

We can also use Greedy Inference with an sequence to sequence model

- How to stop beam search encoder-decoder? Keep going until we have N outputs.
- How to score different output with different length? usually longer output will have higher length, so we need to normalize the score, to get the average probability per word

We will solve this with 'attention'

Give each output step a representation of the input that is most useful for that output decision

Hidden vectors from input $h_1, h_2, h_3, \dots, h_N \in \mathbb{R}^d$

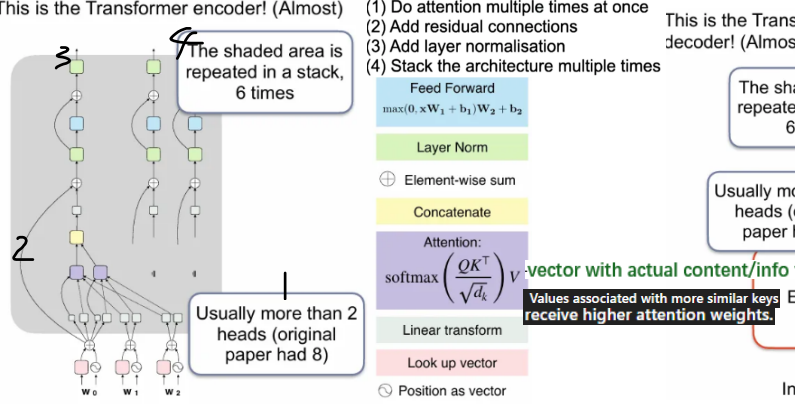
Hidden vector for current output step $s_i \in \mathbb{R}^d$

Calculate attention scores $e_i = [s_i^T h_1, s_i^T h_2, s_i^T h_3, \dots, s_i^T h_N] \in \mathbb{R}^N$

Normalise $a_i = \frac{\exp(e_i)}{\sum_j \exp(e_j)} \in \mathbb{R}^N$

Calculate weighted average $\tilde{h}_i = \sum_j a_j h_j \in \mathbb{R}^d$

This is dot product attention



chrF — compare character ngrams

We compare TP/TN/FP/FN for all n-gram

$FP = FP\text{-}1\text{-gram} + \dots + FP\text{-}4\text{-gram}$

$F_{\beta}\text{-score} = \frac{(1 + \beta^2)TP}{(1 + \beta^2)TP + \beta FN + FP}$

smaller β means closer to Precision

larger β means closer to recall

Robust to typos

BLEU — compare word ngrams

Translation quality

1. Calculate precision for unigrams, bigrams, trigrams, 4grams separately

2. Take the geometric mean

3. Apply a penalty if the output is short

Strongly depends on word order

no 4grams matches means a score of 0!

What if sentence cannot be split on whitespace to get token?

start with small units, then combine units

Byte-Pair Encoding (BPE)

1. Set vocabulary to be all characters

2. Find the two vocabulary items that are adjacent most frequently,

3. Create a new vocabulary item for that pair and update data

4. Check if the vocabulary is size K (e.g., 100,000). If not, go to 2.

WordPiece Consider all possible vocabulary pairs, then Choose the one that will decrease perplexity most if added to the model

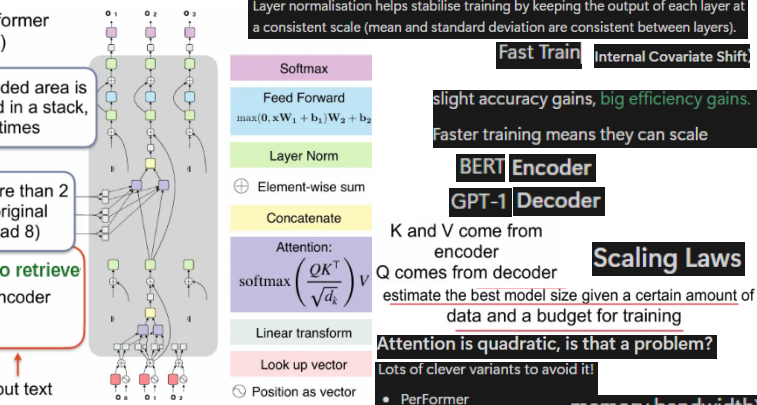
start with big and small units, then delete units

more meaningful subword units

1. Set vocabulary to be all characters and all frequent character sequences, up to full words.

2. Find vocabulary items to remove using a complex method.

3. Repeat until the vocabulary gets down to the desired size.



Encoder — BERT, ELECTRA

Ignore output for unmasked words

masked-token prediction — BERT can do self-attention over the whole sequence

span prediction — SpanBERT mask subword of word or a continuous token set

random-token correction — BERT correct wrong token (model is not told which token was changed)

combination of masked-token + random-token

token edit detection — Edits are not random use a small LLM to generate hard edits

next sentence prediction — BERT judge next sentence (NSP) task.

RoBERTa outperforms BERT by training longer, using more data, and removing the

Encoder-Decoder — T5, BART

masked sequence prediction — BART mask might include multiple token

masked span prediction — T5 mask cannot be empty, include multiple token

deleted sequence prediction — BART or T5 work out where deletion happened, as well as what was deleted

permuted sequence prediction — BART Rearrange sentence

Rotated Sequence Prediction — BART To enhance the model's ability to understand changes in sentence order or structure.

infilling sequence prediction — BART The mask might be empty

Mean Reciprocal Rank (MRR)

earlier the correct answer appears in the ranking, the higher the score.

Decoder — ELMo, GPT

next-token prediction — GPT

Sentiment (Classification) — all

NER (Token predictions) — En & De

Coreference (Structured Output) — Encoder-Decoder

Summarisation (Generation) — Decoder

Translation (Generation, in another language) — Decoder

Reducing Model Size

Efficiency of LLM

Distillation use a big model to train a small model so the small model can perform well

Adding sparsity removing unnecessary features or connections, doesn't help on GPUs because sparse

Increasing training memory efficiency

Parameter-Efficient Fine Tuning (PEFT) — LoRA Reduce weight matrix size

Numerical approximation

Quantization — LLM.int8 Outlier reduced numerical precision version 16-bit

Mixtures of models to less important weights

SMoE divide the model into multiple smaller "expert"

In-Context Learning ICL

Low perplexity prompts are better

Order of examples matters

unsorted examples are best

more number of shots N provide instructions and examples in the input

any label is acceptable, have similar performance in the end

Bad explanations usually lead to the wrong answer!

Url Filtering

Trafilatura text extraction on the raw HTML

FastText LanguageFilter removing any document with en language score lower

Quality filtering

Minhash deduplication deduplicate

PII Formatting anonymize email and public IP addresses

Train a model-based detector — overpredict/overfilter

Instruction Tuning

FLAN-T5 higher generalization

Alpaca LIMA Unclear Careless Ambiguous guidelines? annotators? cases?

0.80-1.0 good

1. Treating retrieval as an LM itself and doing model ensembling

Larger K values generally lead to lower perplexity.

a. kNN-LM Larger dataset size makes relying on external context more reliable.

Smaller λ means more dependence on the LM's internal (in-context) information.

2. Provide the retrieved content in the prompt

Retrieval ICL = ICL + auto-select relevant examples for each input

a. Retrieval ICL Find the most relevant examples in your training data and use them in the prompt.

b. RAG

c. RETRO

d. REALM

e. FLARE

3. Interactive methods (related to Agents, which we will see later)

Risk of dataset

Avoid overfitting into one dataset

minimal edits true answer

Shortcuts / Validity make it wrong!

Does our task have statistical power?

Social Bias

correcting for chance agreement

Is this a task we want to create?

Cohen's Kappa

no clear pattern bigger model benefits more from instruction fine tuning

During instruction fine-tuning: Perplexity may increase instead of decreasing.

Filtered data produce better generation quality high quality

Increase training examples does not improve generation quality

Preference Optimization

RLHF Approach Without making too big change reward model

DPO directly train from preference data

Avoids training a reward model, which may be unreliable

Agent

Reasoning: happened before generating prediction

Self-consistency voting with multiple outputs

Reflection: consider about previous part

Tree of thoughts exploring multiple possible "thought paths"

Acting: doing something beyond the generation: RAG