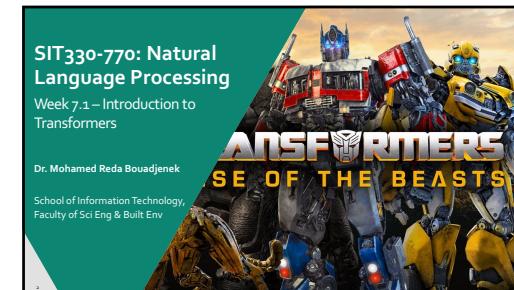
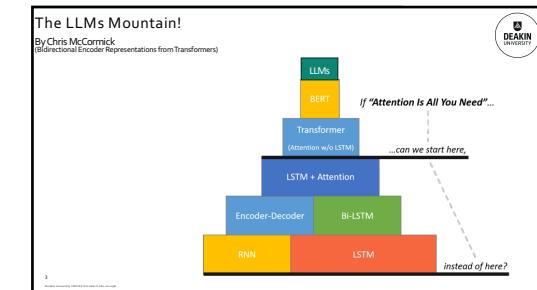




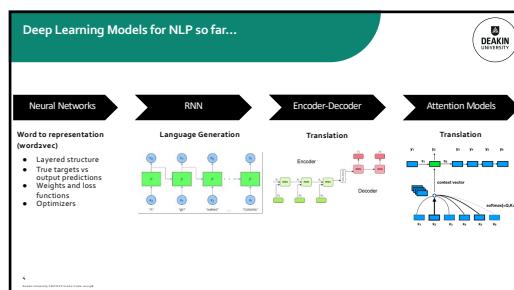
1



2



3



4

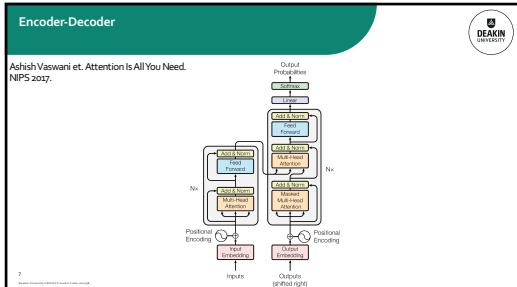
Problem/Motivation

- Encoder-decoder models have largely used RNN and LSTM, but the computation is sequential
 - RNN experiences vanishing gradient
 - Both are slow to train, even with factorization and conditional computation
- ConvS2S uses CNN to compute representations in parallel
 - Computation scales linearly with distance between positions
- Transformer** provides constant computation and uses Multi-headed Attention to negate reduced effectiveness

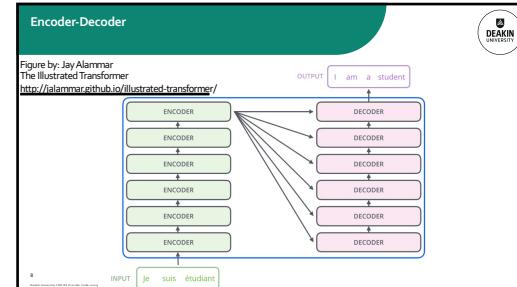
5



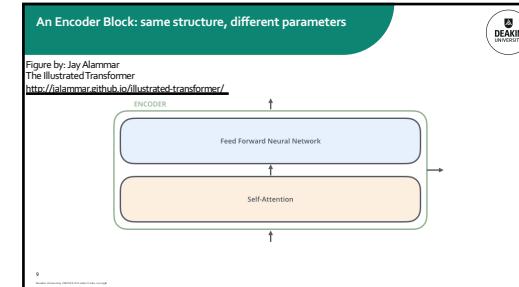
6



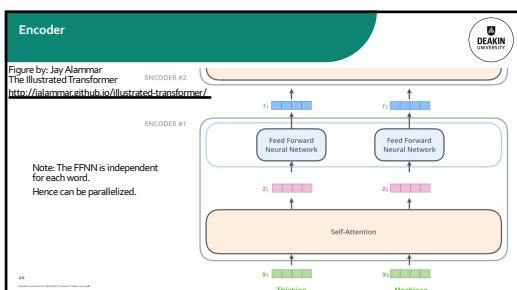
7



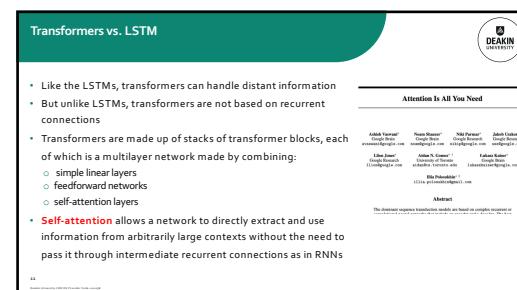
8



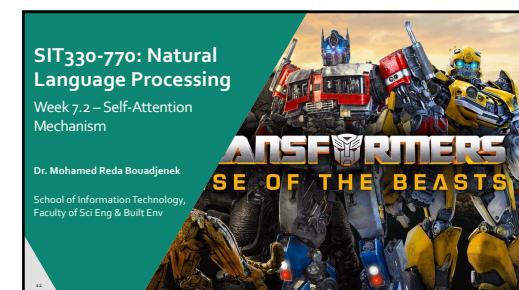
9



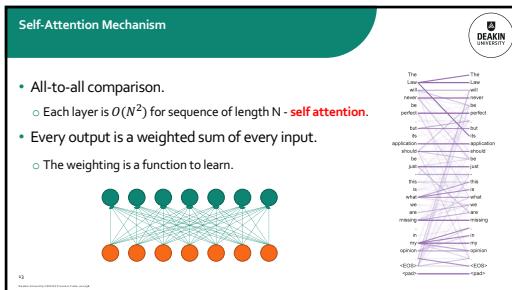
10



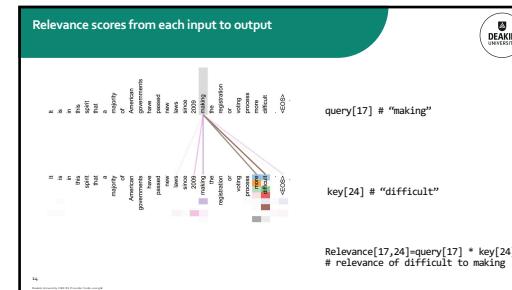
11



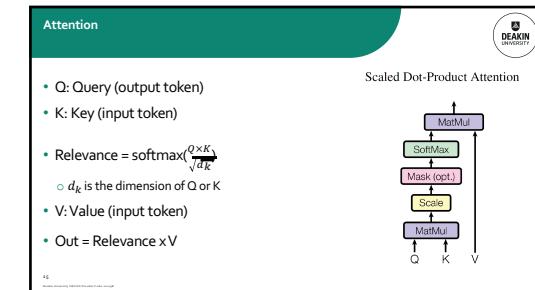
12



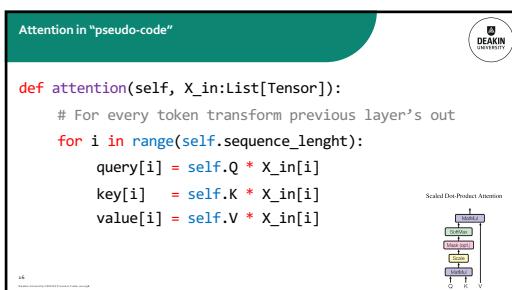
13



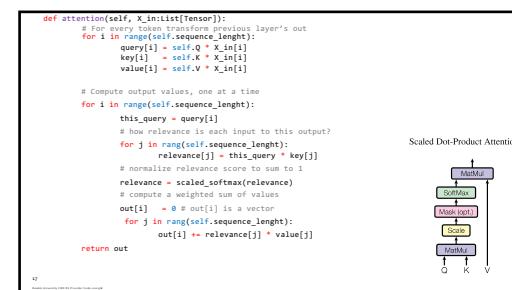
14



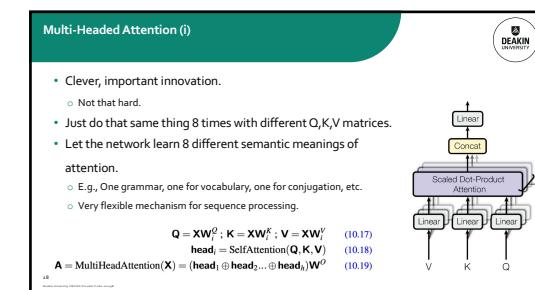
1



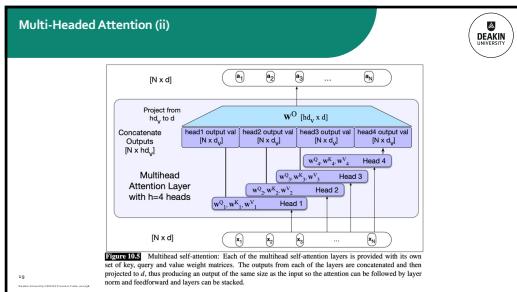
16



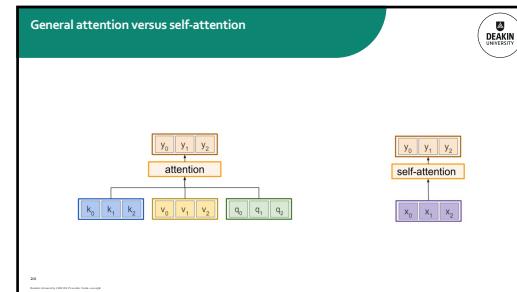
17



1



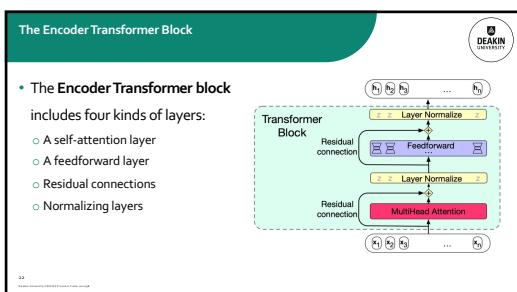
19



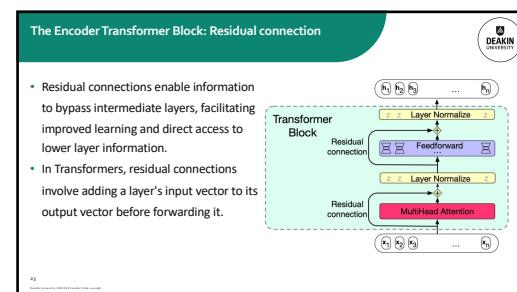
20



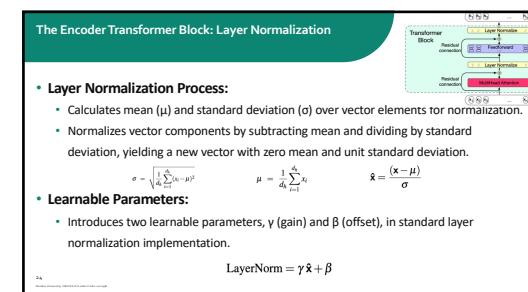
21



22



23



24

The Encoder Transformer Block: feedforward layer

- The feedforward layer consists of N position-wise networks, each positioned independently.
 - Each network is a fully-connected 2-layer neural network, comprising one hidden layer and two weight matrices.
- While the weights remain consistent across positions, the parameters differ from layer to layer.
- Unlike attention mechanisms, the feedforward networks operate independently at each position, enabling parallel computation.

15 Input Layer $\in \mathbb{R}^N$ Hidden Layer $\in \mathbb{R}^N$ Output Layer $\in \mathbb{R}^N$

25

The Encoder Transformer Block: Putting it all together

- The Transformer block includes four kinds of layers:
 - A self-attention layer
 - Residual connections
 - Normalizing layers
 - A feedforward layer

$$\begin{aligned} T^1 &= \text{SelfAttention}(X) \\ T^2 &= X + T^1 \\ T_3 &= \text{LayerNorm}(T^2) \\ T^4 &= \text{FFN}(T_3) \\ T^5 &= T^4 + T^3 \\ H &= \text{LayerNorm}(T^5) \end{aligned}$$

16

26

Encoder-Decoder Transformers

Ashish Vaswani et al. Attention Is All You Need. NIPS 2017.

17

27

SIT330-770: Natural Language Processing

Week 7.4 – The Input: Embeddings for Tokens

Dr. Mohamed Reda Bouadjenek

School of Information Technology,
Faculty of Sci Eng & Built Env

TRANSFORMERS: AGE OF THE BEASTS

18

28

The input X

- Source of Input X :**
 - Originates from a sequence of N tokens, where N represents the number of tokens in the context.
 - Matrix X has dimensions $[N \times d]$ where d is the dimension of each embedding.
- Composition of Embeddings:**
 - The transformer model computes two separate embeddings:
 - Input Token Embedding:** Specific to each token in the sequence.
 - Input Positional Embedding:** Reflects the position of each token within the sequence.

19

29

Input Token Embedding

- Each token's initial representation is a vector of dimension d .
- These embeddings evolve through transformer layers, incorporating contextual nuances.
- Embedding Storage:**
 - Stored in matrix E with shape $[|V| \times d]$, where $|V|$ is the vocabulary size.
- Token Conversion and Indexing:**
 - Tokens first converted to vocabulary indices using techniques like BPE or SentencePiece (discussed in Week 3).
 - Example: "thanks for all" converts to indices [5, 4000, 10532, 2224].
 - These indices are used to fetch the corresponding embeddings from E .

20

One-Hot Vector

- One-Hot Vector Representation:
 - Alternative method using one-hot vectors of size $|V|$.
 - Example: Word "thanks" represented as a vector where only the 5th position is 1, all others are 0.
- Multiplying by a one-hot vector that has only one non-zero element $x_i = 1$ simply selects out the relevant row vector for word i , resulting in the embedding for word i ,



Figure 10.10 Selecting the embedding vector for word V_i by multiplying the embedding matrix E with a one-hot vector with a 1 at index i .

31

One-Hot Vector

- We can extend this idea to represent the entire token sequence as a matrix of onehot vectors, one for each of the N positions in the transformer's context window,

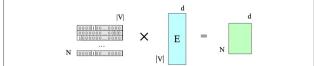


Figure 10.11 Selecting the embedding matrix for the input sequence of token ids W by multiplying a one-hot matrix corresponding to N by the embedding matrix E .

32

SIT330-770: Natural Language Processing

Week 7.5 – The Input: Embeddings for Positions

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env



33

Input Positional Embedding

- How does a transformer model the position of each token in the input sequence?
 - With RNNs, information about the order of the inputs was built into the structure of the model, Not with Transformers
- Solution: Positional**
 - Embeddings Modify the input embeddings by combining them with positional embeddings specific to each position in an input sequence

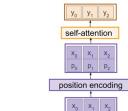
34

Positional Embedding

- Positional encoding assigns a unique representation to each position within a sequence to describe the location or order of entities.
- Limitations of Using Single Numbers for Position:
 - Using index values alone (e.g., sequence position numbers) is problematic for several reasons:
 - Large indices for long sequences can grow unmanageably large.
 - Normalizing indices between 0 and 1 creates inconsistencies across sequences of different lengths.
- Transformers' Approach to Positional Encoding:
 - Instead of single numbers, transformers map each position to a unique vector.
 - This results in a matrix where each row represents an object in the sequence combined with its positional information.

35

Positional encoding



- Desiderata of pos(.) :
 - 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.

36

Positional Encoding Layer in Transformers

The positional encoding is given by sine and cosine functions of varying frequencies:

- k : Position of an object in the input sequence,
- d : Dimension of the output embedding space
- n : User-defined scalar, set to 10,000 by the authors of Attention is All You Need.

Example:

- "I am a robot," with $n=100$ and $d=4$

Positional Encoding Matrix for the response "I am a robot"

37

Coding the Positional Encoding Matrix from Scratch

```

import numpy as np
import matplotlib.pyplot as plt

def getPositionEncoding(seq_len, d, n=10000):
    P = np.zeros((seq_len, d))
    for k in range(seq_len):
        for i in np.arange(int(d/2)):
            denominator = np.power(n, 2*i)
            P[k, 2*i] = np.sin(k/denominator)
            P[k, 2*i+1] = np.cos(k/denominator)
    return P

P = getPositionEncoding(seq_len=4, d=4, n=100)
print(P)
1 [[0.0 1.0 0.0 1.0]
 2 [0.84147098 0.54030231 0.9983342 0.99500417]
 3 [0.90929743 -0.41614684 0.19866933 0.98006658]
 4 [0.14112001 -0.9899925 0.29552021 0.95533649]]

```

38

Positional Embeddings

Figure 10.12 A simple way to model position: add an embedding of the absolute position to the token embedding to produce a new embedding of the same dimensionality.

39

SIT330-770: Natural Language Processing

Week 7.6 – The Task Specific Head

Dr. Mohamed Reda Bouadjenek

School of Information Technology,
Faculty of Sci Eng & Built Env

ANSFORTIERS
SE OF THE BEASTS

40

The Language Modeling Head

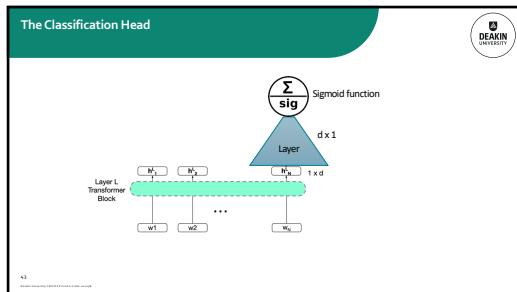
- The Language Modeling Head is an additional neural circuitry integrated with the basic transformer architecture.
- It enables specific tasks such as language modeling by enhancing the transformer's capabilities.
- Given a context of words, a Language Model assigns a probability to each possible next word
- Calculating the probability of the next word "fish" given the context "Thanks for all the": $P(\text{fish}|\text{Thanks for all the})$

41

The Language Modeling Head

Figure 10.13 The language modeling head: the circuit at the top of a transformer that maps from the output embedding for token N from the last transformer layer (h_N^L) to a probability distribution over words in the vocabulary V .

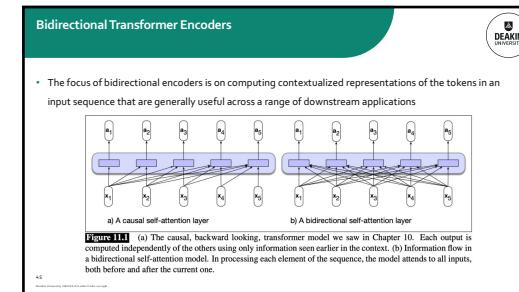
42



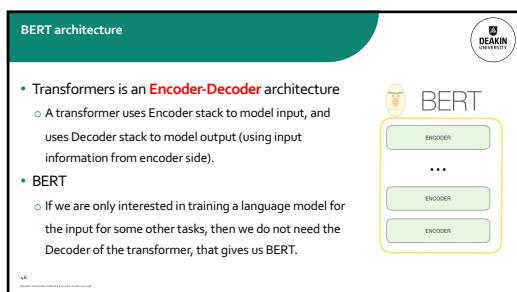
43



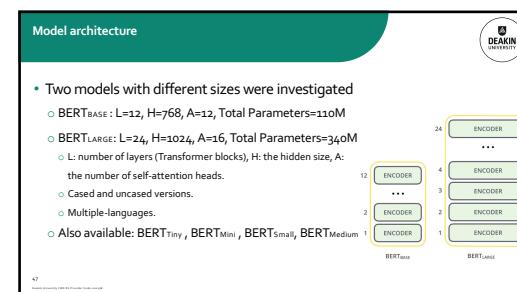
44



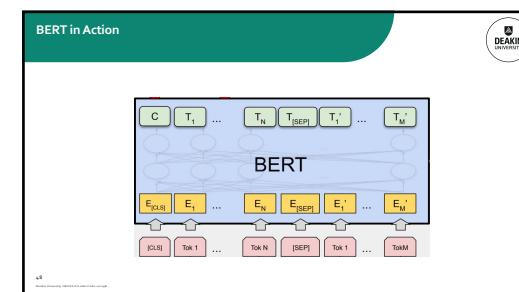
45



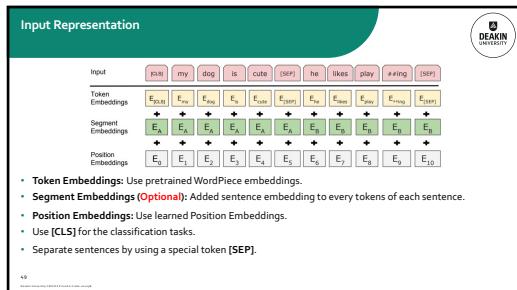
46



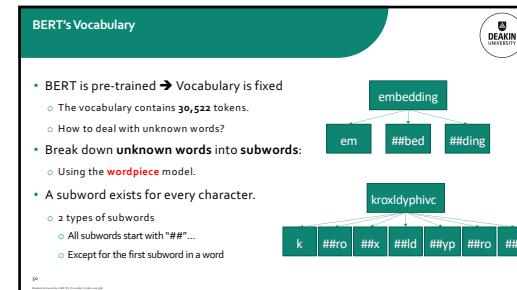
47



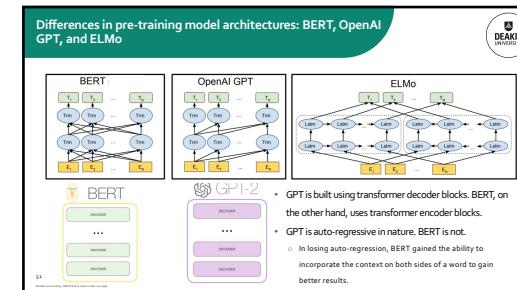
48



49



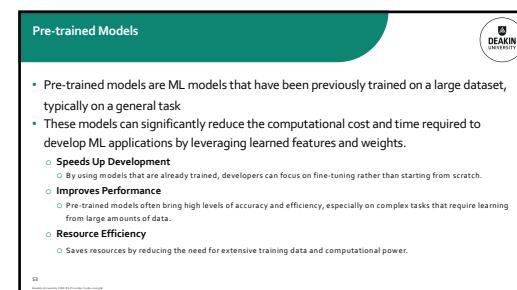
50



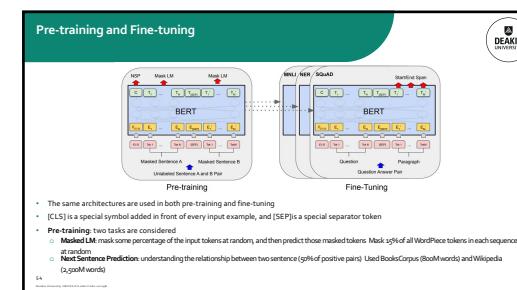
51



52



53



54

Pre-training

Task#1: Masked Language Model (MLM)

- The original approach to training bidirectional encoders is called **Masked Language Modeling (MLM)**
 - MLM uses unannotated text from a large corpus
 - Model is presented with a series of sentences from the training corpus, where a random sample of tokens from each training sequence is selected for use in the learning task. Once chosen, a token is used in one of three ways:
 - It is replaced with the unique vocabulary token [MASK]
 - It is replaced with another token from the vocabulary, randomly sampled based on token unigram probabilities
 - It is left unchanged
- In BERT, 15% of the input tokens in a training sequence are sampled for learning
 - Of these, 80% are replaced with [MASK], 10% are replaced with randomly selected tokens, and the remaining 10% are left unchanged

55

Pre-training

Task#1: Masked Language Model (MLM)

- The MLM training objective is to predict the original inputs for each of the masked tokens using a bidirectional encoder of the kind described in the last section
 - The cross-entropy loss from these predictions drives the training process for all the parameters in the model
 - Note that all the input tokens play a role in the self-attention process, but **only the sampled tokens** are used for learning
- Input:**
 - Original input sequence is first tokenized using a subword model
 - The sampled items which drive the learning process are chosen from among the set of tokenized inputs
 - Word embeddings for all the tokens in the input are retrieved from the word embedding matrix and then combined with positional embeddings to form the input to the transformer

56

Pre-training

Task#1: Masked Language Model (MLM)

Figure 11.4 Masked language model training. In this example, three of the input tokens are selected, two of which are masked and one is not. The masked tokens are passed through the transformer encoder to predict their original values. Note that the other five words in the sequence do not play a role in the training loss.

57

Pre-training

Task#2: Next Sentence Prediction (NSP)

- An important class of applications involves determining the relationship between pairs of sentences
 - paraphrase detection (detecting if two sentences have similar meanings)
 - entailment (detecting if the meanings of two sentences entail or contradict each other)
 - discourse coherence (deciding if two neighboring sentences form a coherent discourse)
- To capture the kind of knowledge required for applications such as these, BERT introduced a second learning objective called Next Sentence Prediction (NSP)
- Training:** The model is presented with pairs of sentences and is asked to predict whether each pair consists of an actual pair of adjacent sentences from the training corpus or a pair of unrelated sentences

58

Pre-training

Task#2: Next Sentence Prediction (NSP)

- In BERT, 50% of the training pairs consisted of positive pairs, and in the other 50% the second sentence of a pair was randomly selected from elsewhere in the corpus
 - The NSP loss is based on how well the model can distinguish true pairs from random pairs
- BERT introduces two new tokens to the input representation
 - After tokenizing the input with the subword model, the token [CLS] is prepended to the input sentence pair, and the token [SEP] is placed between the sentences and after the final token of the second sentence
 - During training, the output vector from the final layer associated with the [CLS] token represents the next sentence prediction

59

Pre-training

Task#2: Next Sentence Prediction (NSP)

Figure 11.5 An example of the NSP loss calculation.

60

Pre-training procedure

- Training data: BooksCorpus (800M words) + English Wikipedia (2.5B words).
- To generate each training input sequences: sample two spans of text (A and B) from the corpus.
 - The combined length is ≤ 500 tokens.
 - 50% B is the actual next sentence that follows A and 50% of the time it is a random sentence from the corpus.
- The training loss is the sum of the mean masked LM likelihood and the mean next sentence prediction likelihood.

61

Under the Hood

ML Architecture Parts	Definition
Parameters	Number of learnable variables/values available for the model
Transformer Layers	Non-linear processing blocks. A transformer block transforms a sequence of word representations to a sequence of context-aware representations (numerical representations)
Hidden Size	Layers of mathematical functions, located between the input and output, that assign weights to words to produce a desired result
Attention Heads	The size of a Transformer block
Processing	Type of processing unit used to train the model
Length of Training	Time it took to train the model

BERT Size & Architecture

Transformer Layers	Hidden Size	Attention Heads	Parameters	Processing	Length of Training	
BERTbase	12	768	12	110M	4 TPUs	4d
BERTlarge	24	1024	16	340M	16 TPUs	4d

62

SIT330-770: Natural Language Processing
Week 7.9 – BERT fine-tuning

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

63

Fine-tuning procedure 1: Classification

- For sequence-level classification task
 - Obtain the representation of the input sequence by using the final hidden state (hidden state at the position of the special token [CLS]) $C \in R^H$
 - Just add a classification layer and use softmax to calculate label probabilities. Parameters $W \in R^{K \times H}$

$$P = \text{softmax}(CW^T)$$

64

Fine-tuning procedure 2: Sentence Pair Classification

- For sentence pair classification task
 - Sentences are separated with a special token [SEP]
 - Obtain the representation of the input sequence by using the final hidden state (hidden state at the position of the special token [CLS]) $C \in R^H$
 - Just add a classification layer and use softmax to calculate label probabilities. Parameters $W \in R^{K \times H}$

$$P = \text{softmax}(CW^T)$$

65

Fine-tuning procedure 3: Named Entity Recognition

- Feed the final hidden representation $T_i \in R^H$ for each token i into a classification layer for the tagset.
- To make the task compatible with WordPiece tokenization
 - Predict the tag for the first sub-token of a word
 - No prediction is made for X

66

Fine-tuning procedure 4: Query Answering 1/2

- Input Question:**
Where do water droplets collide with ice crystals to form precipitation?
- Input Paragraph:**
.... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud.
- Output Answer:**
within a cloud

67

Fine-tuning procedure 4: Query Answering 2/2

- Represent the input question and paragraph as a single packed sequence.
 - The question uses the **A** embedding and the paragraph uses the **B** embedding.
- New parameters to be learned in fine-tuning are start vector $S \in R^H$ and end vector $E \in R^H$.
- Calculate the probability of word i being the start of the answer span:
 $P_{Start} = \text{Softmax}(ST^T)$ and $P_{end} = \text{Softmax}(ET^T)$
- The training objective is the log-likelihood of the correct start and end positions.

68

SIT330-770: Natural Language Processing

Week 7.10 – BERT Performance

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

69

Experiments

- GLUE (General Language Understanding Evaluation) benchmark**
 - Distributes canonical Train, Dev and Test splits
 - Labels for Test set are not provided
- Datasets in GLUE:**
 - MNLI: Multi-Genre Natural Language Inference
 - QQP: Quora Question Pairs
 - ONLU: Question Natural Language Inference
 - SST-2: Stanford Sentiment Treebank
 - CoLA: The corpus of Linguistic Acceptability
 - STS-B: The Semantic Textual Similarity Benchmark
 - MRPC: Microsoft Research Paraphrase Corpus
 - RTE: Recognizing Textual Entailment
 - WNLI: Winograd NLI

70

GLUE Results

System	MNLI-(m/mm)	QQP	ONLU	SST-2	CoLA	STS-B	MRPC	RTE	Average
Pre-OpenAI SOTA	392	363k	108k	67k	8.5k	5.7k	3.3k	2.4k	
Pre-OpenAI SOTA	80,680.1	66.1	62.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76,476.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82,181.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84,693.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86,785.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.⁸ BERT and OpenAI GPT are single-model, single-task. F1 scores are reported for QQP and MRPC. Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

71

SQuAD: The Stanford Question Answering Dataset

System	Dev	Test			
	EM	F1	EM	F1	
Top Leaderboard Systems (Dec 10, 2018)	-	-	-	-	
Human	-	-	91.2	-	
#1 Ensemble - slot	-	-	86.0	91.7	
#2 Ensemble - QANet	-	-	84.5	90.5	
Published	-	-	-	-	
BiDAF+ELMo (Single)	-	-	85.6	88.5	
R.M Reader (Ensemble)	81.2	87.9	82.3	88.5	
Ours	80.8	88.5	-	-	
BERT _{BASE} (Single)	84.1	90.9	-	-	
BERT _{BASE} (Ensemble)	85.8	91.1	-	-	
BERT _{BASE} (Gen+TriviaQA)	85.1	85.1	85.3	91.8	
BERT _{BASE} (Gen+TriviaQA)	86.2	92.2	87.4	93.2	
BERT _{LARGE} (Single)	Ours	78.7	81.9	80.0	83.1

Table 2: SQuAD 1.1 results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

72

SWAG

- The Situations with Adversarial Generations (SWAG)

On stage, a woman takes a seat at the piano.
She ...

- sits on a bench as her sister plays with the doll.
- smiles with someone as the music plays.
- is in the crowd, watching the dancers.
- nervously sets her fingers on the keys.

- The only task-specific parameters is a vector $V \in R^H$
- The probability distribution is the softmax over the four choices

Table 4: SWAG Dev and Test accuracies. ¹Human performance is measured with 100 samples, as reported in the SWAG paper.

73

Conclusions

- Unsupervised pre-training (pre-training language model) is increasingly adopted in many NLP tasks.
- Google Search is applying BERT models for search queries for over 70 languages.
- Major contribution of the paper is to propose a deep bidirectional architecture from Transformer.
- Advance state-of-the-art for many important NLP tasks.
- Cannot do everything in NLP!**

74

SIT330-770: Natural Language Processing

Week 7.11 – Other Models Based on Transformers

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

75

Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks

- A modification of the pretrained BERT network that use Siamese network structures to derive semantically meaningful sentence embeddings that can be compared using cosine-similarity

<https://www.sbert.net/>

76

BERT-based Models

- RoBERTa: Robustly Optimized BERT Approach**
 - Facebook AI research team
 - They used 160 GB of text instead of the 16 GB dataset originally used to train BERT
 - Increased the number of iterations from 100K to 300K and then further to 500K
 - Dynamically changing the masking pattern applied to the training data
 - Removing the next sequence prediction objective from the training procedure
- ALBERT, XLNET, CoBERT

Machine performance on the RACE challenge (SAT-like reading comprehension)

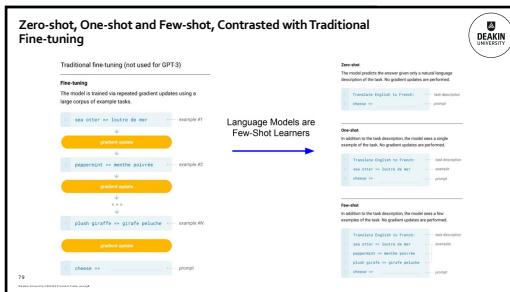
Model	Performance (%)
GPT	44.8
RoBERTa	46.6
BART	52.0
ALBERT	51.8
XLNET	51.2
CoBERT	54.4

77

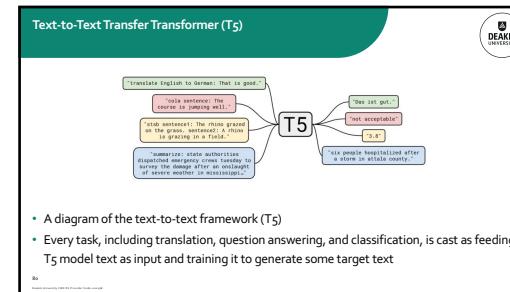
Generative Pre-trained Transformer (GPT)

- GPT** (Generative Pre-trained Transformer) is a series of language generation models developed by OpenAI. These models are based on the Transformer architecture (2018)
- GPT-2** (Generative Pre-trained Transformer 2) was the second model in the GPT series, released in 2019. It was trained on a large corpus of internet text and was designed for language generation tasks such as question answering, and text summarization
- GPT-3** (Generative Pre-trained Transformer 3) Released in 2020, with over 175 billion parameters, and was trained on a much larger and diverse dataset, including web pages, books, and scientific articles
- GPT-4** (Generative Pre-trained Transformer 4) Released in 2023, a multimodal model which can accept image and text inputs and produce text outputs

78



79



80



81

The HuggingFace website features a green header with the text 'More than 5,000 organizations are using Hugging Face'. Below the header, there's a section for 'AI2' partners like Intel, Meta AI, Microsoft, and Google AI. Another section highlights 'Problems solvers' such as audio classification, image classification, object detection, question answering, summarization, text classification, and translation. A note states: 'HuggingFace has a course on NLP: https://huggingface.co/course/chapter0/1?fw=pt_35'. At the bottom, there's a footer with links to 'About' and 'Contact'.

82

The documentation page for HuggingFace includes a green header with the text 'Documentation and Tutorials'. Below the header, there's a list of useful links:

- There is a wide range of examples provided, here are some useful links:
- Notebooks:
 - <https://huggingface.co/transformers/v3.0.2/notebooks.html>
 - <https://github.com/huggingface/transformers/tree/main/notebooks>
- Models:
 - <https://huggingface.co/models>
- Course:
 - <https://huggingface.co/course/chapter0/1?fw=pt>

83

The getting started page for HuggingFace includes a green header with the text 'Getting Started'. Below the header, there's a list of instructions:

- Use pip for the installation, which is the package manager for Python. In notebooks, you can run system commands by preceding them with the ! character, so you can install the Transformers library as follows:

```
!pip install transformers
```
- You can make sure the package was correctly installed by importing it within your Python runtime:

```
import transformers
```

84

Pipelines

- The pipelines are a great and easy way to use models for inference
- These pipelines are objects that abstract most of the complex code from the library, offering a simple API dedicated to several tasks, including Named Entity Recognition, Masked Language Modeling, Sentiment Analysis, Feature Extraction and Question Answering



85

Example (1): Sentiment Analysis

```
from transformers import pipeline
classifier = pipeline("sentiment-analysis")
classifier("I've been waiting for a HuggingFace course my whole life.")
No model was supplied, defaulted to distilbert-base-uncased-finetuned-tweet-classification-v0.0-999999
Using a pipeline without specifying a model name and revision in production is not recommended.
Downloading 100% [██████████] 62MB/62B [00:00<00:00, 28.5kB/s]
Downloading 100% [██████████] 28MB/28B [00:00<00:00, 83.2kB/s]
Downloading 100% [██████████] 48.048.0 [00:00<00:00, 1.9kB/s]
Downloading 100% [██████████] 232K/232K [00:00<00:00, 967kB/s]
[{"label": "POSITIVE", "score": 0.9598849521446228}, {"label": "NEGATIVE", "score": 0.9994558691978455}]
```



86

Example (2): Question Answering

```
from transformers import pipeline
question_answerer = pipeline("question-answering")
question_answerer(
    question="Where do I work?", context="My name is Sylvain and I work at Hugging Face in Brooklyn",
)
{'score': 0.6949767470359802, 'start': 33, 'end': 45, 'answer': 'Hugging Face'}
```

```
from transformers import pipeline
oracle = pipeline(model="deepset/roberta-base-squad2")
oracle(question="Where do I live?", context="My name is Wolfgang and I live in Berlin")
{'score': 0.9190717935562134, 'start': 34, 'end': 40, 'answer': 'Berlin'}
```



87



88

Summary



- Today we learned about:
 - Transformers
 - Attention is All you need
 - BERT



89