



**Transformers**: Encoder: classification, masked lang modelling, named entity recognition e.g BERT, RoBERTa, DeBERTa  
 Transformer decoder: translation, summarisation, free-form QA  
 Transformer Decoder: causal language model, text generation  
 Vision Transformers: visual classification, obj detection, img generation  
 Other Models: img captioning/descriptions, text-to-image generation  
 Transformer contains stack of encoder and decoder layers

#### Encoder

Each encoder layer has 2 sublayers: (multi-head self-attention (MHA)  $\rightarrow$  residual + skip connection (R&N))  $\rightarrow$  (position-wise feed-forward network  $\rightarrow$  R&N)  $\rightarrow$  residual + skip connection (R&N)

First Sublayer input:  $S_{\text{MHA}}$  ( $S_{\text{MHA}}$   $\rightarrow$  R&N ( $S_{\text{MHA}}$ ))

**Self-Attention**: Mechanism that allows each token to attend to all other tokens in the seq, to compute a weight representation of the sequence.

Attention(Q, K, V) = softmax(Q $\cdot$ K $^T$ )/sqrt(d $_h$ ) $V$

To obtain Q, K, V matrices, project encodings through weight matrices  $Q = w \cdot W^T$ ,  $K = w \cdot W'$ ,  $V = w \cdot W''$  (all W matrices are shape  $S_{\text{enc}} \times d_h$ )

MHA performs self-attention head amount of times which allows different heads to attend to parts of the sequence differently EG original Transformer used 8 heads, meaning we run self-attention 8 times over 8 different Q, K, Vs

Final representation concatenates and projects the self-attention outputs

Each head has dimensionality D/heads  $L_N(\hat{x}_n) = \hat{x}_n + \beta$

**Normalization**:  $x_{\text{hat}} = (\text{x} - \text{mean}) \cdot \gamma + \beta$

Layer normalization normalizes the features of one sample in one batch

Two key steps: 1.  $x_{\text{hat}} = \text{normalize}(\text{x})$ . 2. Transform  $x_{\text{hat}}$  with learned parameters gamma, beta

Purpose of gamma and beta is to allow the neural network to learn the optimal scale and shift for each feature after normalization, so NN can adjust the range and mean of each feature to better fit the task at hand

**Residual Connections** help mitigate the vanishing gradient problem by providing a shortcut for information to flow to later layers of a network. The output of an earlier layer is added directly to the output of a later layer so current layer can focus on learning the difference between the two outputs, rather than learning an entirely new transformation

Position-wise Feedforward Network  $FNN(x) = \max(0, xW_1 + b_1 + b_2)$

A position-wise FF network is an NLP  $W_1 \in \mathbb{R}^{d_{\text{enc}} \times d_{\text{ff}}}$ ,  $W_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{enc}}}$ ,  $b_1, b_2 \in \mathbb{R}^{d_{\text{enc}}}$

Position-wise just means it is applying the same transformation to every element in the sequence EG 2 layered network, with ReLU activation

**Positional Encoding (PE)**

Transformers are position invariant by default (diff to RNNs which are seq). PEs are a way to inject position information into the embeddings. Done by adding a pos vector to embedding of x. PE is a function of 2 inputs: word position and model dimensionality

#### Decoder

3 sublayers: (masked MHA  $\rightarrow$  R&N)  $\rightarrow$  (cross-attention  $\rightarrow$  R&N)  $\rightarrow$  (position-wise feedforward  $\rightarrow$  R&N).

During  $I_n$ , we perform auto-regressive generation: 1. Source sentence gets encoded via our encoder 2. We feed SOS token to decoder. It then predicts the first word (i.e. 1). 3. We append prediction to SOS token, and then use this to predict the next word (i.e. 2). Continues until EOS token is predicted i.e. end of the generation.

During training we feed the whole target sequence as input to the decoder (so don't have to run decoder start once).

Reinforcement Learning from Human Feedback (RLHF): Optimise model to maximise expected rewards (as assigned by humans) - Expensive as:

• Can model human preferences as a separate NLP problem, i.e. Train another LM to predict the human score for a given text

Advanced Prompting: Improve task performance using prompting. Chain-of-thought (COT): Ask model to reason about its problem solving.

• Zero shot: "Work through the problem step-by-step"

• N-Shot: "Similar to the these examples of step-by-step reasoning, solve the following problems"

Retrieval-Augmented Generation (RAG): Gets factual knowledge from a different system (e.g. Database or "Googling")

Limitations of instruction fine-tuning: Tasks like open-ended creative generation have no right answer, language modelling penalizes all token-level mistakes equally, but some errors are worse than others.

Reinforcement Learning from Human Feedback (RLHF): Optimise model to maximise expected rewards (as assigned by humans) - Expensive as:

• Can model human preferences as a separate NLP problem, i.e. Train another LM to predict the human score for a given text

• Also Human judgments are noisy and miscalibrated! Instead of asking for the key, value tensors from the last encoder layer, and send them to all the decoder layers. So, query comes from the current decoder layer

Cross attention matrix shape =  $T \times S$ . Note: K and V are from the encoder.

Transformer does 4 things every training loop: 1. Create source and target mask 2. Run encoder 3. Run decoder 4. Output logits for token prediction

Authors use weight tying in the decoder (Embedding matrix and output projection matrix are shared) and use a decaying learning rate

**Pre-Training Models**

Contextual Word Embeddings: Meaning of a word can depend on its context, need to take context into account when constructing word representations

ELMo - Embeddings from Language Models: Take a large corpus of plain text and train two language models:

1. One recurrent language model going forward, left-to-right.  
 2. A second recurrent language model going backward, right-to-left

When we need a vector for a word, combine the representations from both directions (see Machine Translation for more details)

**Model Types**:

• **Encoders**: Are able to access the whole sequence, using context on both sides of each token.

• **Decoders**: Are able to access context on the left. Language models, good for generating text.

• **Encoder-decoder**: Input is processed using an encoder, then output is generated using a decoder.

**BERT: Bidirectional Encoder Representations from Transformers**:

Taken in a sequence of tokens, gives as output a vector for each token using stacked encoder blocks. Trained using the following:

• **Masked Language Modelling (MLM)**: Self-supervised training task that needs no labelled data. Hide  $k\%$  of the input words behind a mask, train model to predict them.

• Too little mask - too expensive to train

• Actual strategy: used Pick 15% of input words as training targets - 80% of those are replaced with [MASK] token, 10% are replaced with a random word, 10% are the left as the original word

• If we masked all chosen words, model wouldn't necessarily learn to construct good representations for non-masked words.

• Second Training Objective: Given two sentences, predict whether they appeared in the original order in the source text, doesn't provide much additional benefit, so it has not been used in newer versions.

Bert Variants: layers: hidden state dim; attention heads: parameters

1. BERT-base: 12 : 768 : 12 : 110 million

2. BERT-large: 24 : 1024 : 16 : 340 million

Trained on: BooksCorpus (800 million words) & English Wikipedia (2.500 million words). Pre-trained with 64 TPU chips for a total of 4 days

Using BERT: BERT-like models give us a representation vector for every input token by removing masked LM head. Represent a whole sequence by vector for end of sequence token.

**Pre-Training Models**

• Multimodal models: Use multiple type of data in a model (text, images, etc)

• BERT-Variants: Roberta, DeBERTa

• Mask contiguous tokens: SPANBERT

• Distil large model into small one (Model Distillation): DistilBERT, ALBERT

• Use sparse Attention: BigBird, LongFormer

• Application specific: ClinicalBERT, MedBERT, PubMedBERT, BEHRT

• Use existing knowledge graphs with entity embeddings: ERNIE

**Pre-training Encoder-Decoder Models**: Initial input sequence can be fully attended over using the encoder, getting the maximum information. Then, a dedicated decoder is used for generating the output sequence.

Popular for machine translation: encoder and decoder can have separate vocabularies and focus on different languages.

Techniques:

• Can't use MLM as no direct correspondence between inputs and outputs

• Prefix Language Modelling: Source: "I like Target: Skidbi toilet

• Sentence permutation decoding: src = "party", good": tgt = "good party".

• Token Masking: "I like Target: skidbi toilet": tgt = "like skidbi toilet"

• SPANBERT-like objective/ span corruption

• Original Text: "Thank you for partying with me"

• Inputs: "Thank  $\langle \text{X} \rangle$  partying  $\langle \text{Y} \rangle$ "

• Targets: " $\langle \text{X} \rangle$  you for  $\langle \text{Y} \rangle$  with me"

encoder-decoder models trained with span corruption found to work better than regular decoders.

**Instruction Training Enc-Dec**:

• Encoder inputs are the "natural language" prompts for a given task

• e.g.: Please answer the following questions in 10 steps."

• T5 (Text-To-Text Transfer Transformer) is trained using the span corruption unsupervised objective, along with a number of different supervised tasks.

• FLAN-T5 (Fine-tuned LAnguage Net) is the same size as T5 but trained on much more data, more languages and 1.8K tasks phrased as instructions.

**Pre-training Decoder Models**: We can train on unlabelled text, optimizing  $p_{\text{theta}}(w_{t+1} | t-1, \dots, t)$ . Great for tasks where the output has the same vocabulary as the pre-training data.

Alternative ways of using pre-trained decoders:

• Full Fine-tuning: Supervised training for particular input-output pairs. Or we can put a new layer on top of fine-tune the model for a desired task such as use "classification head" for a classification task.

• One possibility is to use Language Inference Label pairs of spans as inputs as entailing/contradicting/neutral.

• Zero-shot: Give the model a natural language description of the task, have it generate the answer as a continuation.

• One-shot: In addition to the description of the task, give one example of solving the task. No gradient updates are performed.

• Few-shot: In addition to the task description, give a few examples of the task as input. No gradient updates are performed.

Advanced Prompting: Improve task performance using prompting.

Chain-of-thought (COT): Ask model to reason about its problem solving.

• Zero shot: "Work through the problem step-by-step"

• N-Shot: "Similar to the these examples of step-by-step reasoning, solve the following problems"

Retrieval-Augmented Generation (RAG): Gets factual knowledge from a different system (e.g. Database or "Googling")

Limitations of instruction fine-tuning: Tasks like open-ended creative

generation have no right answer, language modelling penalizes all token-level mistakes equally, but some errors are worse than others.

Reinforcement Learning from Human Feedback (RLHF): Optimise model to maximise expected rewards (as assigned by humans) - Expensive as:

• Can model human preferences as a separate NLP problem, i.e. Train another LM to predict the human score for a given text

• Also Human judgments are noisy and miscalibrated! Instead of asking for the key, value tensors from the last encoder layer, and send them to all the decoder layers. So, query comes from the current decoder layer

Cross attention matrix shape =  $T \times S$ . Note: K and V are from the encoder.

Transformer does 4 things every training loop: 1. Create source and target mask 2. Run encoder 3. Run decoder 4. Output logits for token prediction

Authors use weight tying in the decoder (Embedding matrix and output projection matrix are shared) and use a decaying learning rate

**Pre-Training Models**

Contextual Word Embeddings: Meaning of a word can depend on its context, need to take context into account when constructing word representations

ELMo - Embeddings from Language Models: Take a large corpus of plain text and train two language models:

1. One recurrent language model going forward, left-to-right.

2. A second recurrent language model going backward, right-to-left

When we need a vector for a word, combine the representations from both directions (see Machine Translation for more details)

**Model Types**:

• **Encoders**: Are able to access the whole sequence, using context on both sides of each token.

• **Decoders**: Are able to access context on the left. Language models, good for generating text.

• **Encoder-decoder**: Input is processed using an encoder, then output is generated using a decoder.

**BERT: Bidirectional Encoder Representations from Transformers**:

Taken in a sequence of tokens, gives as output a vector for each token using stacked encoder blocks. Trained using the following:

• **Masked Language Modelling (MLM)**: Self-supervised training task that needs no labelled data. Hide  $k\%$  of the input words behind a mask, train model to predict them.

• Too little mask - too expensive to train

• Actual strategy: used Pick 15% of input words as training targets - 80% of those are replaced with [MASK] token, 10% are replaced with a random word, 10% are the left as the original word

• If we masked all chosen words, model wouldn't necessarily learn to construct good representations for non-masked words.

• Second Training Objective: Given two sentences, predict whether they appeared in the original order in the source text, doesn't provide much additional benefit, so it has not been used in newer versions.

Bert Variants: layers: hidden state dim; attention heads: parameters

1. BERT-base: 12 : 768 : 12 : 110 million

2. BERT-large: 24 : 1024 : 16 : 340 million

Trained on: BooksCorpus (800 million words) & English Wikipedia (2.500 million words). Pre-trained with 64 TPU chips for a total of 4 days

Using BERT: BERT-like models give us a representation vector for every input token by removing masked LM head. Represent a whole sequence by vector for end of sequence token.

**Pos (Part of Speech) Tagging**

Tagger Tag (example) | ADJ (adj) ADV (adv) INT (int) PUNCT (punct) (cat) | PNP (pnp) UH (uh) RP (rp) (SYM) (%)(%) (X) (other) PREP (prep) (on) AUX (aux) (CCONJ (cc) (fwd) (DEP) (dep) (the) (num) (one) (PART) (part) (S) (pron) (he) (SCONJ (if) (poS) uses Name Entity Recognition (NER), preprocessing to select ADJs, adverbs, lemmas, neural syntactic/semantic parsing, dictation analysis i.e. if ILM generated, keyword/phrase recognition for indexing + searching, spam detection + filtering

Solved for mainstream languages, e.g. Spacy, NLTK  $\hat{T} = \text{argmax}_T P(T|W)$

**Probabilistic POS Tagging**

Given word sequence  $W = w_1, \dots, w_n$ , estimate tag sequence  $T = t_1, \dots, t_n$  Compute  $P(T|W)$ , to many-to-many classification

Generative (Bayes)  $P(T|W) = P(W|T)P(T)$   $P(W|T) = P(w_1|t_1)P(w_2|t_2) \dots P(w_n|t_n)$

1. Chain rule - Markov assumption (bigram):  $P(T|W) = P(t_1|w_1)P(t_2|w_2, t_1) \dots P(t_n|w_n, t_{n-1})$

2. Word depends only on its tag:  $P(W|T) = P(w_1|t_1)P(w_2|t_2) \dots P(w_n|t_n)$

Combine to  $P(T|W) = P(t_1|w_1)P(w_2|t_2) \dots P(t_n|w_n, t_{n-1})P(w_n|t_n)$

$t_1 = \text{Cat}(t_1, w_1)$   $t_2 = \text{Cat}(t_2, w_2, t_1)$  ...  $t_n = \text{Cat}(t_n, w_n, t_{n-1})$

Table of  $P(w_i|t_i)$ ,  $(P(w_i|t_i))$  given another tag  $t_i$  is cols. Let  $\hat{s} \leftarrow$  be a  $t_i$ -1.

To tag, sequence words as cols,  $\text{P}(t_i|t_{i-1})P(w_i|t_i)$  as entries. Can pick maxes

$t_i = \text{argmax}_T P(T|W)$

$\approx \text{argmax}_{t_1, t_2, \dots, t_n} \prod_{i=1}^n P(t_i|t_{i-1})P(w_i|t_i)$

**Hidden Markov Model (HMM) tagger**

• State sequence  $(t_0, t_1, \dots, t_n)$

• Transition probability matrix  $a_{ij}$  is prob. state  $i \rightarrow$  state  $j$ .  $A = [a_{ij}]_{n \times n}$

• Observation sequence  $(w_0, w_1, \dots, w_n)$

• State probabilities  $b_i(w_i)$  (likelihood of observation  $w_i$  from state  $i$ )  $B = [b_i(w_i)]_{n \times |V|}$

Start:  $t_0 = [\text{root}]$ ,  $\beta = 1$

1. Shift:  $\alpha_i = \sum_j a_{ij} \beta$

2. Left-Arc:  $\alpha_i = \sum_j a_{ij} \alpha_j$

3. Right-Arc:  $\alpha_i = \sum_j a_{ij} \alpha_j$

Finish:  $\alpha_n = [\text{end}]$ ,  $\beta = \alpha_n$

**Dependency Parsing**: Another form of extracting syntactic structure.

• Connect words in sentence to indicate dependencies. Directional arrow from heads to dependents. Can be typed (arrows are labelled).

• Allows: Ability to deal with morphologically rich languages, head-dependent relationships provide approximation to semantic relationship between predicates and arguments, can be used to analyse complexity of a language

**Dependency structure**: A directed graph  $G(V, A)$

•  $V$  = set of vertices (words, punctuation, ROOT)

•  $A$  = set of ordered pairs of vertices (i.e. arcs)

**Dependency Tree**: Acyclic dependency structure with a single ROOT node with no incoming arcs. Each vertex has exactly one incoming arc.

**Projective vs Non-Projective**: A dependency structure where there are no crossing dependency arcs when the words are laid out in their linear order.

**Automated Techniques**:

• Dynamic programming is cubic time / not accurate

• Shift: From left to right, fast (linear) but slightly less accurate, MultiParser

• Spanning Tree (graph-based, constraint satisfaction)

• Calculate full tree at once, slightly more accurate, slower, MSTParser

• MailParser: Greedy choice of attachment for each word in order, guided by ML classifier, operates as a stack machine, formally 3 data structure

• Stack  $O$ : starts with ROOT

• Buffer  $B$ : starts with all words in sentence

• Arc Set  $A$ : starts empty

With operations: shift/ left arc/ right arc. Optionally, set of dep. labels for left and right arc actions ("40").

Start:  $O = [\text{root}]$ ,  $\beta = \emptyset$

1. Shift:  $O, (\text{w} | \beta), \beta \rightarrow O \cup \{\text{w}\}, \beta$

2. Left-Arc:  $(\text{o} | \text{w} | \beta), \beta \rightarrow (\text{o} | \text{w}), \beta$

3. Right-Arc:  $(\text{o} | \text{w} | \text{w}), \beta \rightarrow (\text{o} | \text{w}), \beta, \text{A} \cup \{\text{r}(\text{w}, \text{w})\}$

Finish:  $\alpha = [\text{end}]$ ,  $\beta = \emptyset$

**Action Decisions**: Actions are predicted using a discriminative classifier over each move. Traditionally uses lots of engineered features such as: top of stack word, word length, top of stack word, top of stack word, etc.

**Evaluation**: Treebanks exist for dependency parsing as well. Can use standard metrics between reference and hypothesis (labelled arcs).

**Neural Parsing**: Do dependency parsing as a SEQ2SEQ task where the output sequences are a linearized tree (using brackets). Can use many architectures mentioned previously (RNNs, LSTMs, transformers). Train with cross-entropy. Evaluate as translation BLEU or parsing tasks (parseval)

**Advanced Neural Parsing**: Use graph-based models.

**Exam/Tutorial Ons**

13) The dimensions for word embeddings and GRU states are  $E$  and  $H$ , respectively. The vocabulary size is  $V$ . Provide the size ( $R$ ? ) of  $W$ ,  $Wh$  and  $W_{\text{in}}$  and  $R$  in  $R = \text{Cat}(x, t, 1)$ .  $t = \text{H}(E, H)$   $Wh = \text{H}(H, t)$   $br = \text{H}(1, t)$

14) Doing optimize  $P(w|t)$  vs  $P(w|t, E)$  because not differentiable & extrinsic vs intrinsic

15) NMT loss is the sum of log-probabilities for every token prediction. How many log-probability terms would we have for the above sample batch with three training examples? batch size \* (sum of all target lengths)