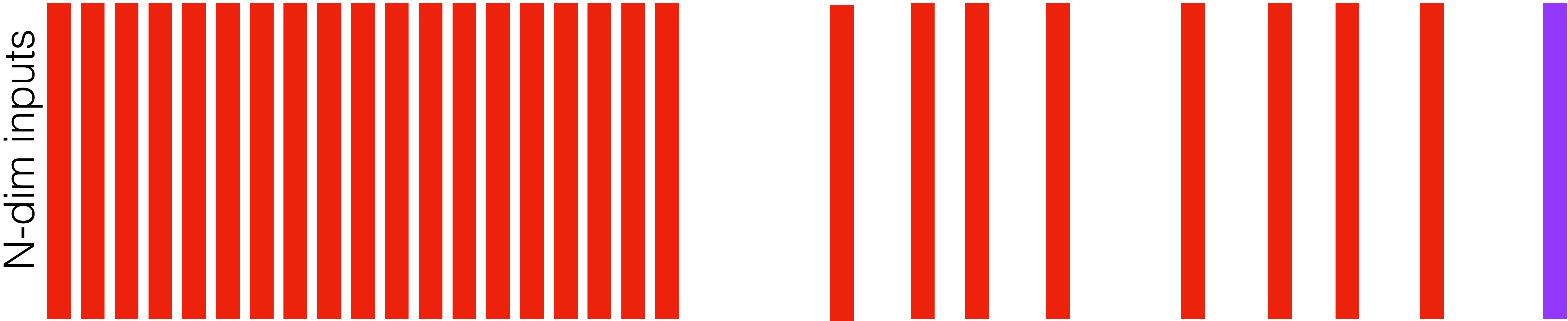
x

(N x S)

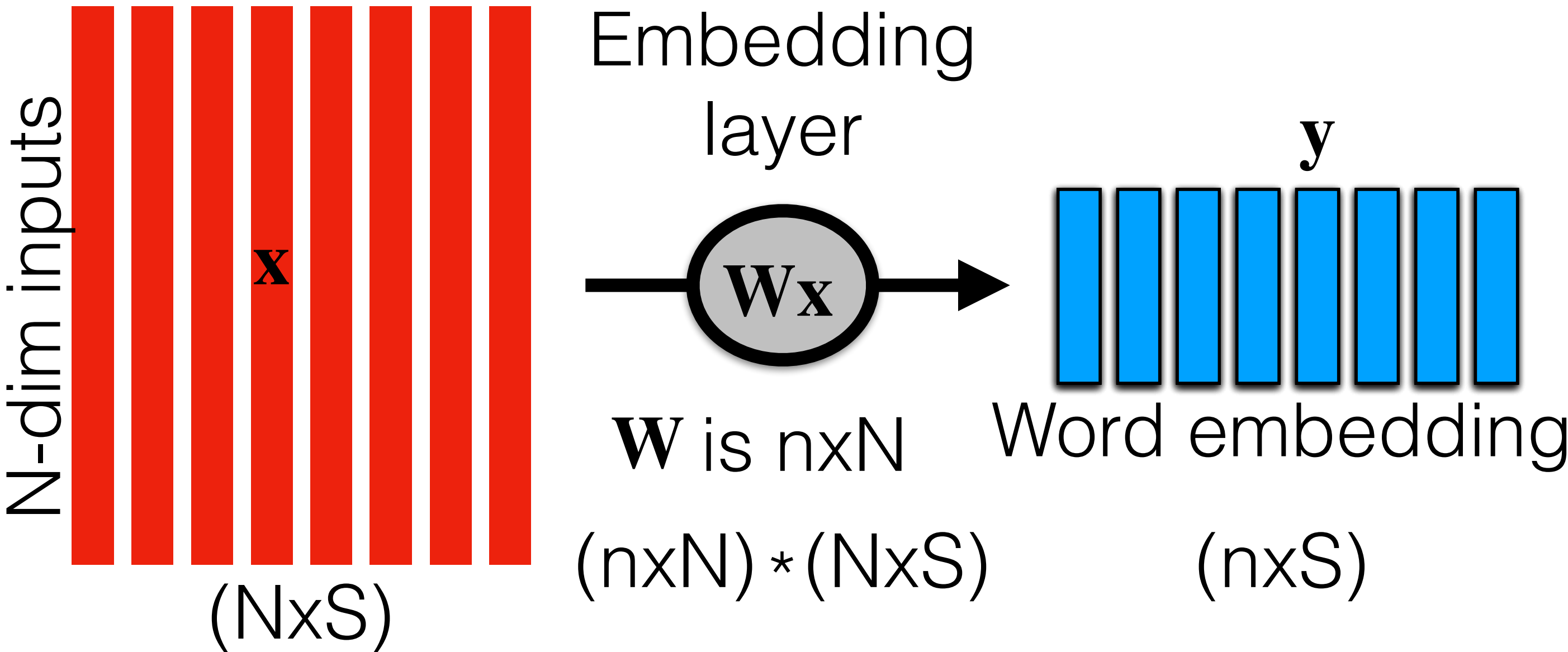
Word2vec represents words as low-dimensional continuous vectors

N-word vocabulary (one-hot enc.)

Karel is the best teacher in the whole world



S input words



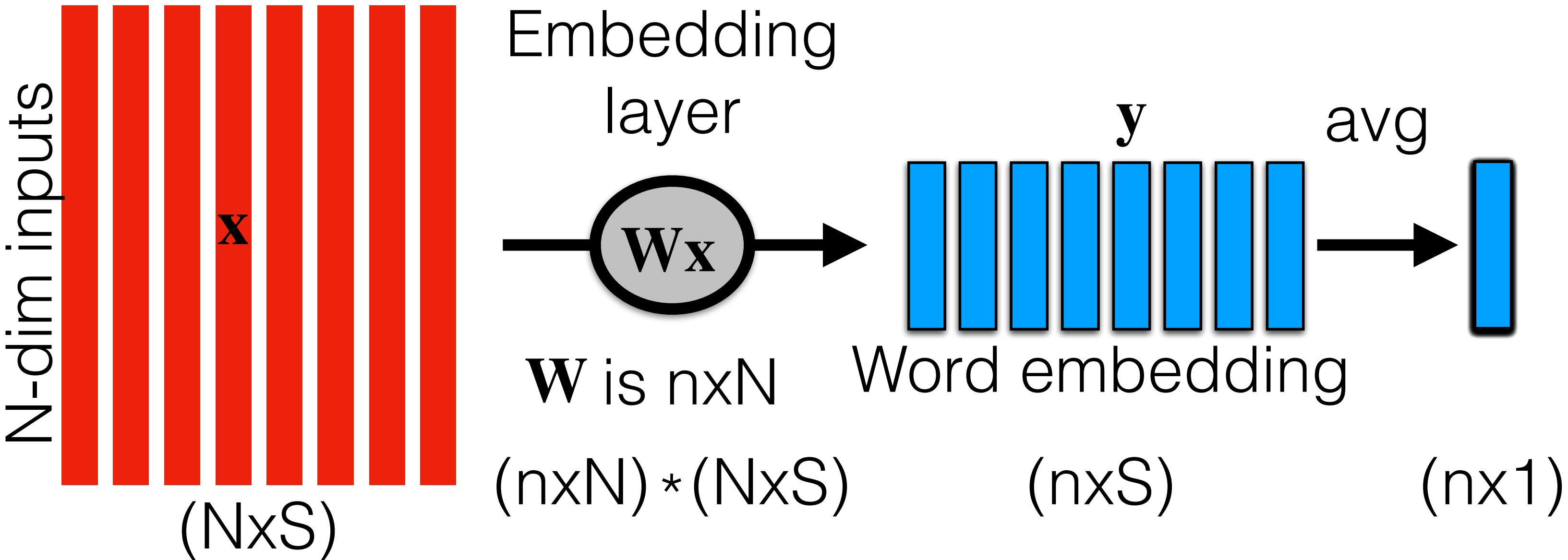
Word2vec represents words as low-dimensional continuous vectors

N-word vocabulary (one-hot enc.)

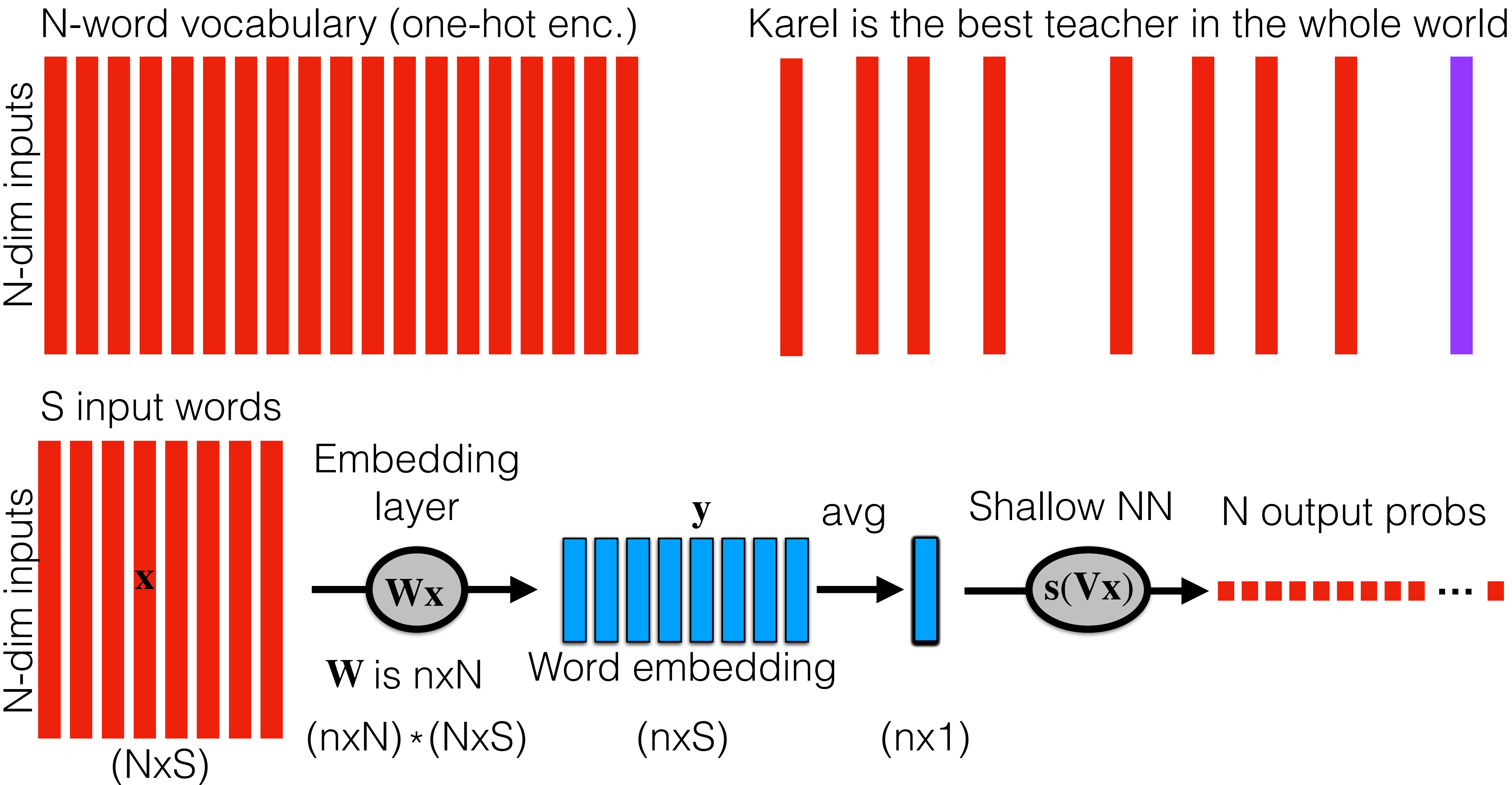
Karel is the best teacher in the whole world



S input words

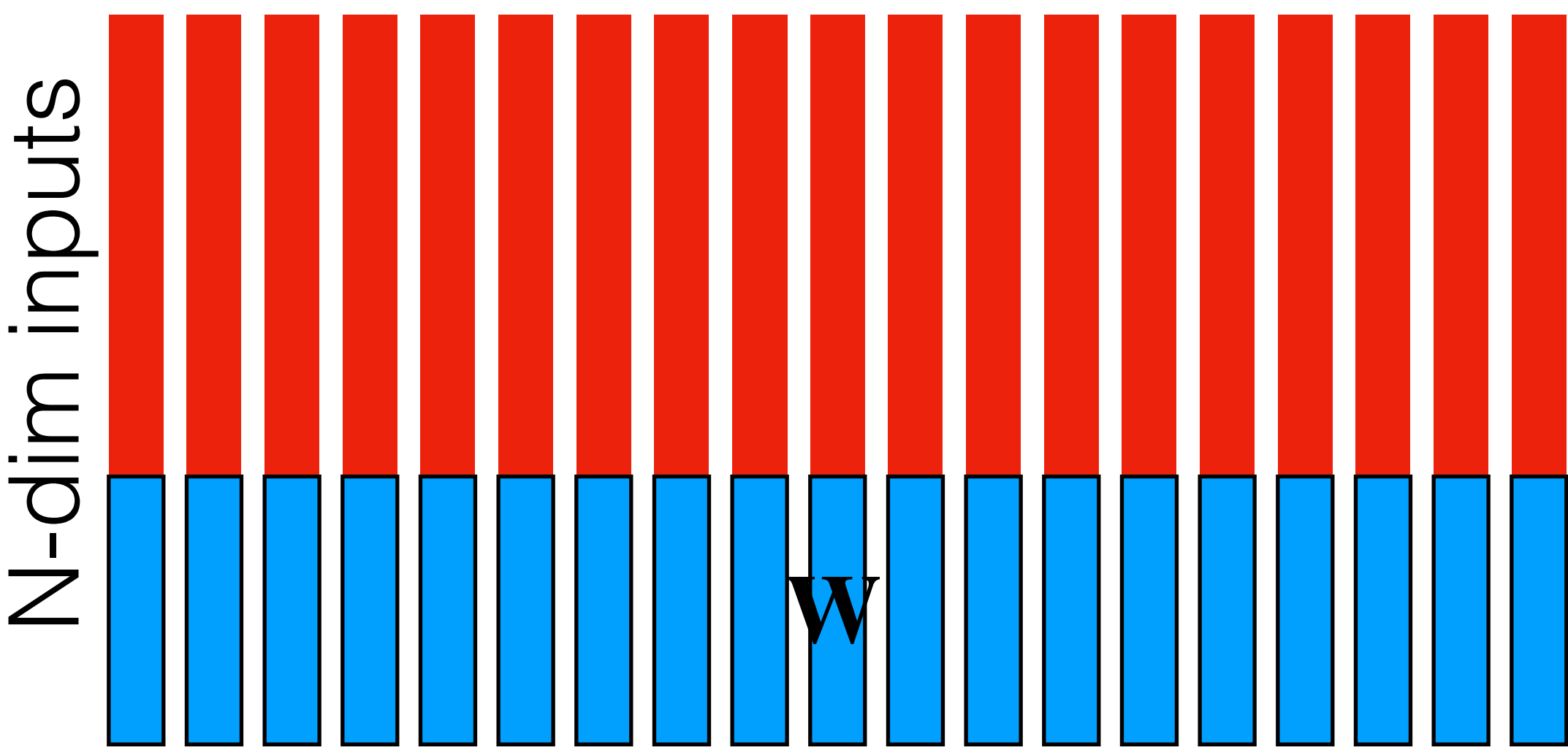


Word2vec represents words as low-dimensional continuous vectors

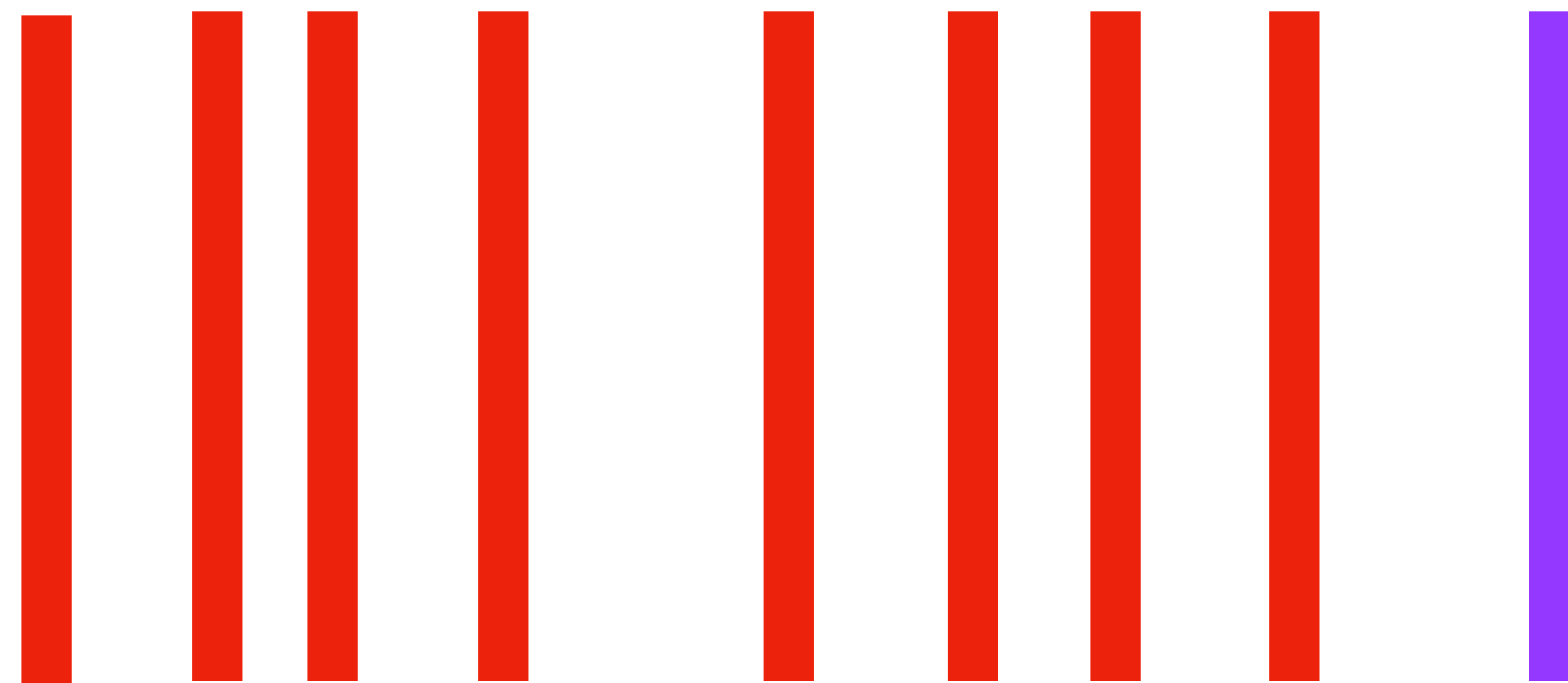


Word2vec represents words as low-dimensional continuous vectors

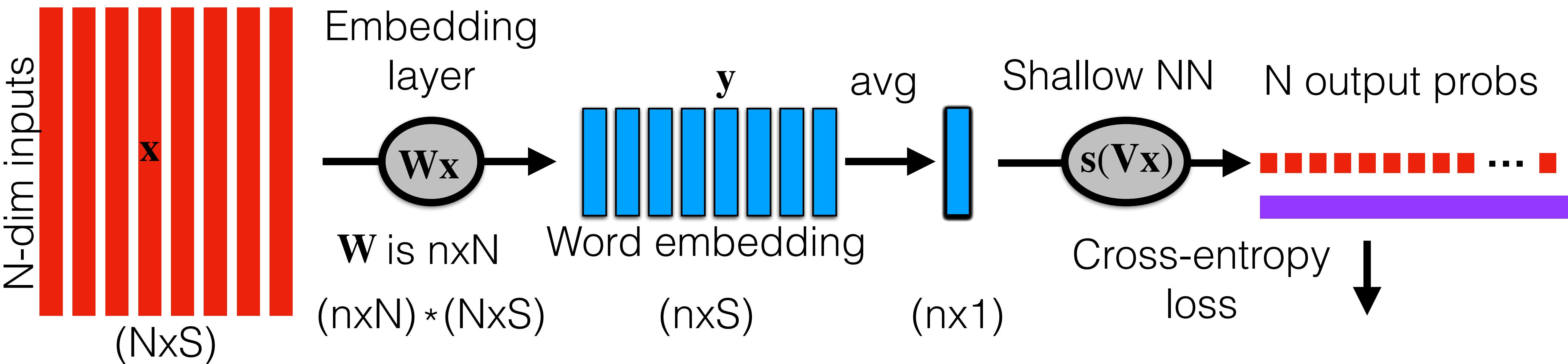
N-word vocabulary (one-hot enc.)



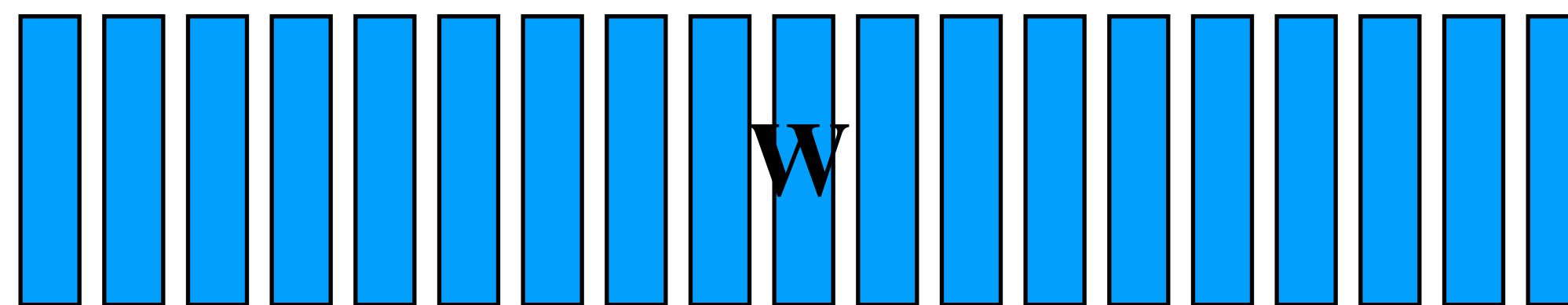
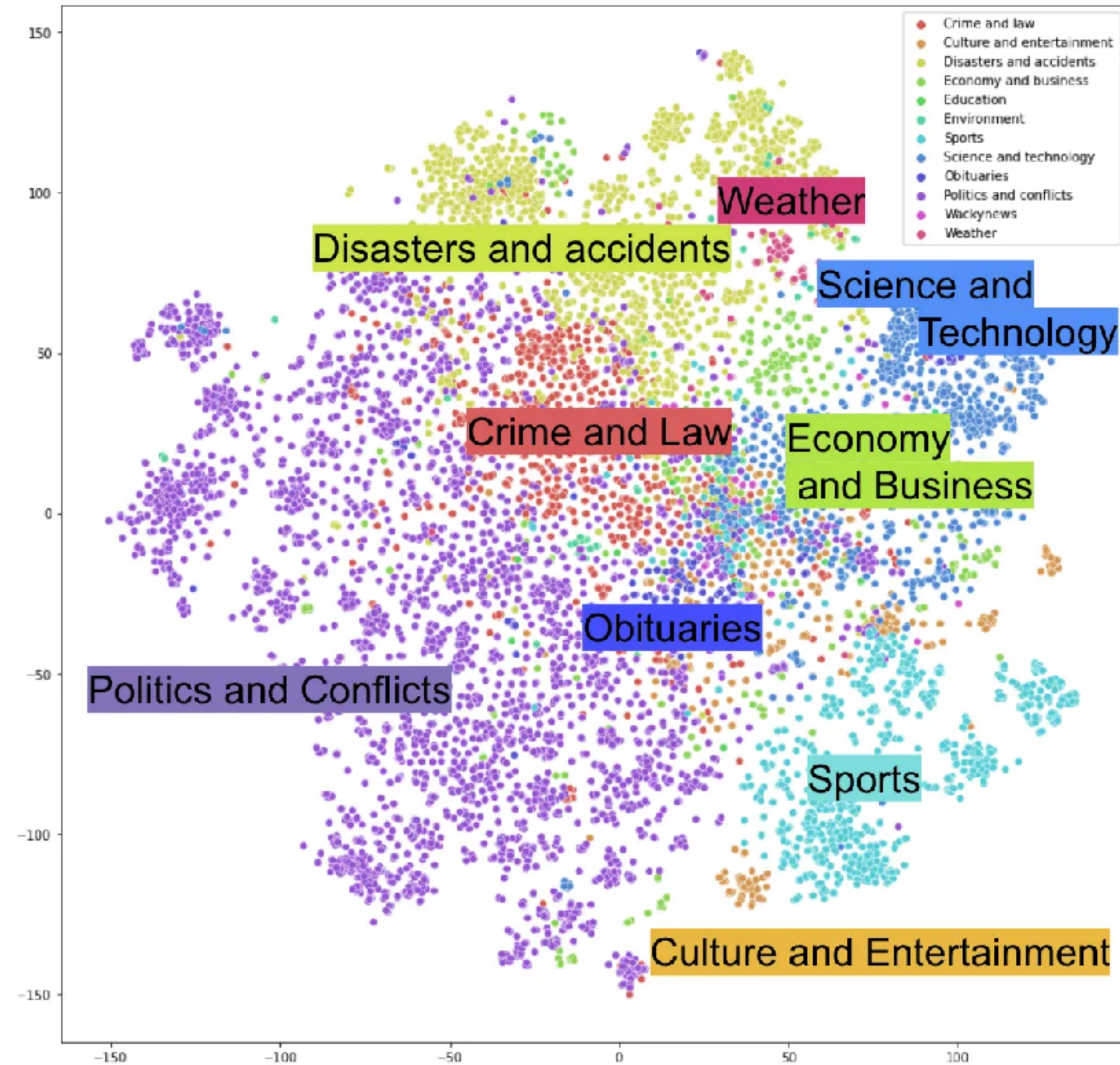
Karel is the best teacher in the whole world



S input words



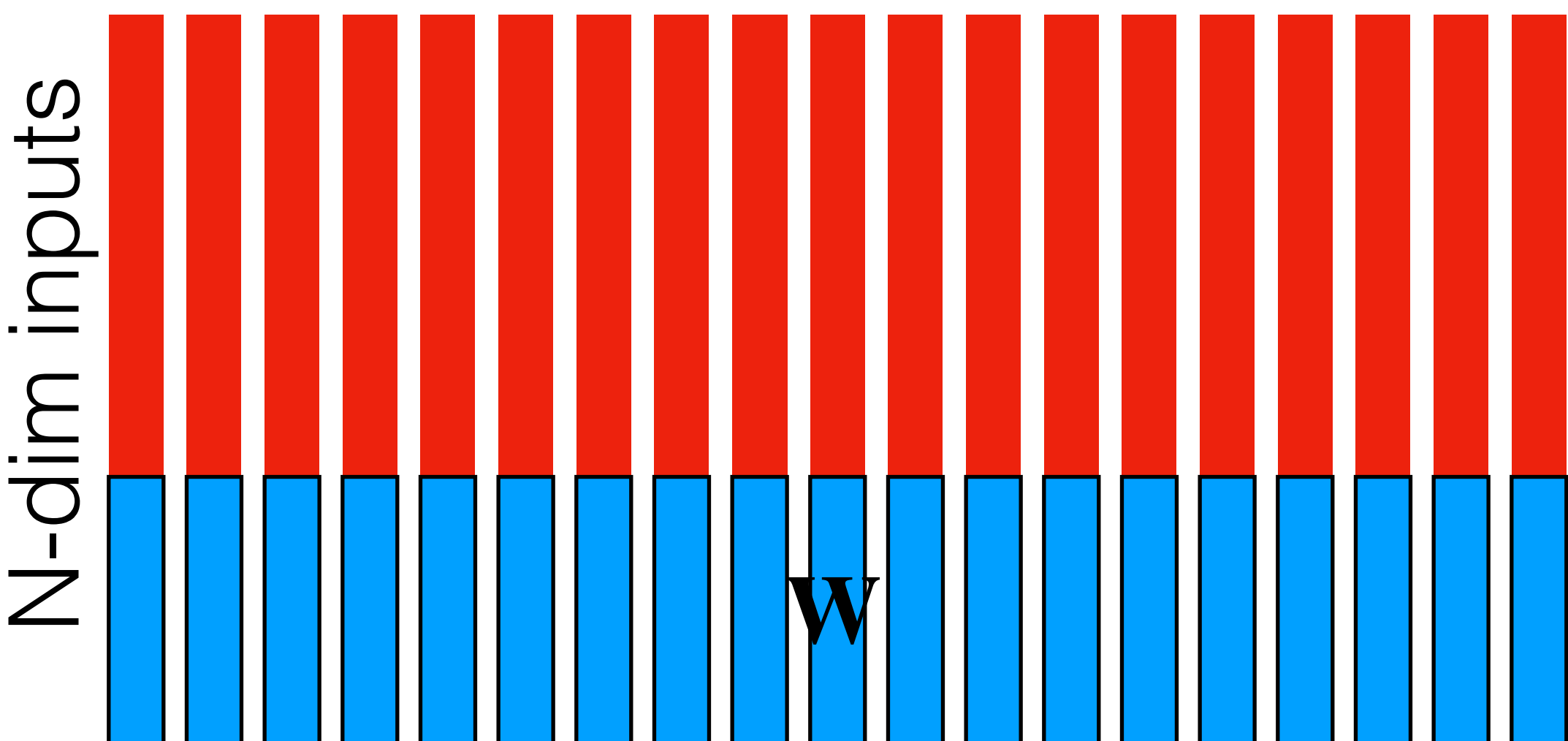
Word2vec represents words as low-dimensional continuous vectors



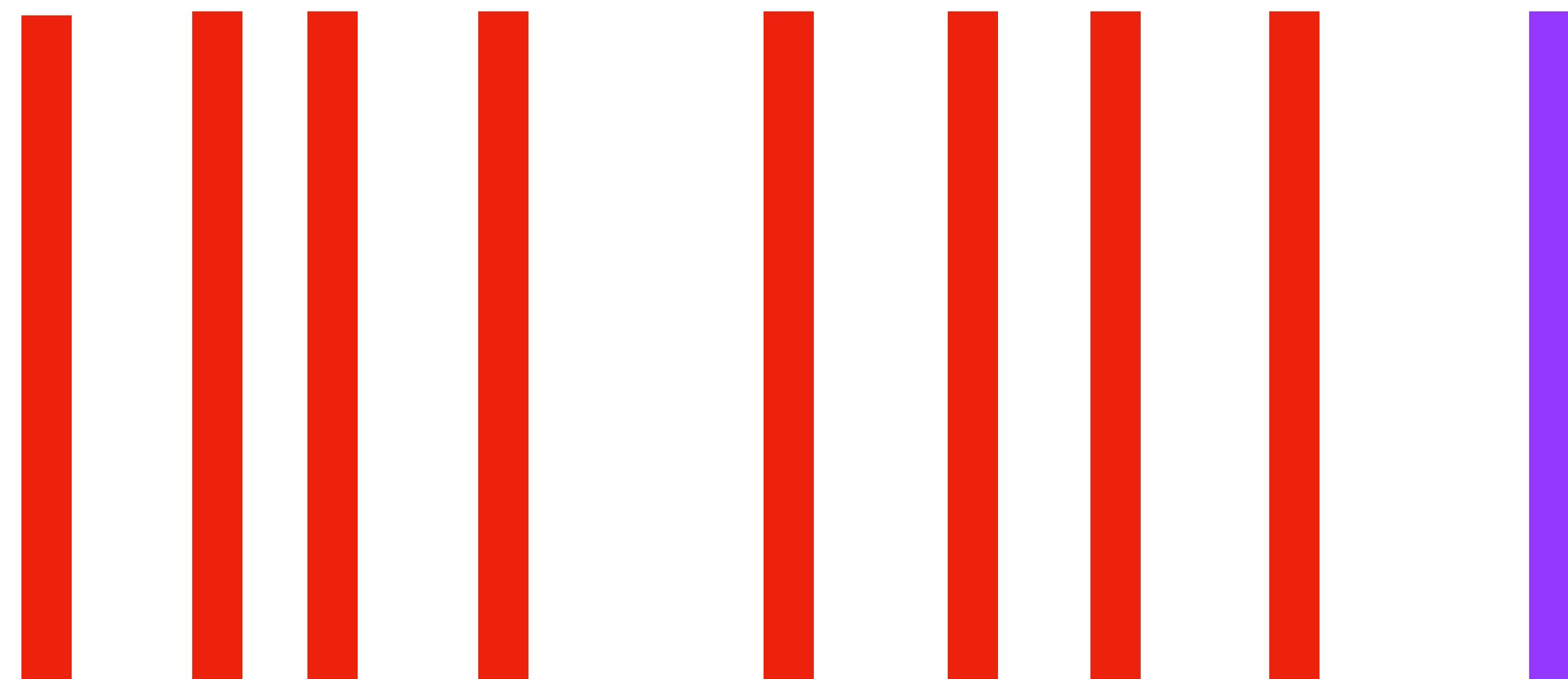
Word embedding

Word2vec represents words as low-dimensional continuous vectors

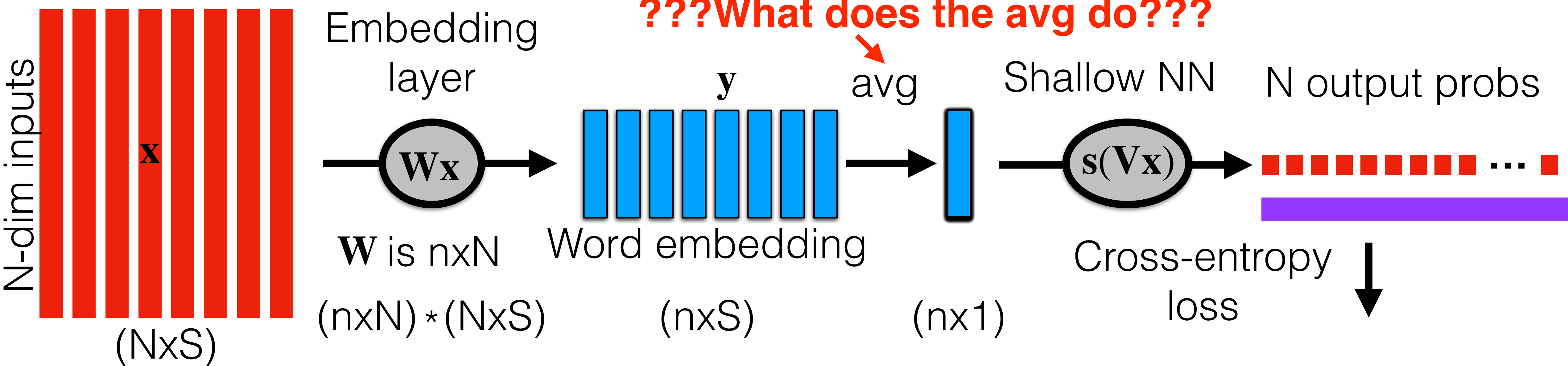
N-word vocabulary (one-hot enc.)



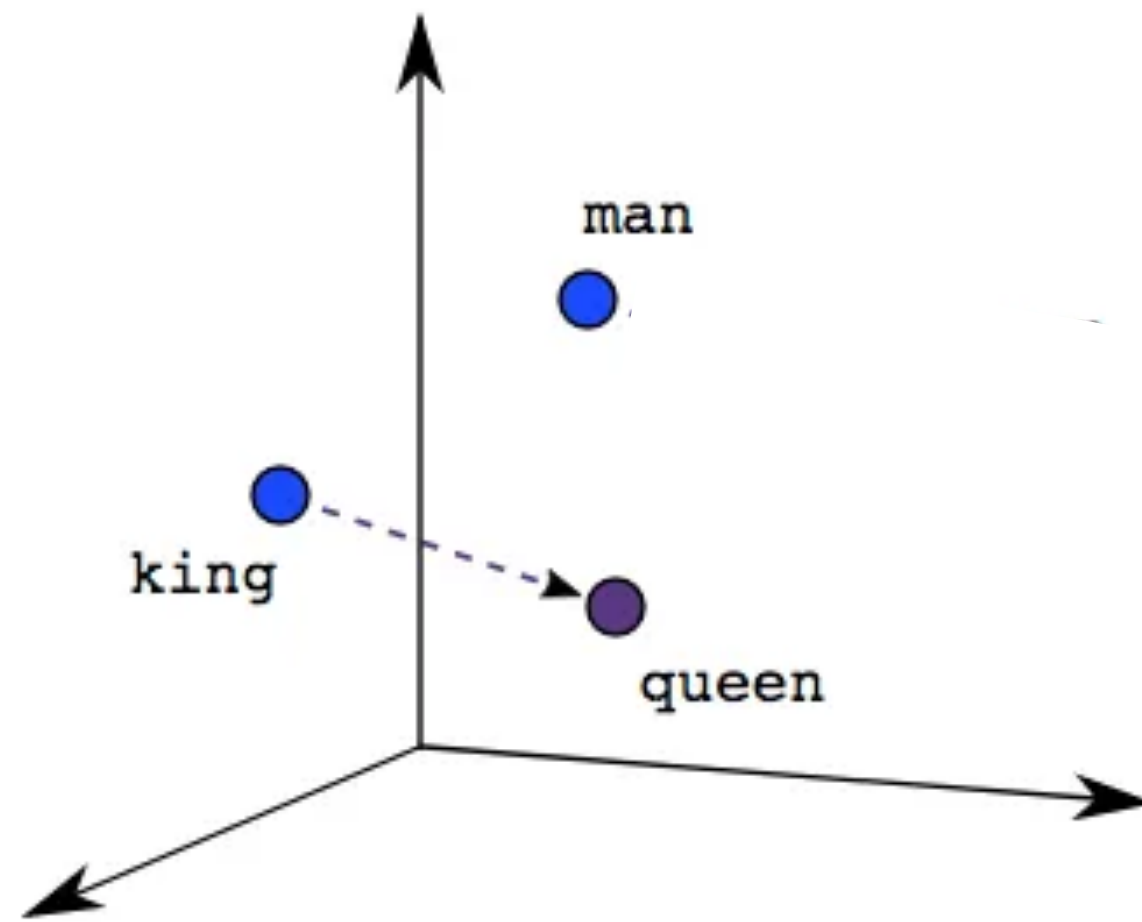
Karel is the best teacher in the whole world



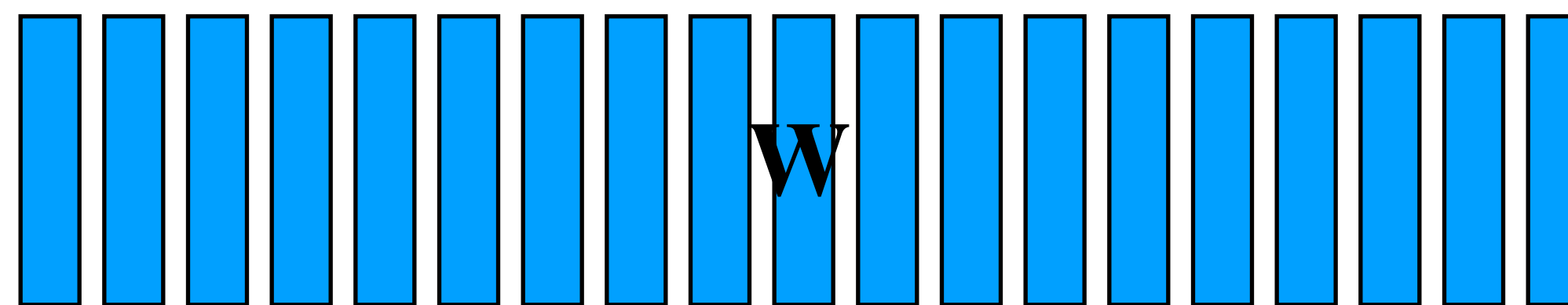
S input words



Word2vec represents words as low-dimensional continuous vectors

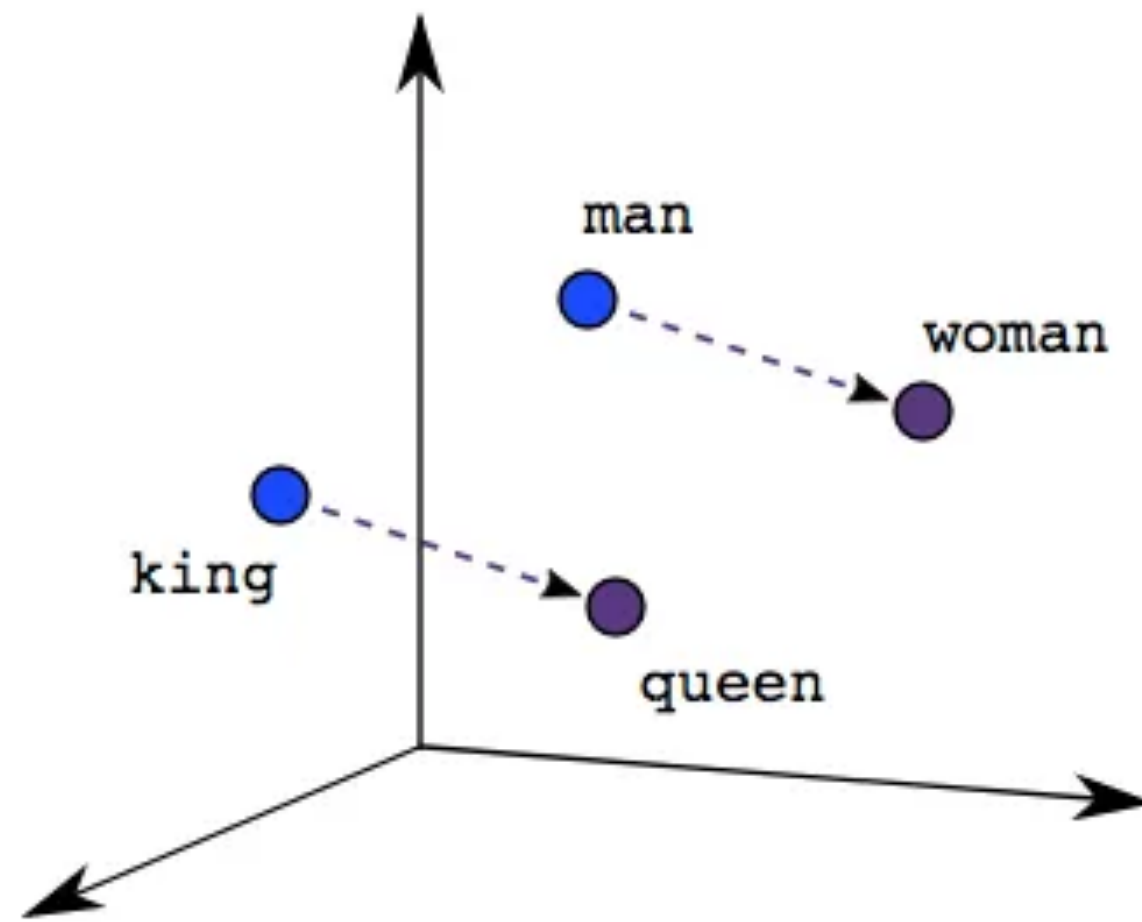


Male-Female



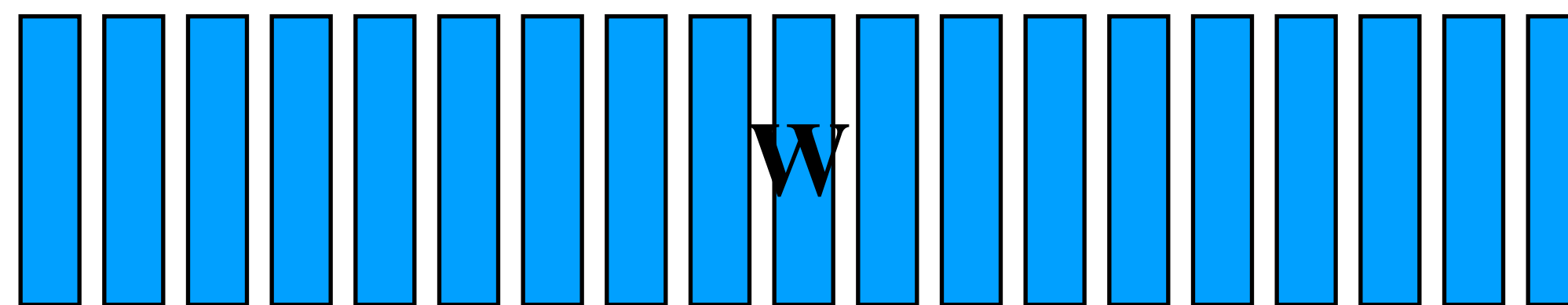
Word embedding

Word2vec represents words as low-dimensional continuous vectors



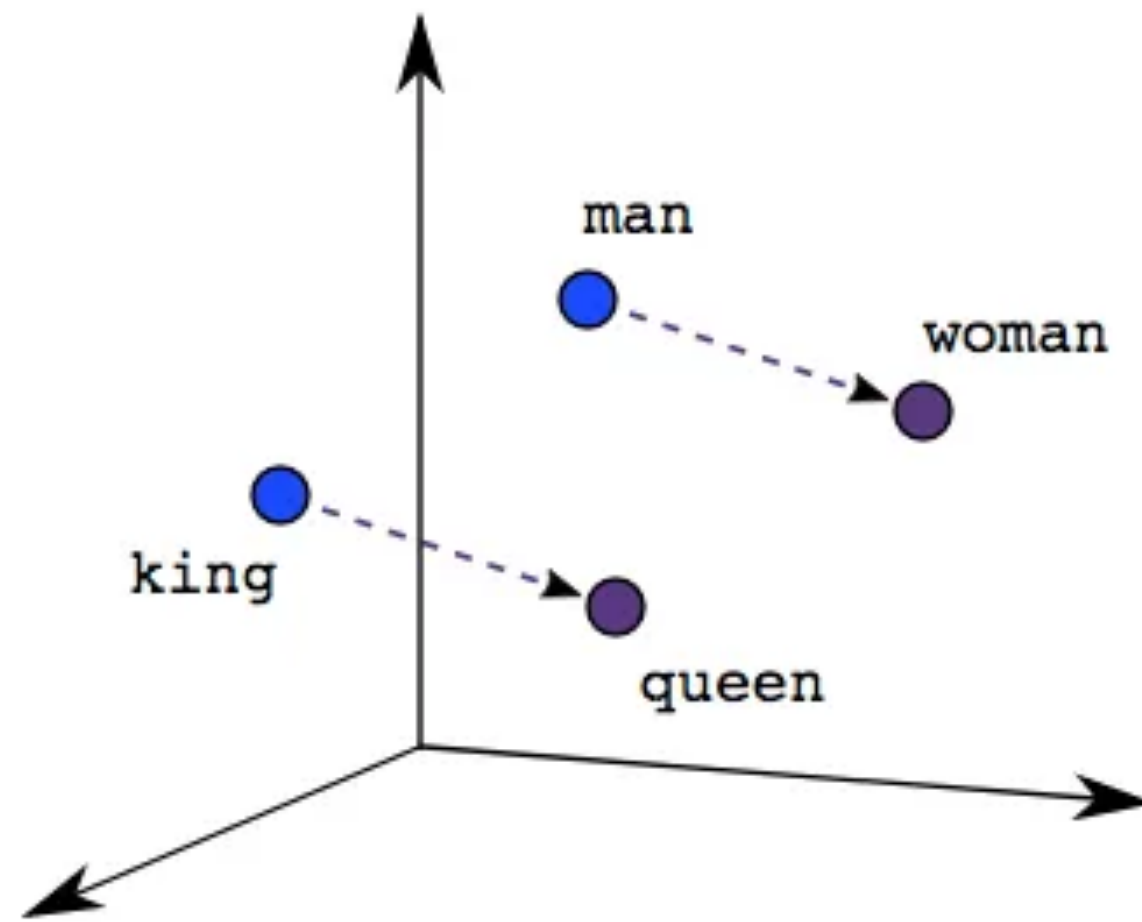
Male-Female

Word algebra: $\text{king} - \text{man} + \text{woman} = \text{queen}$

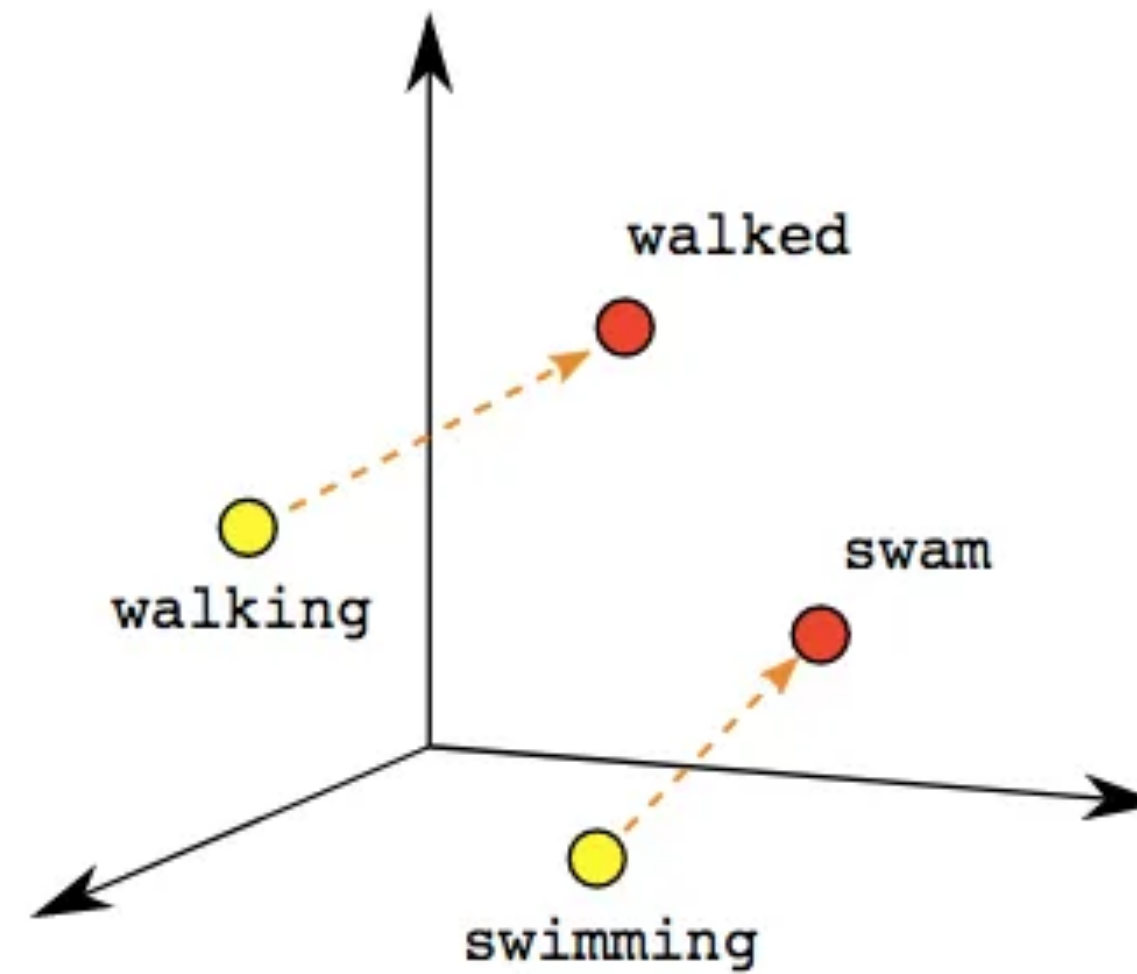


Word embedding

Word2vec represents words as low-dimensional continuous vectors

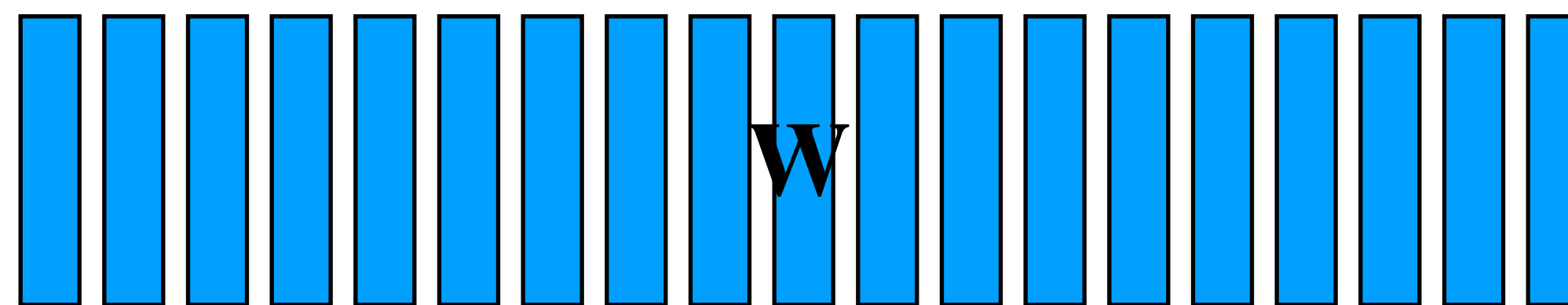


Male-Female



Verb Tense

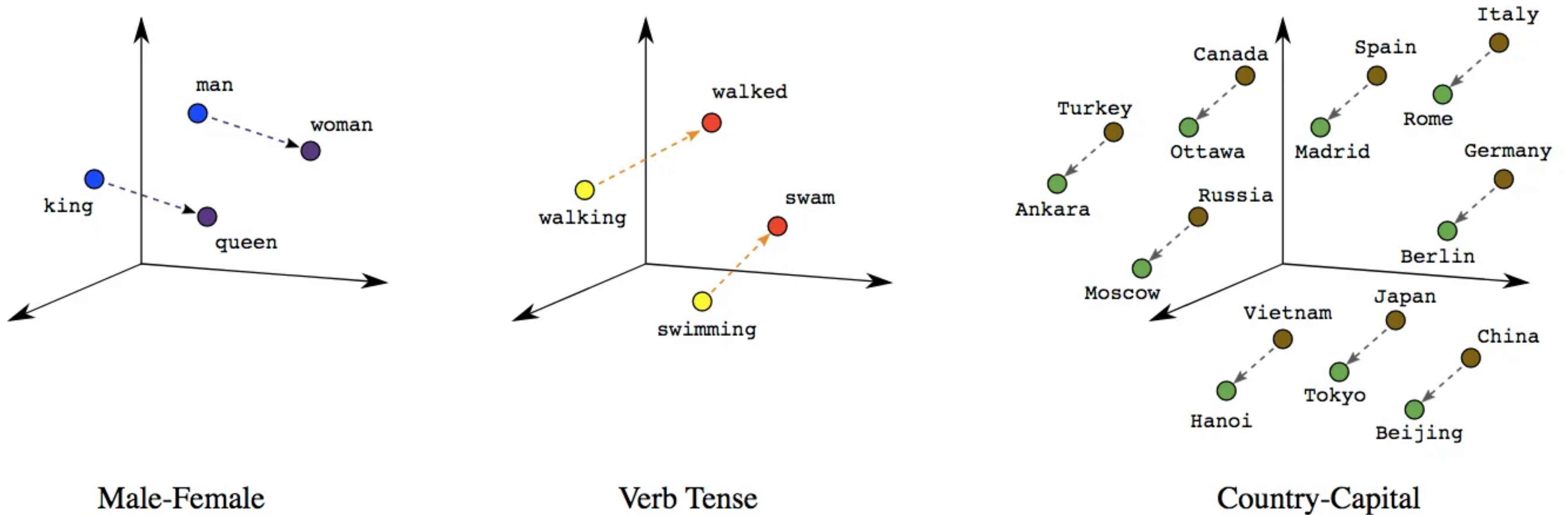
Word algebra: $\text{king} - \text{man} + \text{woman} = \text{queen}$



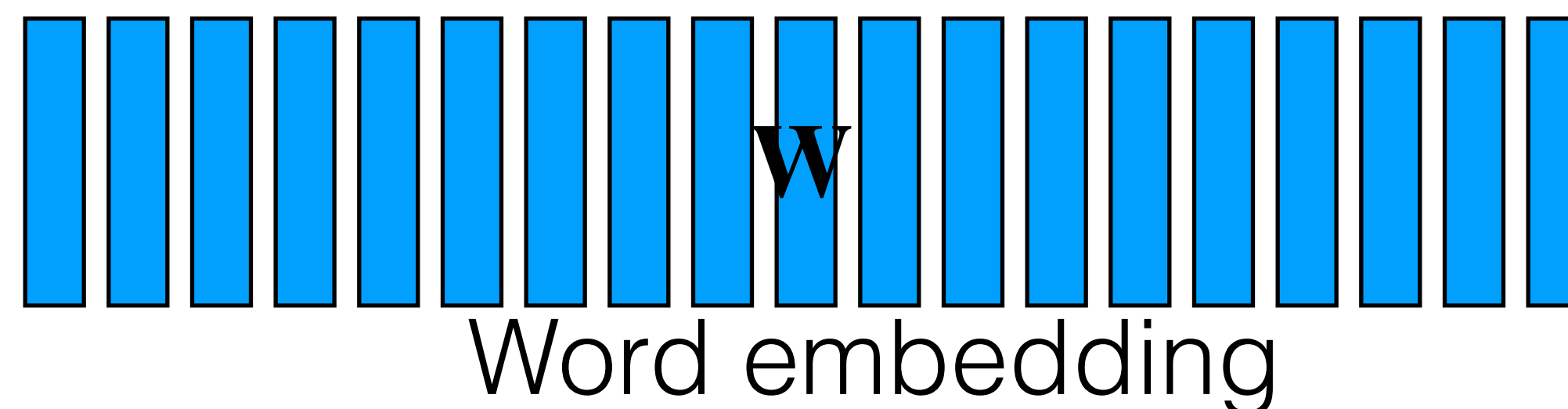
Word embedding

Word2vec represents words as low-dimensional continuous vectors

<https://dash.gallery/dash-word-arithmetic/>

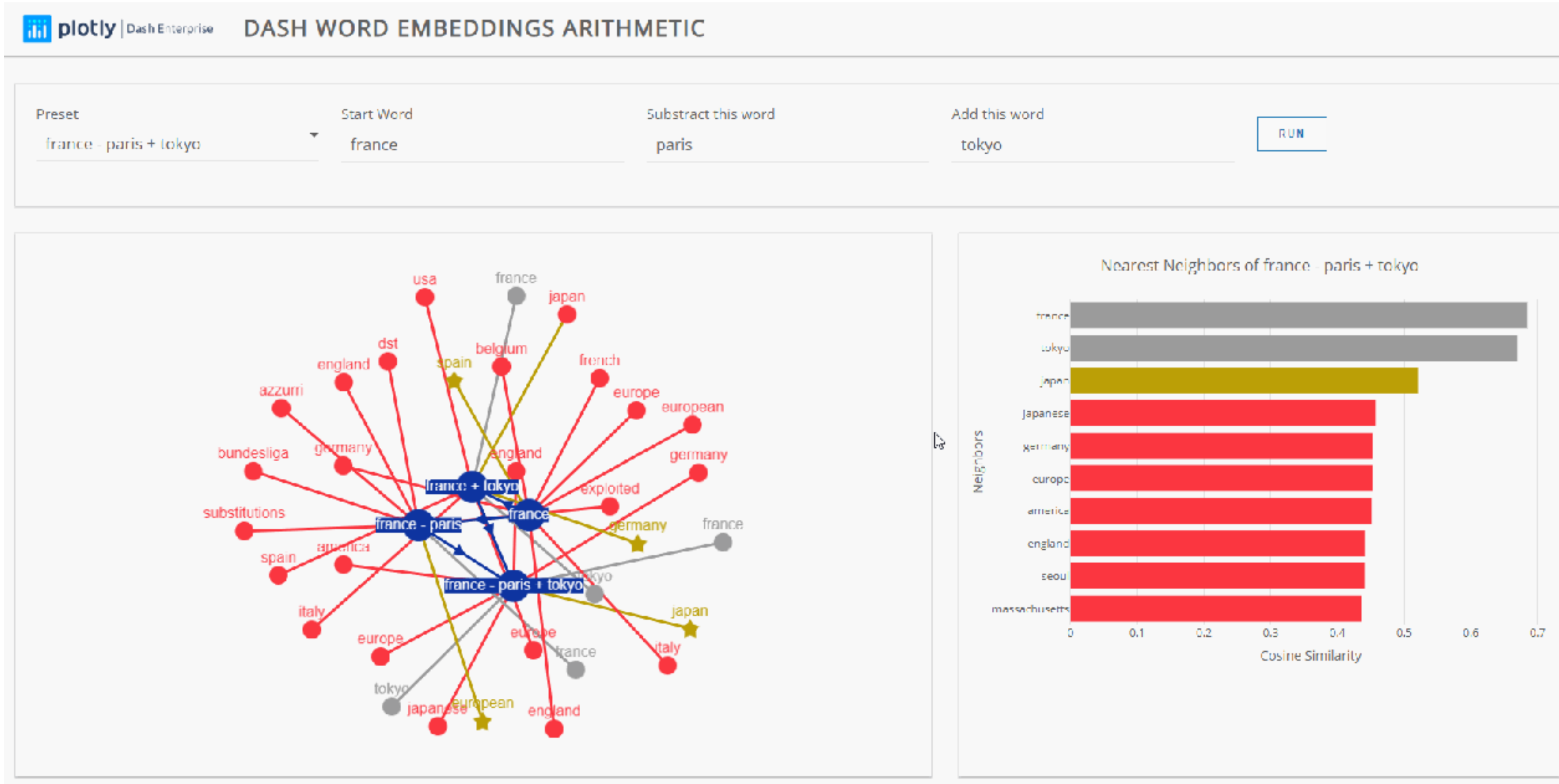


Word algebra: $\text{king} - \text{man} + \text{woman} = \text{queen}$

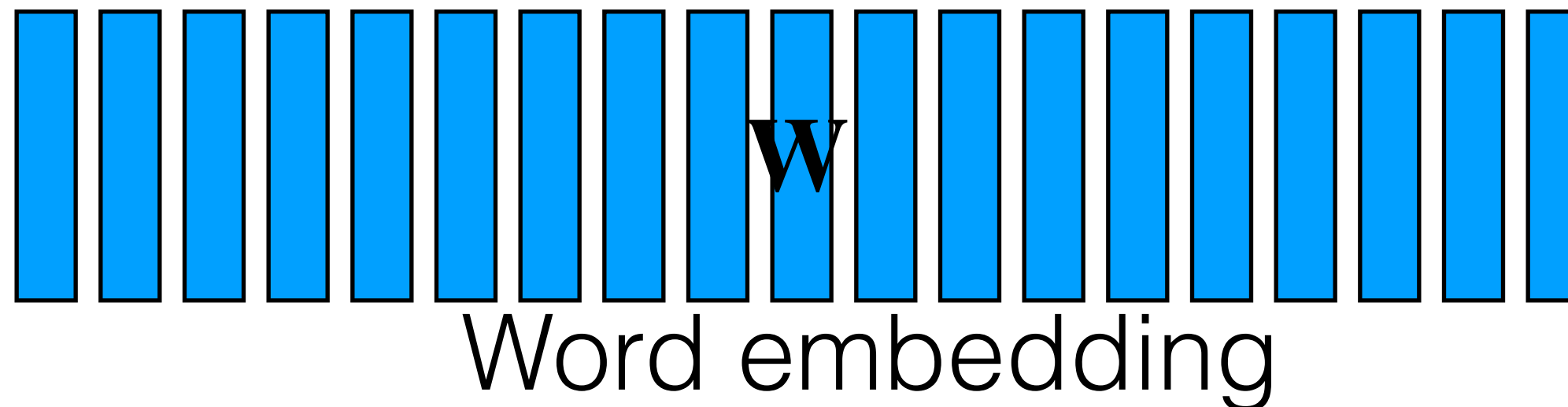


Word2vec represents words as low-dimensional continuous vectors

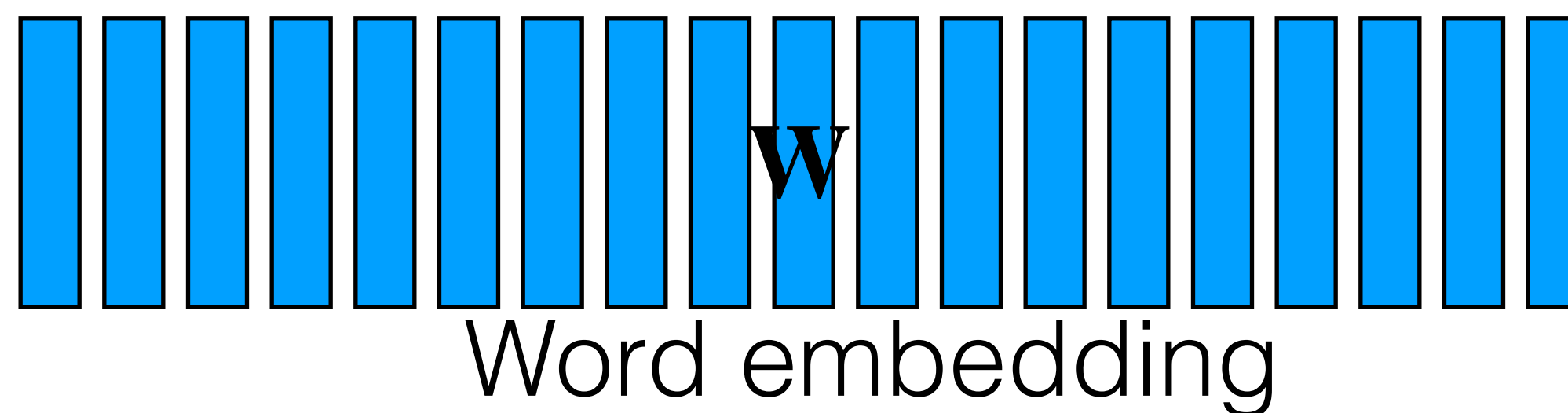
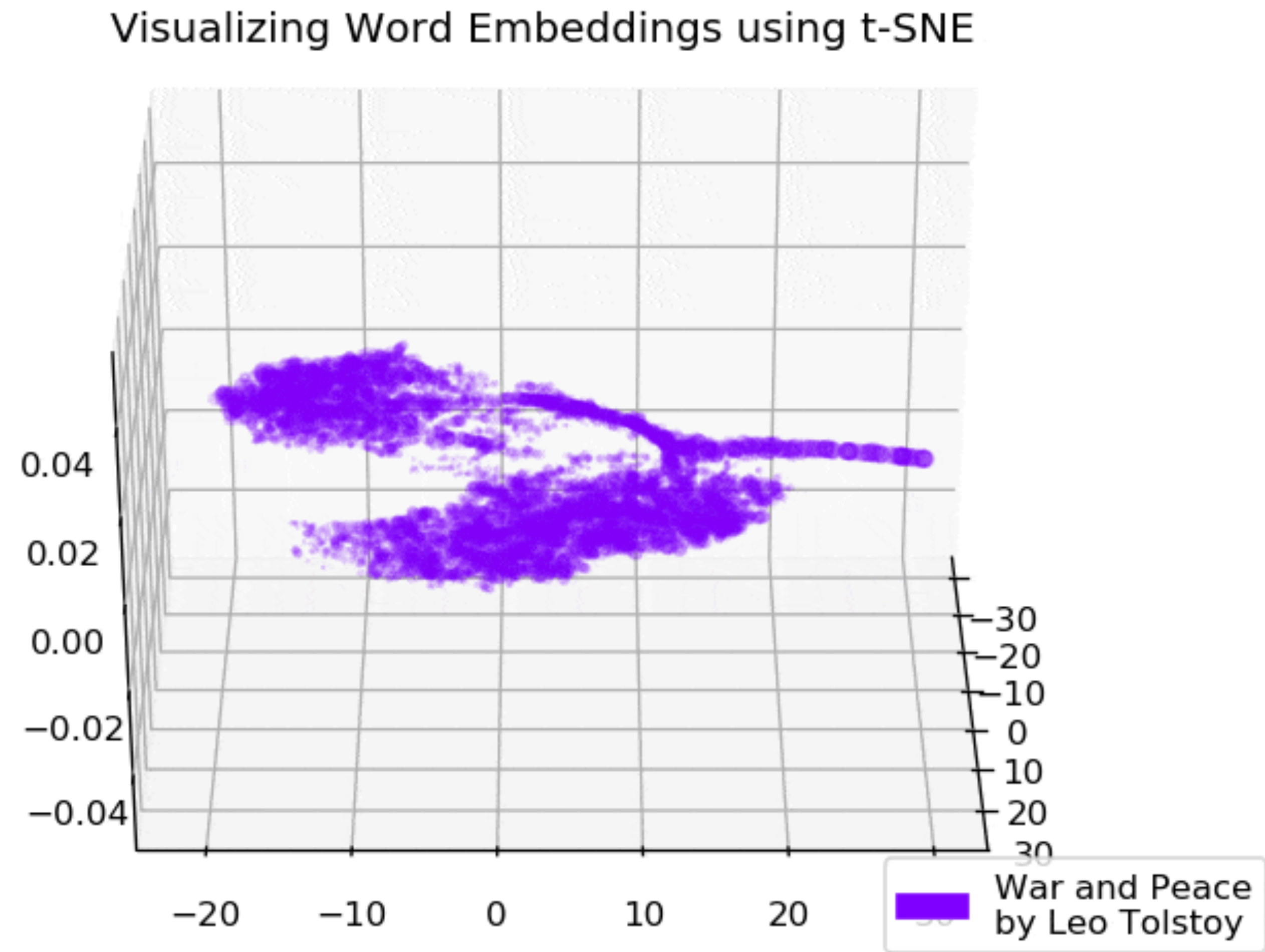
<https://dash.gallery/dash-word-arithmetic/>



Word algebra: king - man + woman = queen

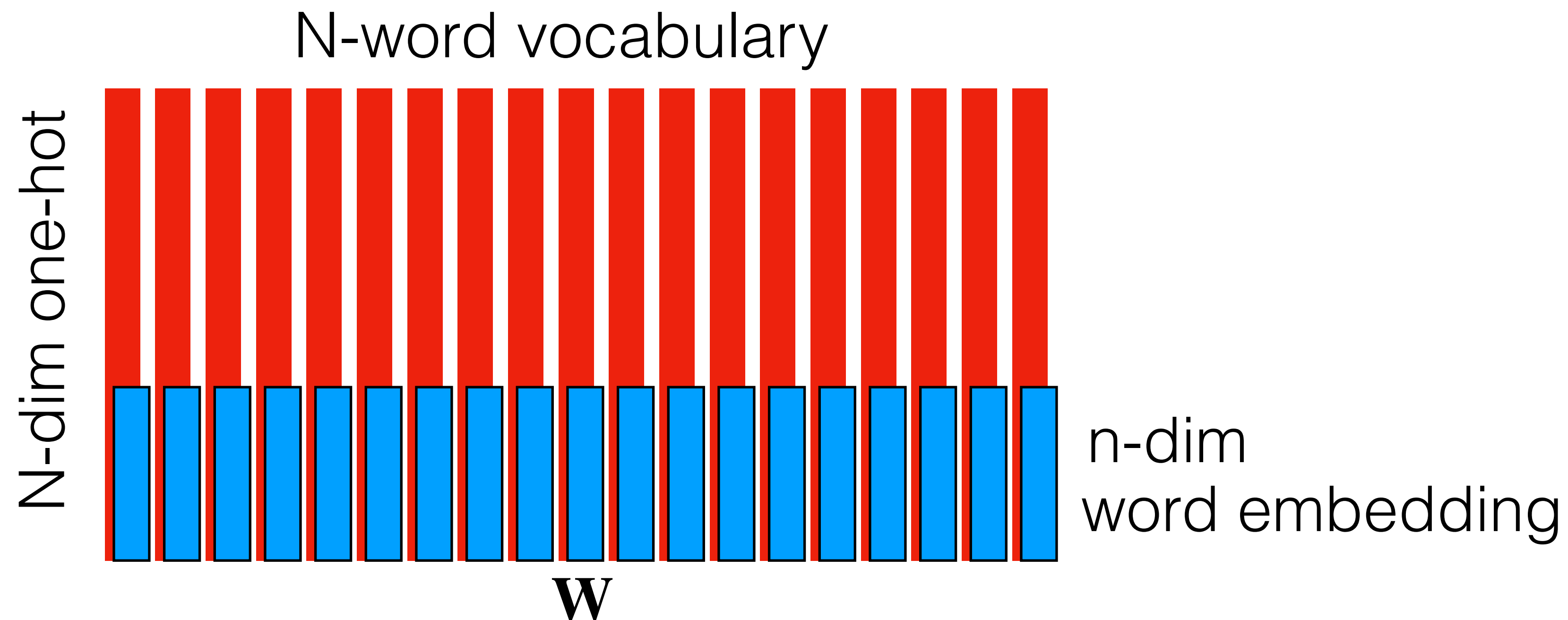


Word2vec represents words as low-dimensional continuous vectors



Word2vec represents words as low-dimensional continuous vectors

- **N N-dim orthonormal vectors** projected into **n-dim** space where **$n \ll N$** (large-scale models BERT, GPT has $n=768-1024$)
- How many orthogonal vectors in “n”-dim space?
- Can I represent only “n” independent concepts?

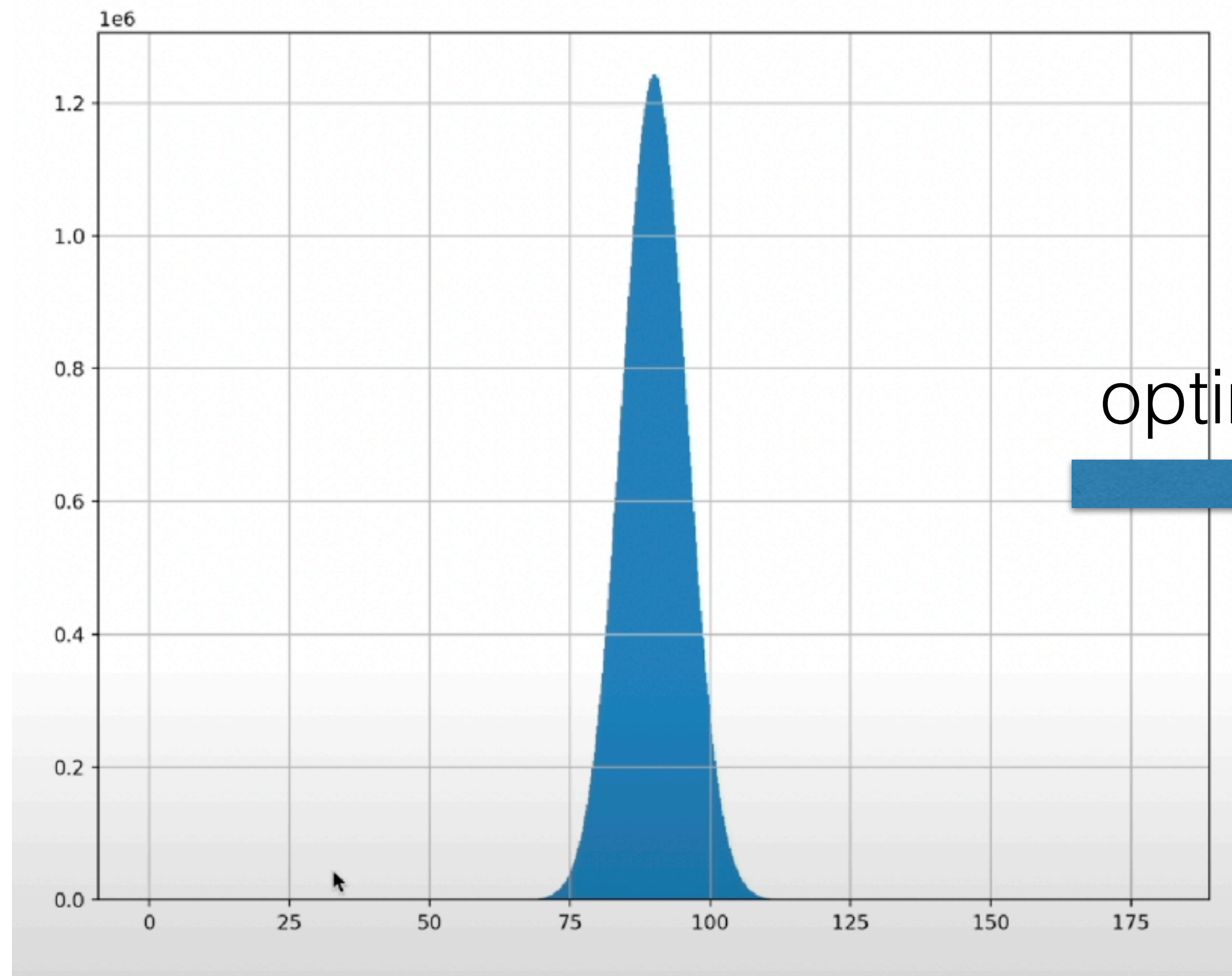


How many independent concepts do you fit in N-dimensional space?

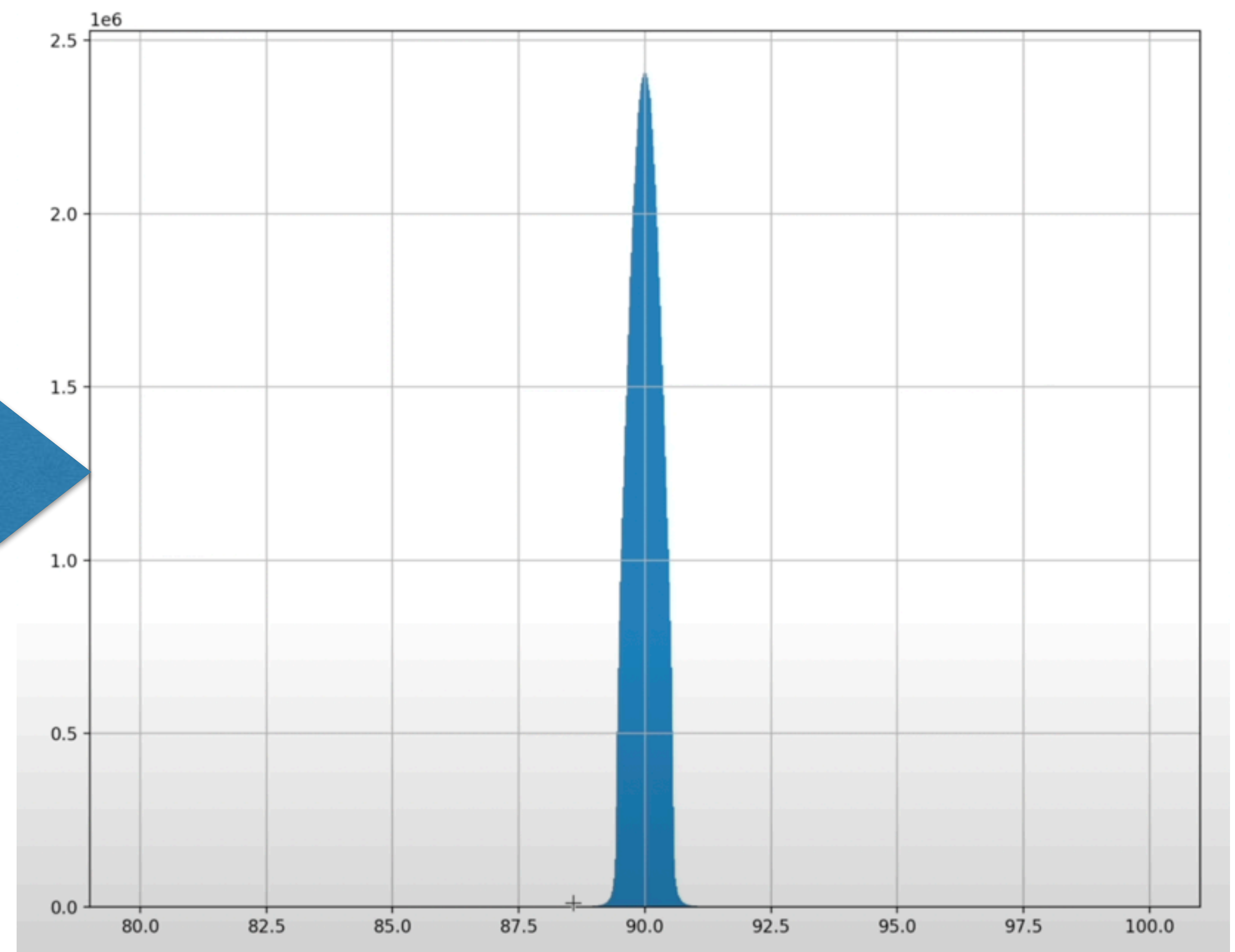
- How many orthogonal vectors in “n”-dim space?

Angles between randomly generated 10000 vectors in 100-dimensional space

Optimized vectors are almost independent

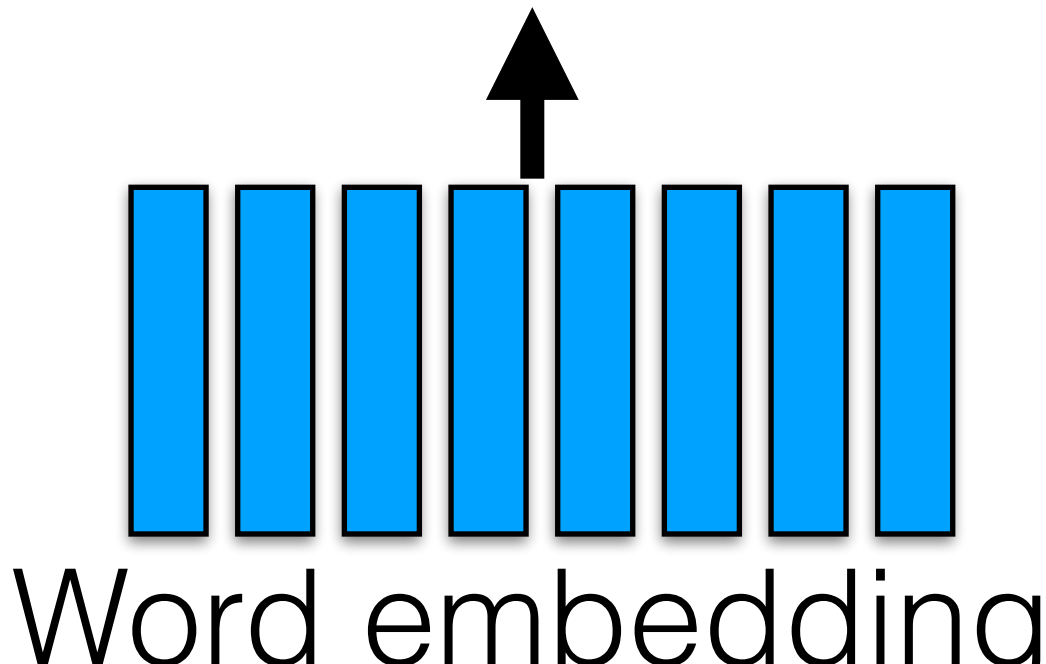


optimize

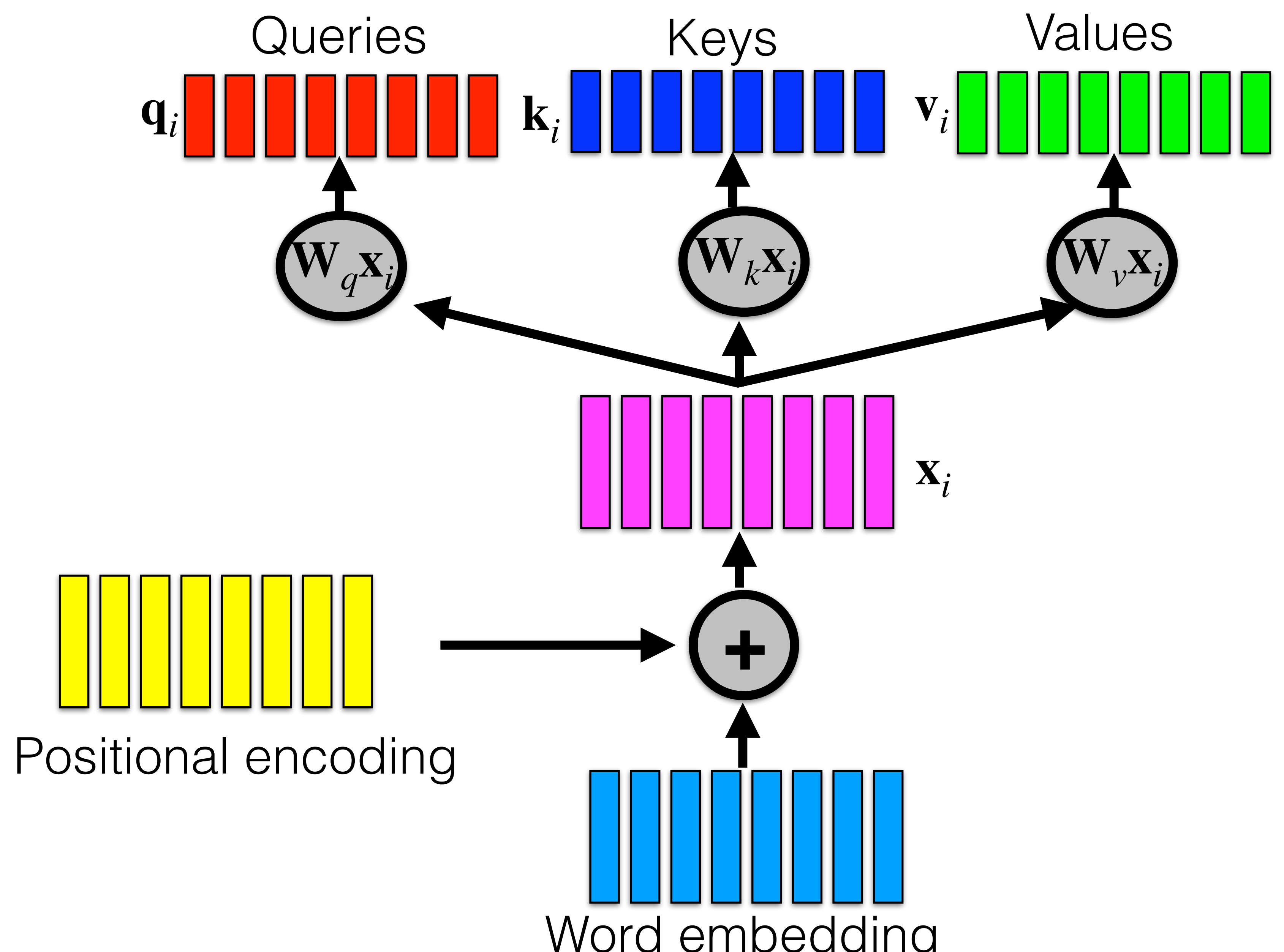


consequence of Johnson-Lindenstrauss Lemma

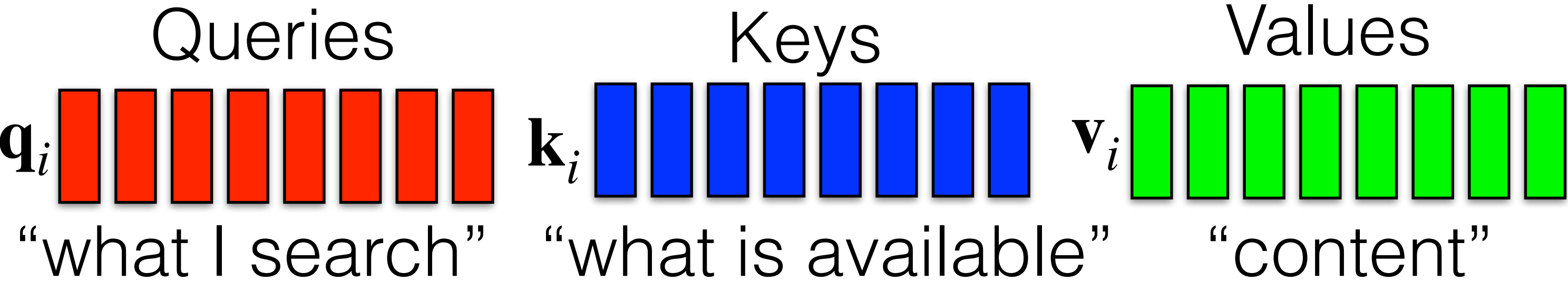
encoder



encoder

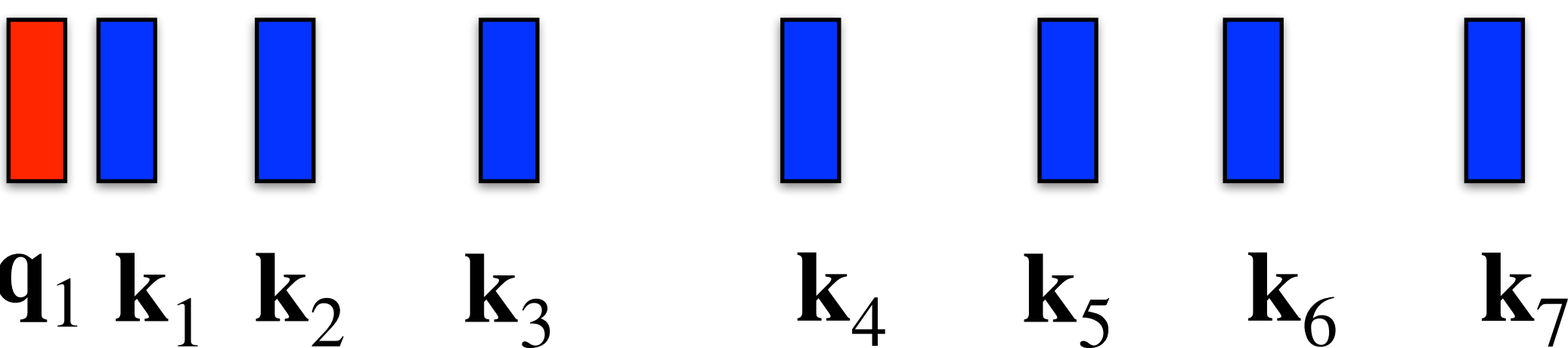





encoder



Karel is teacher and Mario is plumber.

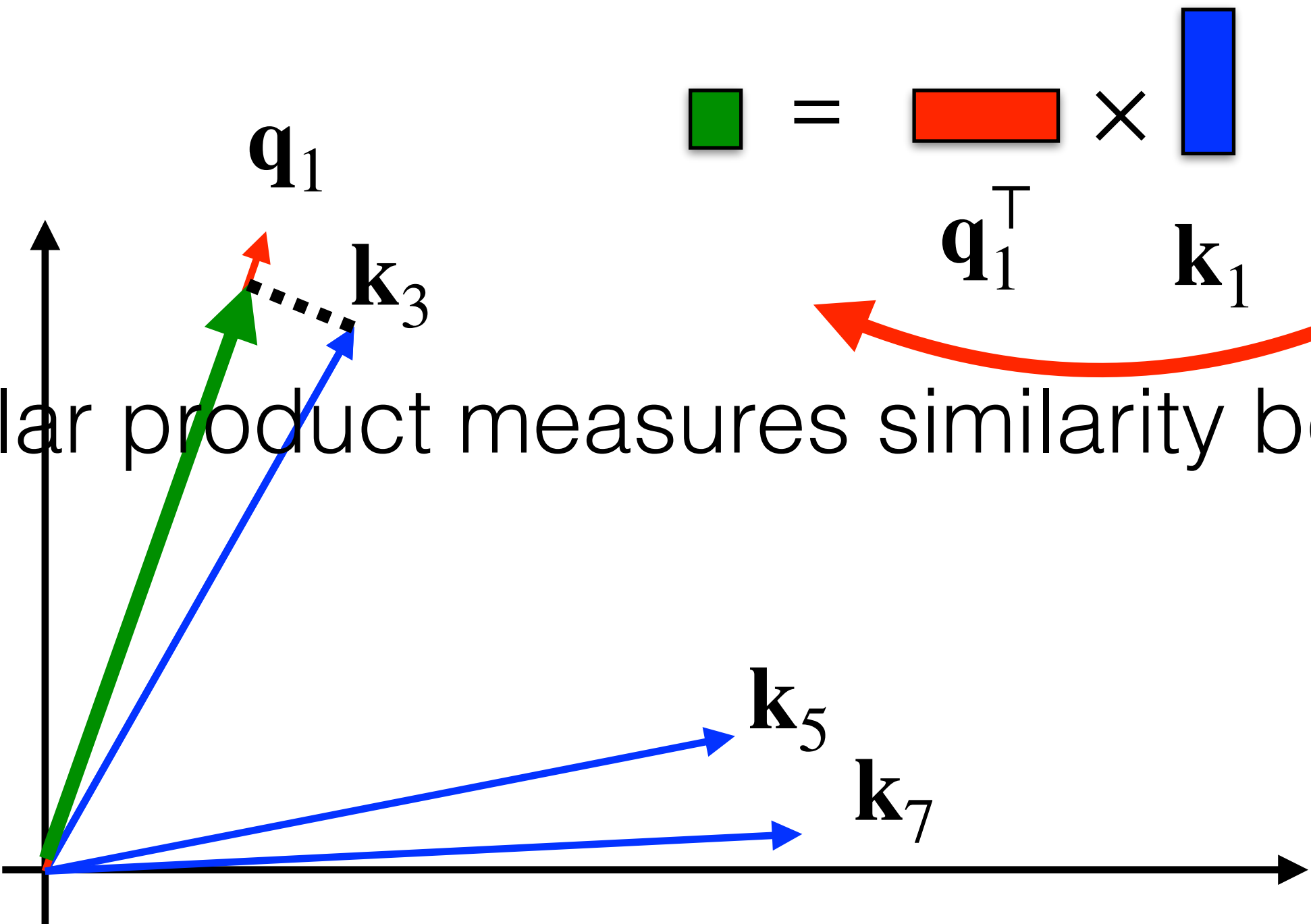
Which words contributes
to meaning of Karel?



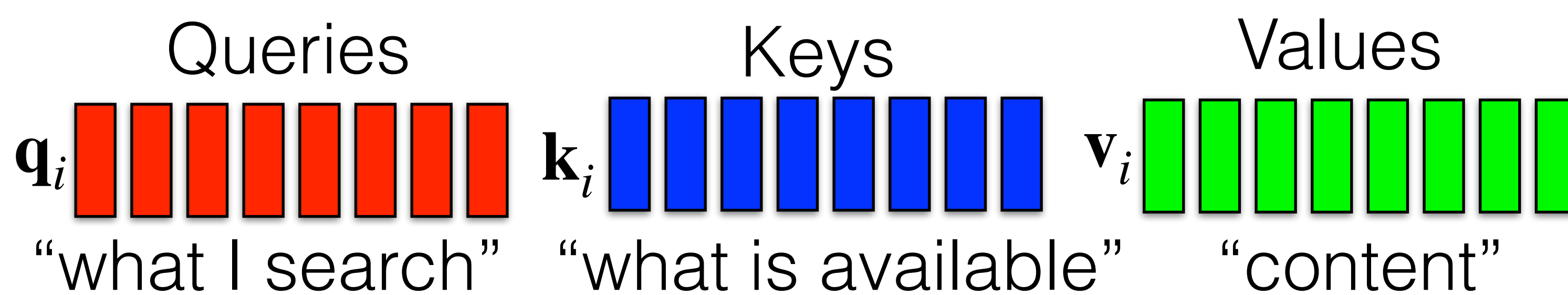
 =  \times 

\mathbf{q}_1^\top \mathbf{k}_1

Scalar product measures similarity between vectors.



encoder

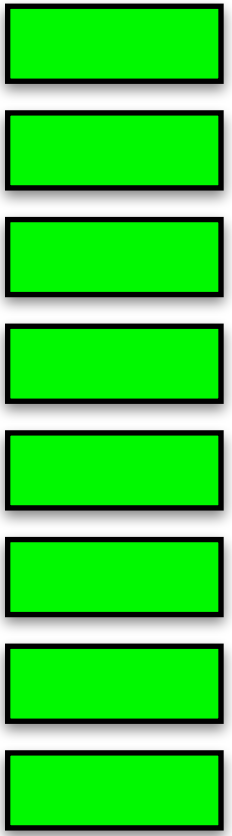


Scalar product measures similarity between vectors.

Get attention weights

$$s \left(\begin{array}{c} \text{red bar} \\ \mathbf{q}_i^\top \end{array} \times \begin{array}{c} \text{blue bars} \\ \mathbf{k}_1 \quad \dots \quad \mathbf{k}_n \end{array} \right) = \begin{array}{c} \text{orange bars} \\ \mathbf{a}_1 \quad \dots \quad \mathbf{a}_n \end{array}$$

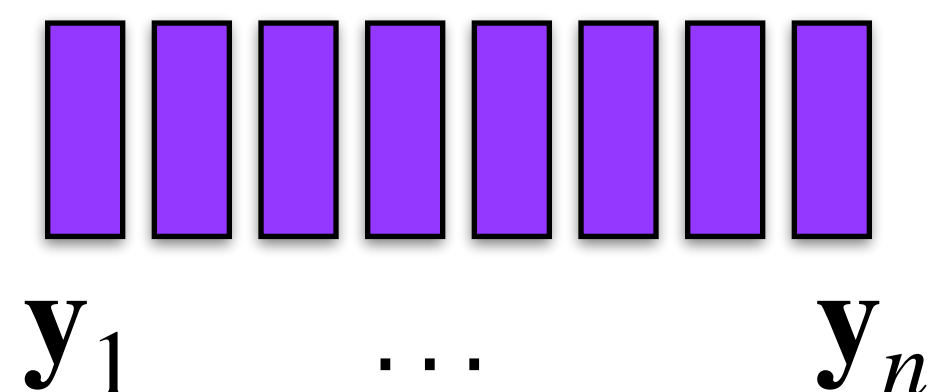
Attention-weighted sum of values

$\mathbf{a}_1 \quad \dots \quad \mathbf{a}_n$  $\mathbf{v}_1 \dots \mathbf{v}_n$

I am adding “teacher vector” to “Karel vector”

$$\mathbf{a}_1 \mathbf{v}_1 + \mathbf{a}_2 \mathbf{v}_2 + \dots + \mathbf{a}_n \mathbf{v}_n = \mathbf{y}_1$$

Outputs:



Avoid for-loop by smart matrix multiplication

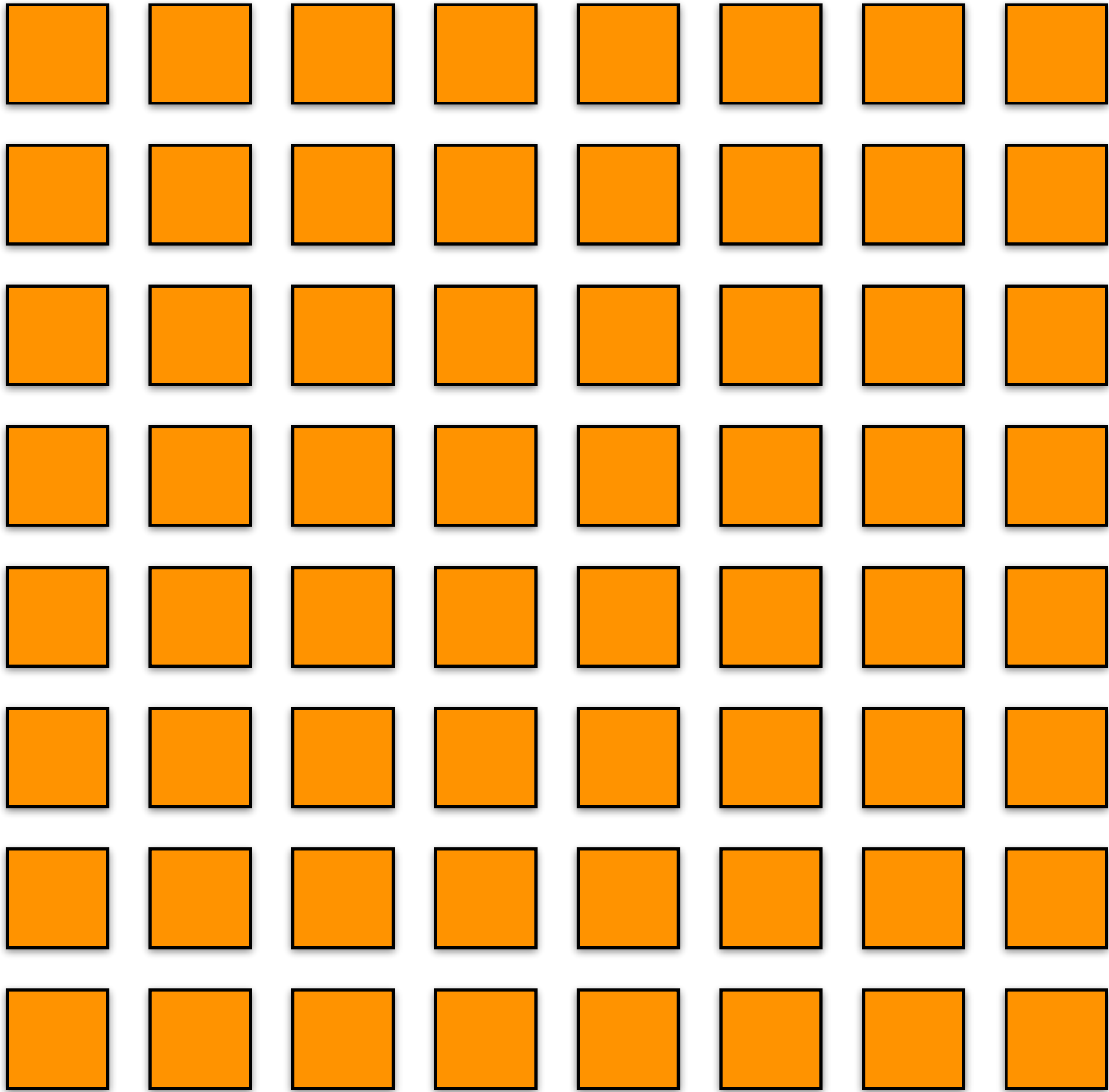
The diagram illustrates the triplet loss components. It shows three horizontal rows of colored rectangles representing embeddings. The first row, labeled 'Queries' above and \mathbf{q}_i to the left, contains eight red rectangles. Below it is the text '“what I search”'. The second row, labeled 'Keys' above and \mathbf{k}_i to the left, contains eight blue rectangles. Below it is the text '“video captions”'. The third row, labeled 'Values' above and \mathbf{v}_i to the left, contains eight green rectangles. Below it is the text '“video files”'.

The diagram illustrates the matrix multiplication $Q^T \times K = A$. Matrix Q^T is a 10x3 matrix of red rectangles. Matrix K is a 3x8 matrix of blue rectangles. Matrix A is a 10x8 matrix of orange rectangles, representing the result of the multiplication.

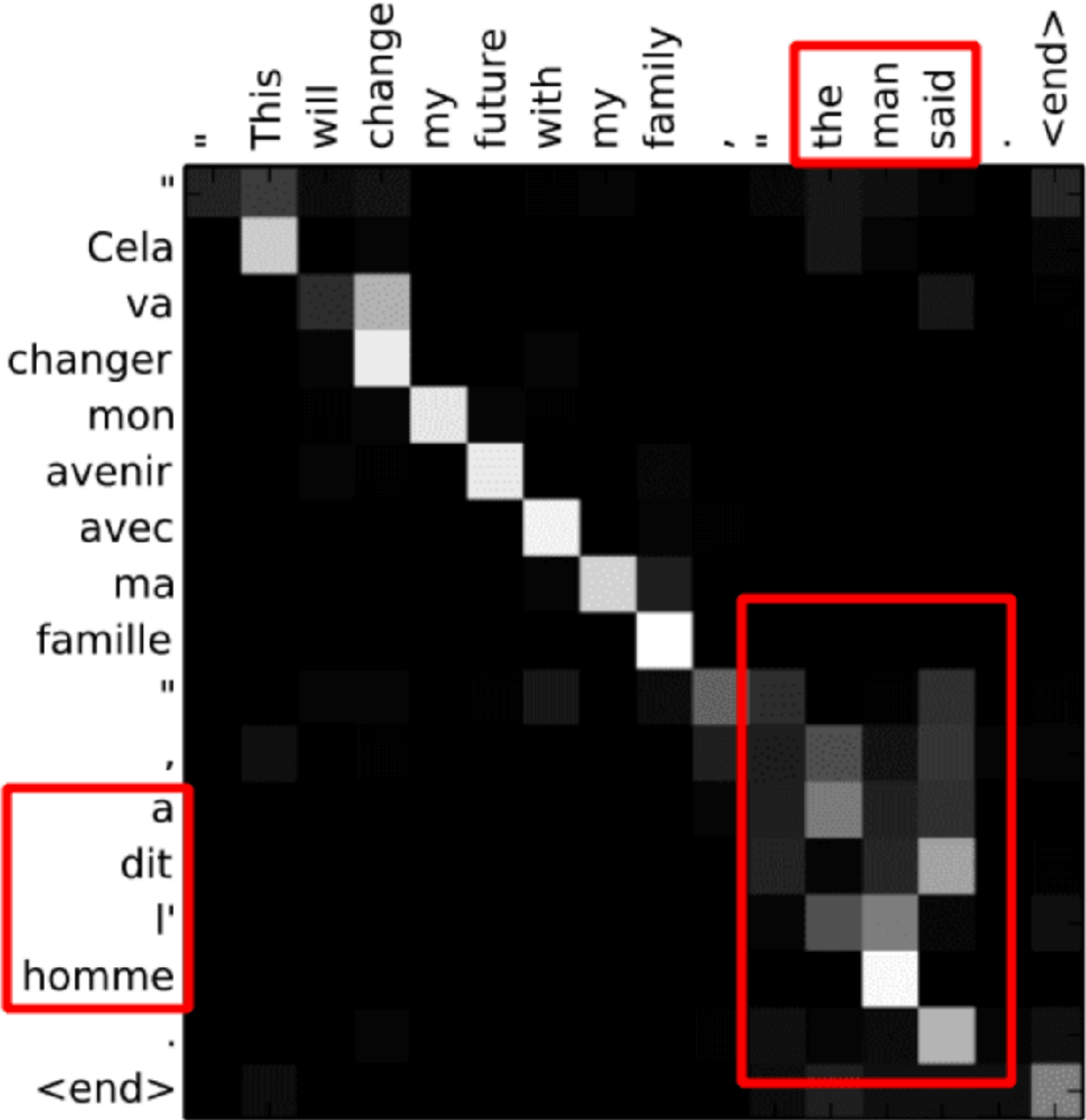
Diagram illustrating the matrix multiplication $\mathbf{A} \times \mathbf{V}^T = \mathbf{s}(\mathbf{Q}^T \mathbf{K}) \mathbf{V}^T = \mathbf{Y}$.

- \mathbf{A} : An 8x8 matrix represented by orange squares.
- \mathbf{V}^T : A column vector represented by 8 green rectangles.
- $\mathbf{s}(\mathbf{Q}^T \mathbf{K}) \mathbf{V}^T$: The result of the multiplication, shown as a row of 8 purple rectangles.
- \mathbf{Y} : The final output, labeled "Output", represented by 8 purple rectangles.

encoder



=



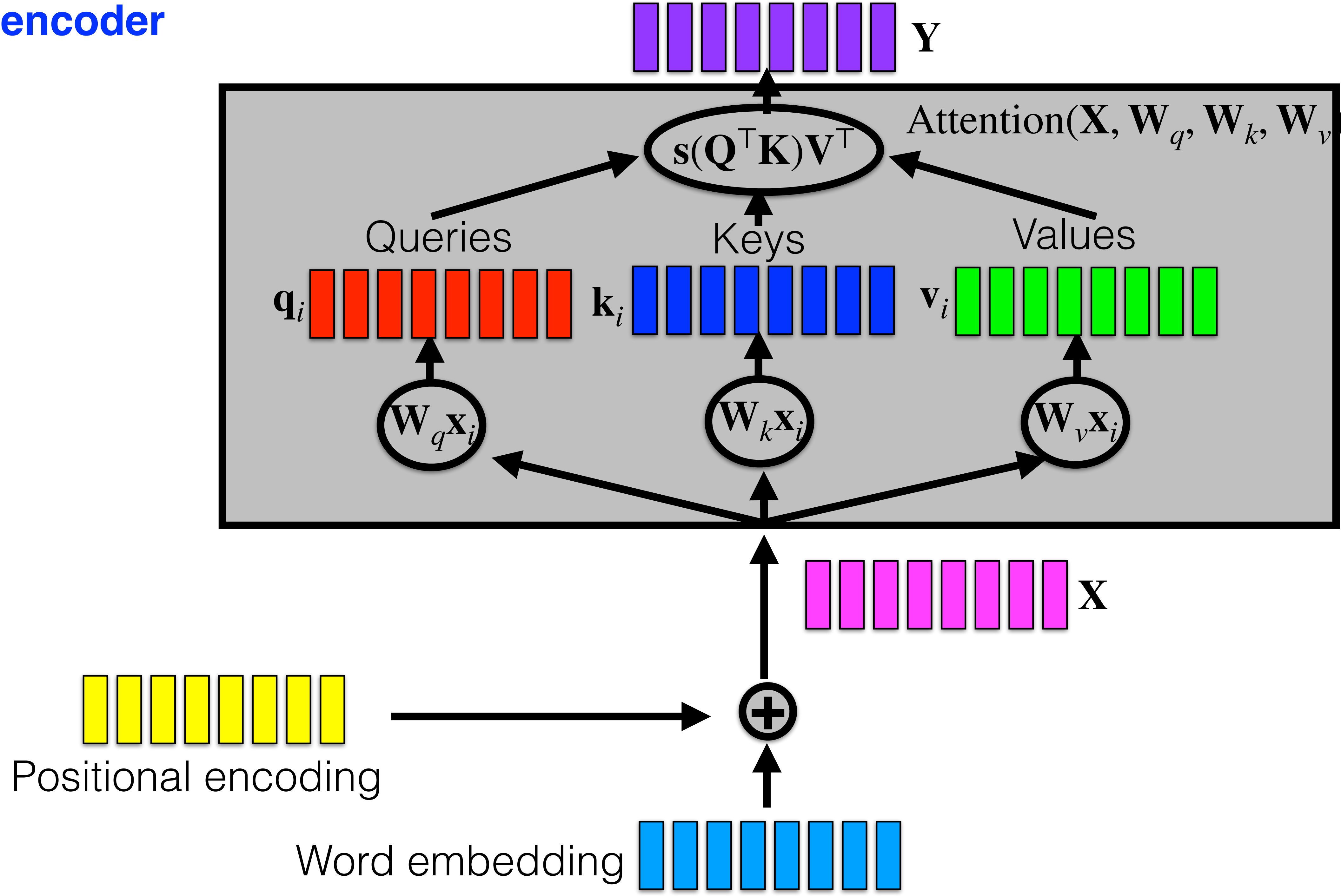
The diagram illustrates the triplet loss components. It shows three horizontal rows of colored rectangles representing embeddings. The first row, labeled 'Queries' above and \mathbf{q}_i to the left, consists of eight red rectangles with the text 'what I search' below. The second row, labeled 'Keys' above and \mathbf{k}_i to the left, consists of eight blue rectangles with the text 'video captions' below. The third row, labeled 'Values' above and \mathbf{v}_i to the left, consists of eight green rectangles with the text 'video files' below.

The diagram illustrates the matrix multiplication $Q^T \times K = S$. On the left, Q^T is represented by a vertical stack of 10 red rectangles, and K is represented by a horizontal row of 8 blue rectangles. These are enclosed in large parentheses with a multiplication symbol between them. An equals sign follows, leading to the result matrix S , which is a grid of 10 rows and 8 columns of orange rectangles.

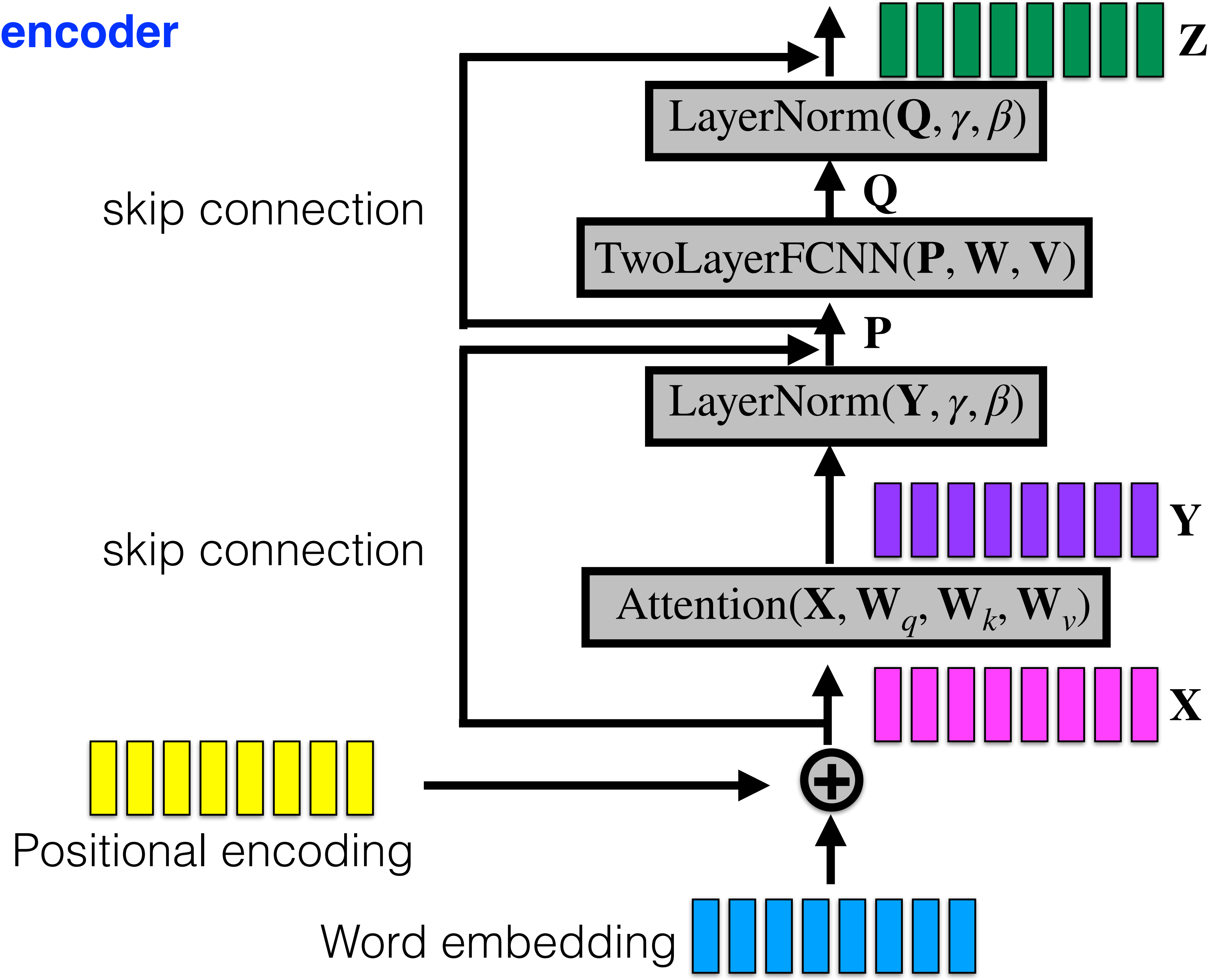
Diagram illustrating the matrix multiplication $\mathbf{A} \times \mathbf{V}^T = \mathbf{s}(\mathbf{Q}^T \mathbf{K}) \mathbf{V}^T = \mathbf{Y}$.

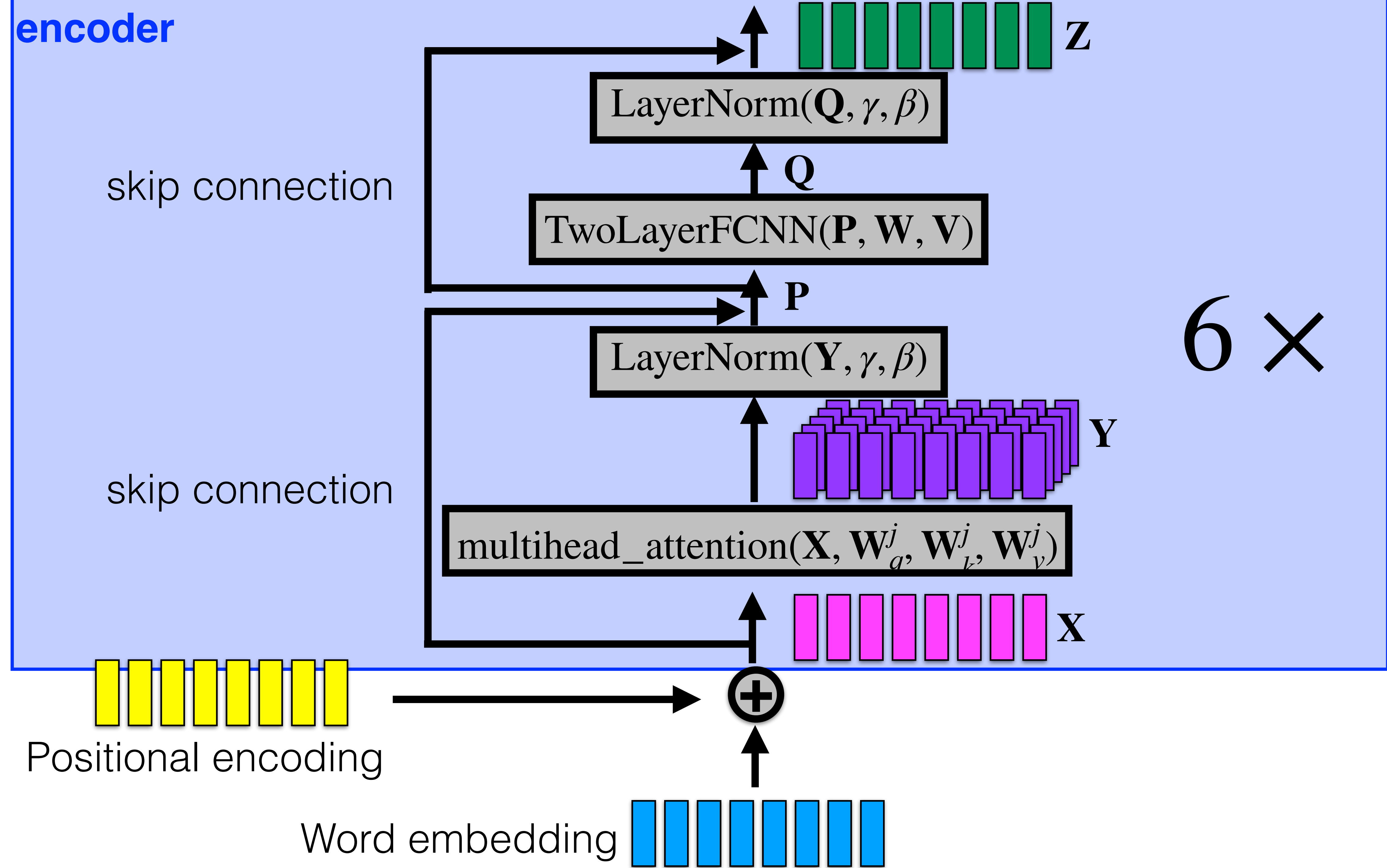
- \mathbf{A} : An 8x8 matrix represented by orange squares.
- \mathbf{V}^T : A vertical vector represented by 8 green rectangles.
- $\mathbf{s}(\mathbf{Q}^T \mathbf{K}) \mathbf{V}^T$: The result of the multiplication, shown as a horizontal vector of 8 purple rectangles.
- \mathbf{Y} : The final output vector, labeled "Output".

encoder



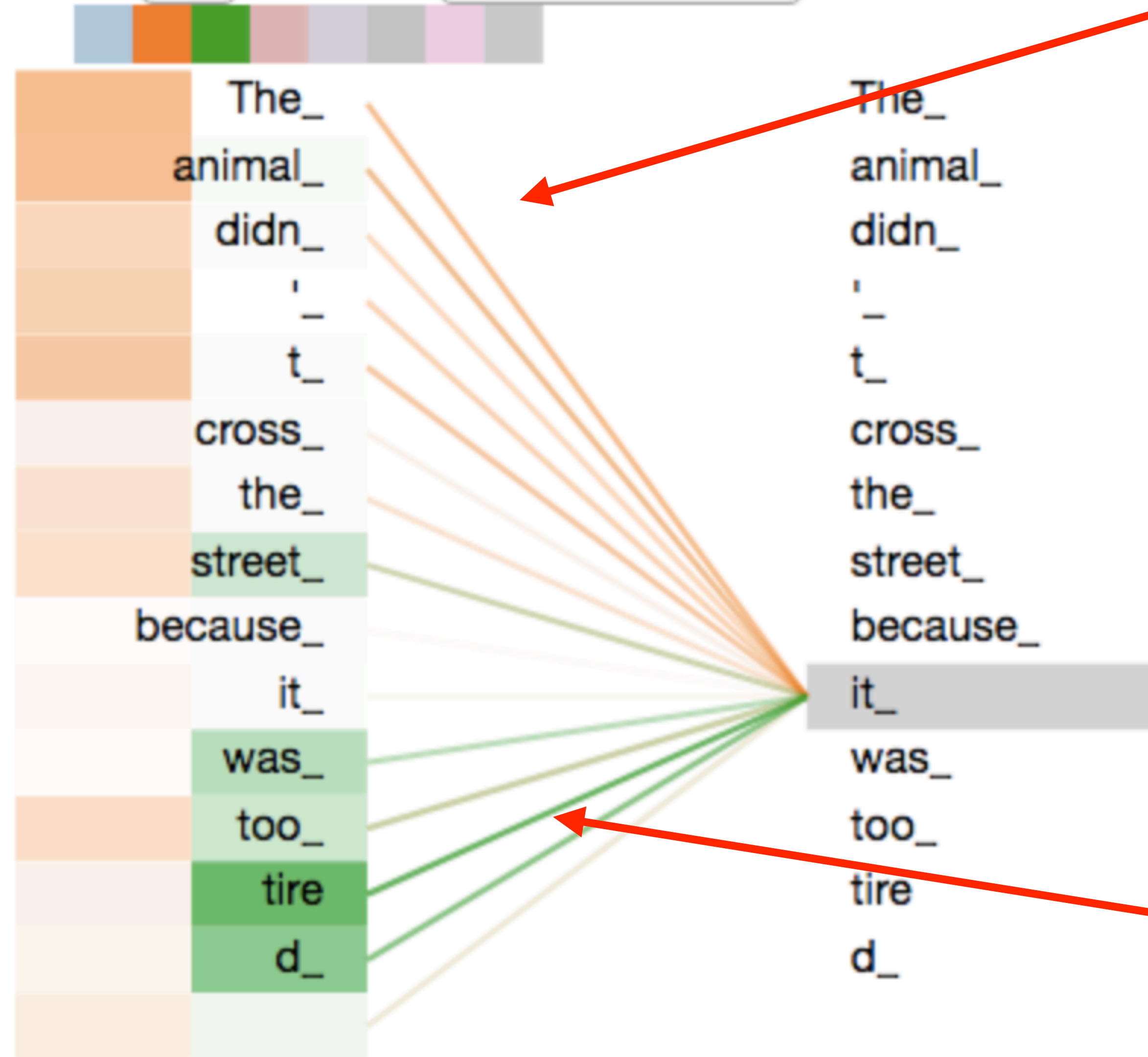
encoder





BertViz

Layer: 5 Attention: Input - Input



attention weights of
orange attention head

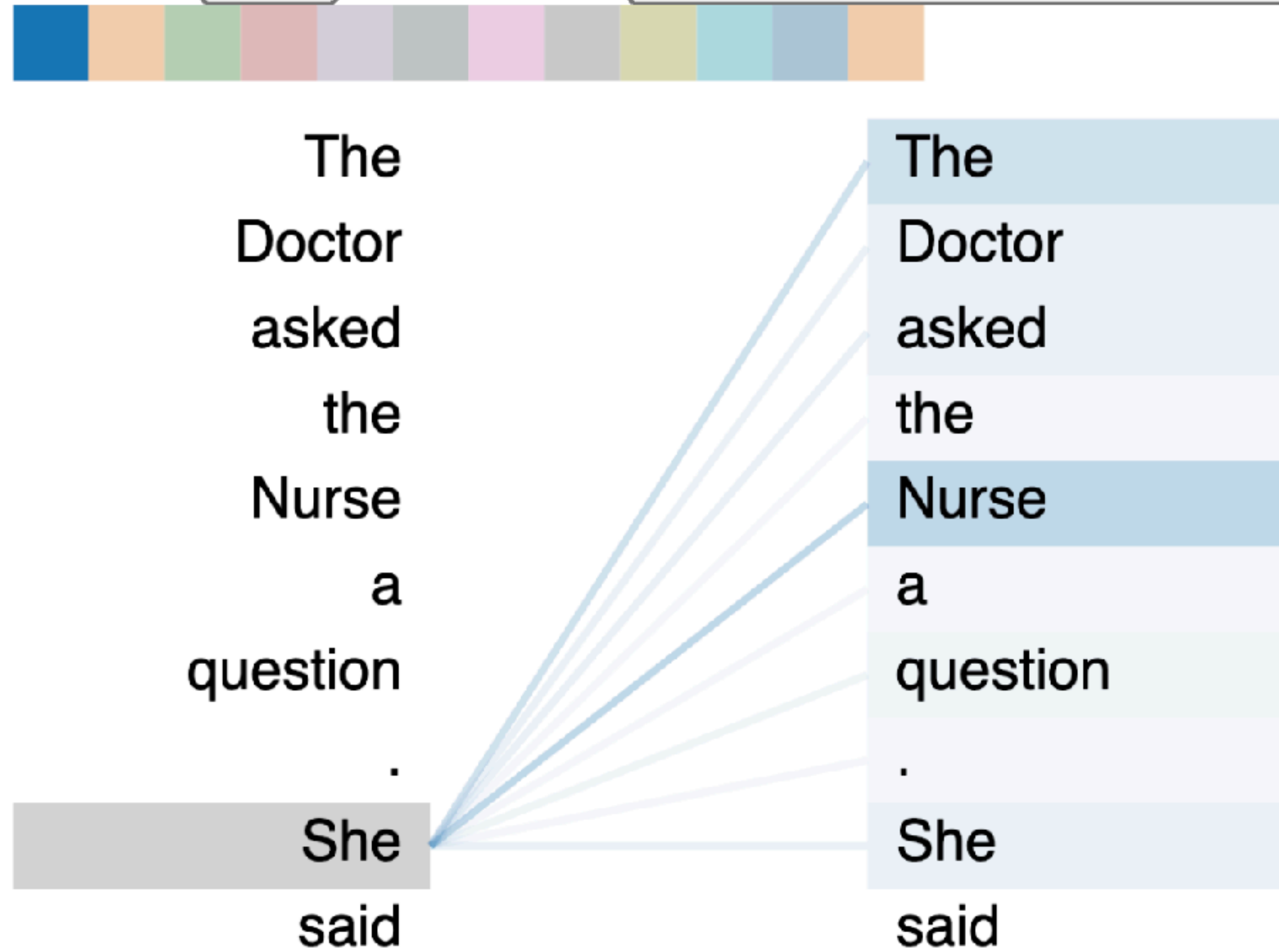
“it=animal” vs “it=street”??

attention weights of
green attention head

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

BertViz

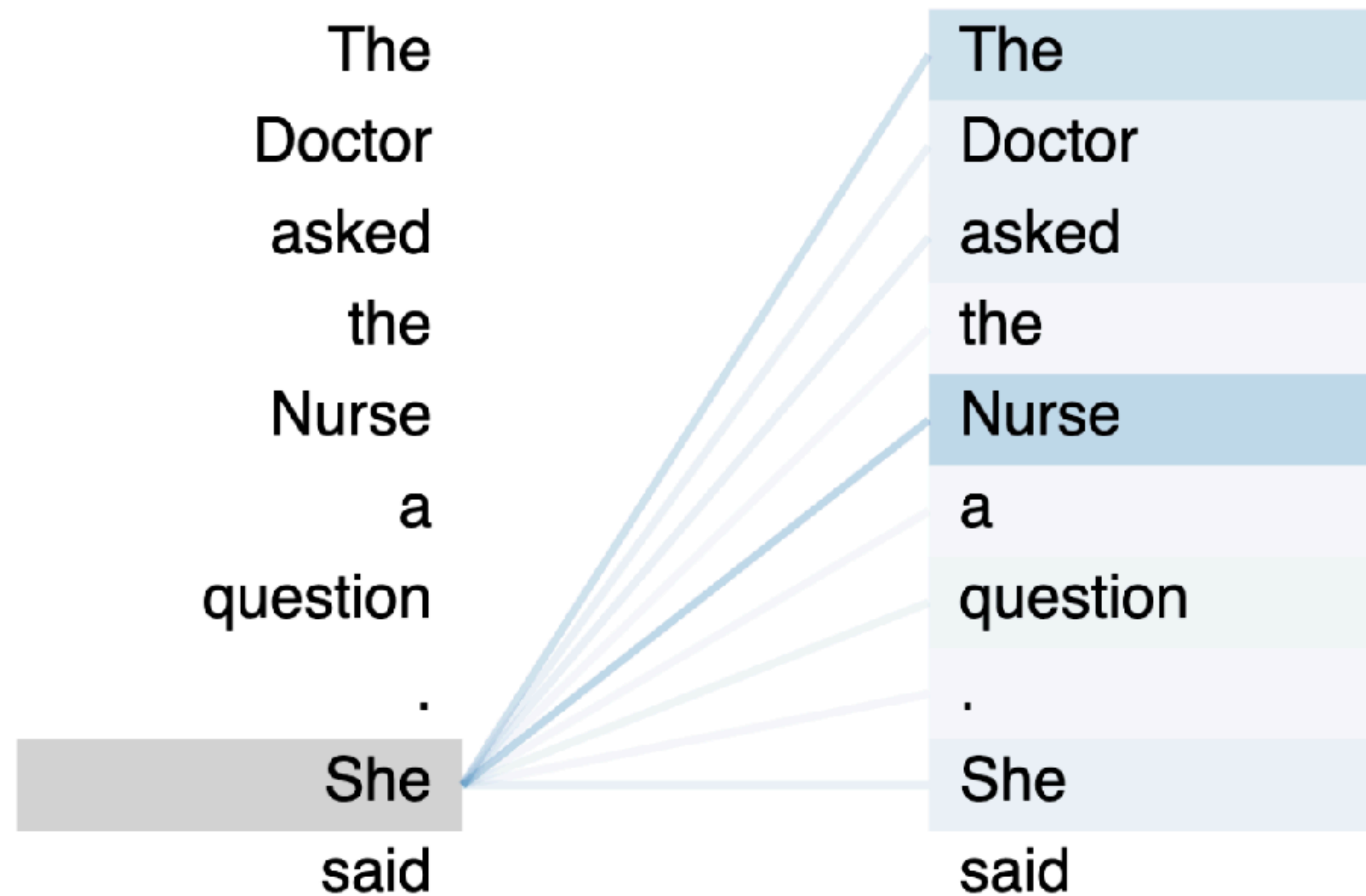
Layer: 0 ▾ Attention: All ▾



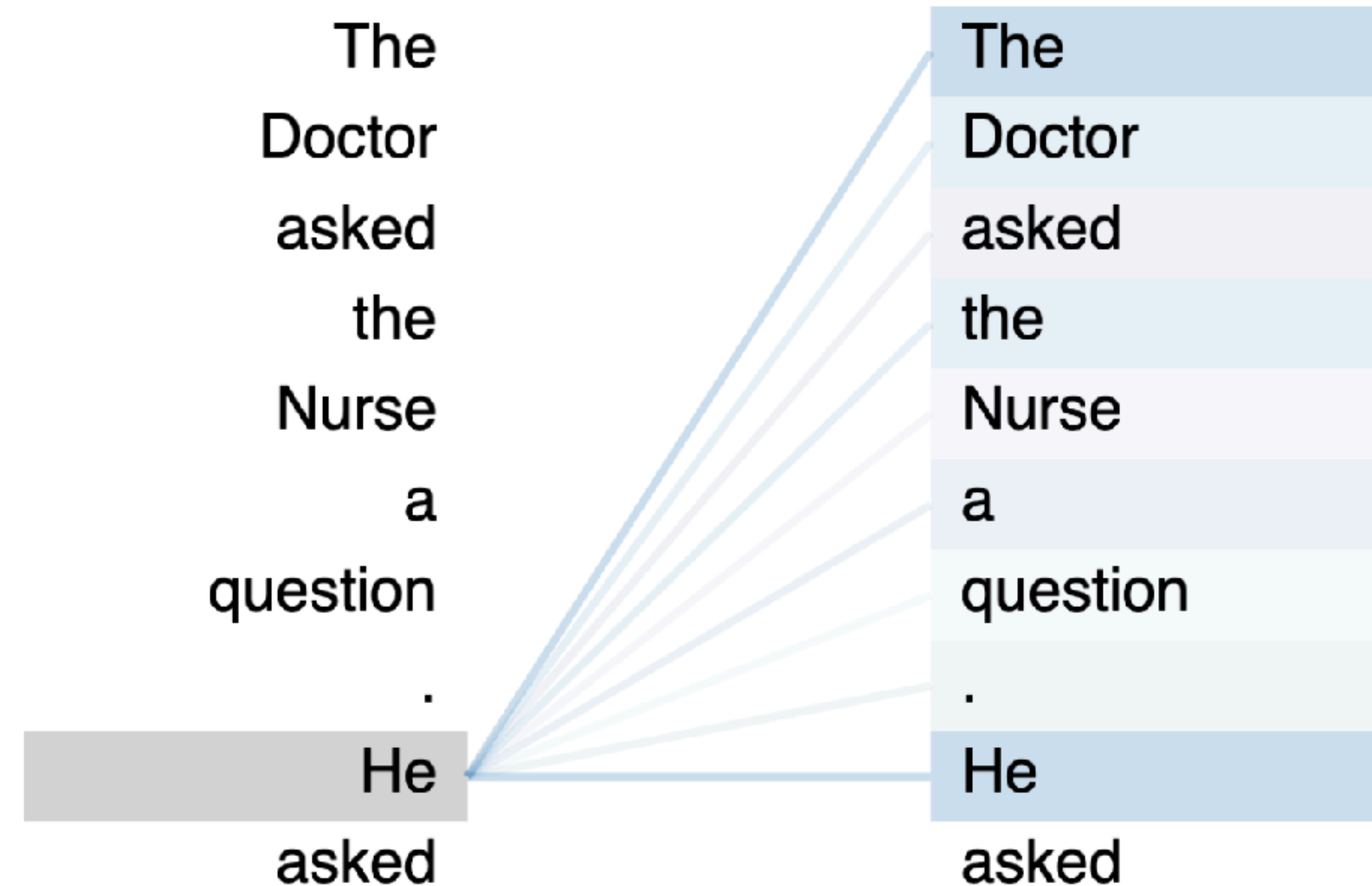
Model assumes "she=nurse"

<https://www.comet.com/site/blog/explainable-ai-for-transformers/>

BertViz (GPT2 model)

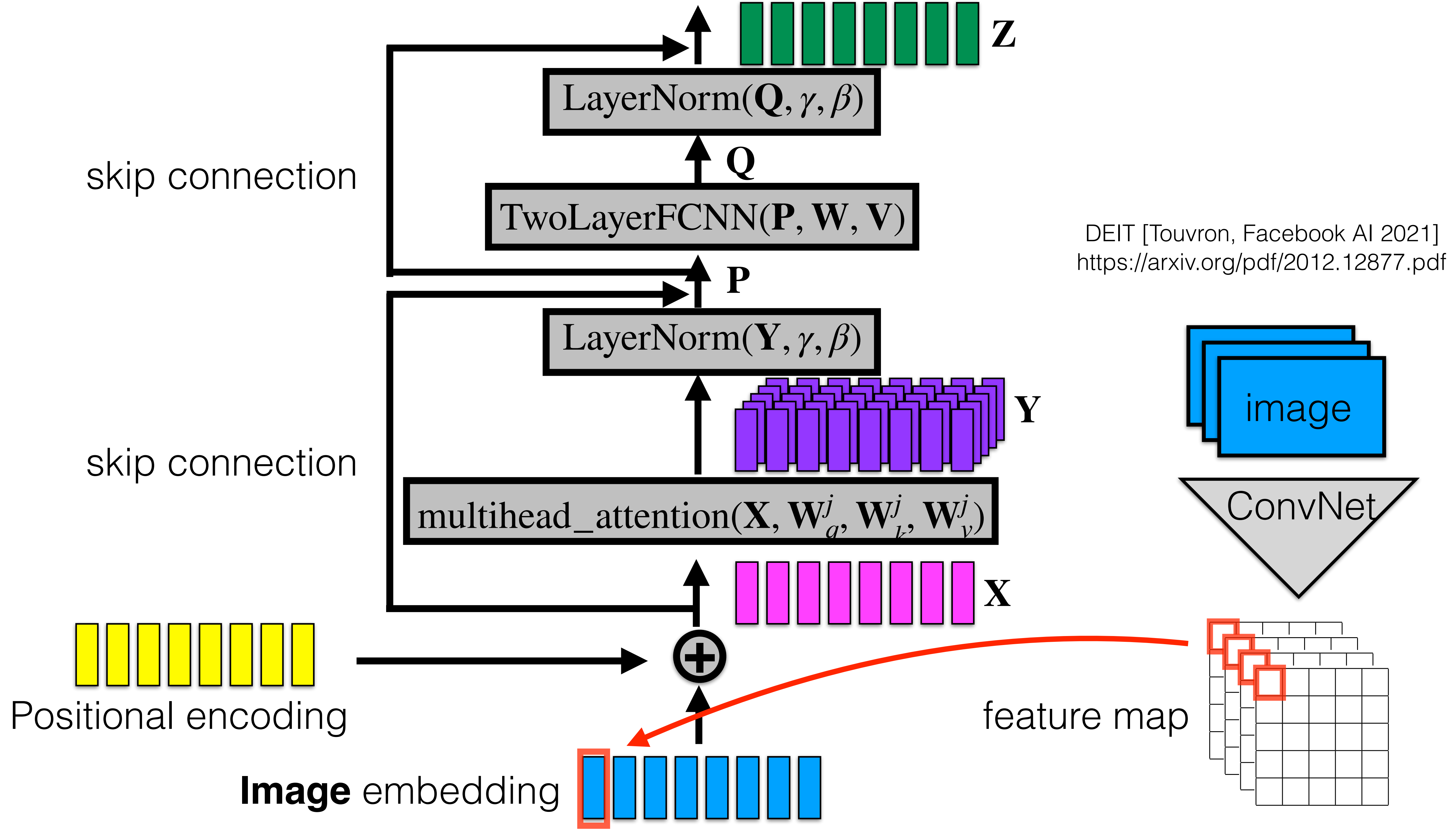


Model assumes “she=nurse”



Model assumes “he=doctor”

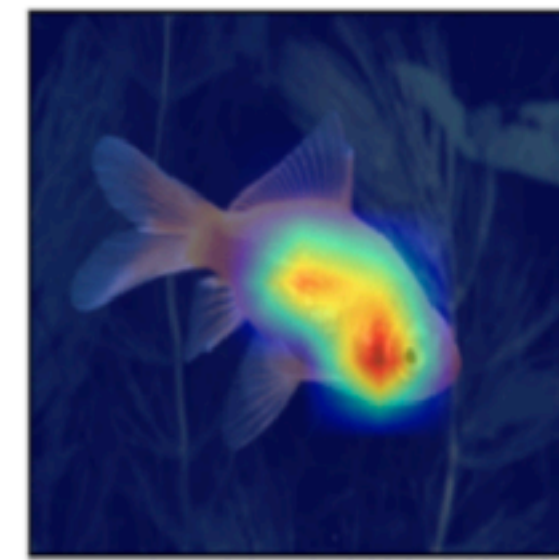
Transformers in images



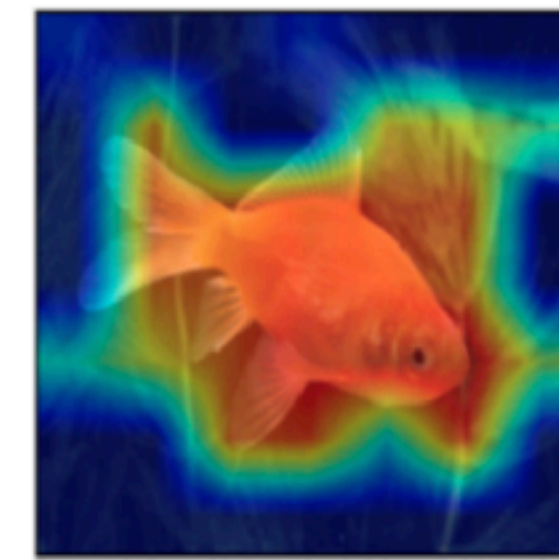
Attention in different Vision transformers (visu by gradCAM)



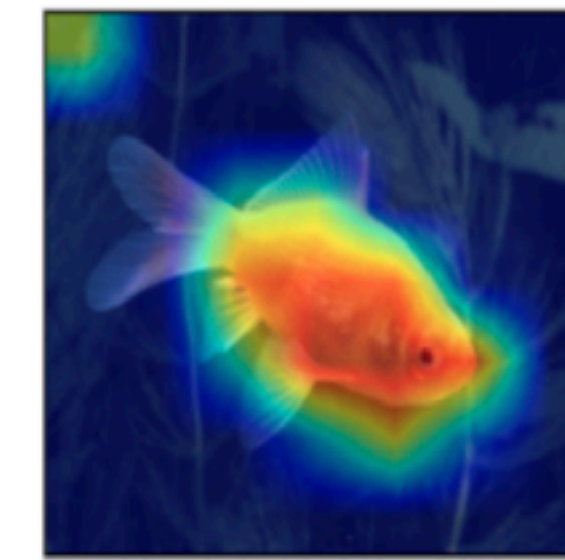
Input Image



Swin Transformer



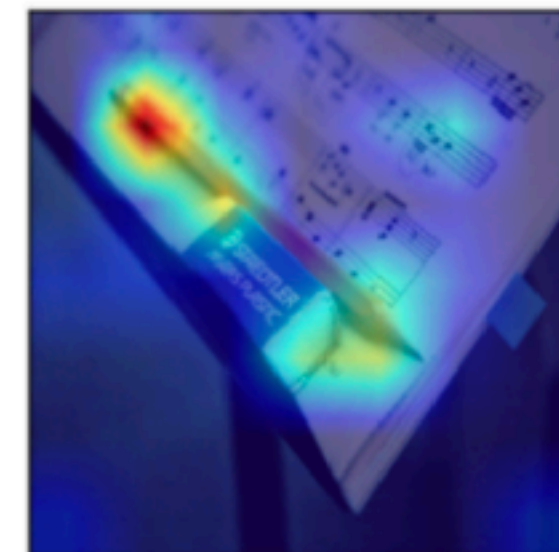
MaxViT



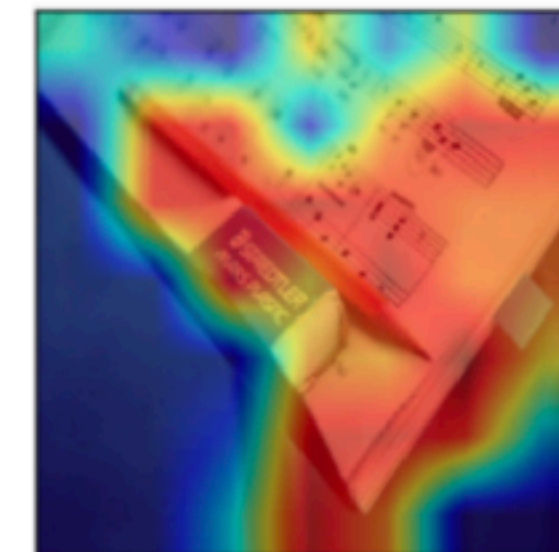
ACC-ViT



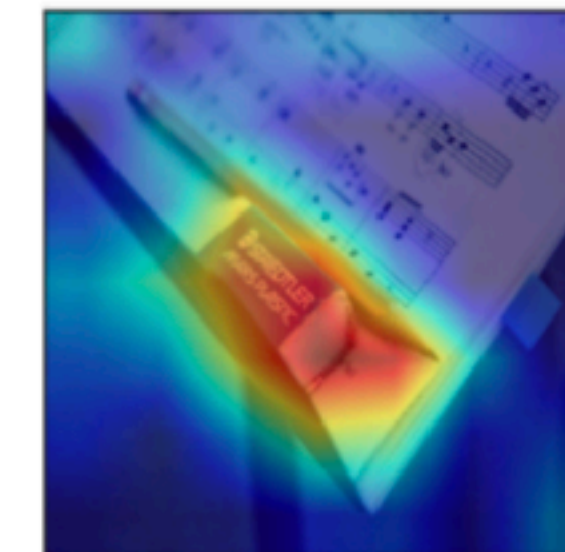
Input Image



Swin Transformer



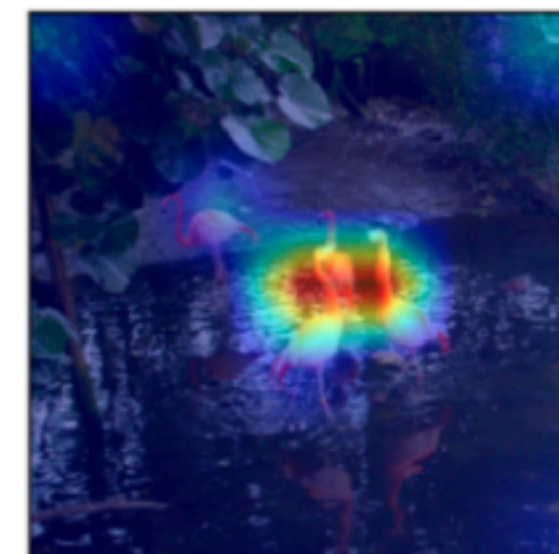
MaxViT



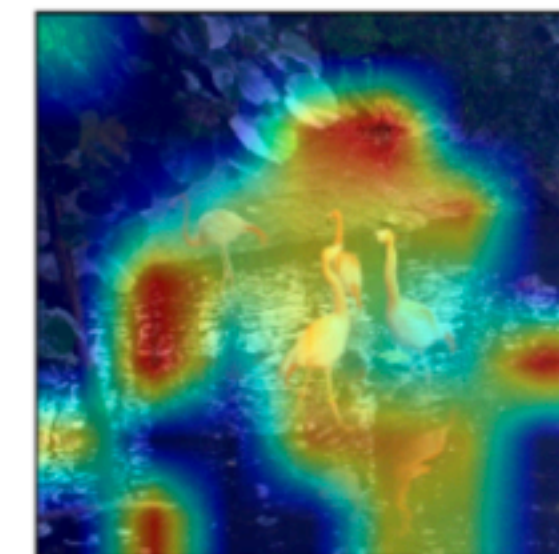
ACC-ViT



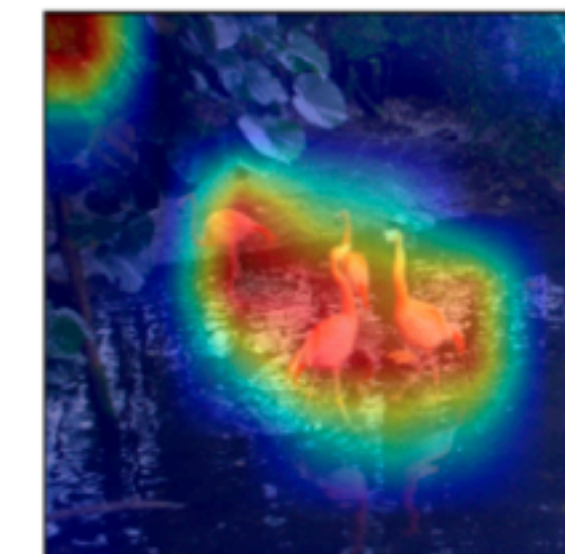
Input Image



Swin Transformer



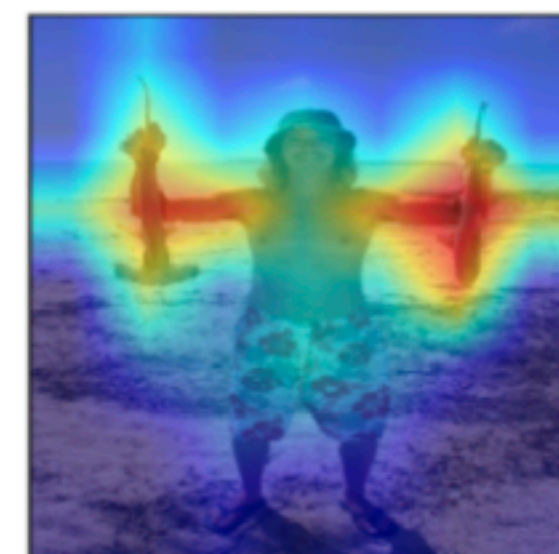
MaxViT



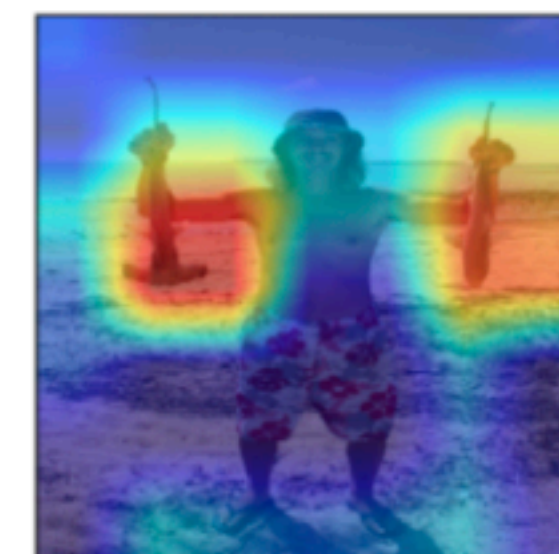
ACC-ViT



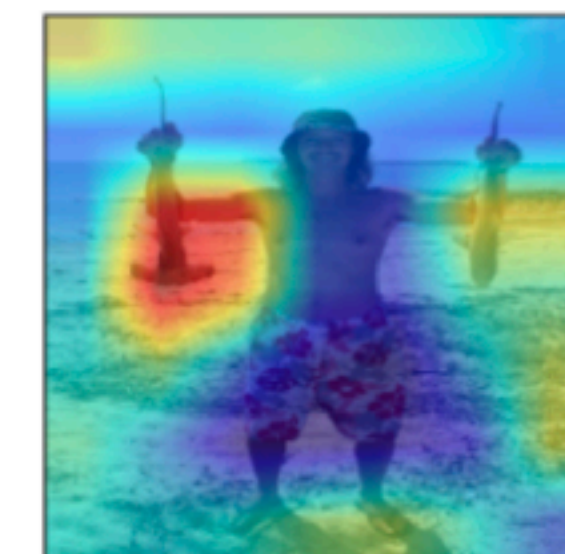
Input Image



Swin Transformer

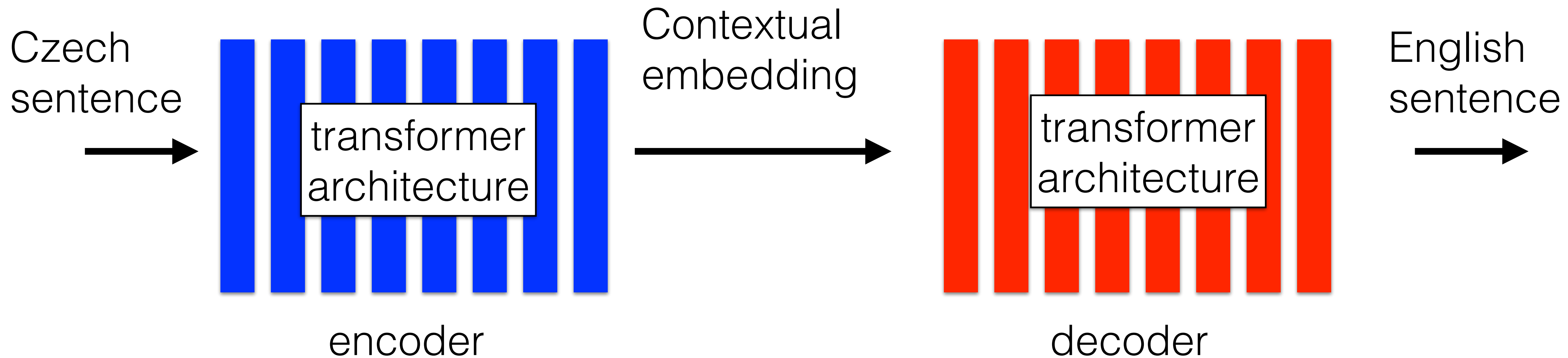


MaxViT



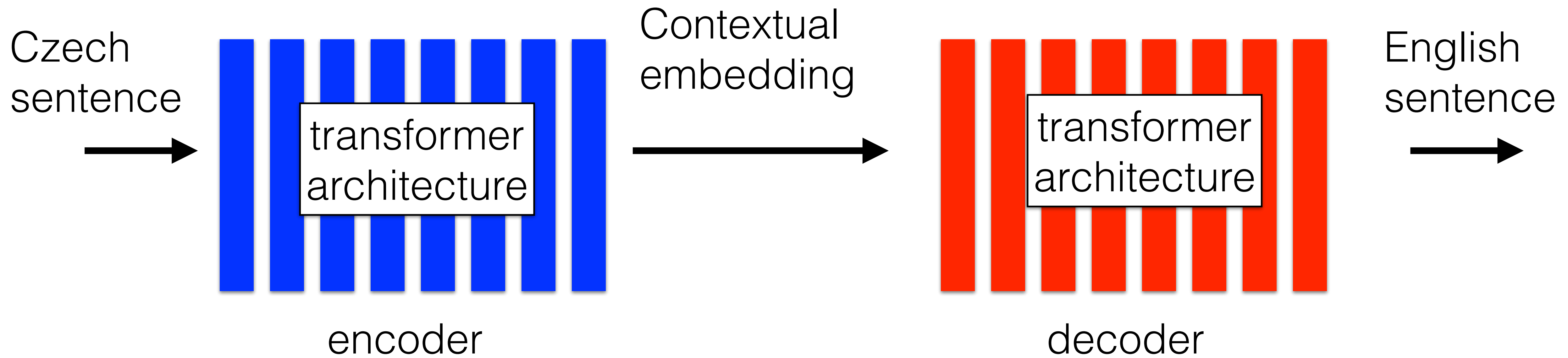
ACC-ViT

Machine translation



"Vy"
"jste"
"dobří"
"studenti"

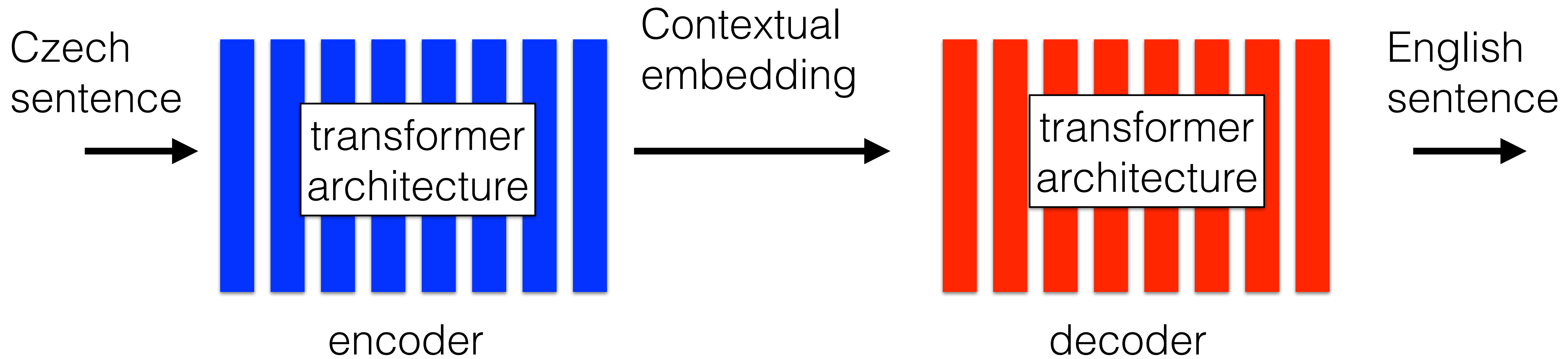
- Encoder is standard transformer with self-attention
- Decoder auto-regressively generates output sentence
- Decoder requires special attentions



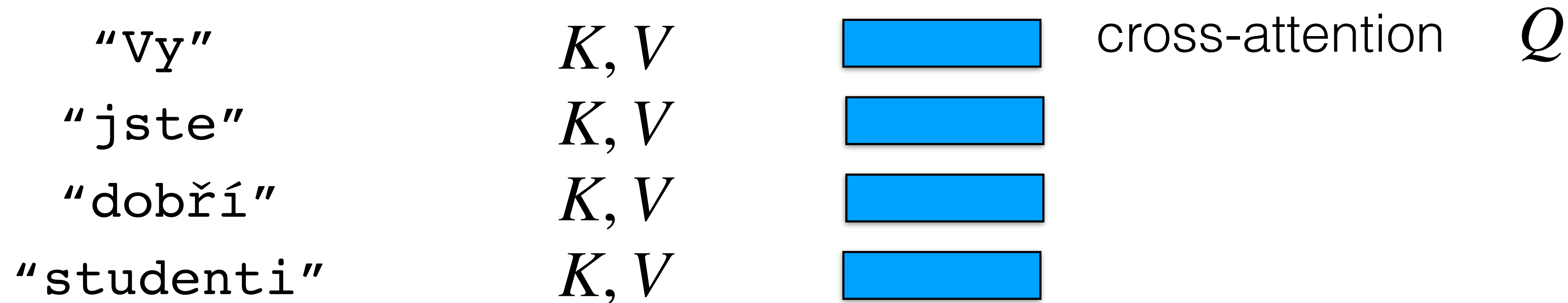
masked self-att. "<SOS>"

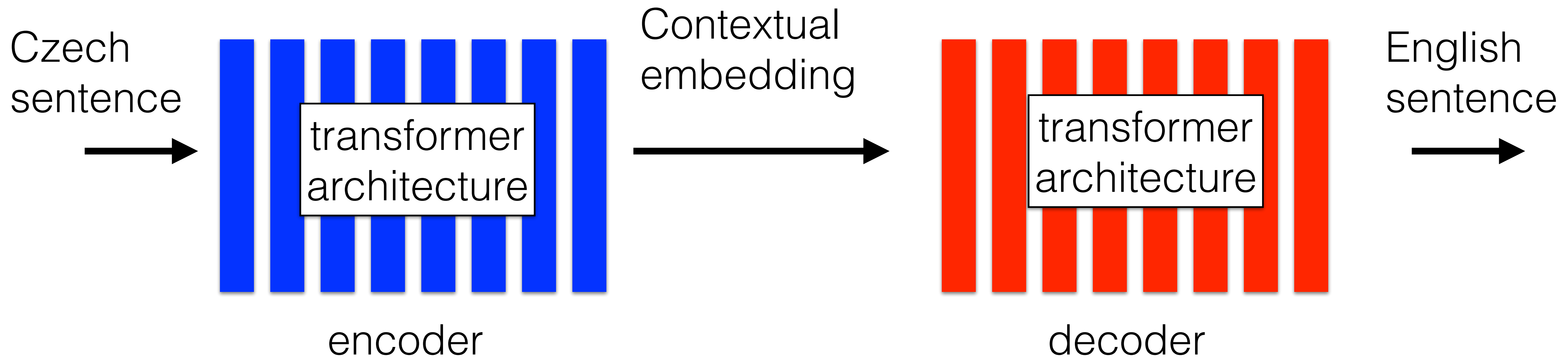
"Vy"
"jste"
"dobří"
"studenti"





masked self-att. " <sos> " K, V, Q



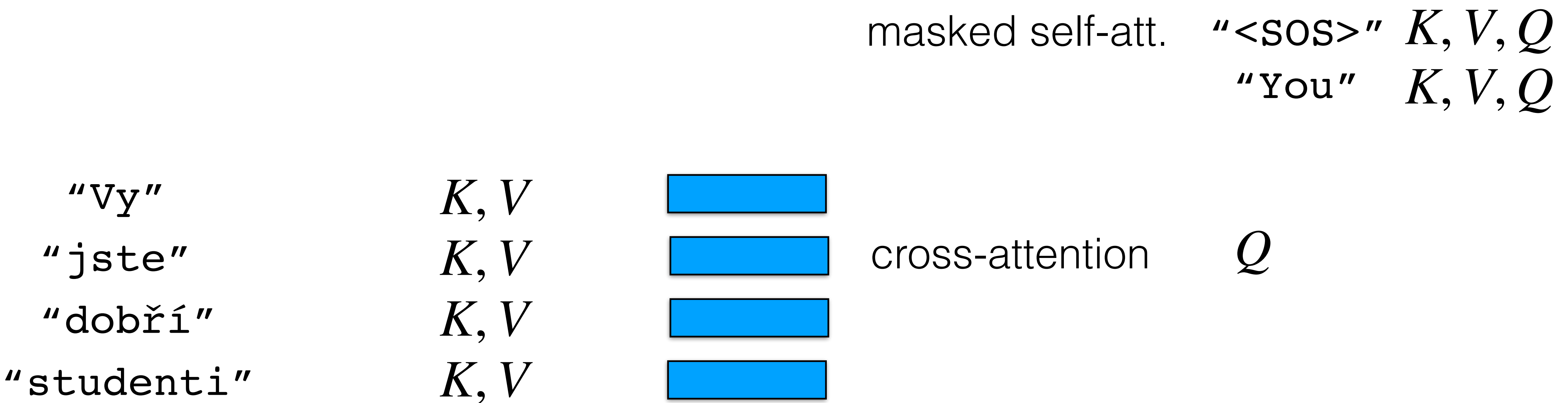
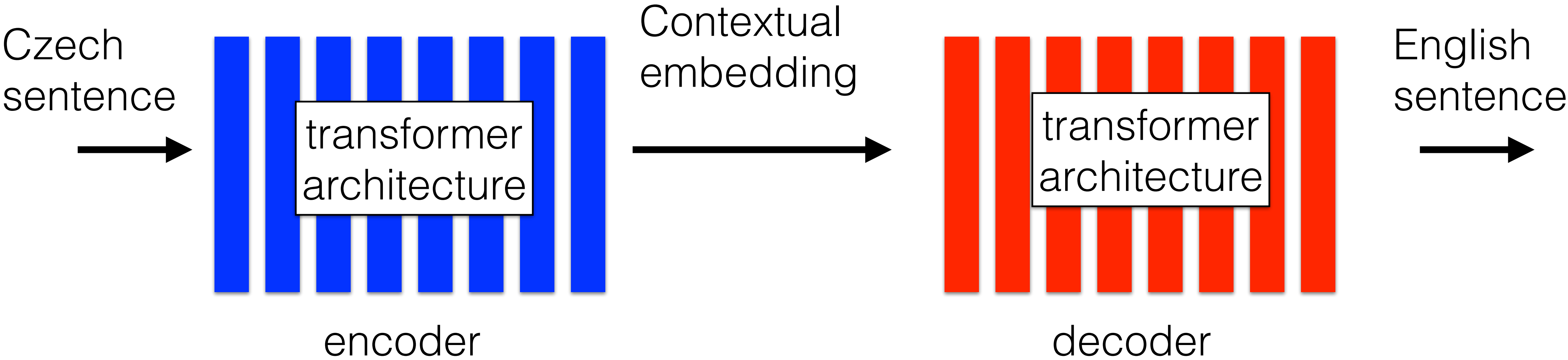


masked self-att. "<SOS>"

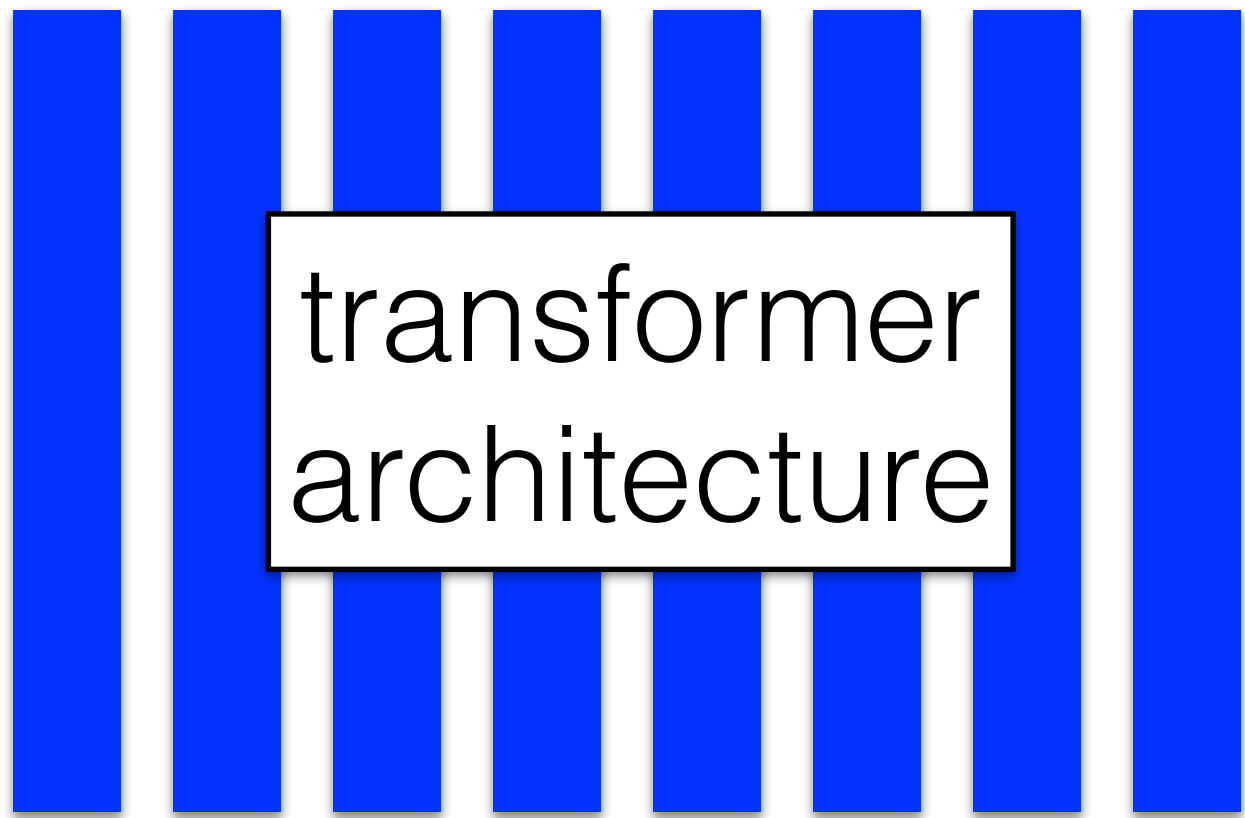
"Vy"
"jste"
"dobří"
"studenti"



cross-attention "You"



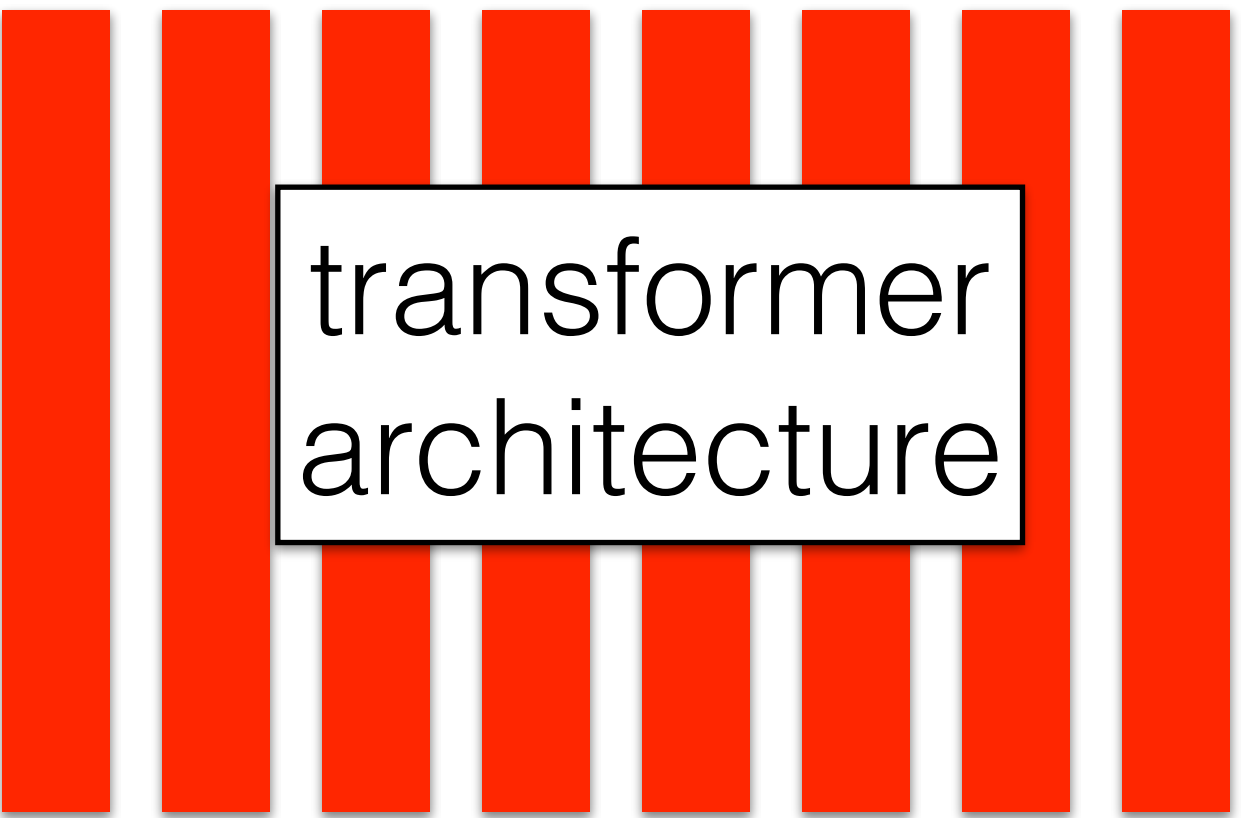
Czech
sentence



transformer
architecture

encoder

Contextual
embedding



transformer
architecture

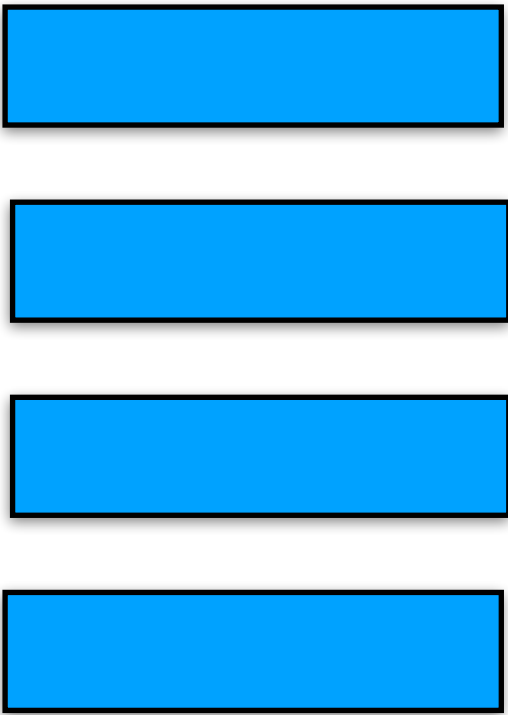
decoder

English
sentence

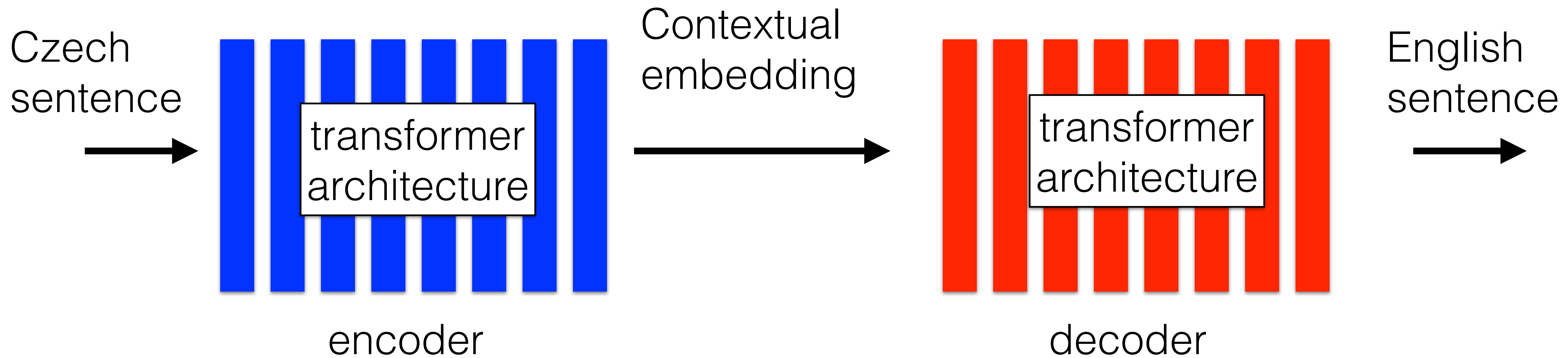


masked self-att. "<SOS>"
 "You"

"Vy"
"jste"
"dobří"
"studenti"

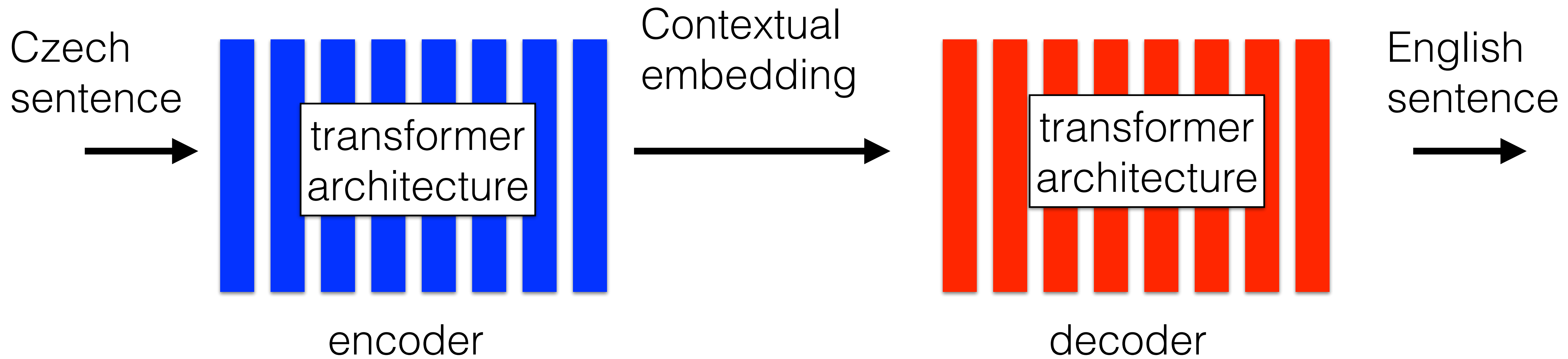


cross-attention "are"



masked self-att. " <SOS> " K, V, Q
 " You " K, V, Q
 " are " K, V, Q

"Vy"	K, V	<div style="background-color: #00bfff; width: 100px; height: 20px;"></div>	cross-attention Q
"jste"	K, V	<div style="background-color: #00bfff; width: 100px; height: 20px;"></div>	
"dobří"	K, V	<div style="background-color: #00bfff; width: 100px; height: 20px;"></div>	
"studenti"	K, V	<div style="background-color: #00bfff; width: 100px; height: 20px;"></div>	



"Vy"
"jste"
"dobří"
"studenti"



masked self-att.

"<SOS>"
"You"
"are"
"good"
"students"

cross-attention

Self-attention

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

Fill in **gaps** (<unknown> words tokens) in sentences (BERT)

Masked self-attention

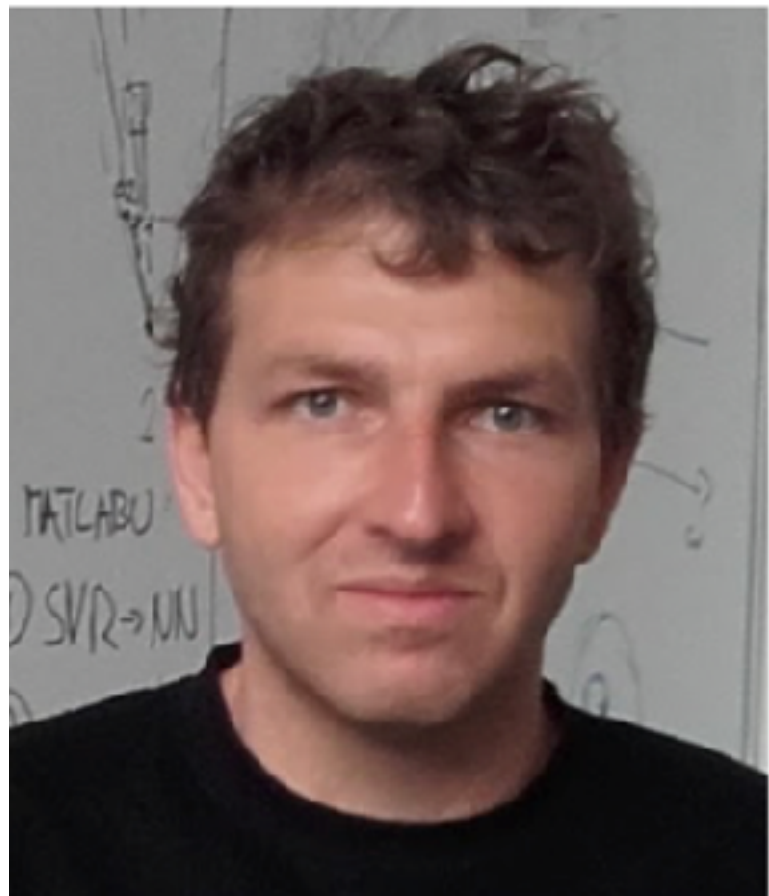
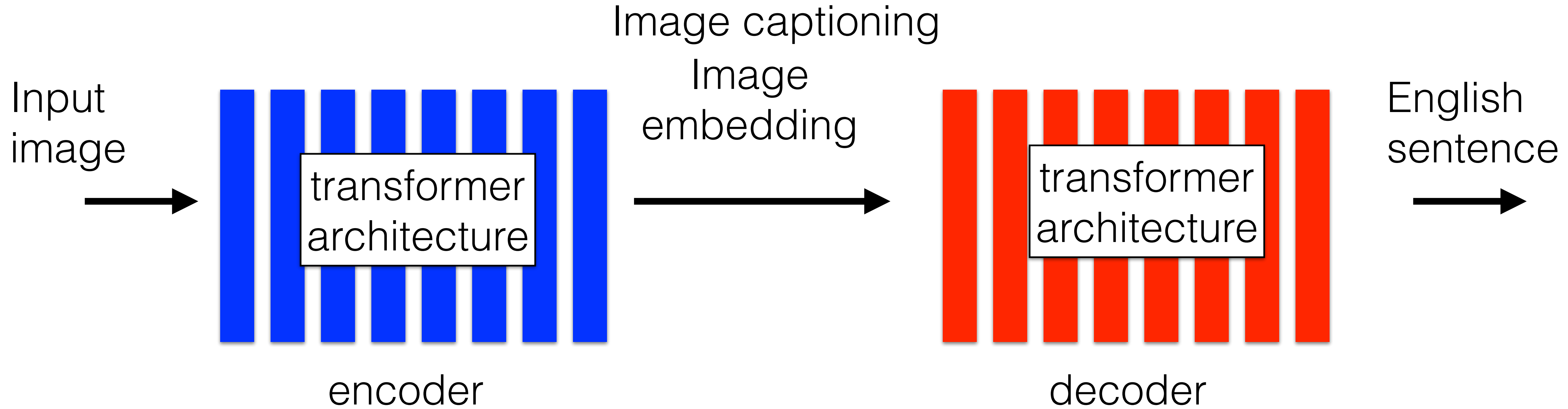
$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^\top}{\sqrt{d_k}} + \text{mask} \right) V$$

$$\text{Mask} = \begin{bmatrix} 0 & -\infty & -\infty \\ 0 & 0 & -\infty \\ 0 & 0 & 0 \end{bmatrix}$$

- Assures temporal coherence without creating dependence on the correct number of <unknown> words in the input.
- Assures better **parallelization** and **generalization** in autoregressive text generation (GPT)

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

Q ... from decoder
 K, V ... from encoder



masked self-att.

"<SOS>"

"Best"

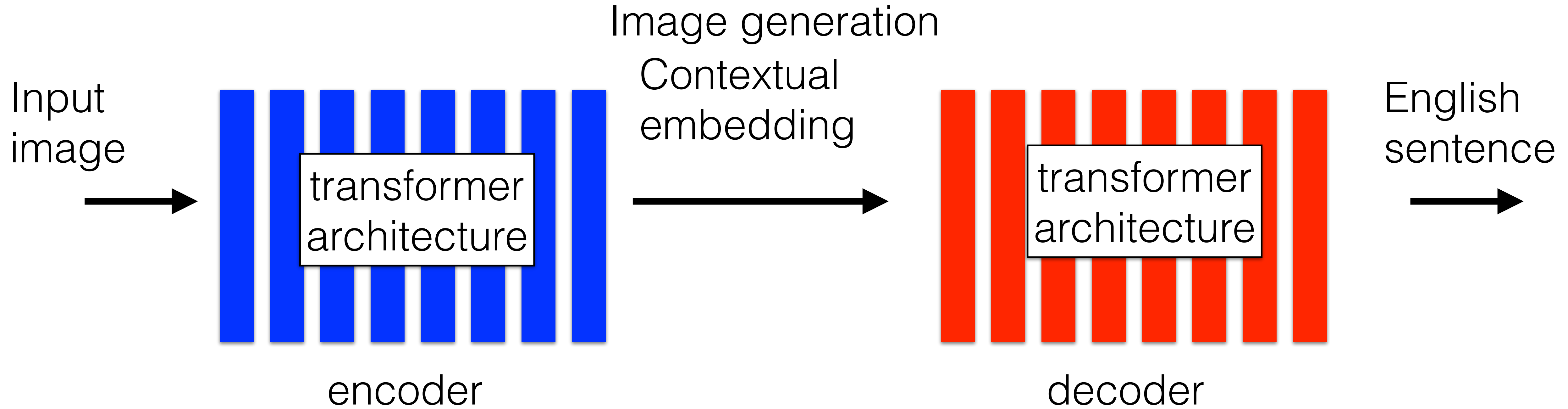
"UROB"

"teacher"

"ever"



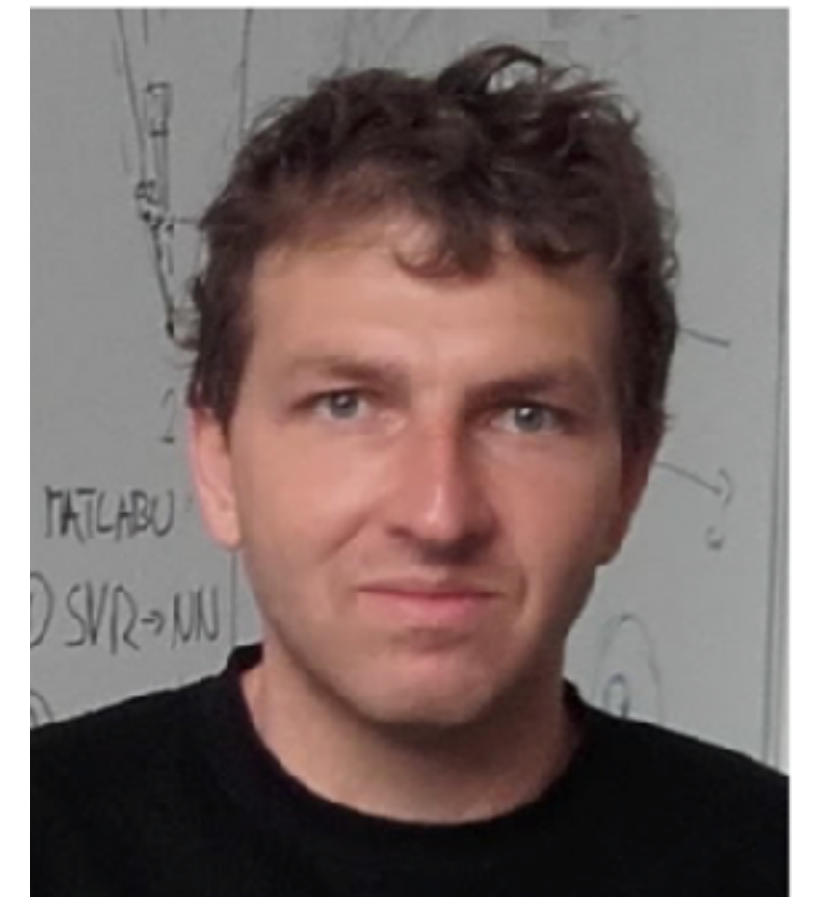
cross-attention

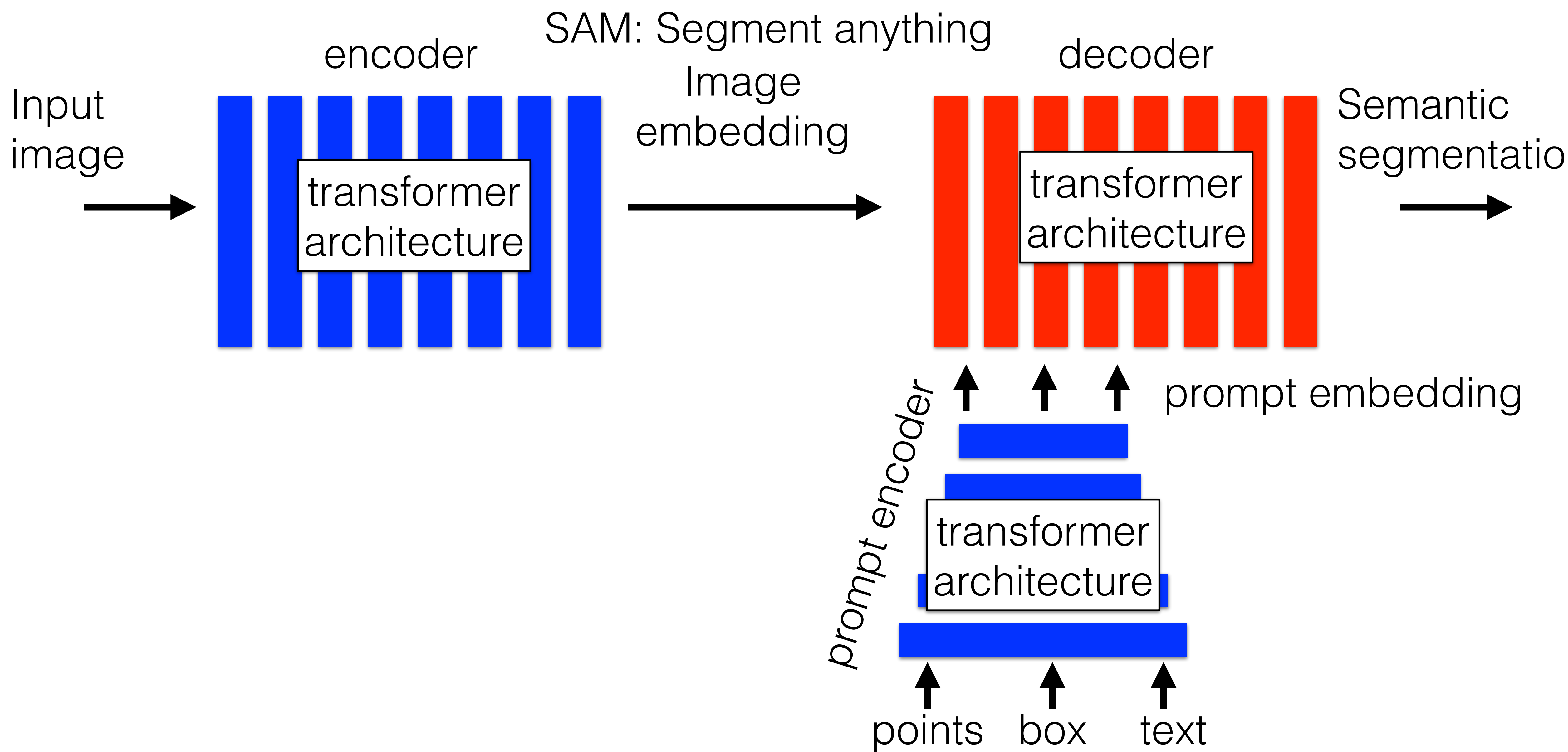


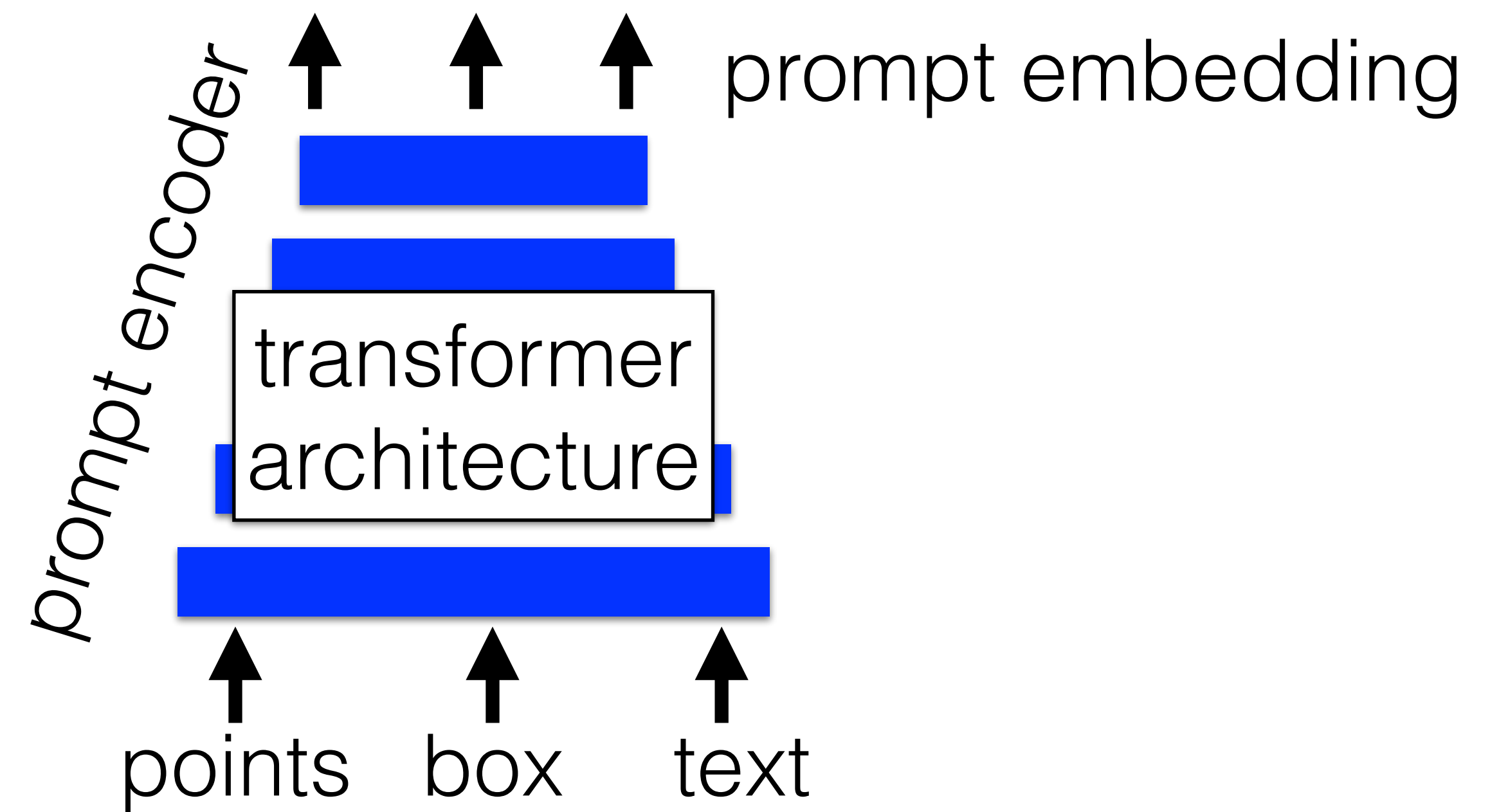
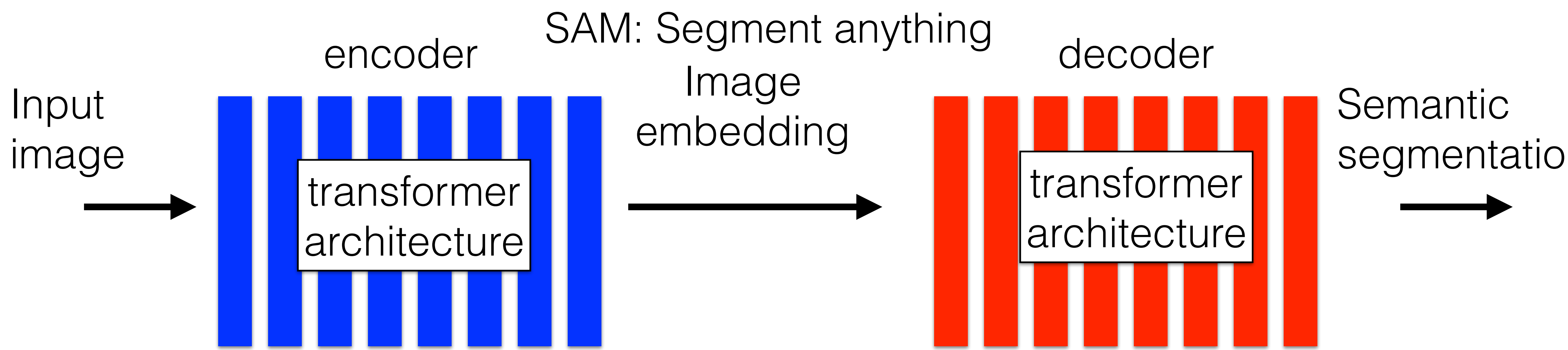
Text encoder
(BERT, GPT)

Generative model
(Diffusion, VQ-VAE, GAN)

"<SOS>"
"Best"
"UROB"
"teacher"
"ever"



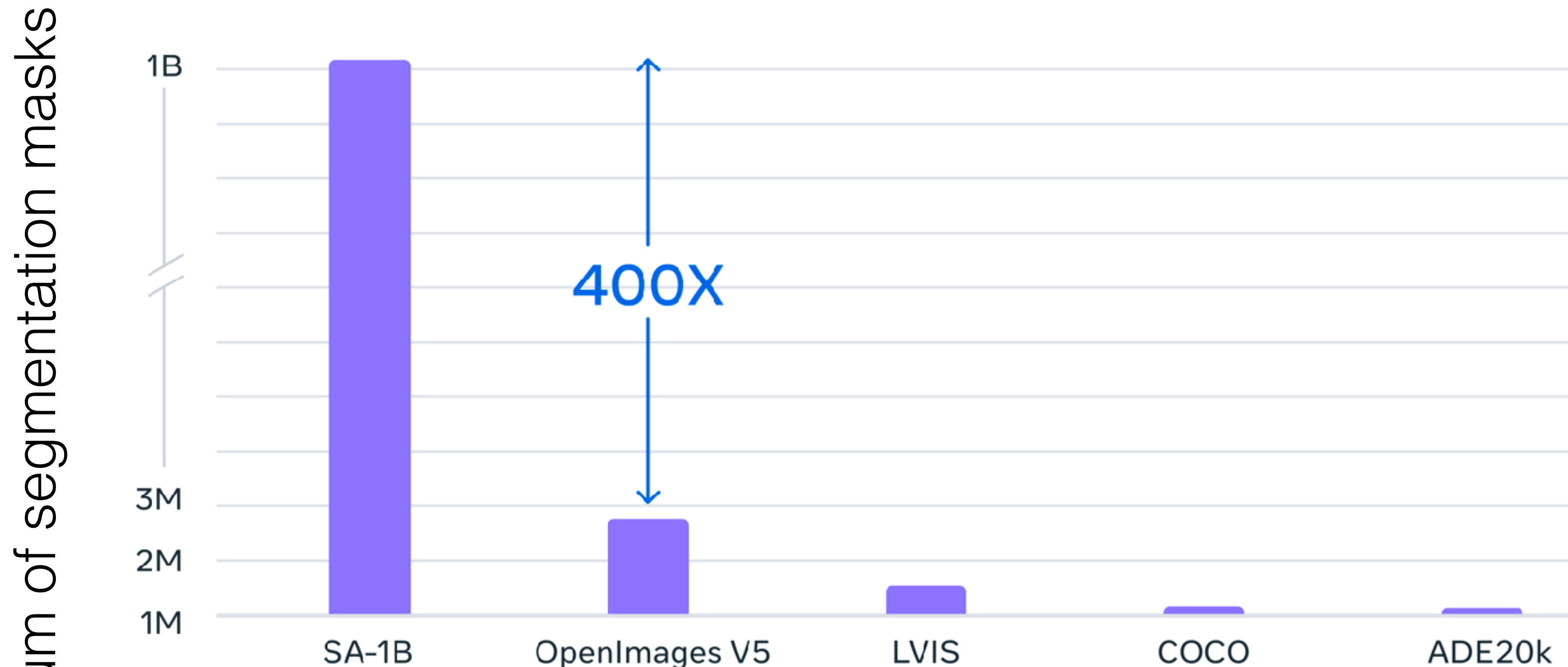




SAM: Segment anything

Convolution is actually **structurally-enforced local attention**.

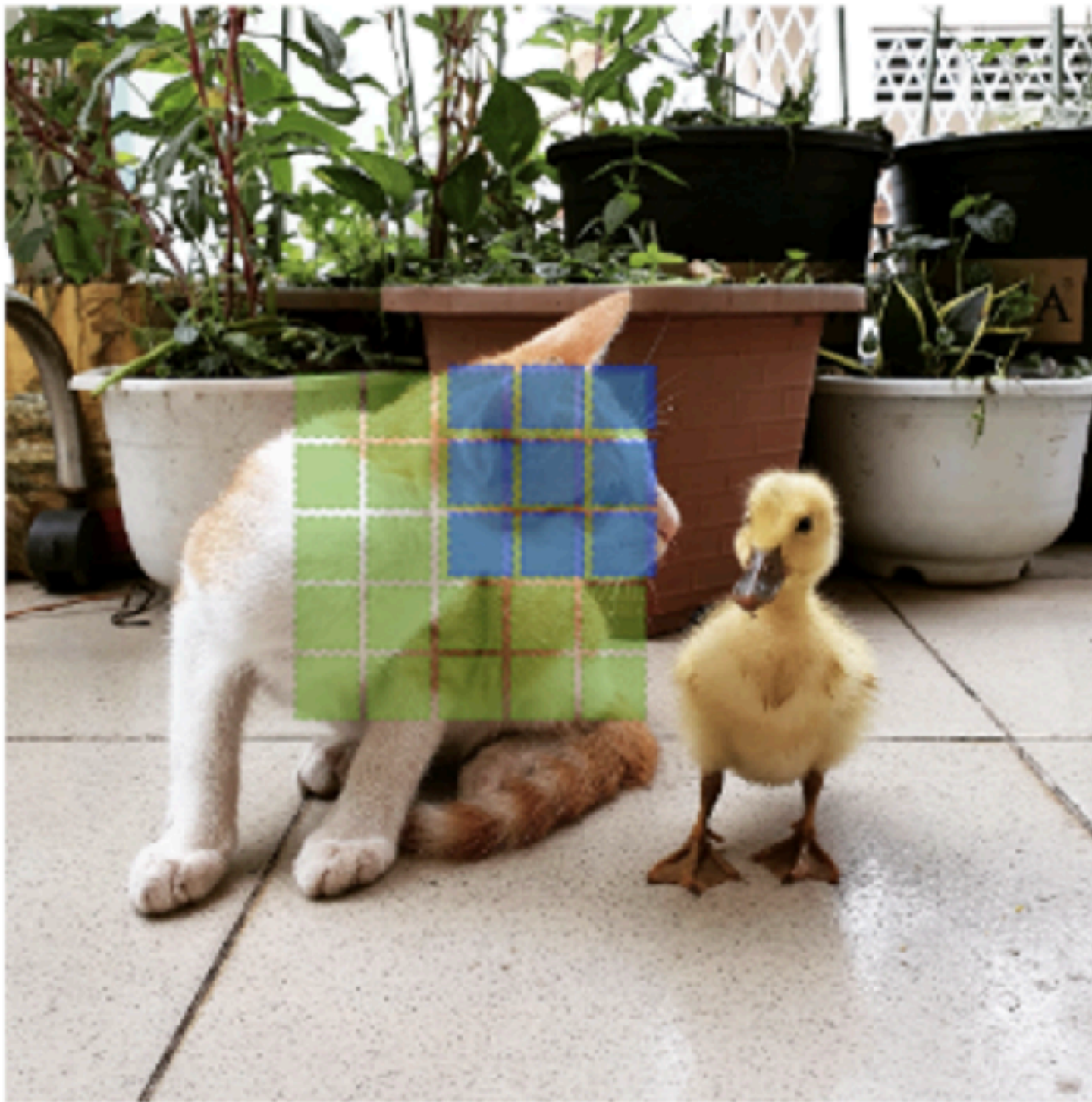
Transformers allow **global attention** and have to **learn** it **from data**.



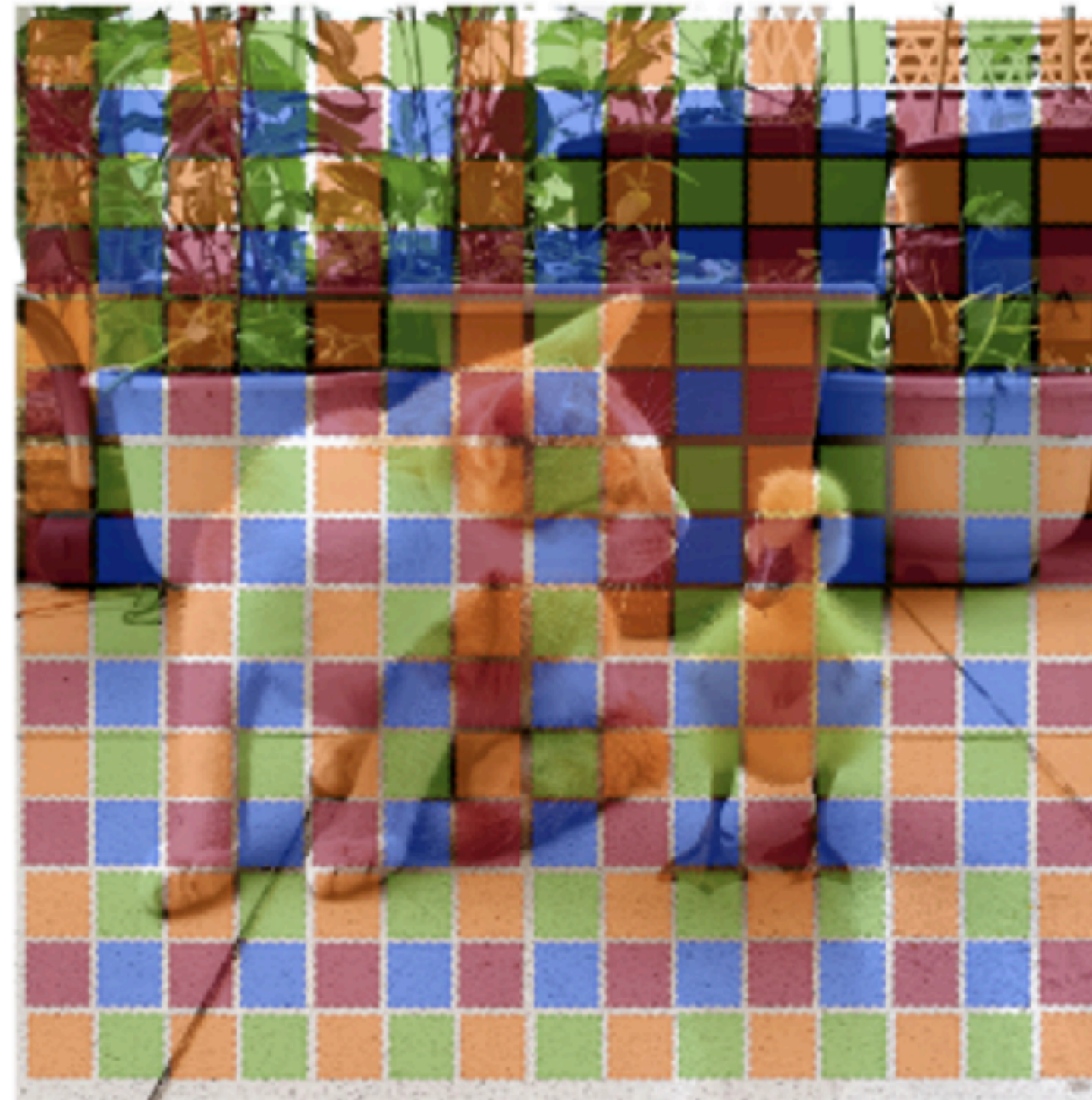
[Facebook 2023] <https://arxiv.org/pdf/2304.02643.pdf>

Attention used for images

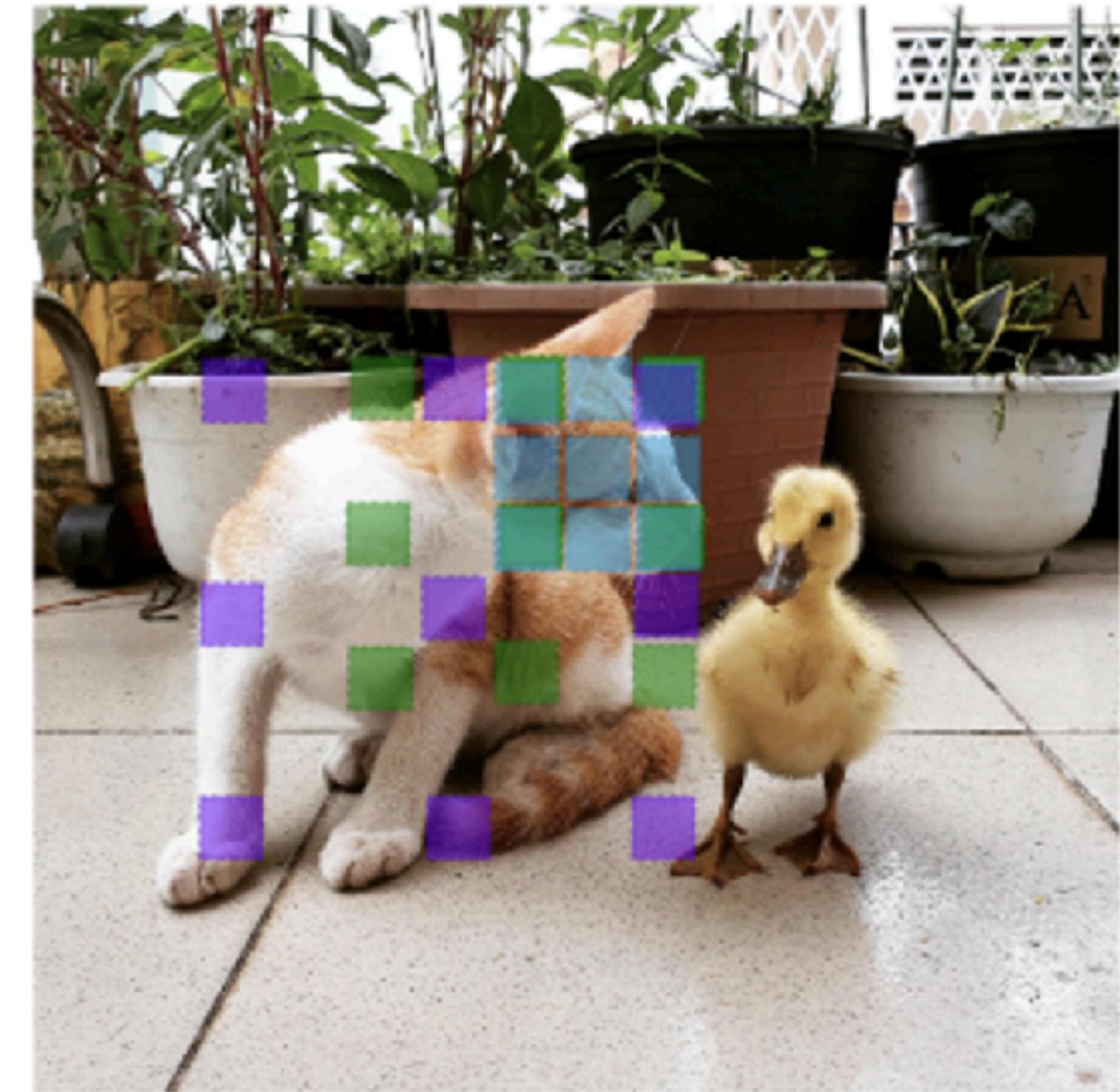
Global attention in early layers can be replaced by **local** attention



(a) Regional Attention



(b) Sparse Attention



(c) Atrous Attention

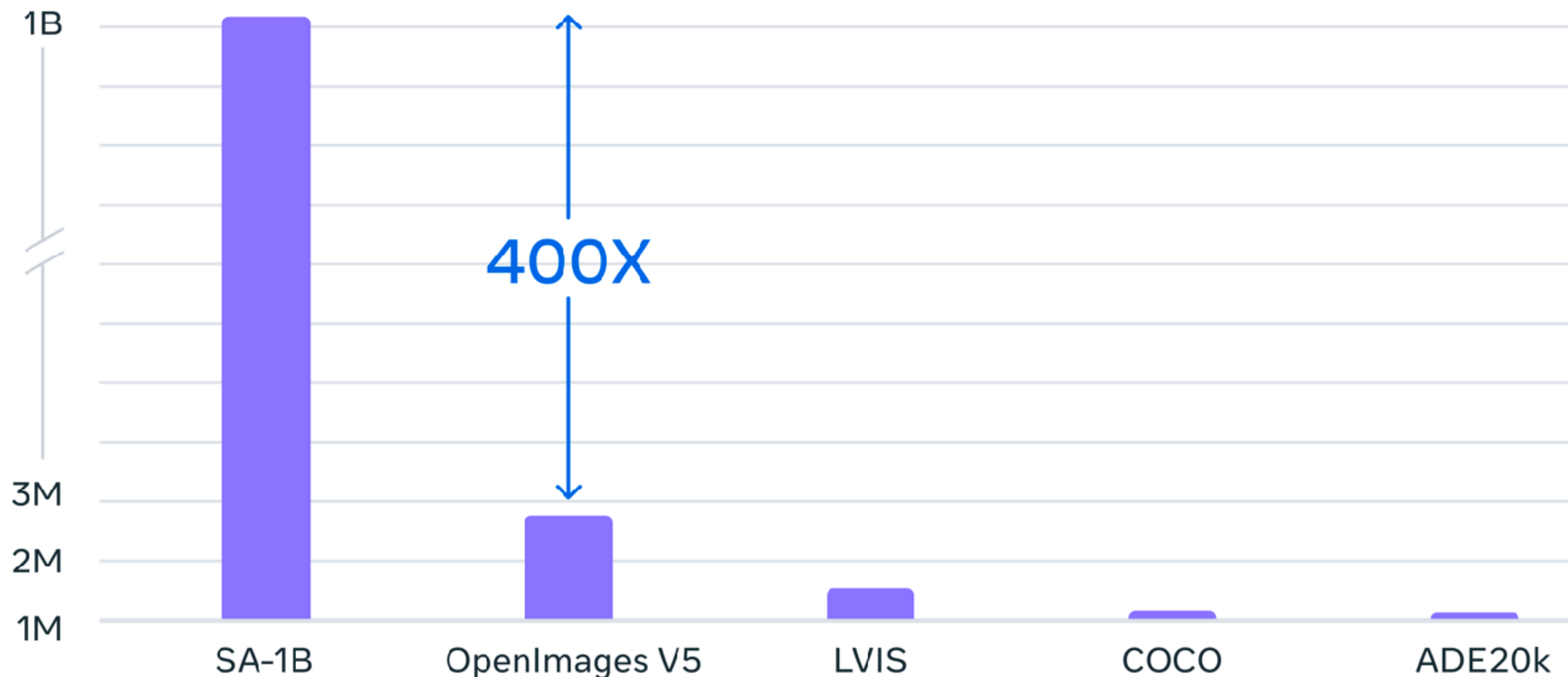
SAM: Segment anything

Convolution is actually **structurally-enforced local attention**.

Transformers allow **global attention** and have to **learn** it **from data**.

Change of paradigm:

- **small** datasets => use **simple** models (strong inductive prior such as convolution)
- **huge** dataset with cheap or free training data => **complex** model **learn everything**
- **semi-supervision** / **self-supervision**



[Facebook 2023] <https://arxiv.org/pdf/2304.02643.pdf>

Image inpainting

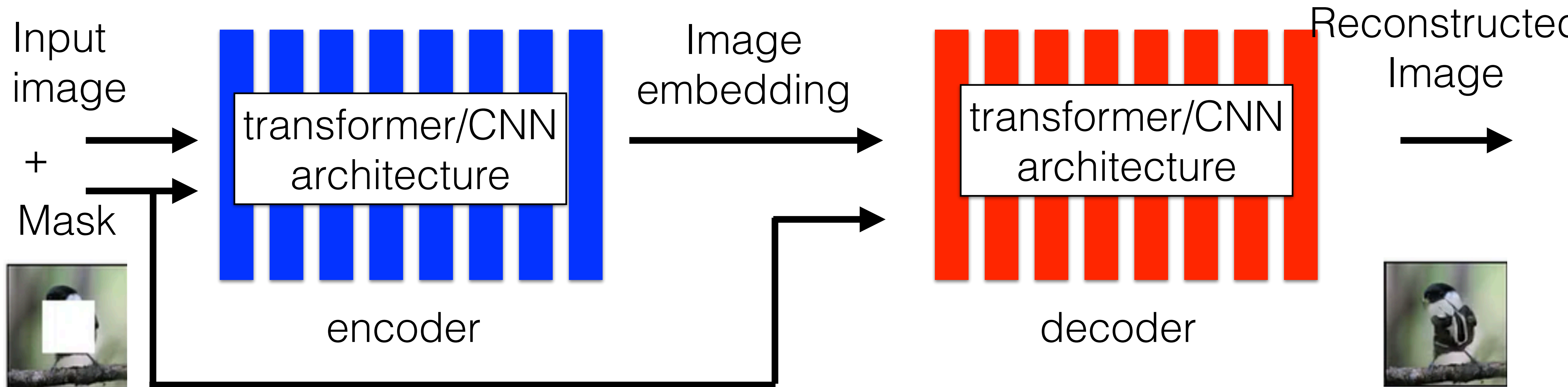
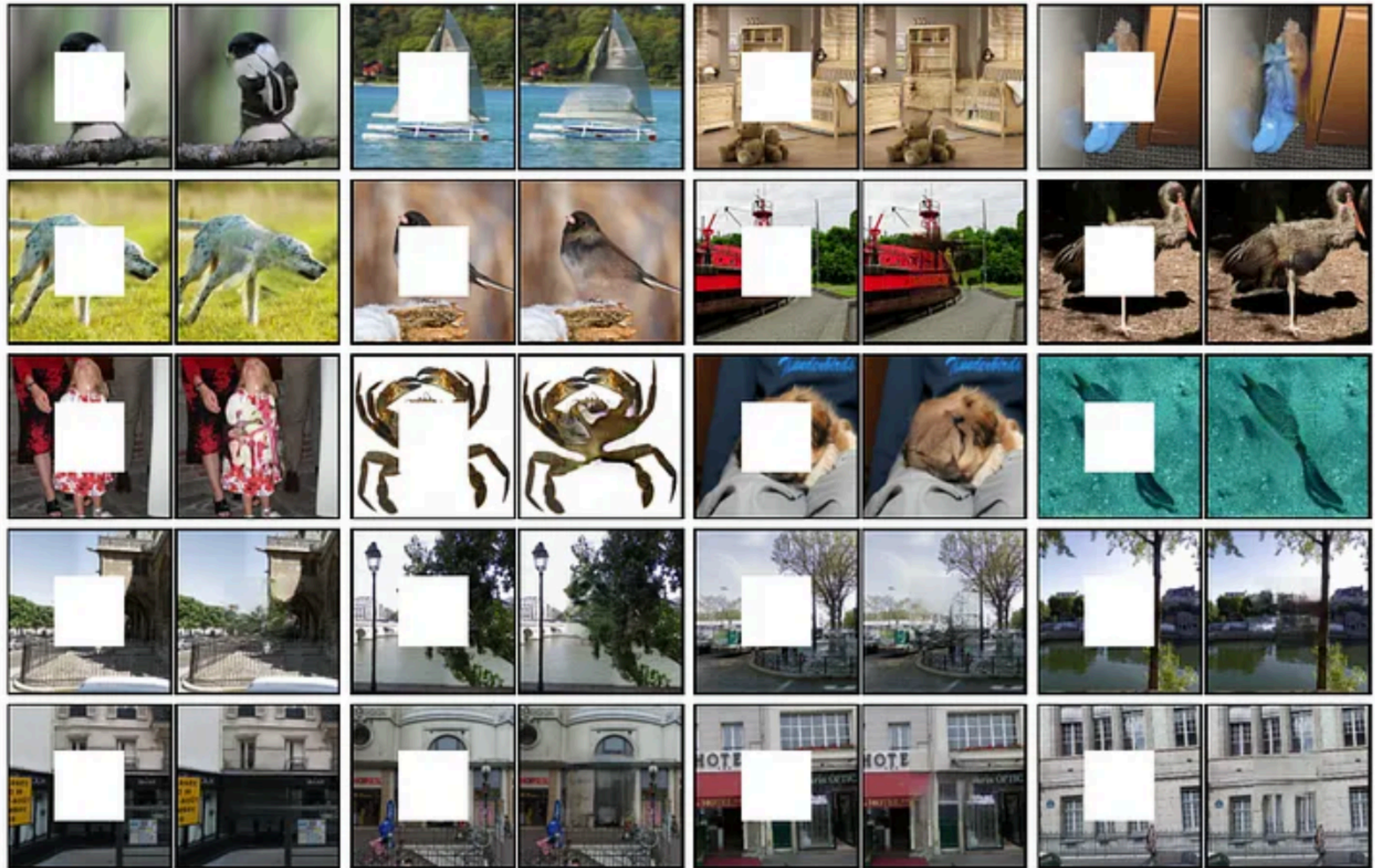


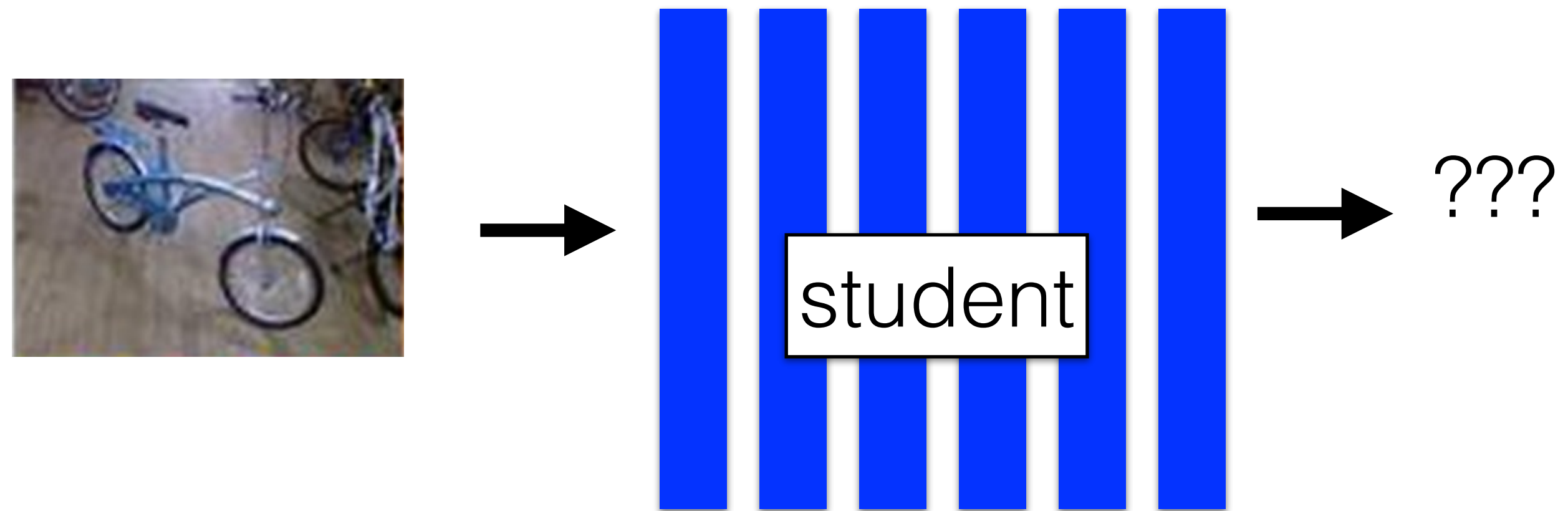
Image inpainting



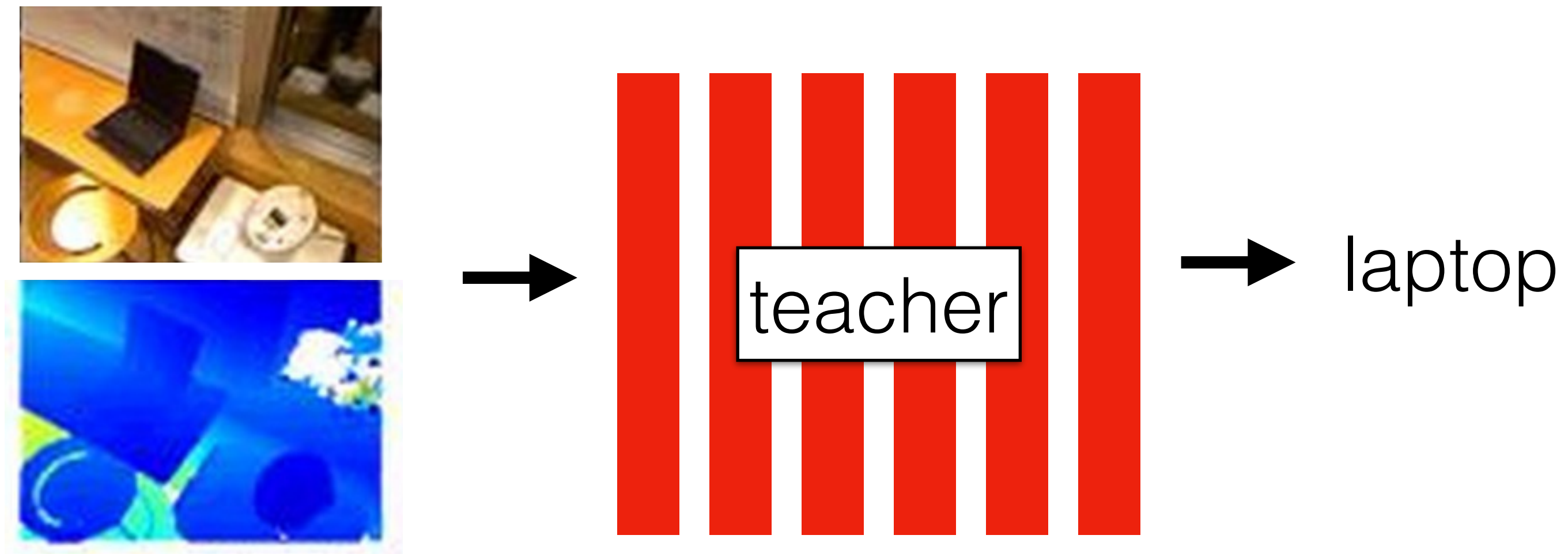
Context encoders CVPR 2016 <https://arxiv.org/abs/1604.07379>

LUPI

learning using privileged information



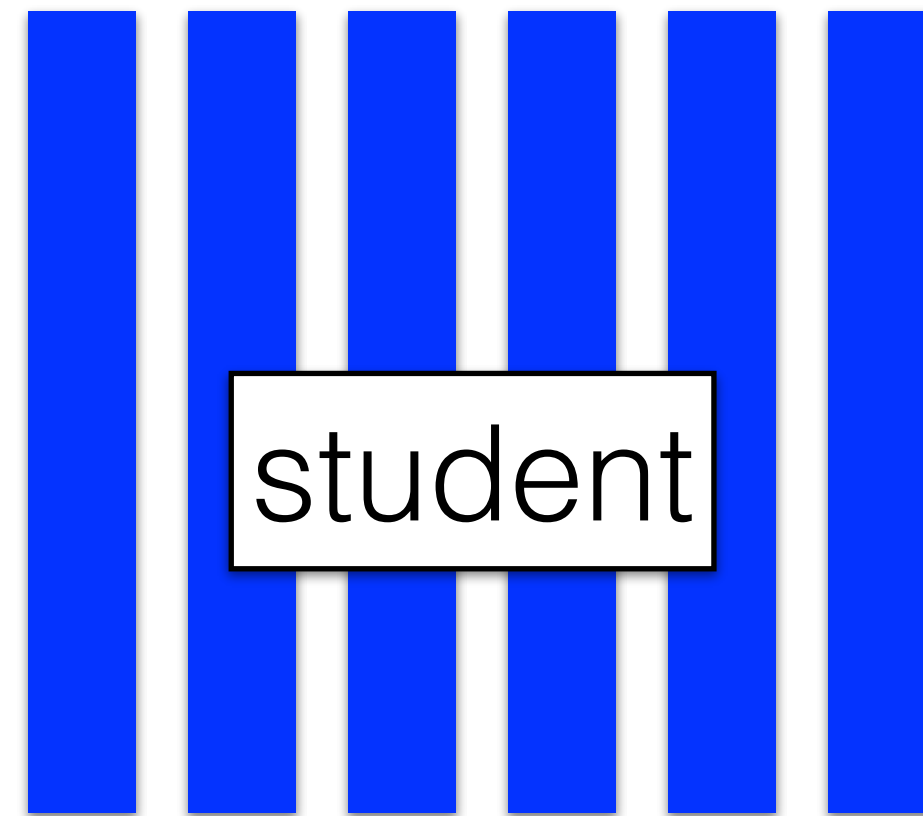
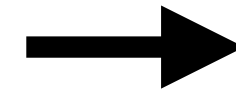
Annotated data



Train **teacher** on easier task
with the access to **privileged**
information

LUPI

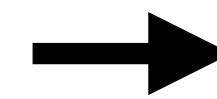
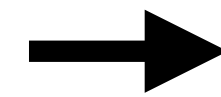
learning using privileged information



bike

Train **student** on teacher's outputs on **not annotated** data

Unannotated data



bike

Use teacher to classify unannotated data.



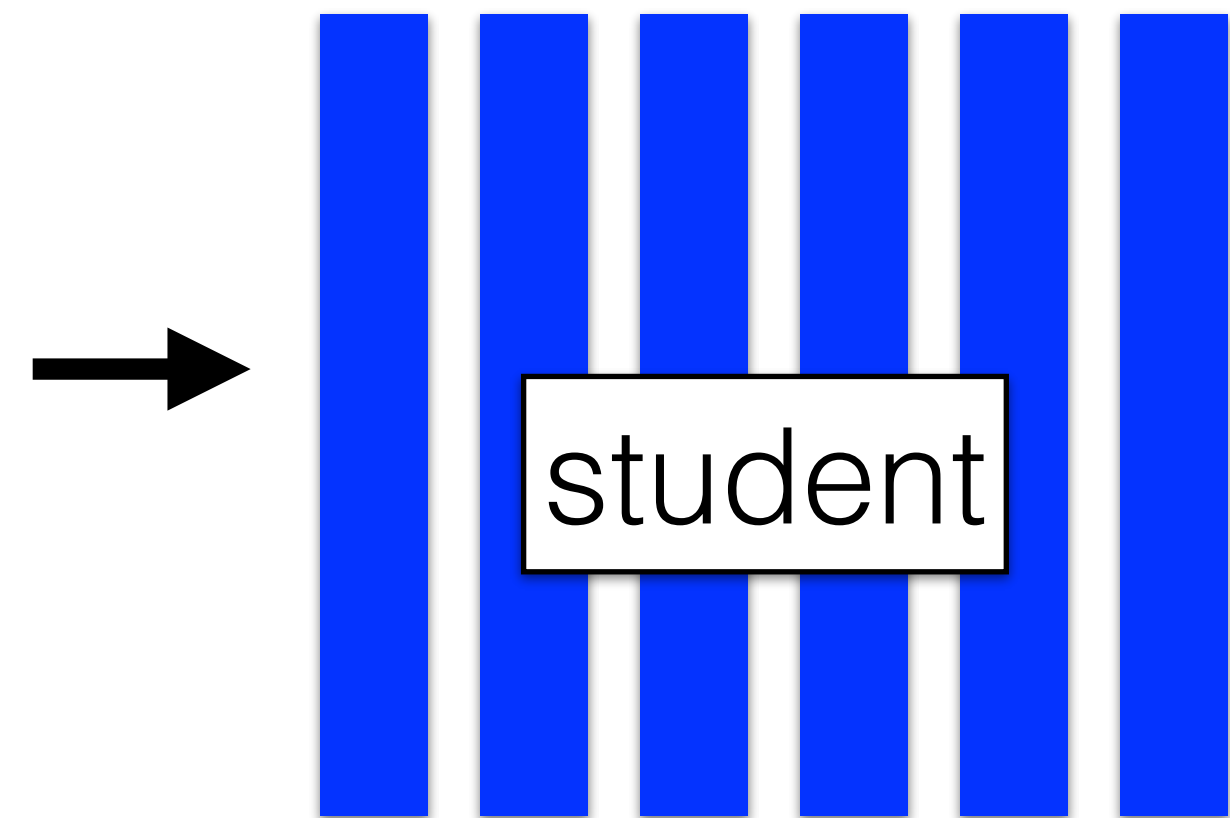


DINO

learning class-level features through the contrastive learning objective applied to the class token



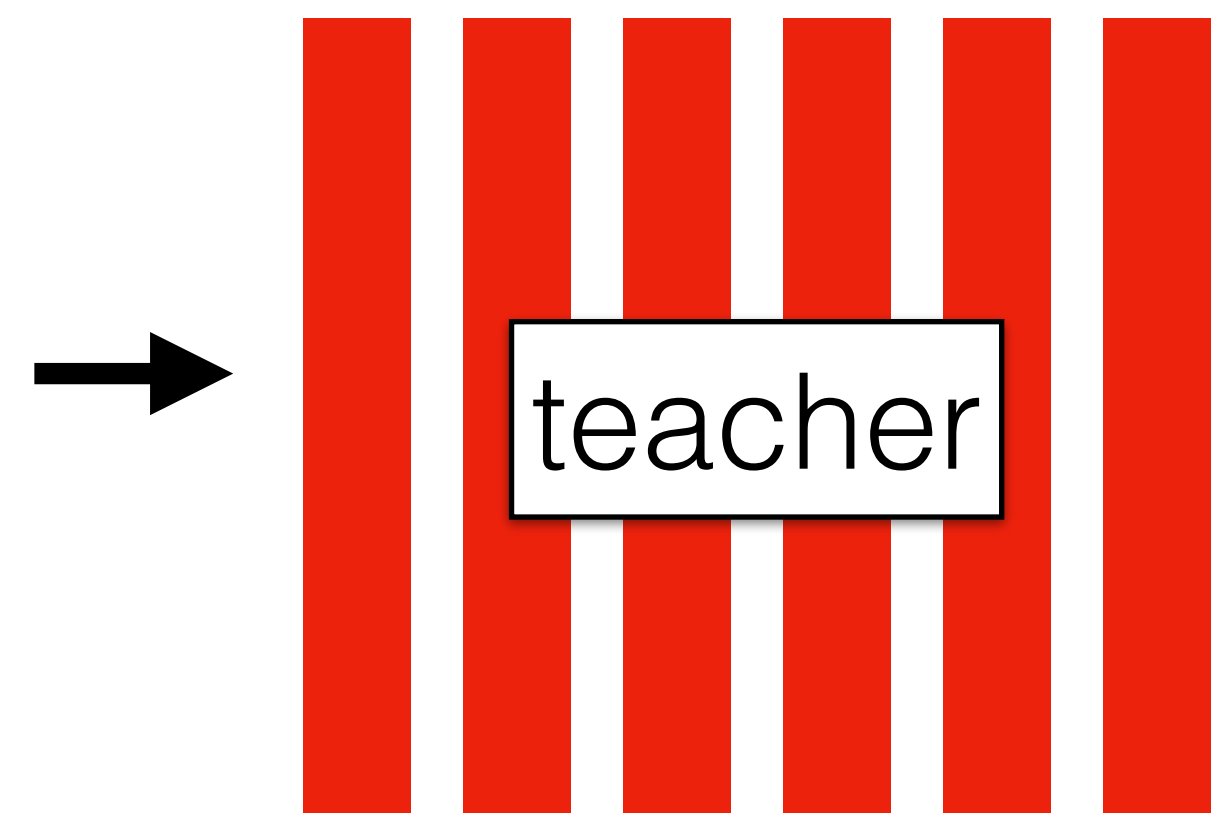
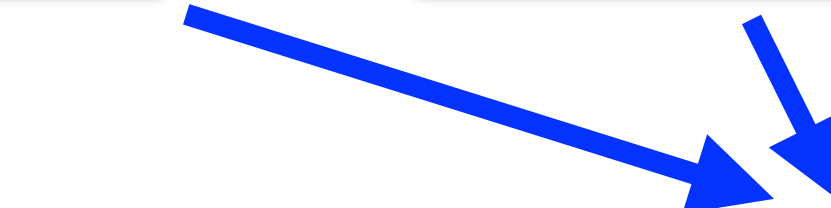
Cropped patches



\mathbf{s}_1

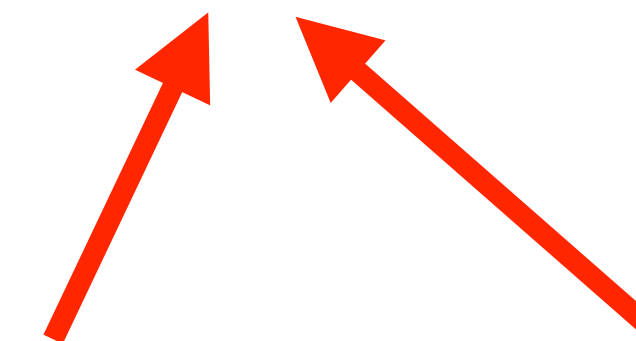
\mathbf{s}_2

$$\mathcal{L} = \sum_i \text{KL}(\text{softmax}(\mathbf{t}_i) \parallel \text{softmax}(\mathbf{s}_i))$$



\mathbf{t}_1

\mathbf{t}_2



Forces the student to predict context-aware embedding

DINO

learning class-level features through the contrastive learning objective applied to the class token

Collapse avoided by:

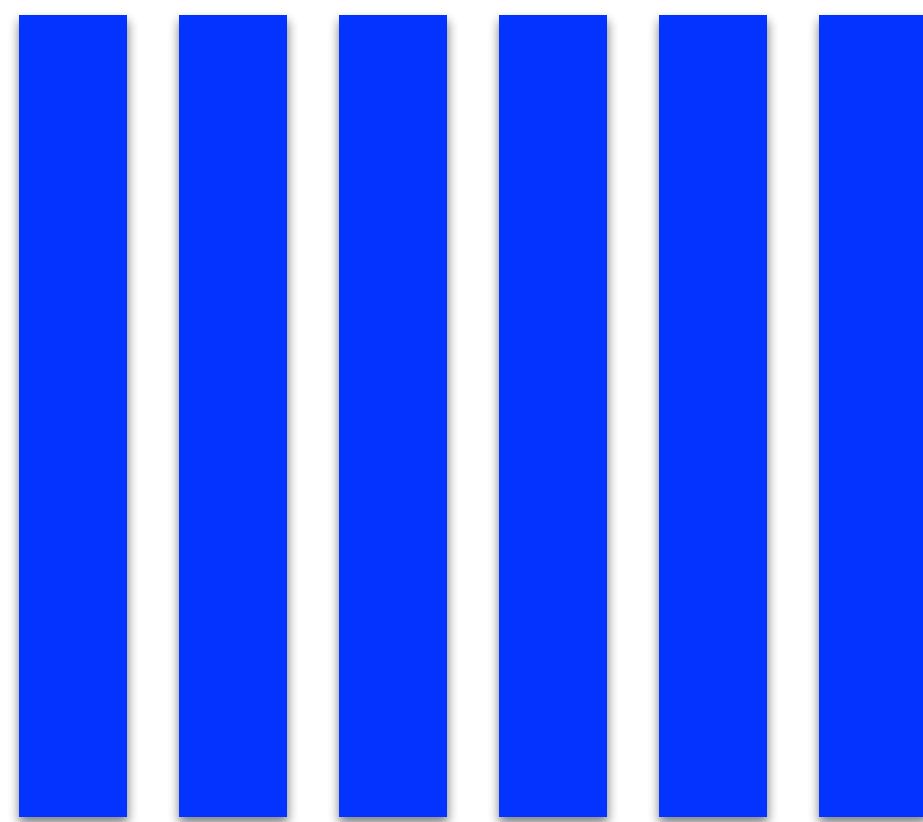
- Sharpening $p_k^{(\mathbf{t}_i)} = \frac{\exp(\mathbf{t}_{i,k}/\tau_t)}{\sum_j \exp(\mathbf{t}_{i,j}/\tau_t)}$
- Centering $\mathbf{t}'_i = \mathbf{t}_i - \mathbf{c}$

Pair of stereo images

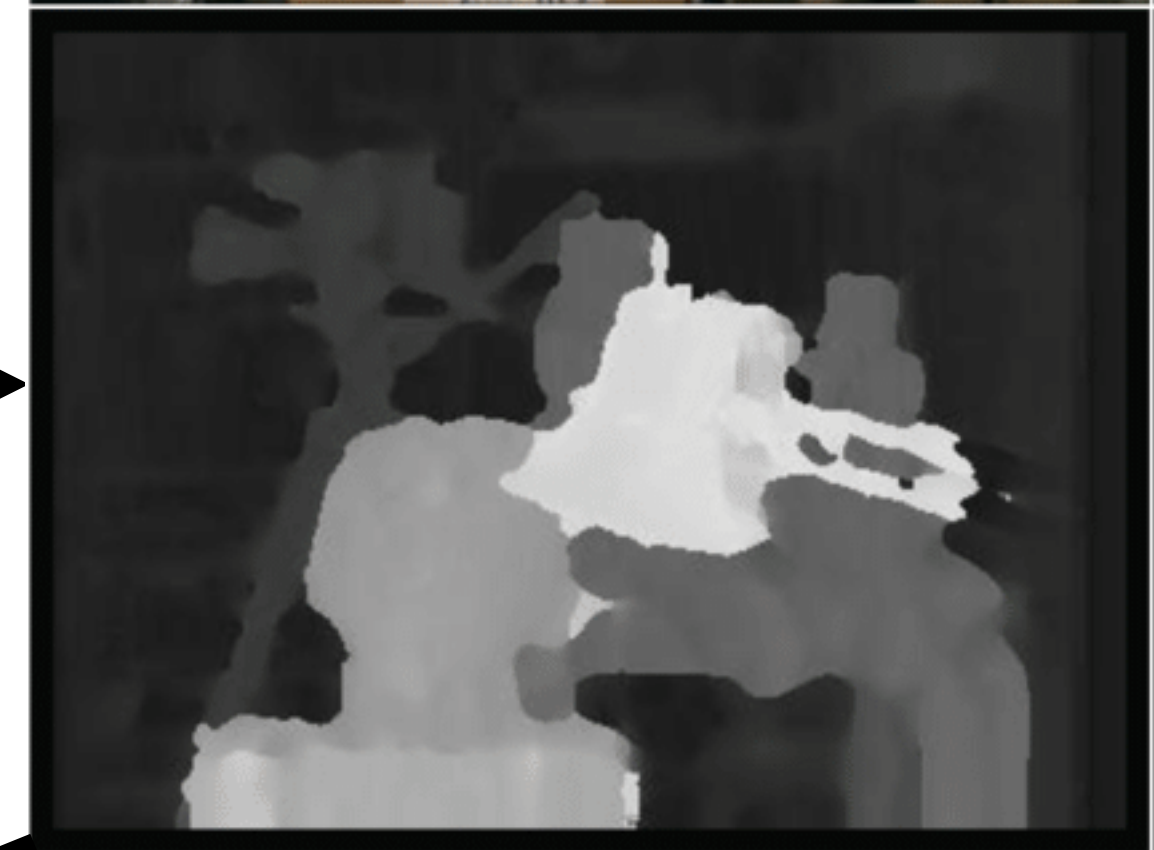
Left image



MonoDepth



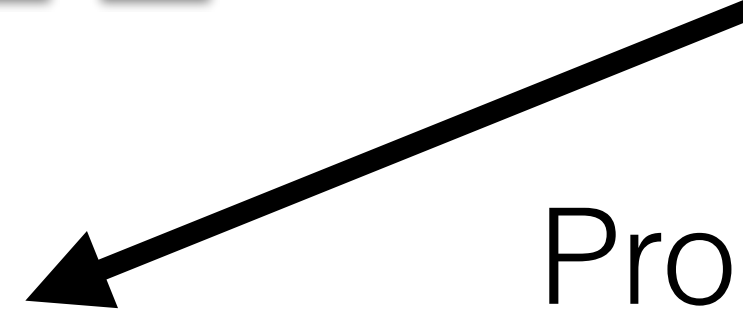
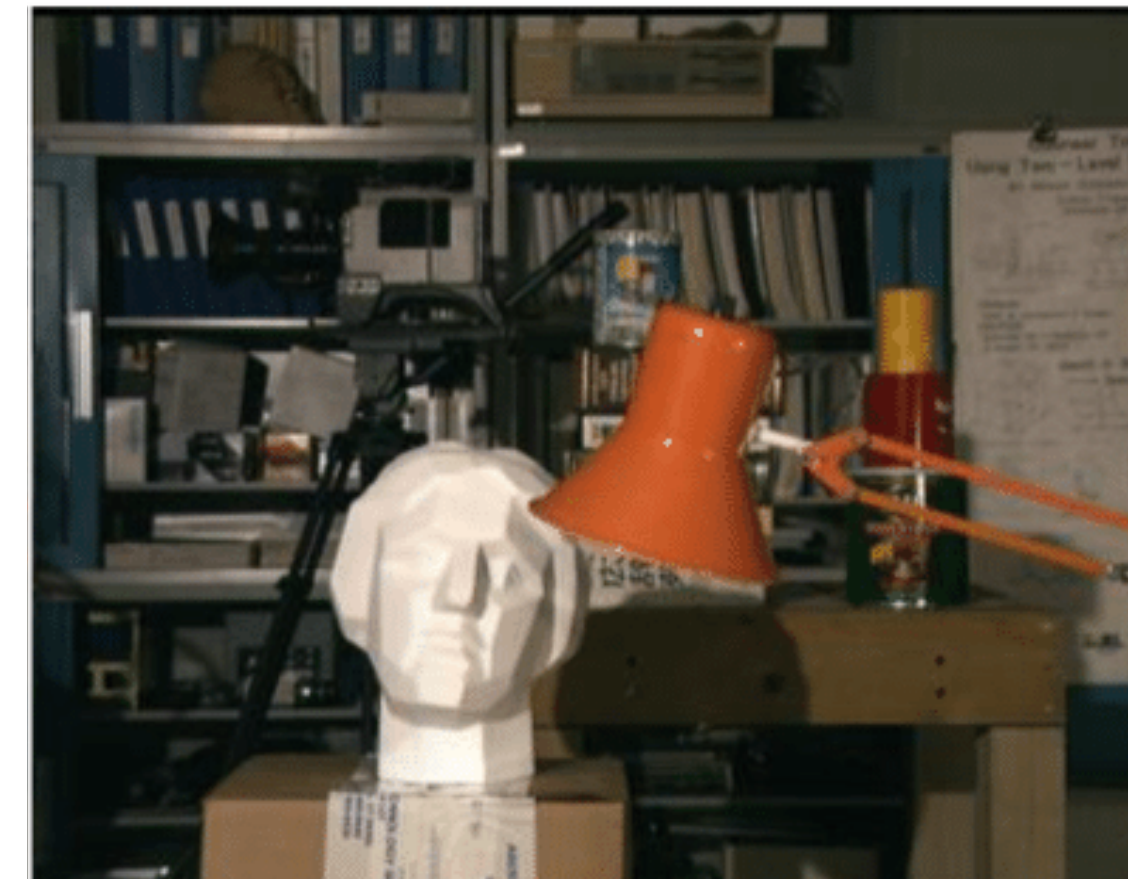
Predicted depth



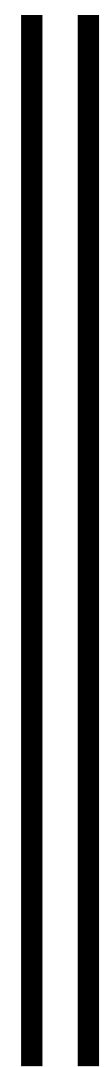
Right image



Projected
Right image



2



2

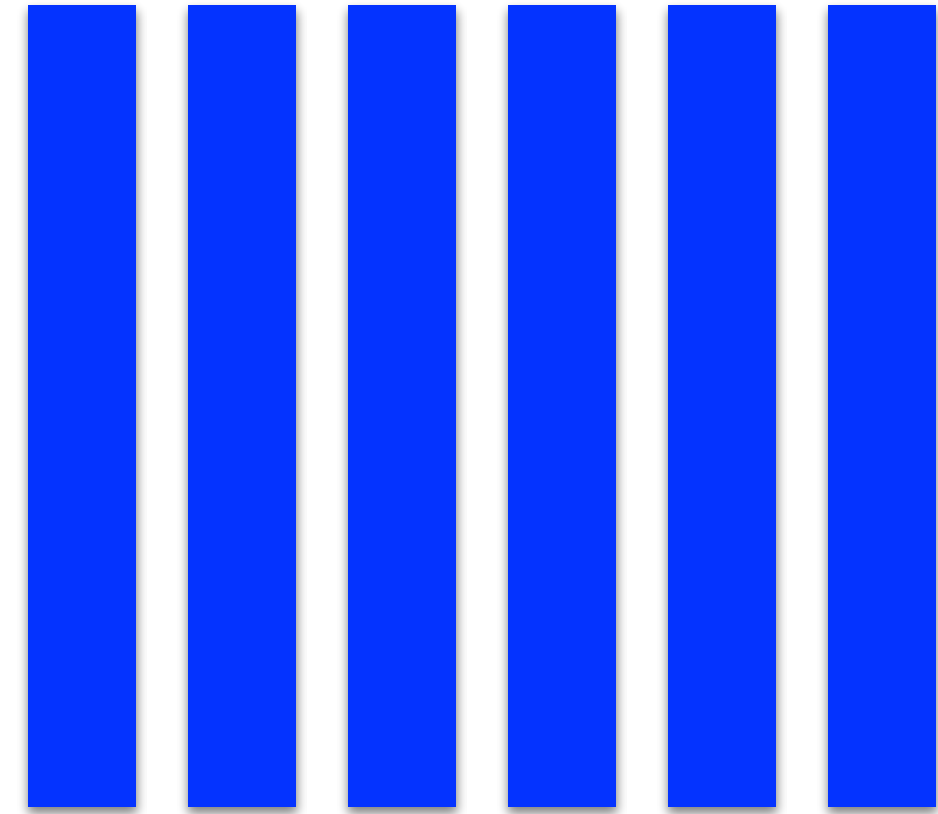
Project **left** image
through the **depth**
to **right** image using
known geometry

Minimize color inconsistency

Two cameras
looking at the
same scene

Co-learning

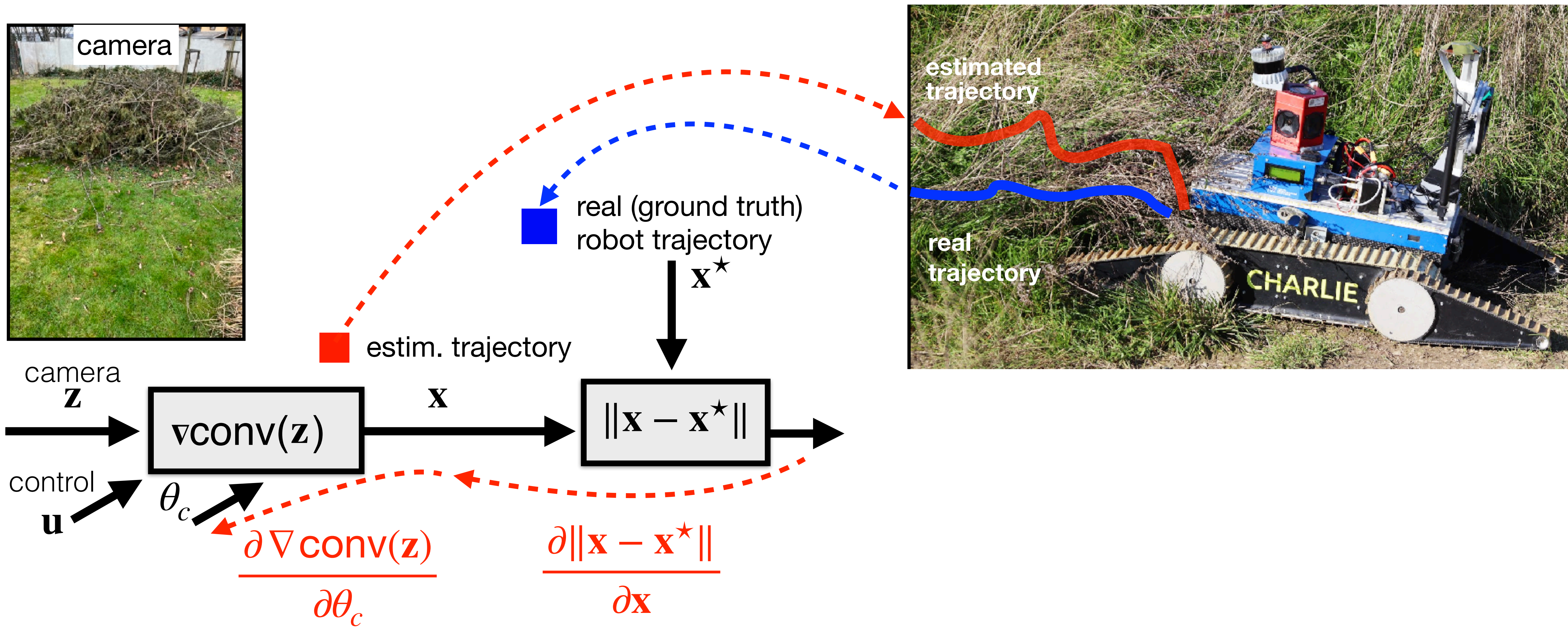
Detected humans



Project
detections
from **red** domain
to **blue** domain



Learning to mimic other sensor



Conclusions

- **ConvNets**

- enforced **local** attention
- data-independent **fixed attention**
- data-independent **fixed** kernel **weights**

- **Transformers**

- learned **global** attention
- data-dependent **dynamic attention** (different for different content from **Q, K**)
- data-dependent **dynamic weights** (different for different content from **V**)

- Big **foundation models** (such as SAM) delivered for various modalities images, depthmaps, pointclouds, text, speech

- In order to deliver billions of training data **self-supervision** is required