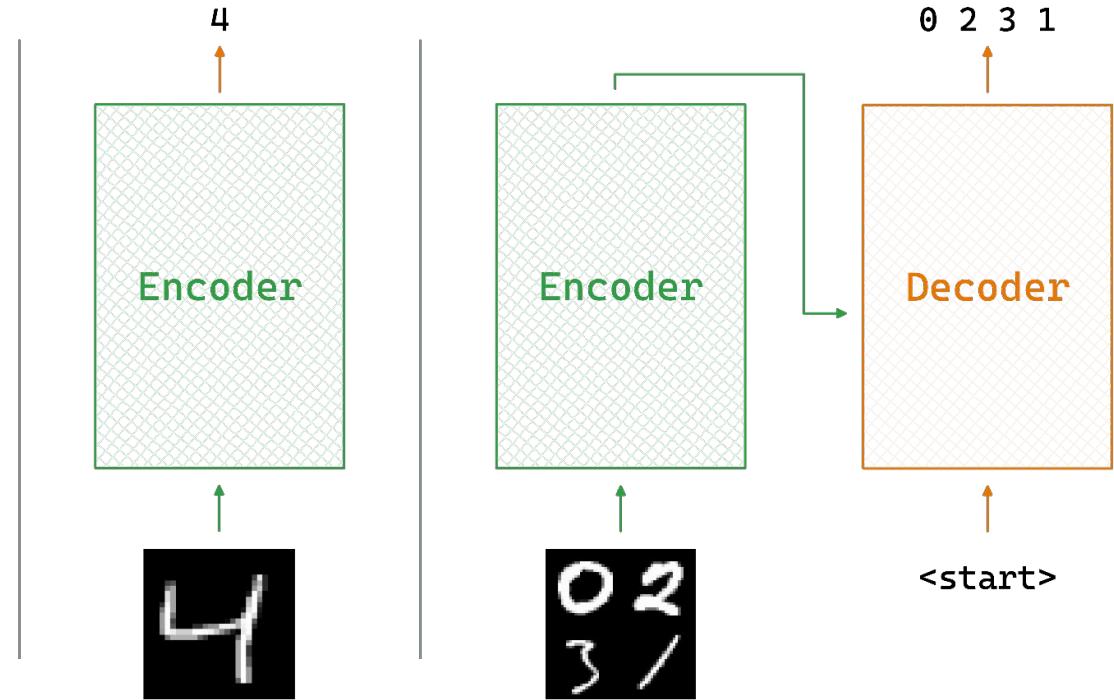


Attention Is All You Need

A bit of a Q&A

Task

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |



Paper

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

Attention Is All You Need

Ashish Vaswani*
 Google Brain
 avaswani@google.com

Noam Shazeer*
 Google Brain
 noam@google.com

Niki Parmar*
 Google Research
 nixp@google.com

Jakob Uszkoreit*
 Google Research
 usz@google.com

Llion Jones*
 Google Research
 llion@google.com

Aidan N. Gomez[†]
 University of Toronto
 aidan@cs.toronto.edu

Lukasz Kaiser^{*}
 Google Brain
 lukasz.kaiser@google.com

Illia Polosukhin[‡]
 illia.polosukhin@gmail.com

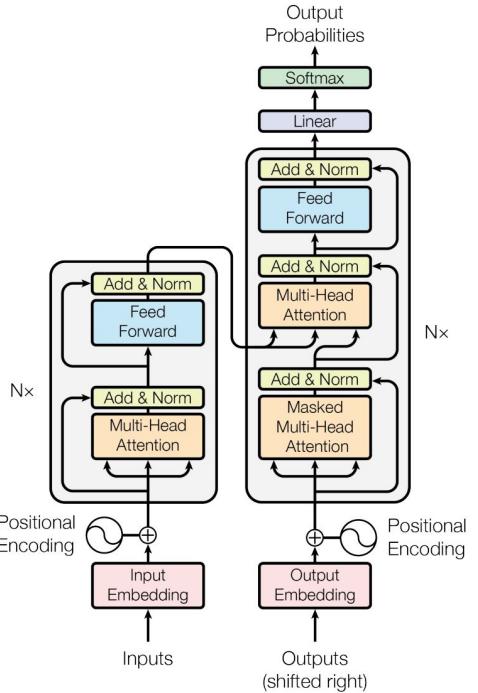
Abstract

The dominant sequence-to-sequence translation models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, without recurrent layers or convolutional layers entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the previous best using neural language ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training time required by the best prior work in literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

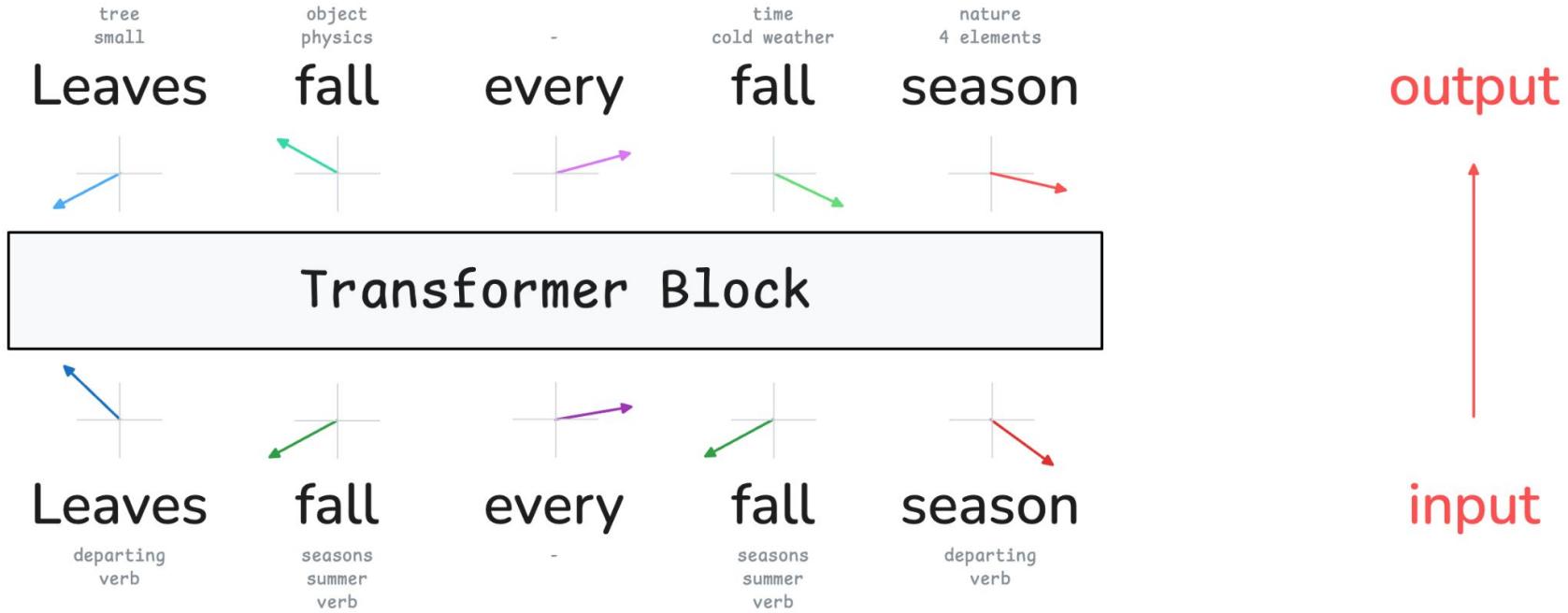
*Equal contribution. [†]Using older codebase. [‡]Work performed while at Google Brain.
[†]Work performed while at Google Research.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

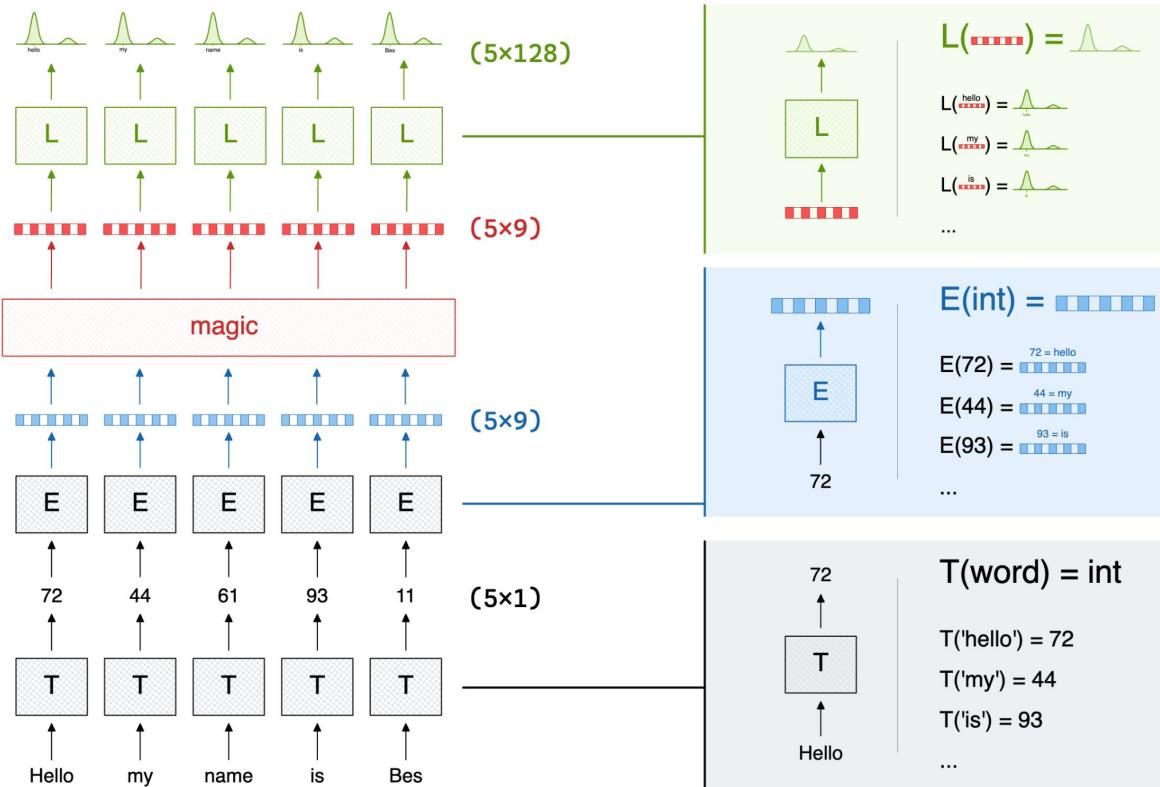
Vaswani et al. 2017



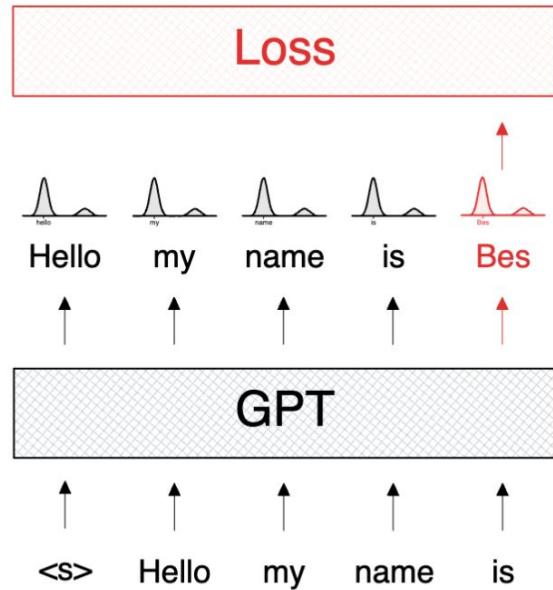
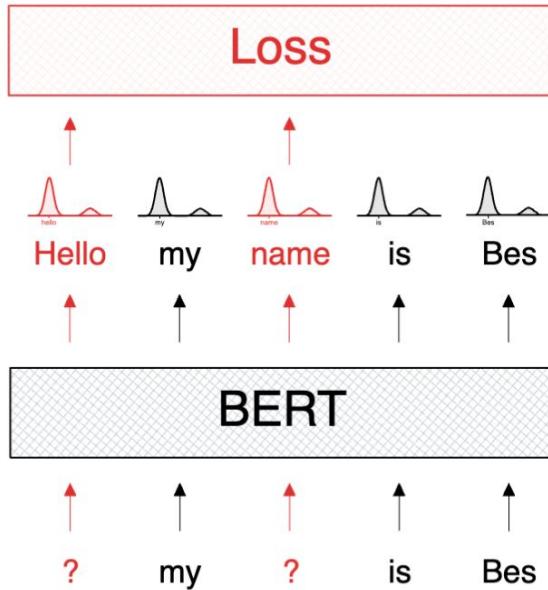
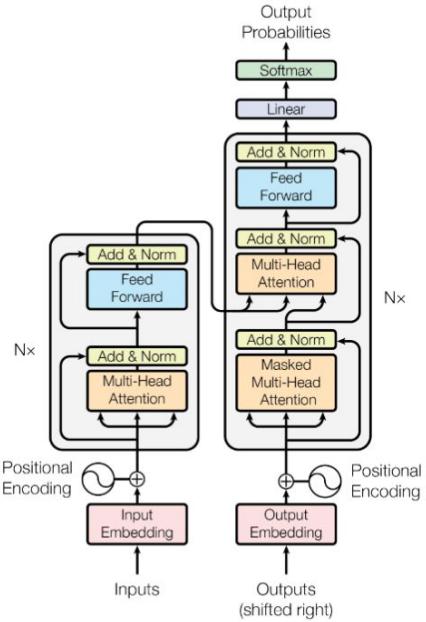
Why “transformer”?



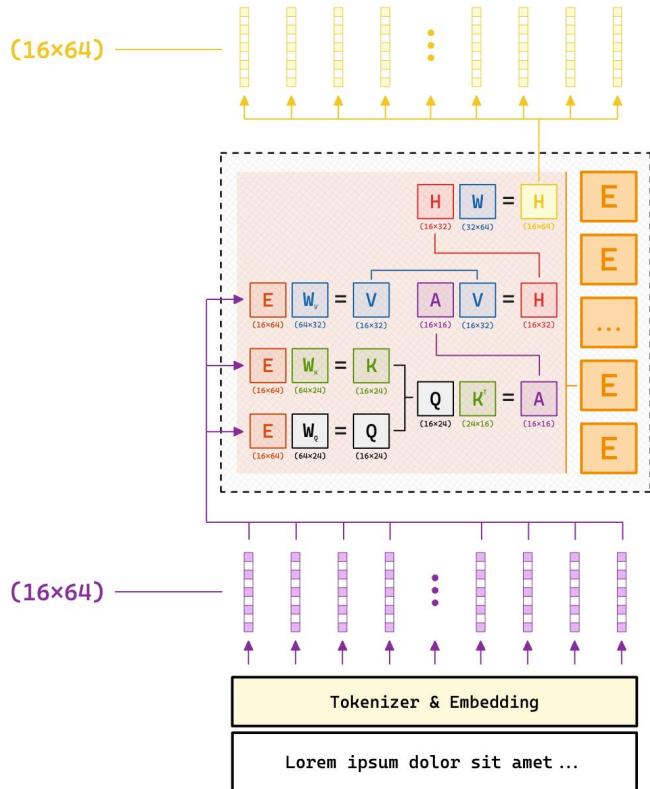
word2vec → transformer



Two Flavours



Attention



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

<https://github.com/besarthoxhaj/attention>

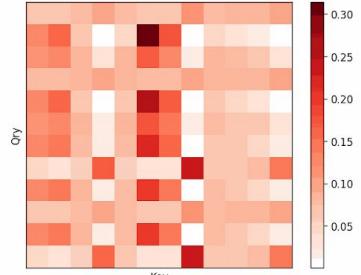
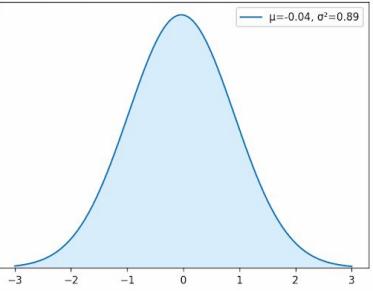
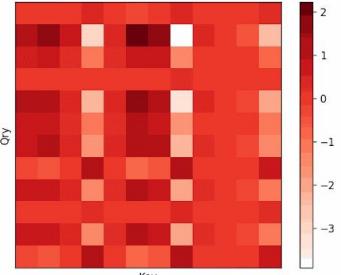
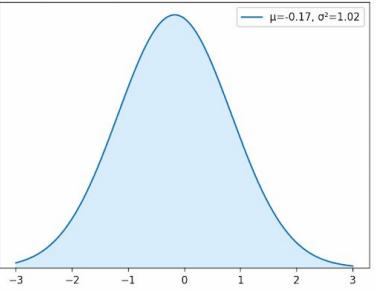


Let me show you :)

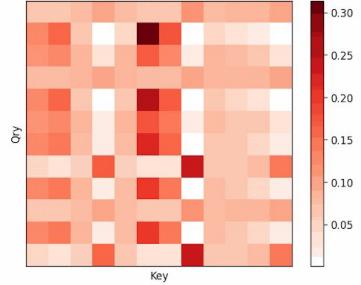
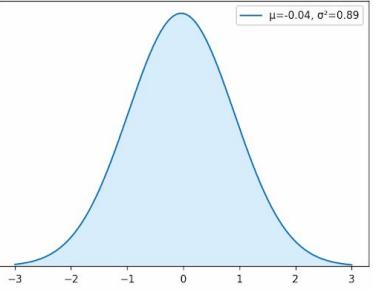
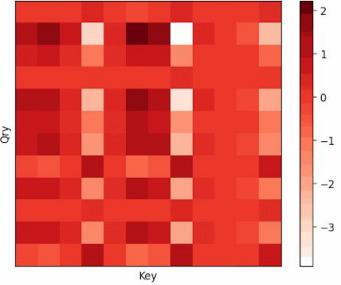
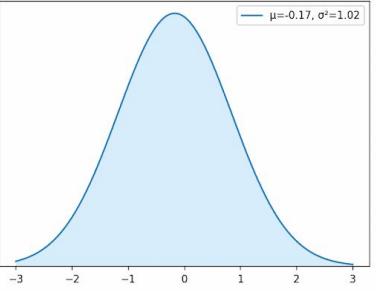
Animation

Embedding size: 1

$\text{softmax}(QK^T)$



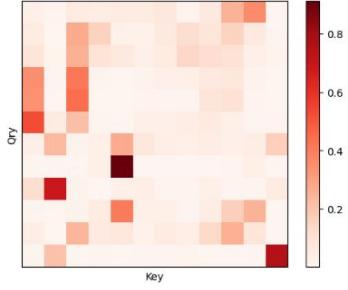
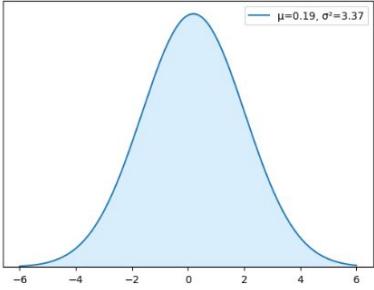
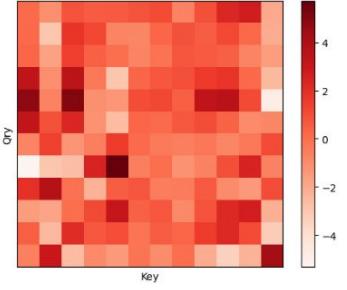
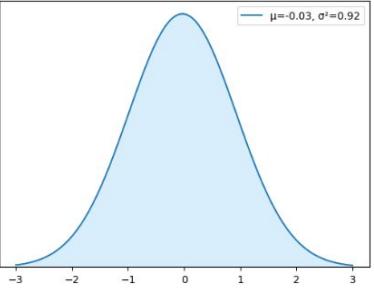
$\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$



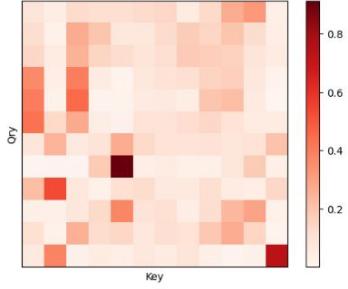
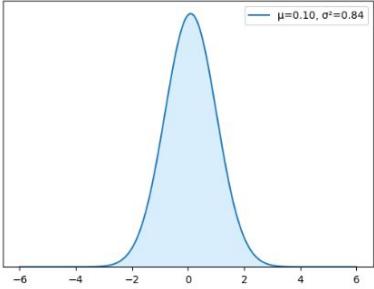
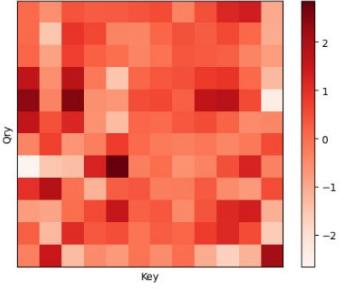
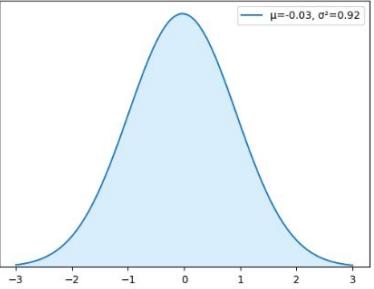
Toy

Embedding size: 4

$\text{softmax}(QK^T)$



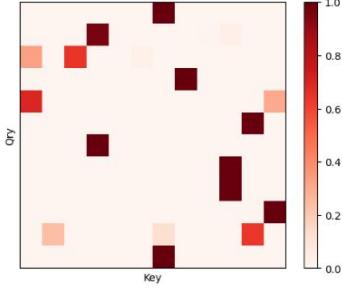
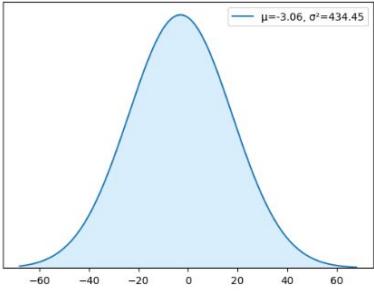
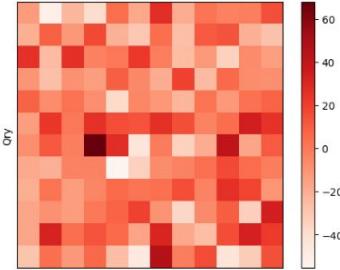
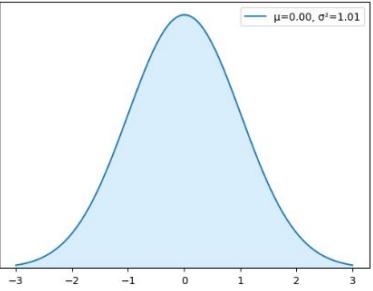
$\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$



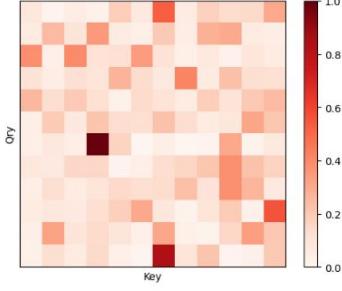
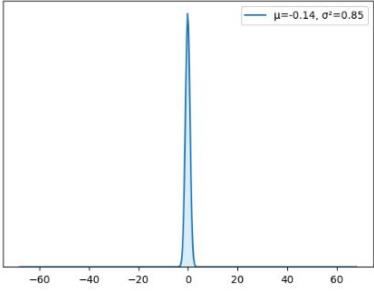
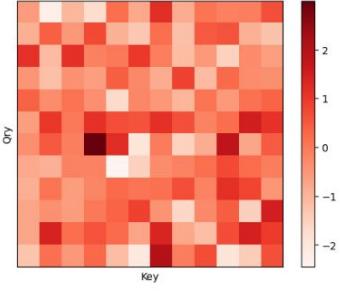
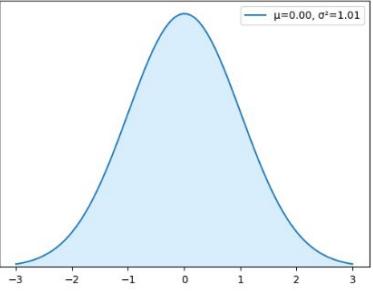
Production

Embedding size: 512

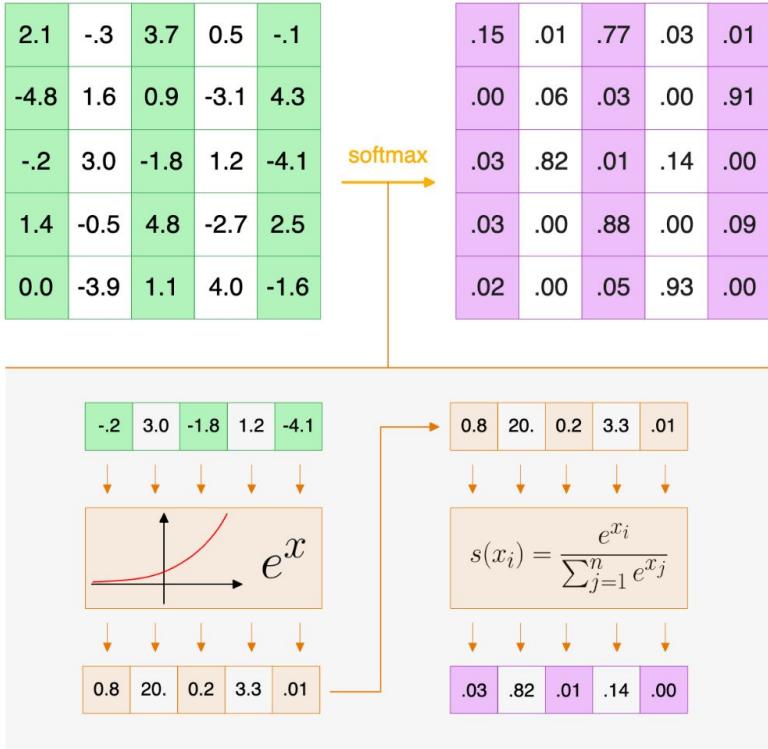
$\text{softmax}(QK^T)$



$\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$



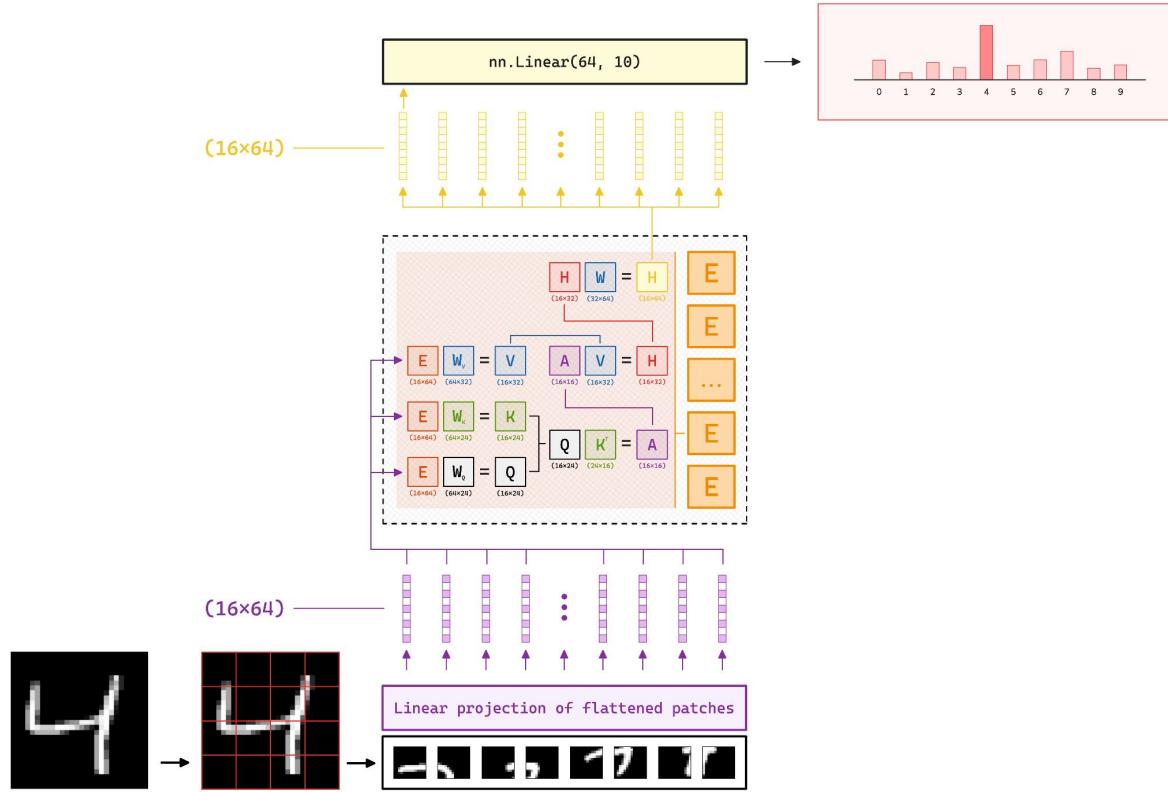
Attention Is Relative



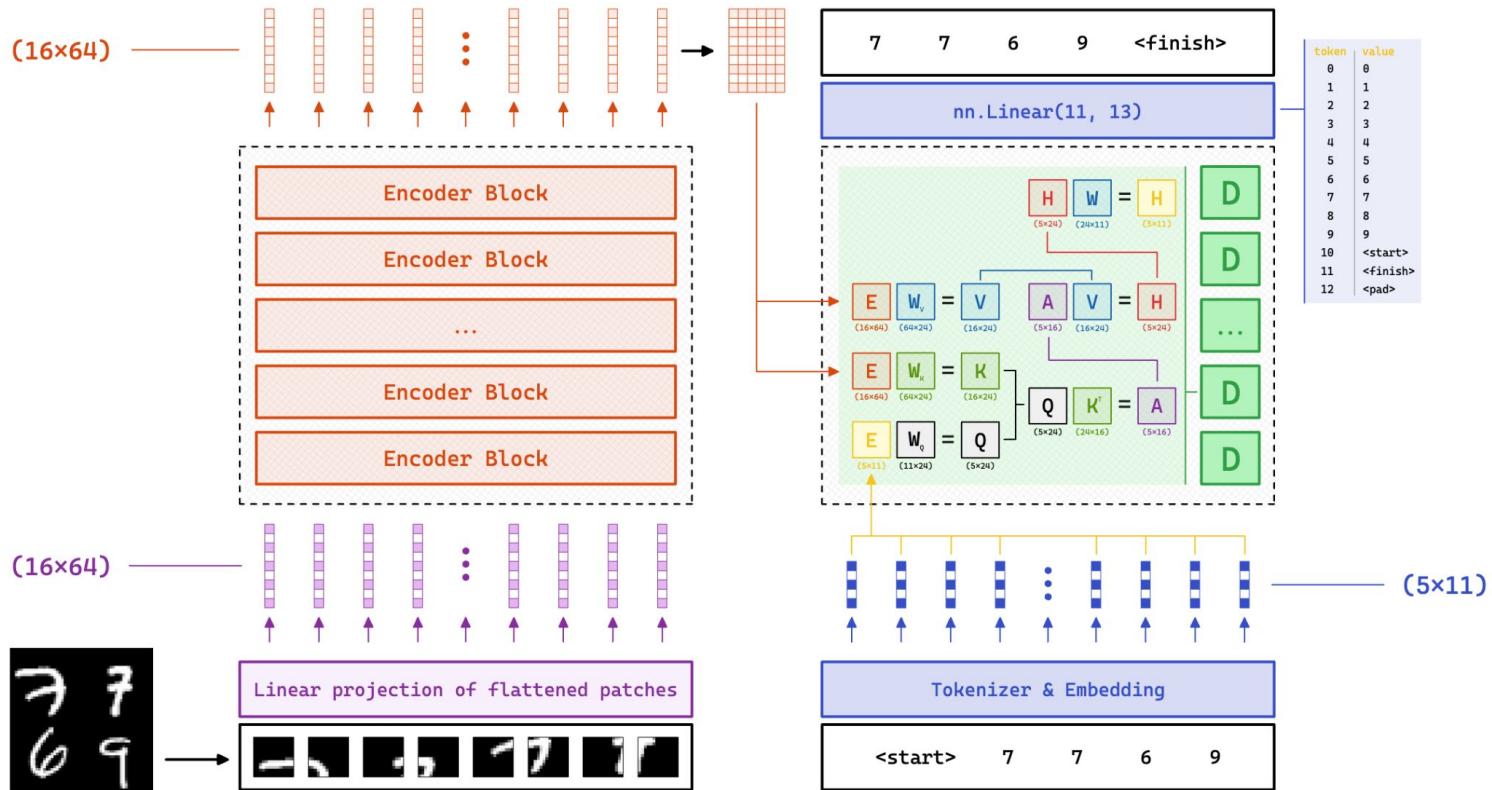
$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

- Guarantees strictly positive outputs
- Translational invariance
- Linear differences \mapsto multiplicative odds
- Maximum-entropy derivation
- Smooth, everywhere-differentiable, and strongly convex
- “Explosive” growth is tunable
- Numerical stability tricks
- Could we use something else?

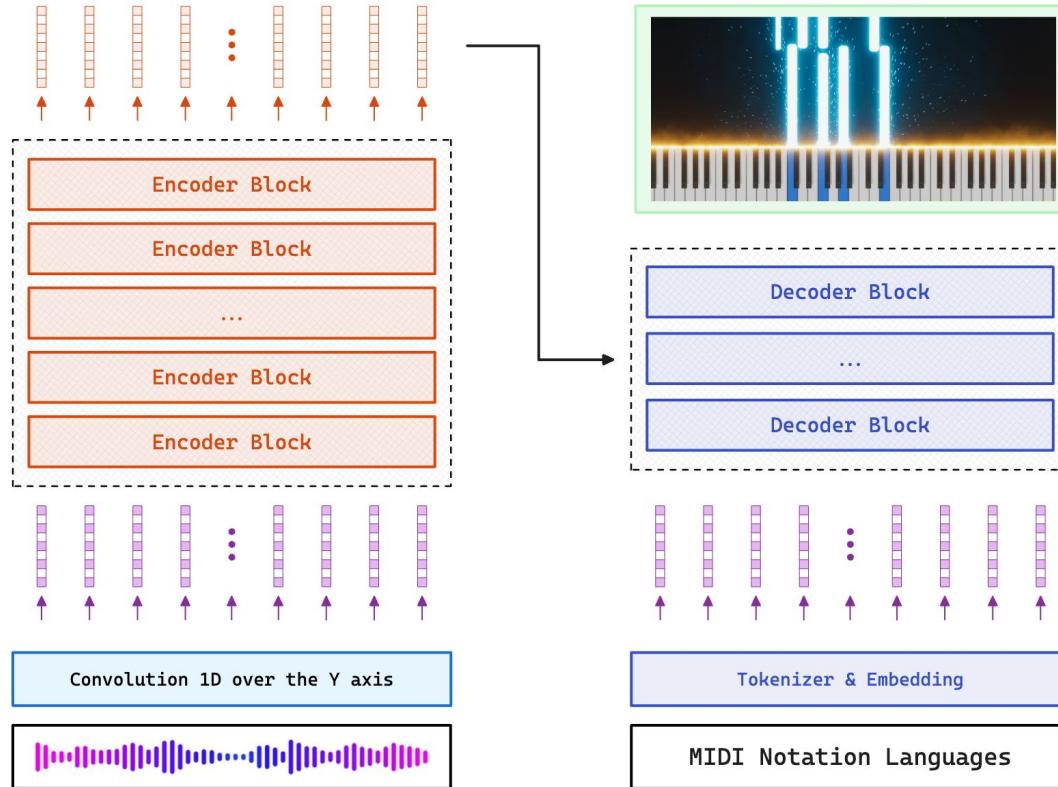
Encoder



Transformer



Audio





Thank you!