# Lecture 5: Attention, transformer architecture

Leonid Sanochkin

# RNN drawbacks

# Self-attention

Self-attention

input #1

| 1 | 0 | 1 | 0 |
|---|---|---|---|

input #2

| 0 | 2 | 0 | 2 |
|---|---|---|---|

input #3

| 1 | 1 | 1 | 1 |
|---|---|---|---|

# Self-attention

Self-attention

input #1

| 1 | 0 | 1 | 0 |

input #2

| 0 | 2 | 0 | 2 |

input #3

| 1 | 1 | 1 | 1 |

# Self-attention

# Self-attention

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

$$= Z$$

# Self-attention



Self-attention

output #1

addition    | 2.0 | 7.0 | 1.5 |

| key | | value | | key | | value | | key | | value | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 2 | 3 | 4 | 4 | 0 | 2 | 8 | 0 |

key: 0 1 1    value: 1 2 3    key: 4 4 0    value: 2 8 0    key: 2 3 1    value: 2 6 3

input #1    1 0 1 0

input #2    0 2 0 2

input #3    1 1 1 1

# Self-attention



walk by river bank

Input embeddings

Queries

Scalar product

Keys

Scaling/
Softmax

Linear combination

Values

Contextualized embeddings
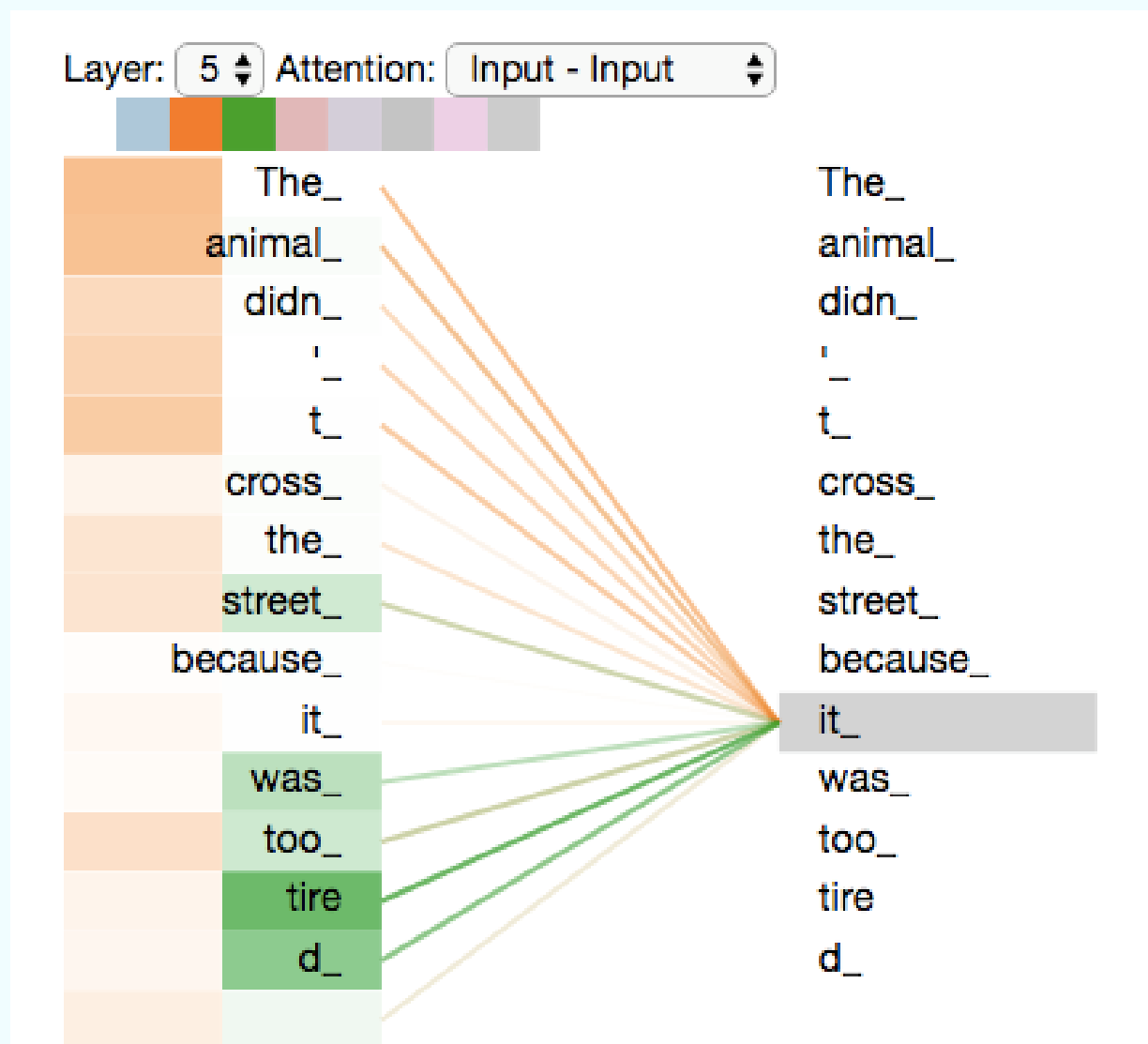
# Multihead self-attention

# Multihead self-attention

# Positional encodings

$$\text{PE}_{pos,2i} = \sin(pos/10000^{2i/d_{model}}),$$

$$\text{PE}_{pos,2i+1} = \cos(pos/10000^{2i/d_{model}}),$$

# Positional encodings

# Transformer architecture



Residual connections
and layer normalization

Feed-forward network:
after taking information from
other tokens, take a moment to
think and process this information

Feed-forward network:
after taking information from
other tokens, take a moment to
think and process this information

Encoder self-attention:
tokens look at each other

queries, keys, values
are computed from
encoder states

Decoder-encoder attention:
target token looks at the source

queries – from decoder states; keys
and values from encoder states

Decoder self-attention (masked):
tokens look at the previous tokens

queries, keys, values are computed
from decoder states

walk | by | the | river | bank

| -0.02 | 0.02 | -0.00 | -0.04 | -0.01 |
| :---: | :---: | :---: | :---: | :---: |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| -0.03 | -0.07 | 0.03 | -0.03 | -0.04 |

Input embedding

Linear projections

Keys

Queries

Values

Layer 1

**Multi-head attention**

64

64

12

Add Norm

**BERT**

Linear projections

Add Norm

Feed-forward layer

x11 layers

Output

| 0.51 | 0.30 | -0.13 | 0.29 | 0.29 |
| :---: | :---: | :---: | :---: | :---: |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| -0.03 | 0.08 | 0.15 | 0.13 | 0.10 |

≈

waterside

AIRI

AIRI

airi.net