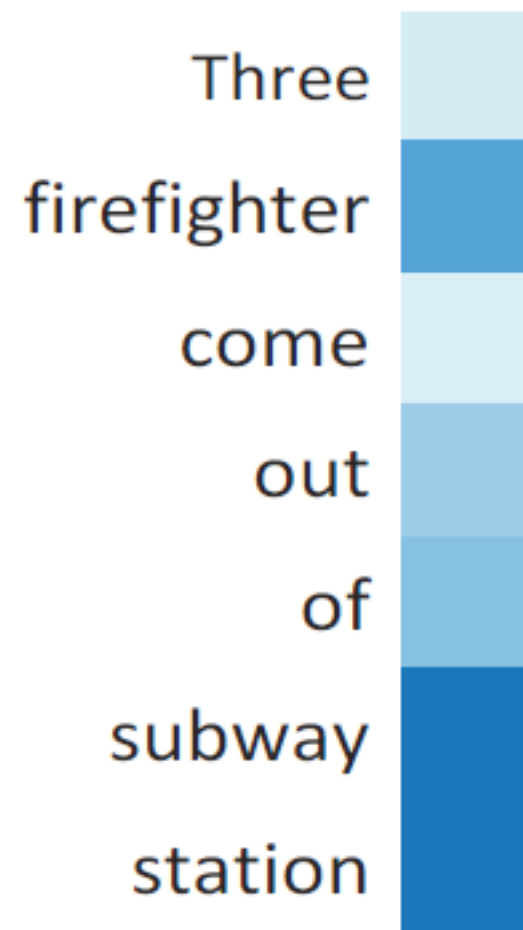# Attention

# Attention

# Visual attention



Feature before mask — Soft attention mask — Feature after mask

High-level part feature
Balloon instance mask → Classification

# SQuAD1

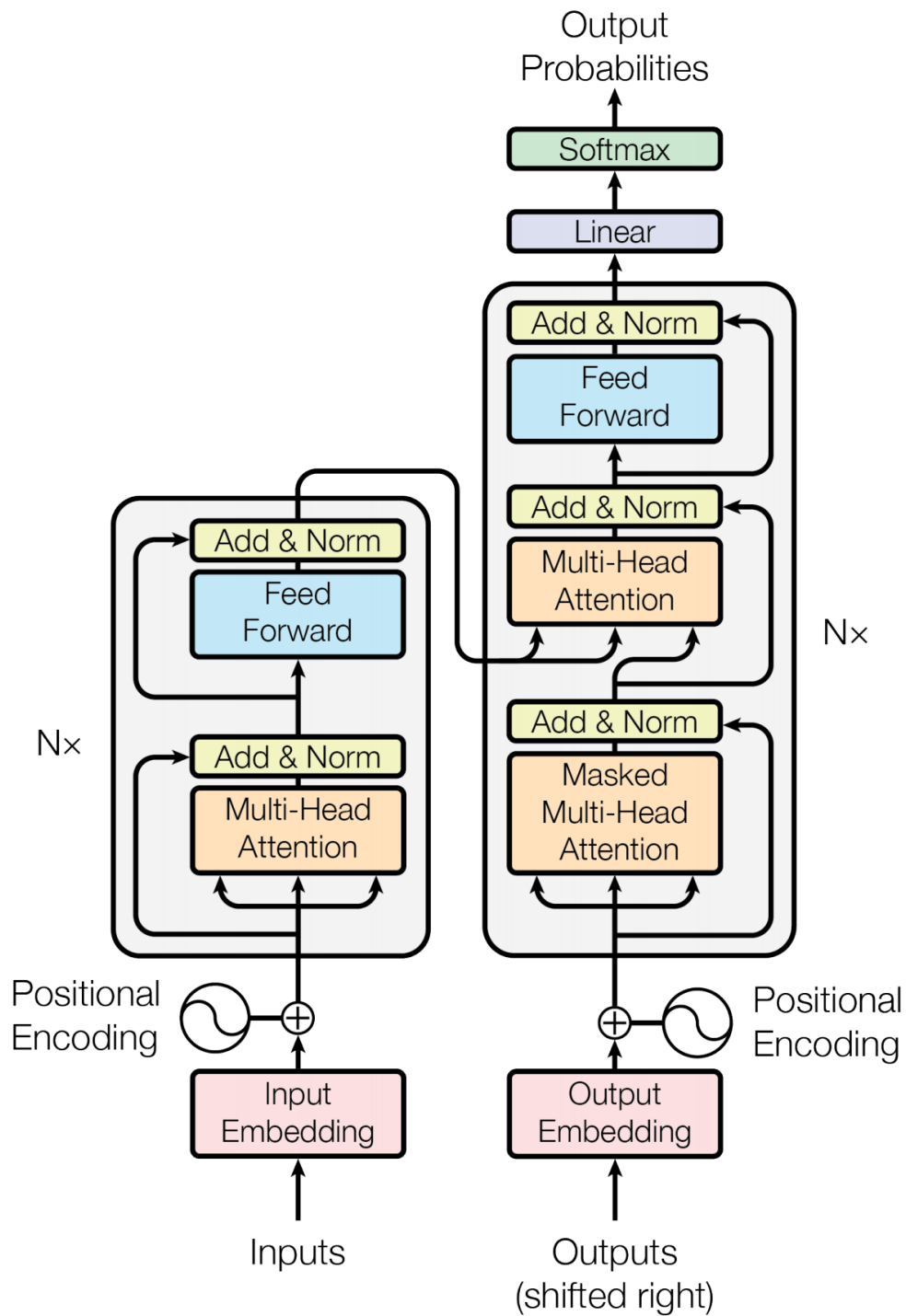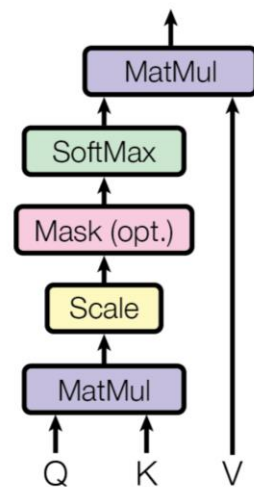| Rank | Model | EM | F1 |
|:---:|:---:|:---:|:---:|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Sep 13, 2018 | nlnet (single model)<br>*Microsoft Research Asia* | **74.238** | **77.022** |
| 2<br>Sep 17, 2018 | Unet (ensemble)<br>*Fudan University & Liulishuo Lab* | 71.553 | 75.011 |
| 2<br>Aug 15, 2018 | Reinforced Mnemonic Reader + Answer Verifier (single model)<br>*NUDT*<br>https://arxiv.org/abs/1808.05759 | 71.699 | 74.238 |
| 2<br>Aug 28, 2018 | SLQA+ (single model)<br>*Alibaba DAMO NLP*<br>http://www.aclweb.org/anthology/P18-1158 | 71.451 | 74.422 |
| 3<br>Sep 14, 2018 | SAN (ensemble model)<br>*Microsoft Business Applications Research Group*<br>https://arxiv.org/abs/1712.03556 | 71.282 | 73.658 |

# SQuAD2

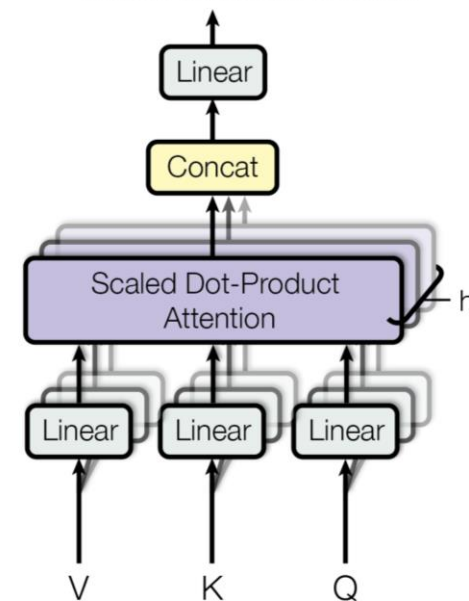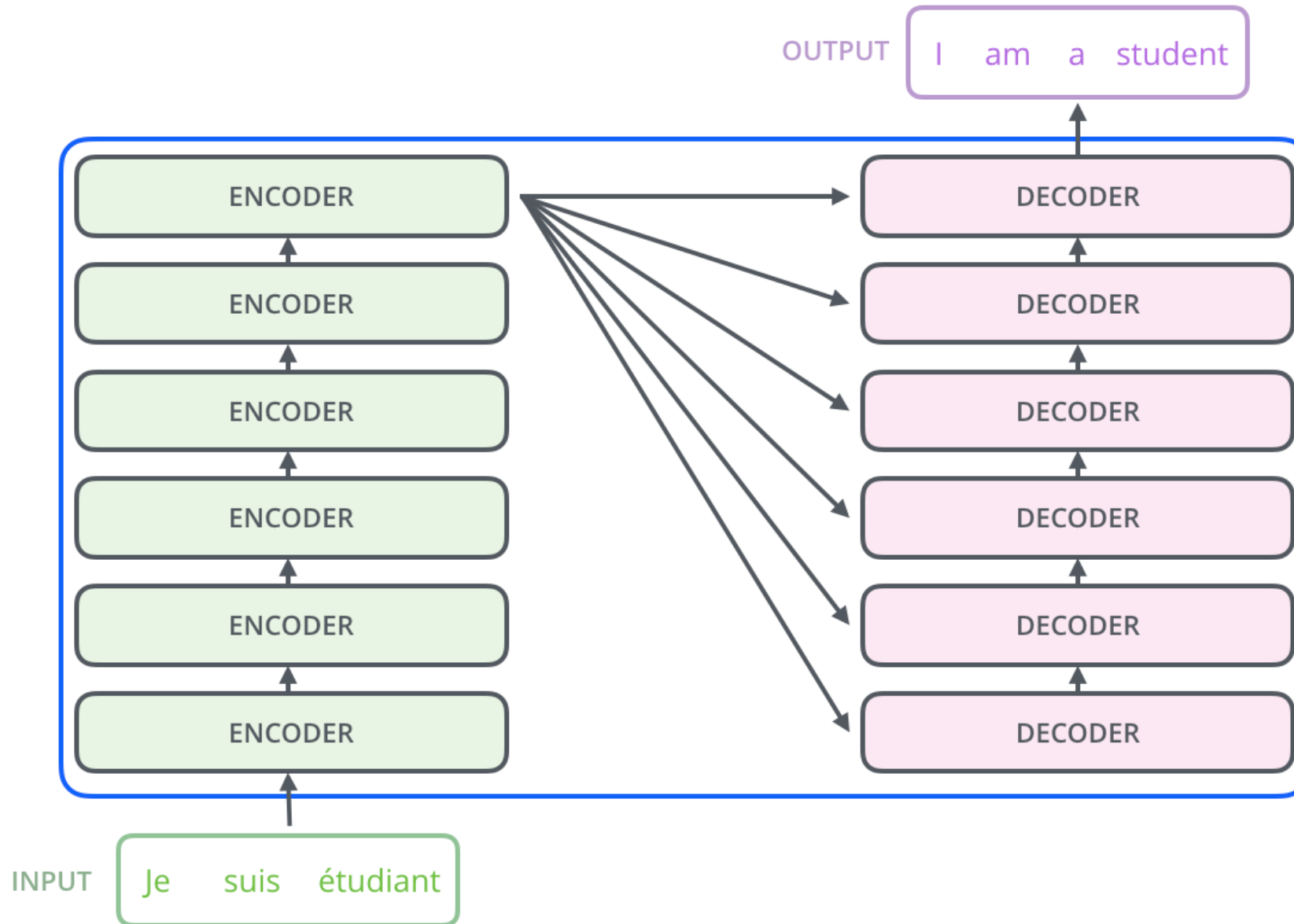| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Jul 22, 2019 | XLNet + DAAF + Verifier (ensemble)<br>*PINGAN Omni-Sinitic* | **88.592** | **90.859** |
| 2<br>Jul 19, 2019 | XLNet + SG-Net Verifier (ensemble)<br>*Shanghai Jiao Tong University & CloudWalk* | 88.050 | 90.645 |
| 3<br>Jul 23, 2019 | XLNet + SG-Net Verifier (single model)<br>*Shanghai Jiao Tong University & CloudWalk* | 87.046 | 89.899 |
| 3<br>Mar 20, 2019 | BERT + DAE + AoA (ensemble)<br>*Joint Laboratory of HIT and iFLYTEK Research* | 87.147 | 89.474 |
| 3<br>Jul 20, 2019 | RoBERTa (single model)<br>*Facebook AI* | 86.820 | 89.795 |
| 4<br>Mar 15, 2019 | BERT + ConvLSTM + MTL + Verifier (ensemble)<br>*Layer 6 AI* | 86.730 | 89.286 |
| 5<br>Mar 05, 2019 | BERT + N-Gram Masking + Synthetic Self-Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 86.673 | 89.147 |

5

# Multi-Head Attention Transformer

# The Illustrated Transformer

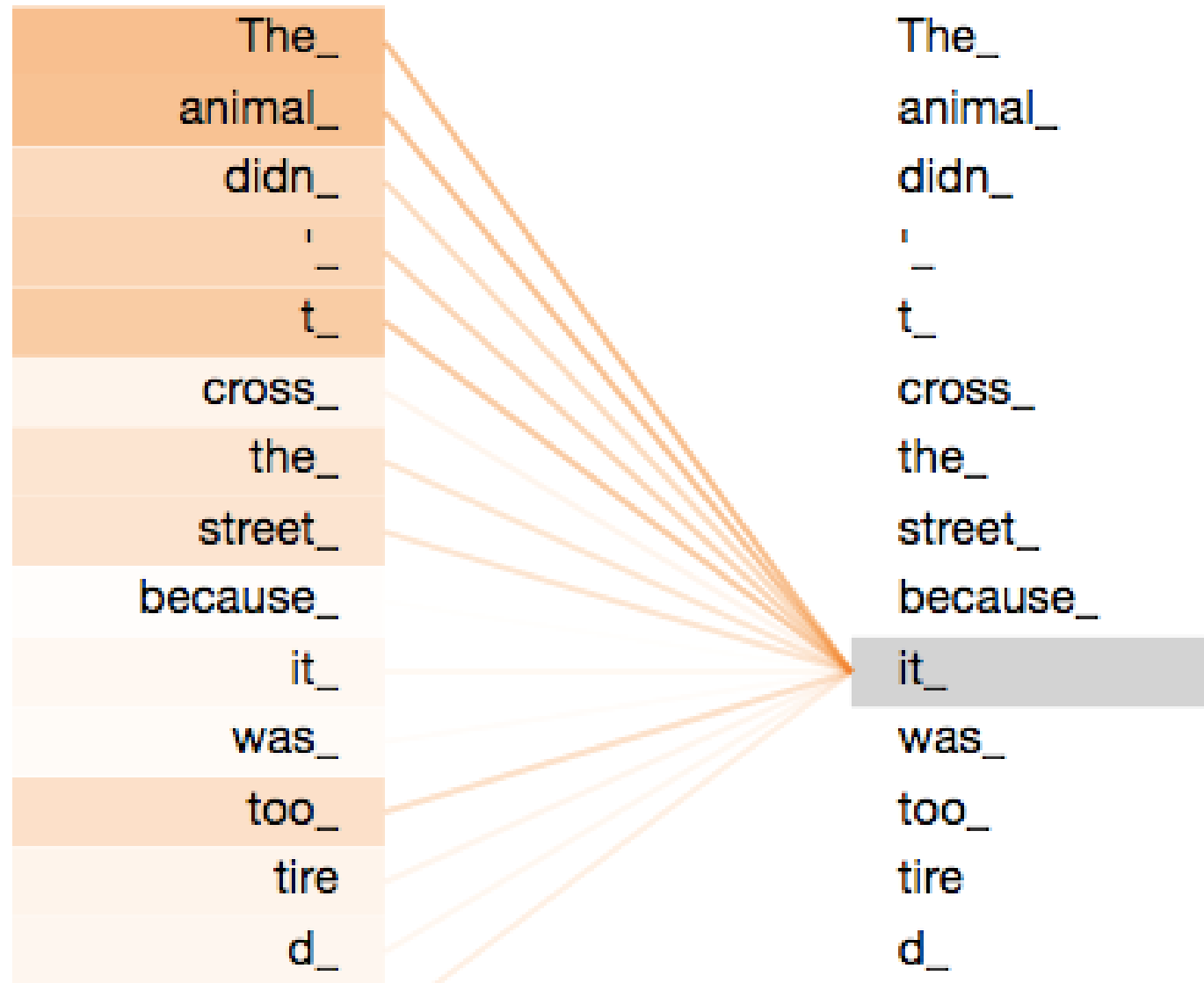http://jalammar.github.io/illustrated-transformer/
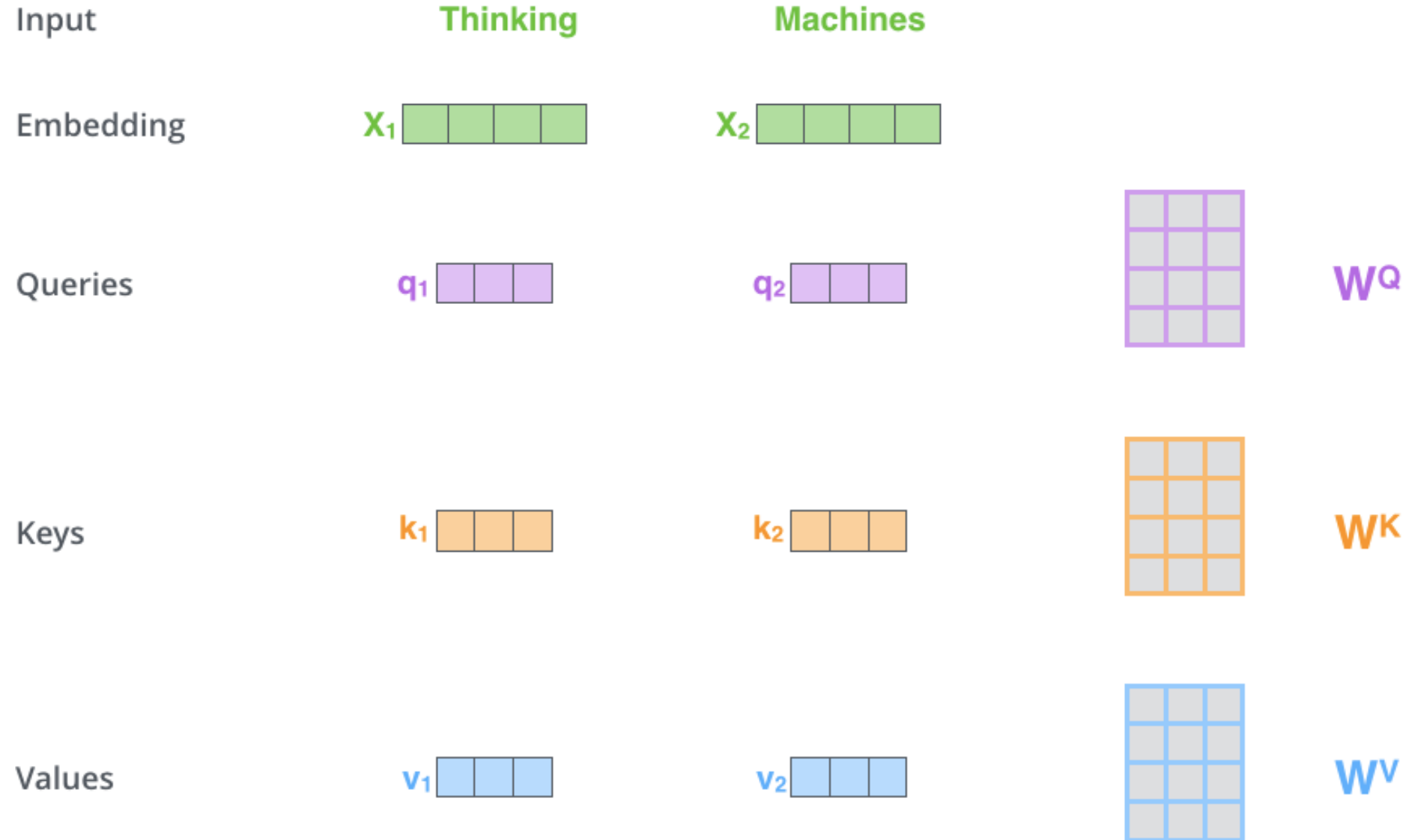
# Transformer

# Transformer

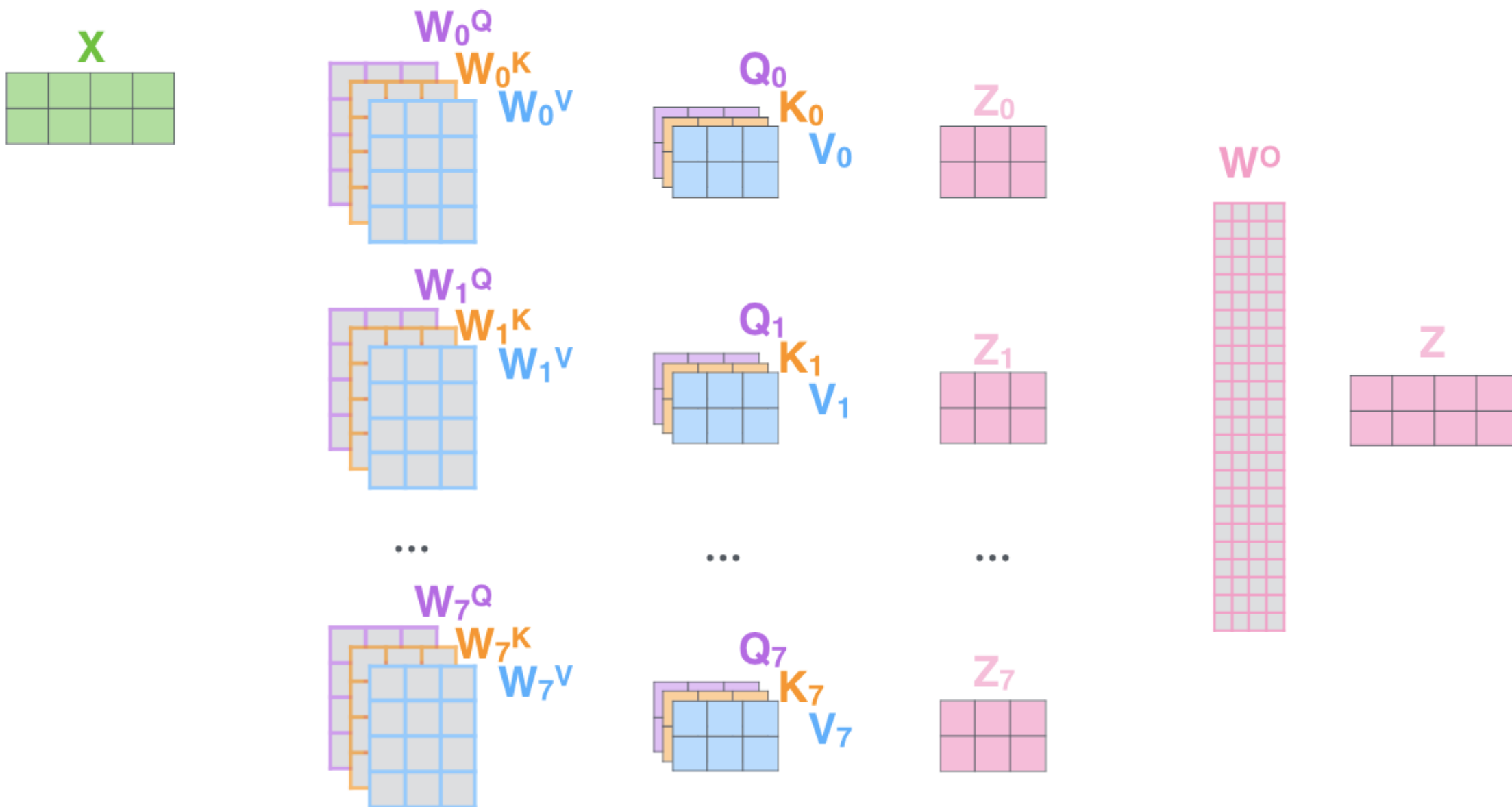# Self-attention

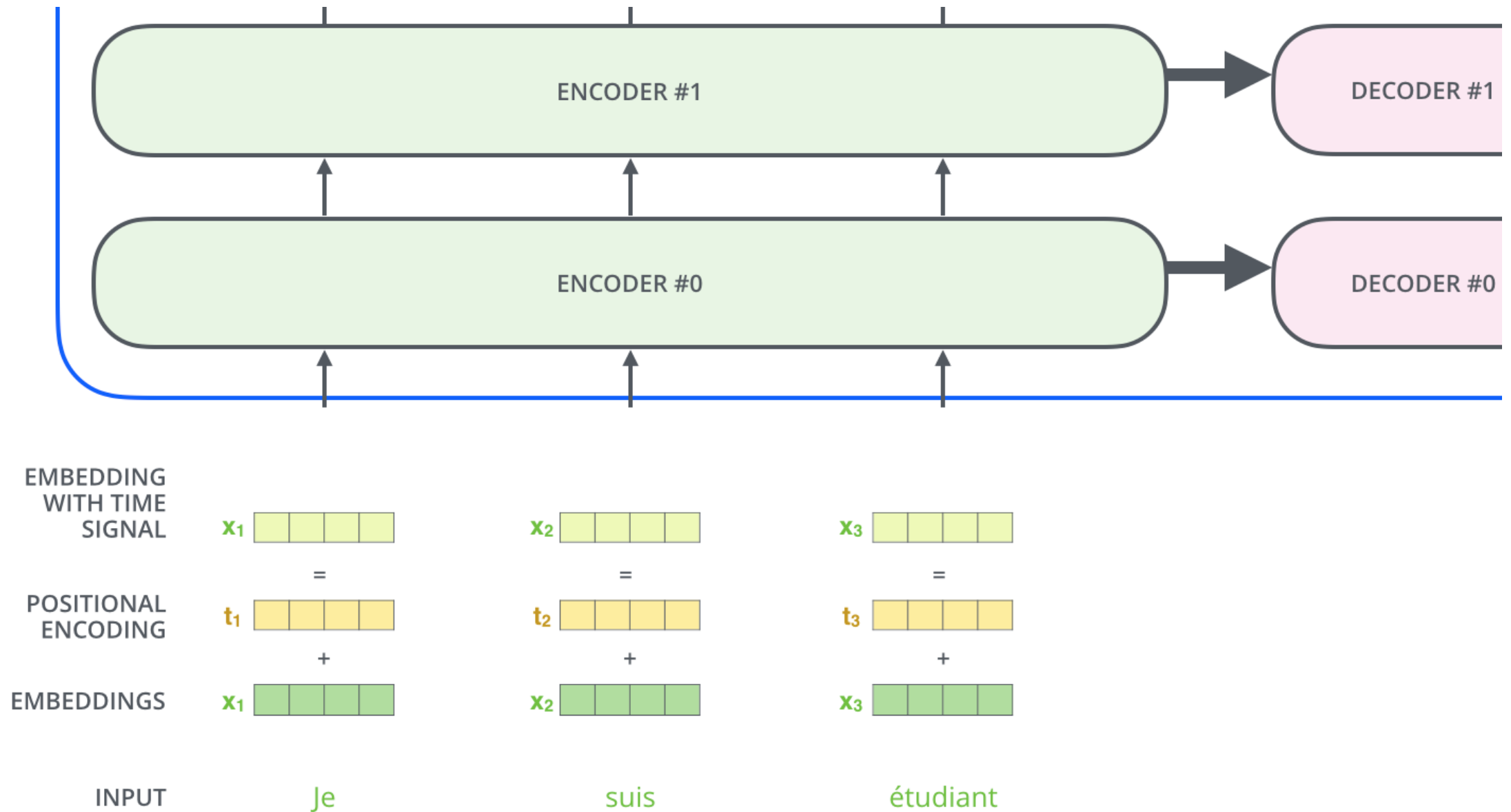# Queries, Keys and Values

# Self attention

# Transformer formula

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

# Multi-headed attention

# Positional encoding

# Positional encoding

- It should output a unique encoding for each time-step (word's position in a sentence)
- Distance between any two time-steps should be consistent across sentences with different lengths.
- Our model should generalize to longer sentences without any efforts. Its values should be bounded.
- It must be deterministic.

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

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

# Positional encoding

# Encoder - decoder

# Bidirectional Encoder Representations from Transformers (**BERT**)
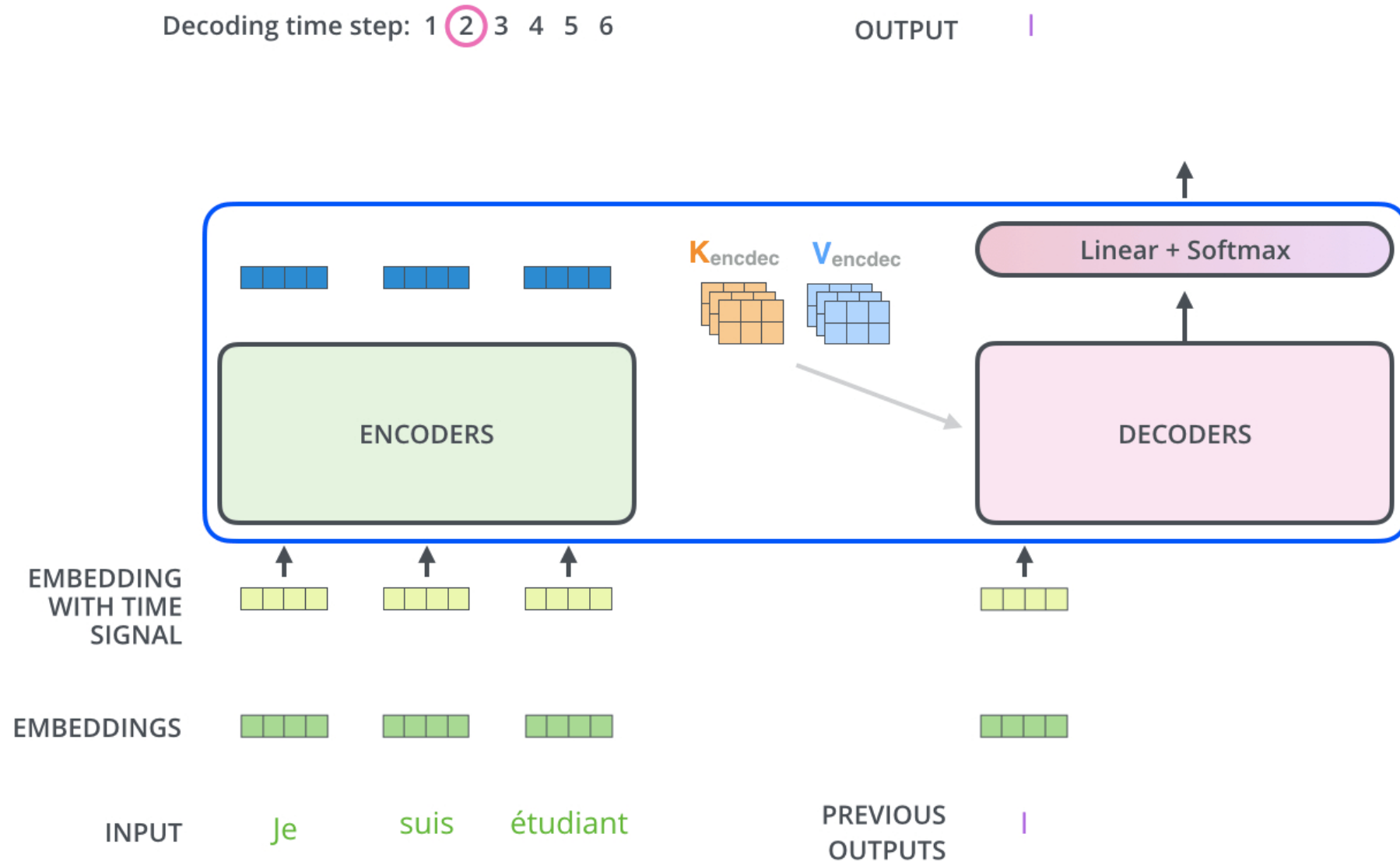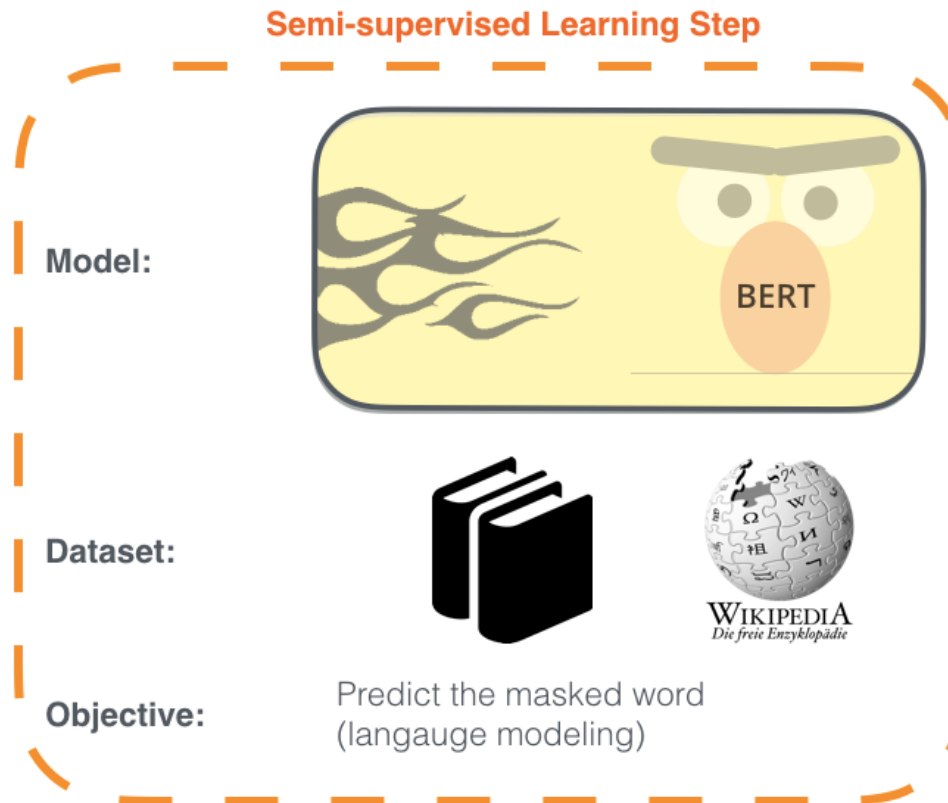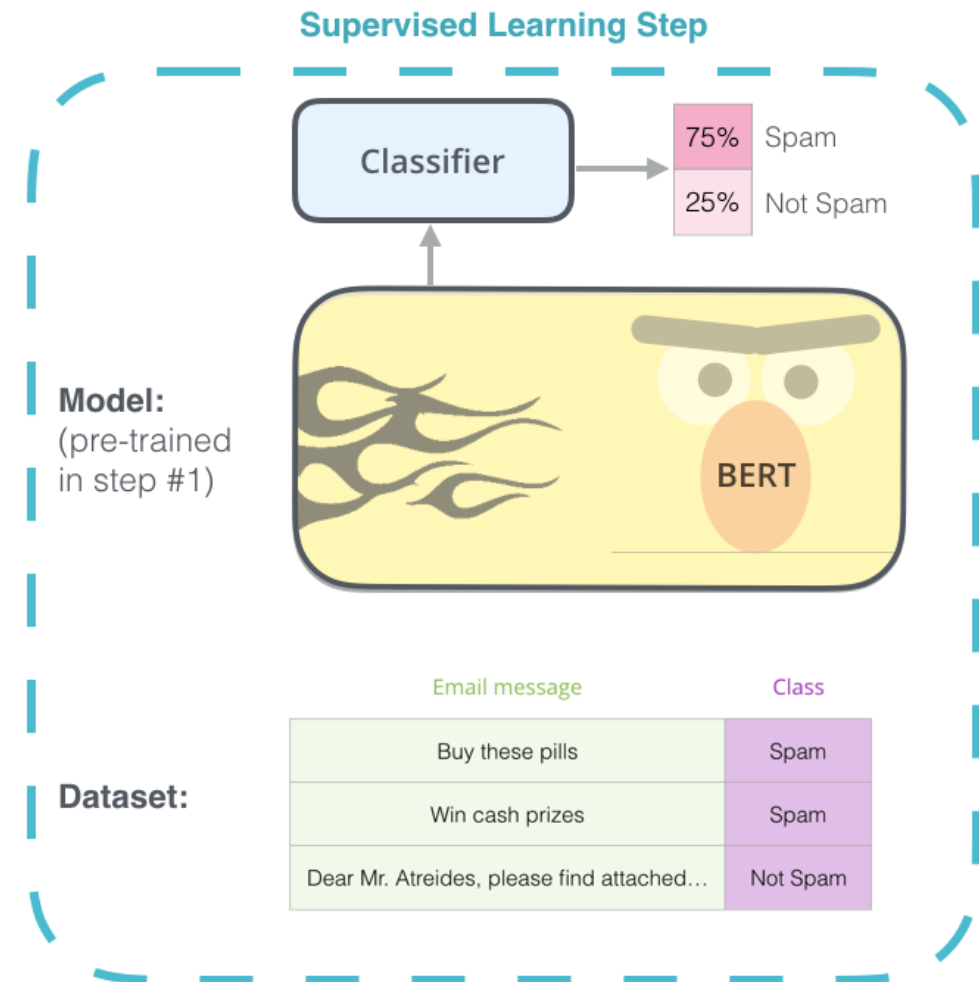
# Bidirectional Encoder Representations from Transformers (**BERT**)

# ChatGPT



**RLHF**

| Low quality data | High quality data | Human feedback | |
| --- | --- | --- | --- |
| **Text** e.g. Internet data | **Demonstration data** | **Comparison data** | **Prompts** |
| Optimized for text completion | Finetuned for dialogue | Trained to give a scalar score for (prompt, response) | Optimized to generate responses that maximize scores by reward model |
| Language modeling | Supervised finetuning | Classification | Reinforcement Learning |
| Pretrained LLM | SFT model | Reward model | Final model |

**Scale**
May '23

| **>1 trillion** tokens | **10K - 100K** (prompt, response) | **100K - 1M** comparisons (prompt, winning_response, losing_response) | **10K - 100K** prompts |

**Examples**
**Bolded**: open sourced

| GPT-x, Gopher, **Falcon**, LLaMa, **Pythia**, **Bloom**, **StableLM** | **Dolly-v2, Falcon-Instruct** | | InstructGPT, ChatGPT, Claude, **StableVicuna** |