

Technical details of explicharr

Kuan Yu¹

September 11, 2018

¹kuanyu@uni-potsdam.de

explicharr³

- ▶ sentence simplification with
- ▶ character-level
- ▶ transformer²

“It is located in Potsdam .” \mapsto “It is in Potsdam .”

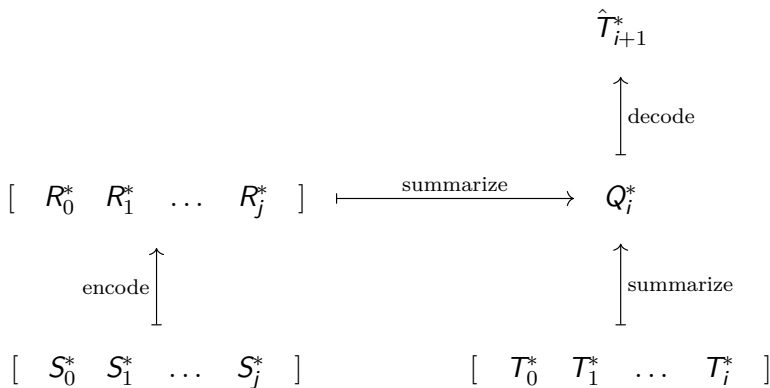
model: $S^* \rightarrow T^*$ where

- ▶ S = the source alphabet
- ▶ T = the target alphabet

²<https://arxiv.org/abs/1706.03762>

³<https://github.com/srewai/explicharr>

encoder-decoder, seq-to-seq, autoregressive



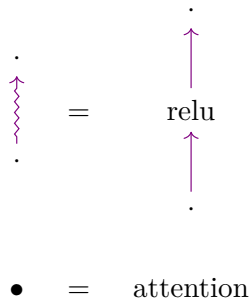
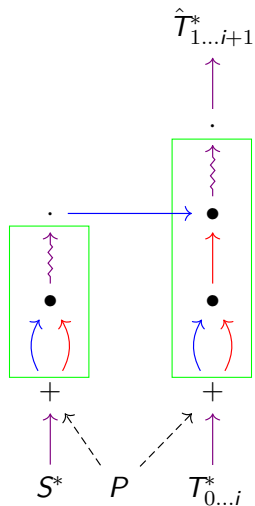
soft attention

given a **query** vector and multiple **value** vectors

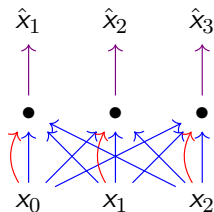
attention: $\downarrow \bullet \downarrow \downarrow \dots \downarrow \mapsto \downarrow$

- ▶ compute a weight for each value, according to the query
- ▶ normalize the weights with softmax
- ▶ take the weighted sum of the values

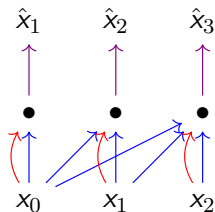
transformer



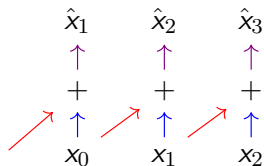
self-attention⁴



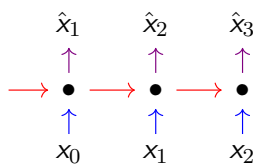
encoder self-attention



decoder self-attention



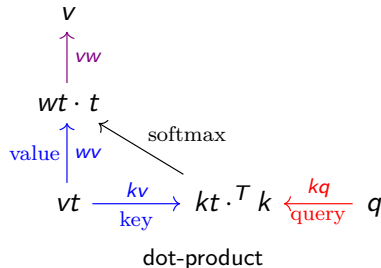
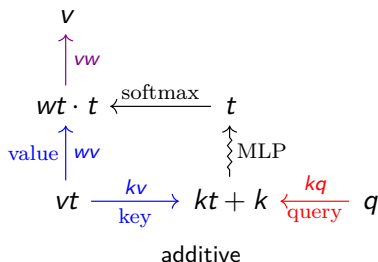
convolutional



recurrent

⁴<https://arxiv.org/abs/1606.01933>

attention cells: additive⁵ vs dot-product⁶, key-value⁷



dimensions: **t**ime, **q**uery, **k**ey, **v**alue, **w** intermediate

$$A \cdot B := AB$$

$$A \cdot^T B := A^T B$$

⁵<https://arxiv.org/abs/1409.0473>

⁶<https://arxiv.org/abs/1508.04025>

⁷<https://arxiv.org/abs/1702.04521>

transformer attention

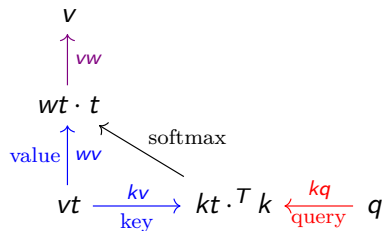
scaled dot-product

- ▶ divide weights by \sqrt{k} before applying softmax
- ▶ raise temperature
- ▶ lower variance

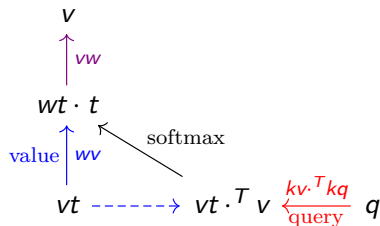
multi-head attention

- ▶ split spaces (query, value, key) into disjoint subspaces (subquery, subvalue, subkey)
- ▶ one attention head for each split
- ▶ concatenate the resulting subvectors

key transformation

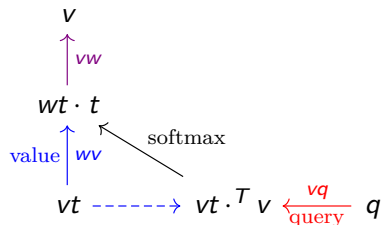


linear

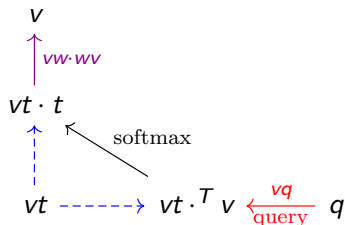


identity

value transformation

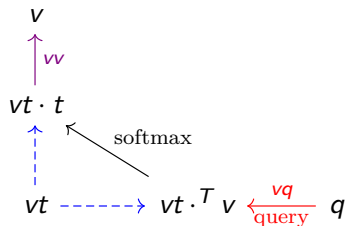


linear

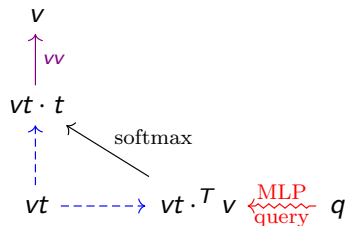


identity

query transformation

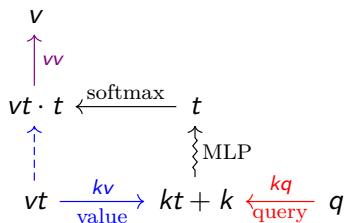


linear

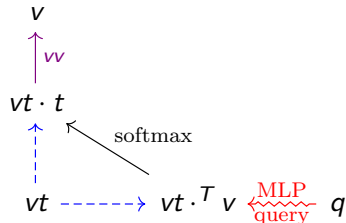


non-linear

additive vs dot-product with non-linear query

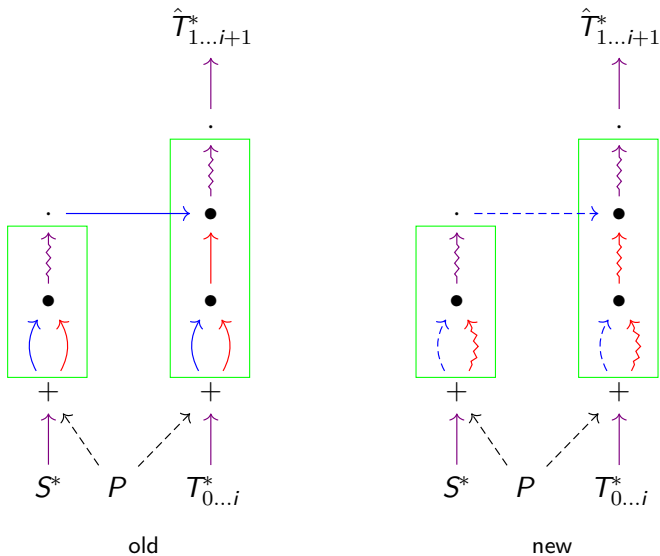


additive



dot-product

transformer



architecture

- ▶ 2 encoder layers, 2 decoder layers
- ▶ 2 input embedding layers, 1 output softmax layer
- ▶ 256 representation dimension, 512 relu in MLPs
- ▶ single-head scaled dot-product attention
- ▶ dropout⁸, residual connection⁹, layer normalization¹⁰ after each attention or MLP sublayer

⁸<http://jmlr.org/papers/v15/srivastava14a.html>

⁹<https://arxiv.org/abs/1512.03385>

¹⁰<https://arxiv.org/abs/1607.06450>

training

- ▶ cross entropy loss with label smoothing¹¹
- ▶ teacher forcing
- ▶ batch size 64
- ▶ ~6 minutes per epoch (~223k instances)
- ▶ for 180 epochs

¹¹<https://arxiv.org/abs/1512.00567>

introspection

- ▶ greedy autoregressive decoding
- ▶ attention weight matrix

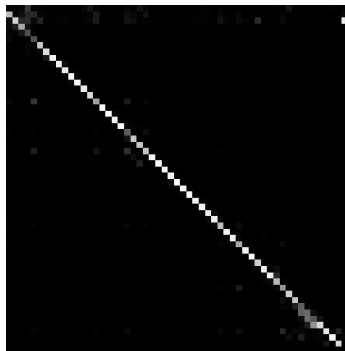
self-attention

- ▶ always a diagonal line
- ▶ encoder layer 1 and decoder layer 2 slightly fuzzy

introspection: normal

The enshrined kami is Isonotakeru no mikoto (五十猛命?)

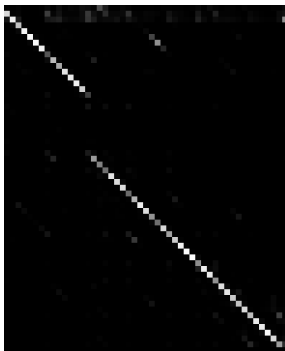
The enshrined kami is Isonotakeru no mikoto (____ ?)



introspection: skip

Sodium iodate (NaIO_3) is the sodium salt of iodic acid .

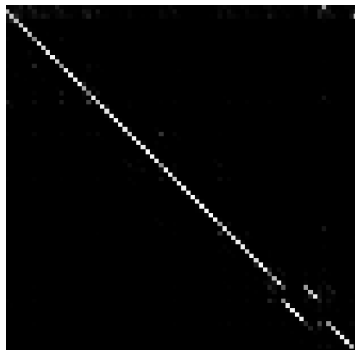
Sodium iodate is the sodium salt of iodic acid .



introspection: swap

Their first child , Prince George of Cambridge , was born on 22 July 2013 .

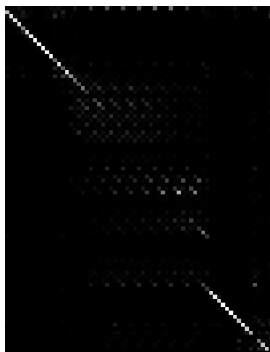
Their first child , Prince George of Cambridge , was born on July 22 , 2013 .



introspection: confused

For example , $2 + 3 = 3 + 2 = 5$ and $2 \cdot 3 = 3 \cdot 2 = 6$, but $2^3 = 8$, whereas $3^2 = 9$.

For example , $2 + 3 = 3 = 3 = 3 = 3 = 3 = 3 = 3$, whereas $3^2 = 9$.



introspection: really confused

1 and - 1 divide (are divisors of) every integer .

In the 19th century , the state of the state of the state
of the state of the state of the state of the state .

