

SIT330-770: Natural Language Processing

Week 10 - Transformers and Pretrained LMs

Dr. Mohamed Reda Bouadjenek

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

reda.bouadjenek@deakin.edu.au

1

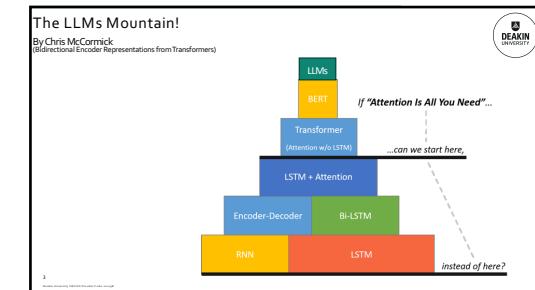
SIT330-770: Natural Language Processing

Week 10.1 – Introduction to Transformers

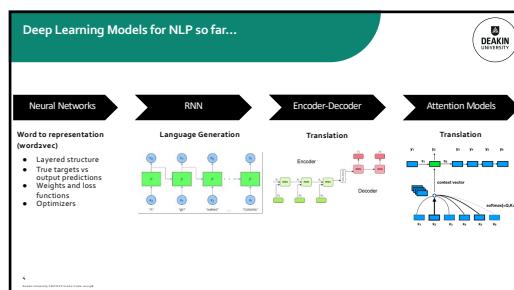
Dr. Mohamed Reda Bouadjenek

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

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

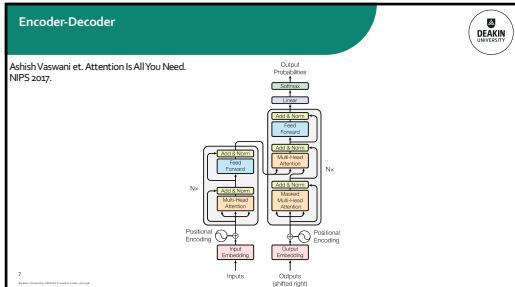
Attention Is All You Need

Ashish Vaswani¹ Google Brain arxiv.org/abs/1706.03762 Noam Shazeer² Google Brain arxiv.org/abs/1706.03762 Niki Parmar³ Google Research arxiv.org/abs/1706.03762 Jakob Uszkoreit⁴ Google Research arxiv.org/abs/1706.03762 Llion Jones⁵ Google Research arxiv.org/abs/1706.03762 Adam N. Gao⁶ University of Toronto arxiv.org/abs/1706.03762 Łukasz Kaiser⁷ Google Brain arxiv.org/abs/1706.03762 Illia Polosukhin⁸ arxiv.org/abs/1706.03762

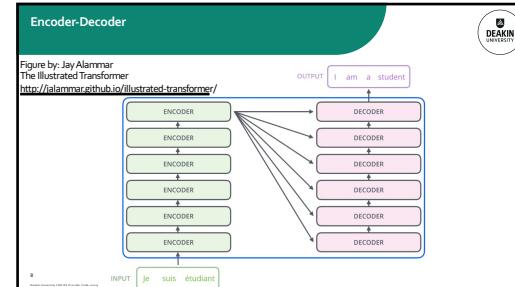
Abstract

The dominant sequence-to-sequence models are based on complex recurrent or

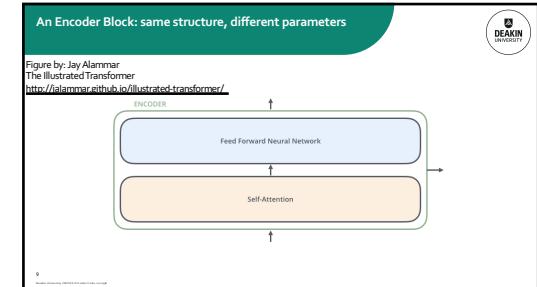
6



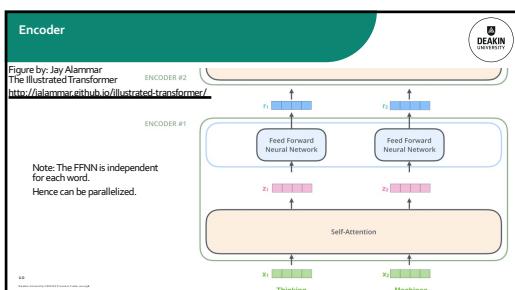
7



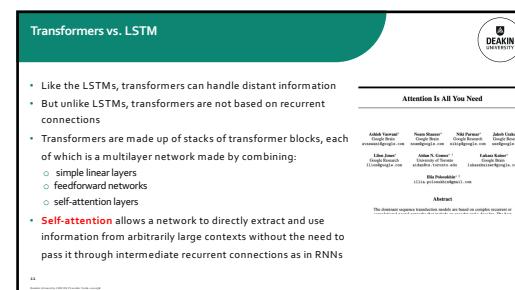
8



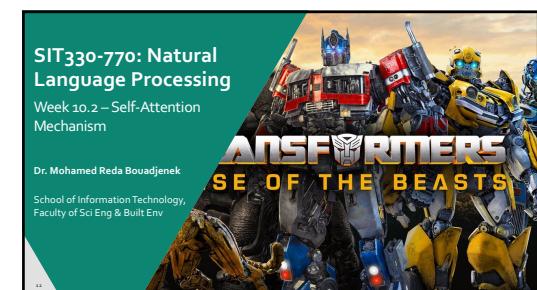
9



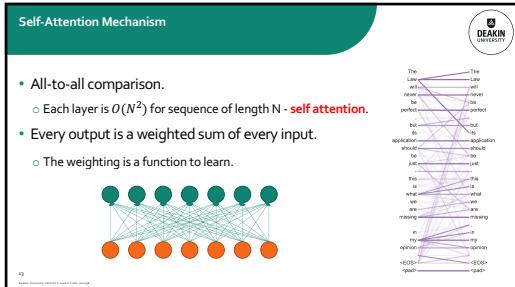
10



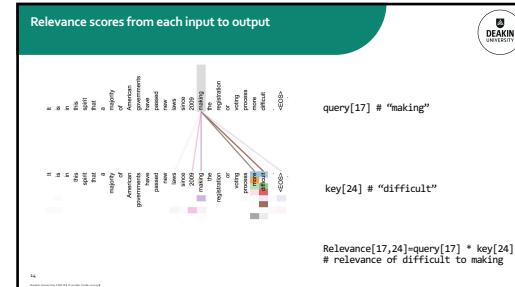
11



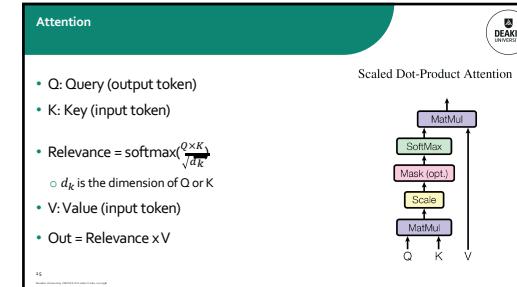
12



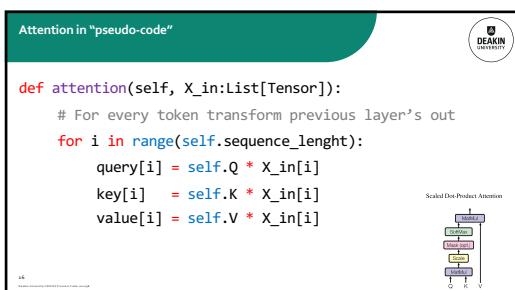
13



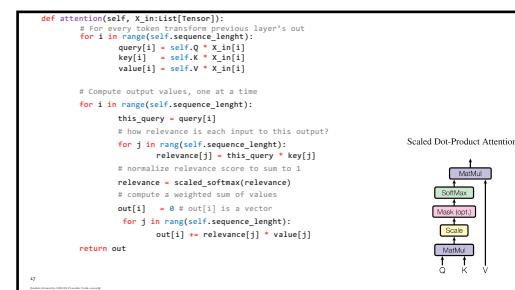
14



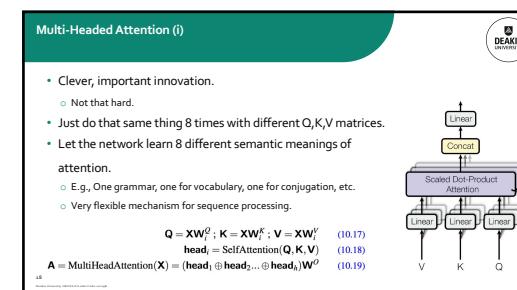
15



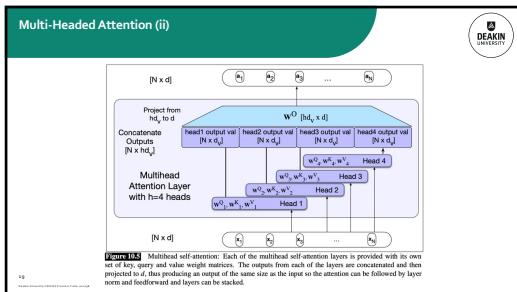
16



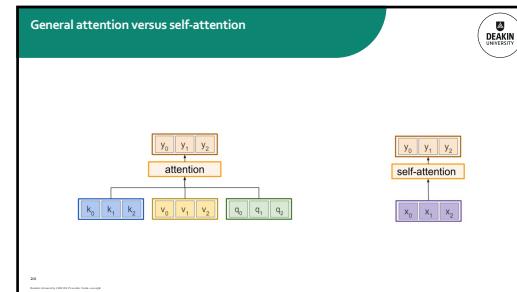
17



18



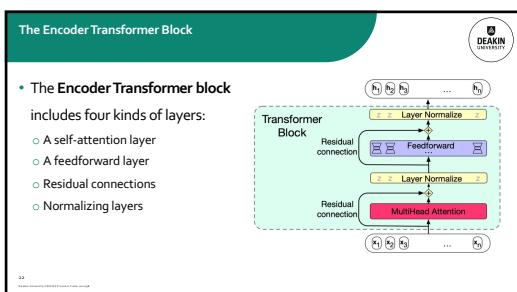
19



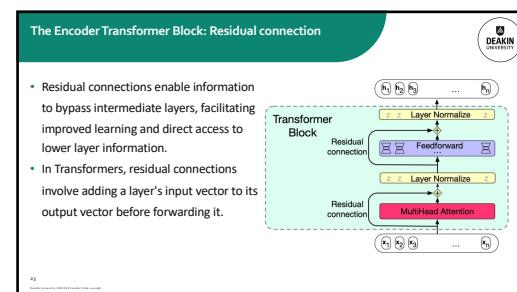
20



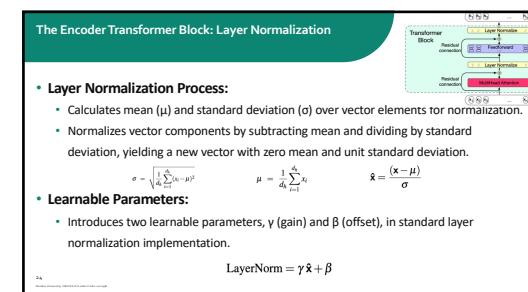
21



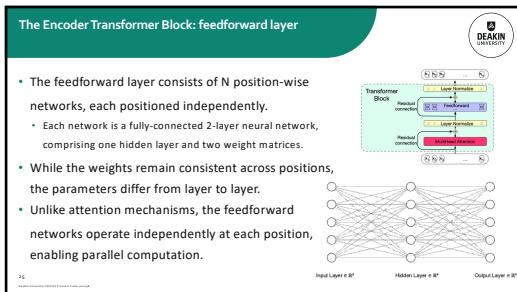
22



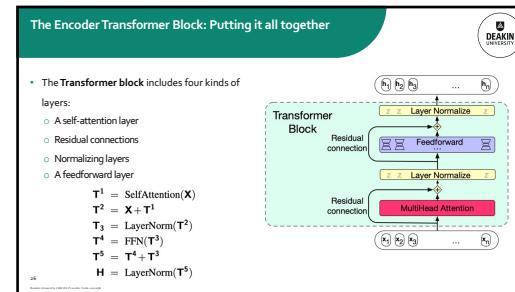
23



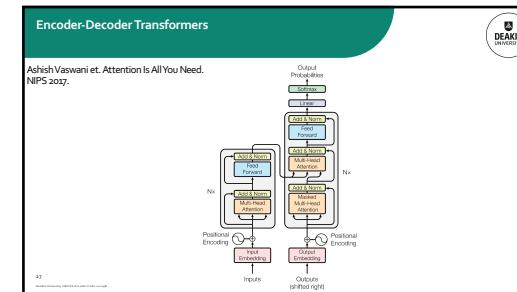
24



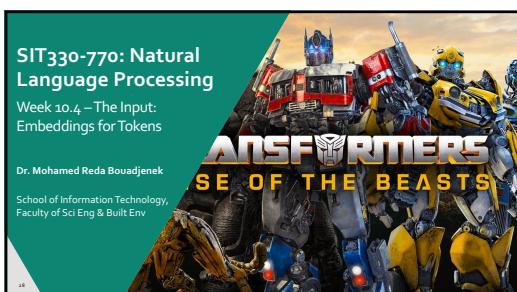
25



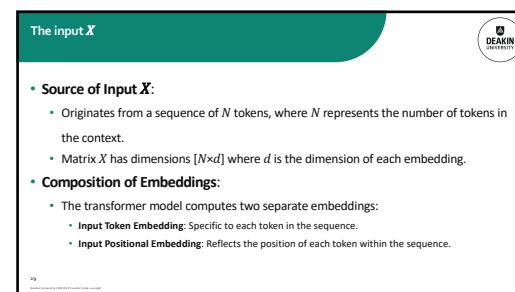
26



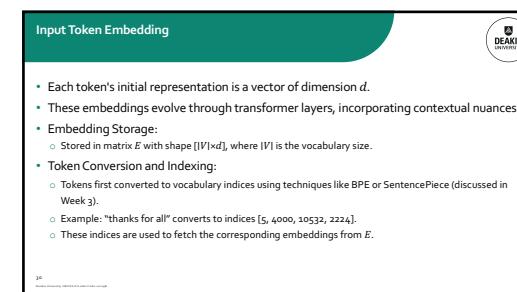
27



28



29



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_5 by multiplying the embedding matrix E with a one-hot vector with a 1 at index 5.

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,

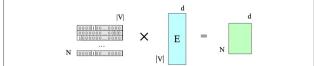


Figure 10.11 Selecting the embedding matrix for the input sequence of token ids W by multiplying a $N \times |V|$ one-hot matrix corresponding to W by the embedding matrix E .

32

SIT330-770: Natural Language Processing

Week 10.5 – 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

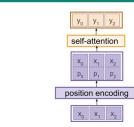
34

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.
- Concatenate/add special positional encoding p to each input vector x
- We use a function pos: $N \rightarrow d$ to process the position j of the vector into a d -dimensional vector
- So, $p_j = pos(j)$

36

Positional Encoding Layer in Transformers

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

- 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 10.6 – The Task Specific Head

Dr. Mohamed Reda Bouadjenek

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

40

The Language Modeling Head

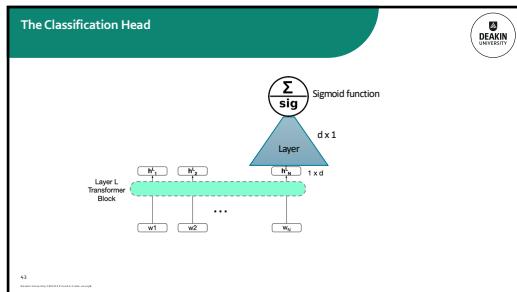
- The Language Modeling Head is an additional neural circuitry integrated with the basic transformer architecture.
- It enables specific tasks such as language modeling by enhancing the transformer's capabilities.
- Given a context of words, a Language Model assigns a probability to each possible next 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 (\vec{h}_N) to a probability distribution over words in the vocabulary V .

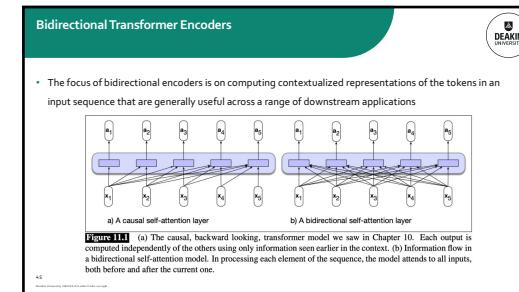
42



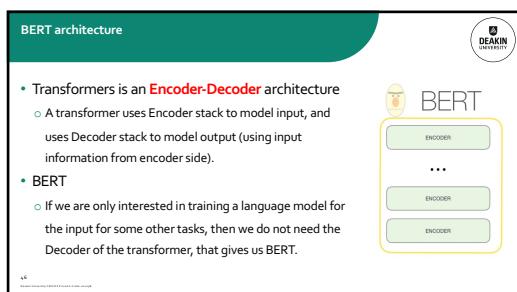
43



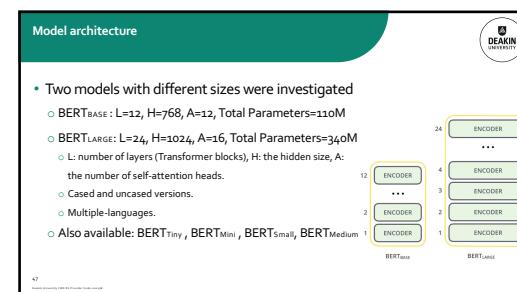
44



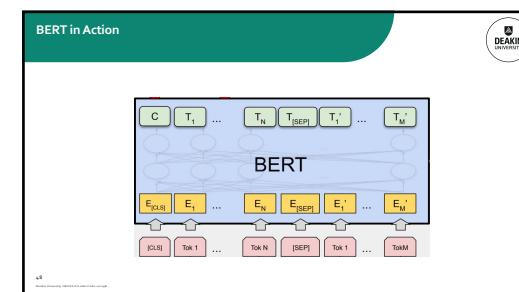
45



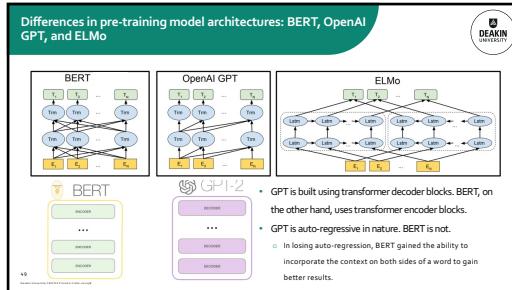
46



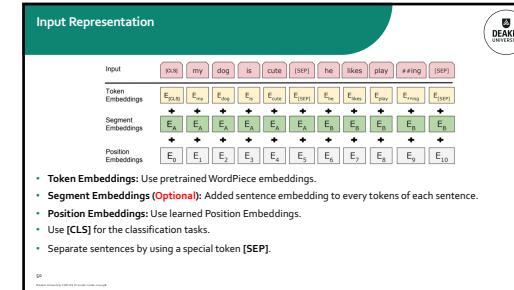
47



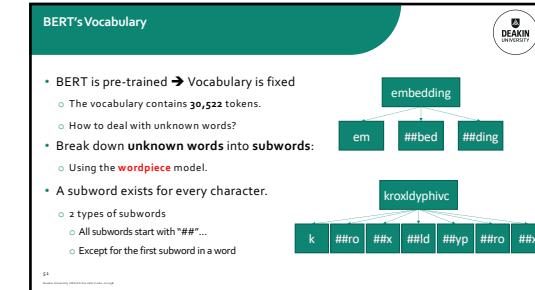
48



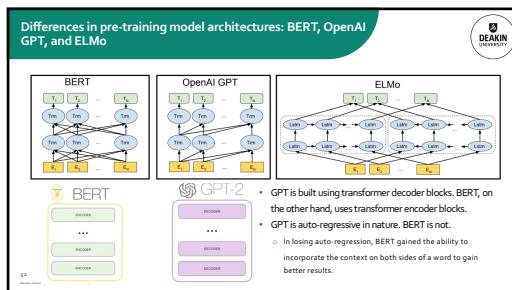
49



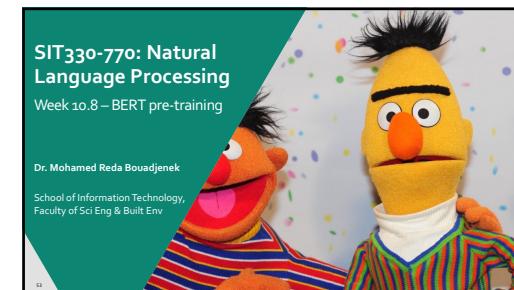
50



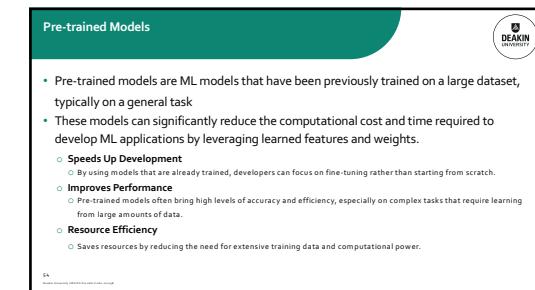
51



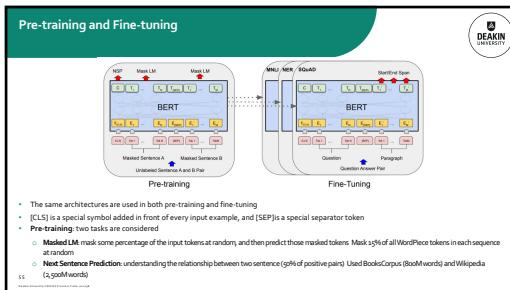
52



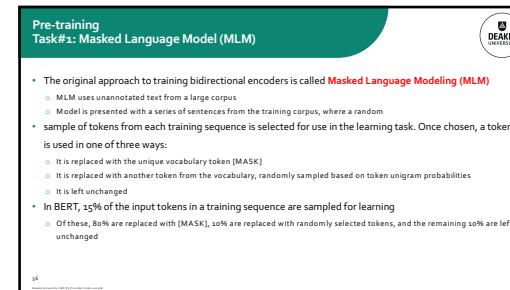
53



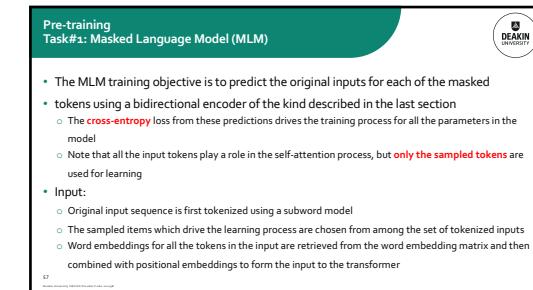
54



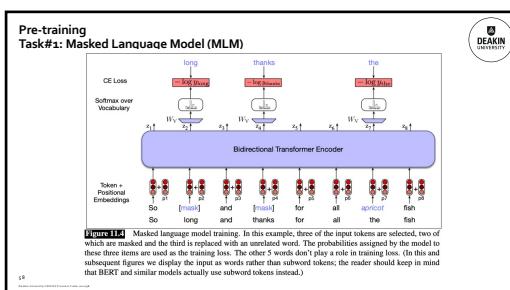
55



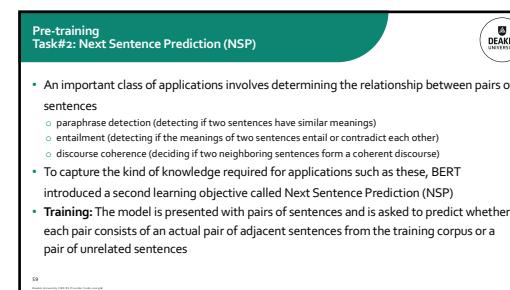
56



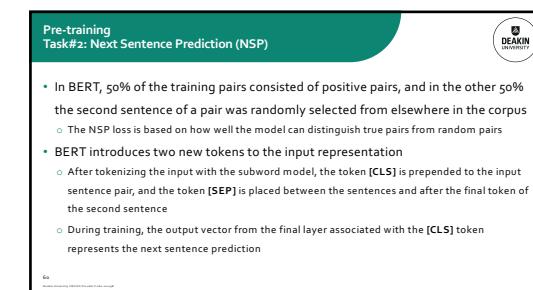
57



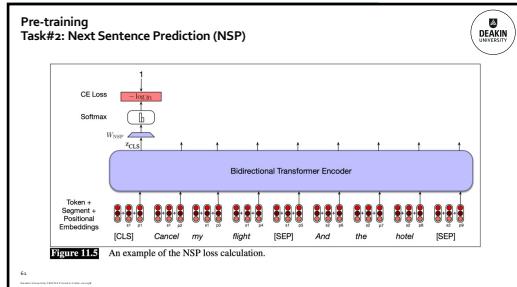
58



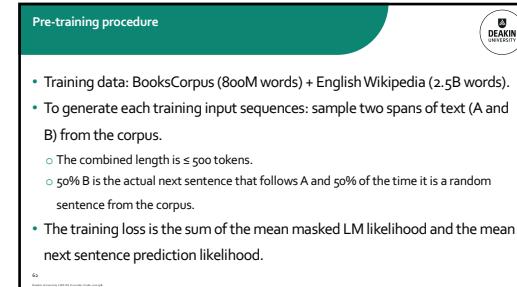
59



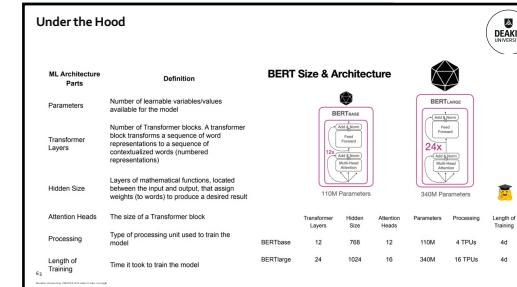
60



61



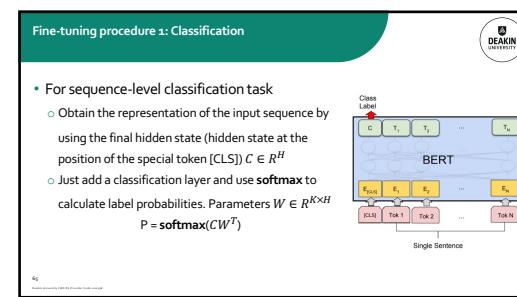
62



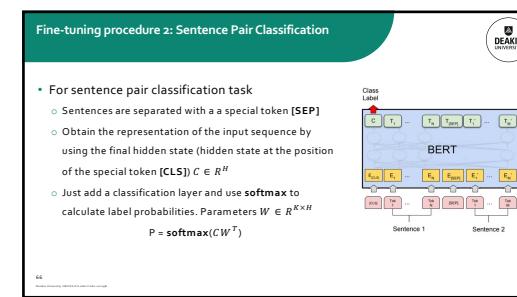
63



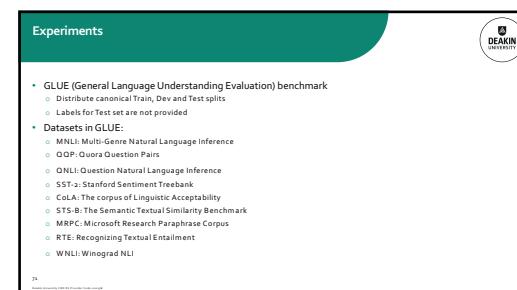
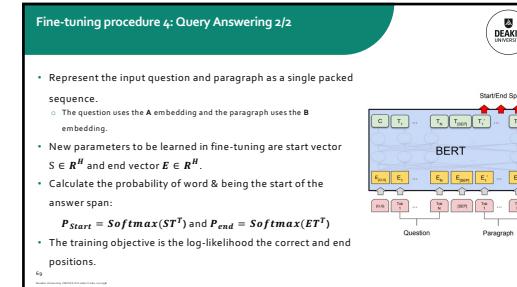
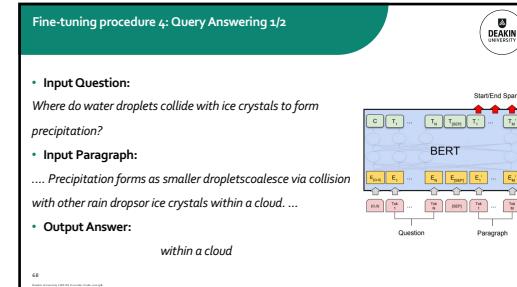
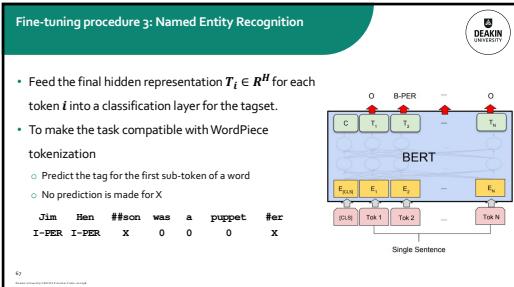
64



65



66



GLUE Results

System	MNLI (em/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
Pre-OpenAI SOTA	39.2	36.3	108.8	35.0	81.0	86.0	61.7	74.0	
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{base}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{large}	86.7/85.9	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. 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.

72

SQuAD: The Stanford Question Answering Dataset

System	Dev	Test	EM	F1	EM	F1
Human	-	-	82.3	91.2		
Top Leaderboard Systems (Dec 10th, 2018)	-	-	84.7	90.5		
#1 Ensemble - albert	-	-	84.5	90.5		
#2 Ensemble - QANet	-	-	84.5	90.5		
Ours	80.8	88.5	-	-		
BERT4LMo (Single)	81.2	87.9	82.3	88.5		
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5		
BERT4LMo (Single) + Ours	80.8	88.5	-	-		
BERT4LMo (Single) + R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5		
BERT4LMo (Single) + SLQA+ (Ensemble)	83.8	91.8	-	-		
BERT4LMo (Ensemble)	83.8	91.8	-	-		
BERT4LMo (Sgl+TriviaQA)	84.2	91.1	85.1	91.8		
BERT4LMo (Em+TriviaQA)	86.2	92.2	87.4	93.2		

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

73

SWAG

System	Dev	Test	EM	F1	EM	F1
Human	96.3	89.0	95.9	95.9	-	-
#1 Single - MIR-MRC (F-Net)	-	-	74.8	78.0		
#2 Single - sheet	-	-	74.2	77.1		
Ours	-	-	71.4	74.9		
SLQA+ (Single)	-	-	71.4	74.9		
BERT4LARGE (Single) + Ours	78.7	81.9	80.0	83.1		

Table 3: SQuAD 2.0 results. We exclude entries that use BERT as one of their components.

• The Situations with Adversarial Generations (SWAG)
 On stage, a woman takes a seat at the piano.
 She –
 a) sits on a bench as her sister plays with the doll.
 b) smiles with someone as the music plays.
 c) is in the crowd, watching the dancers.
 d) nervously sets her fingers on the keys.

• The only task-specific parameters is a vector $V \in \mathbb{R}^H$

• The probability distribution is the softmax over the four choices

74

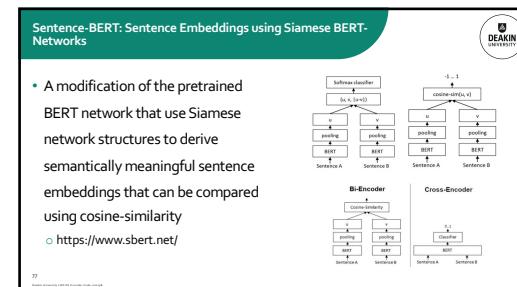
Conclusions

- 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!**

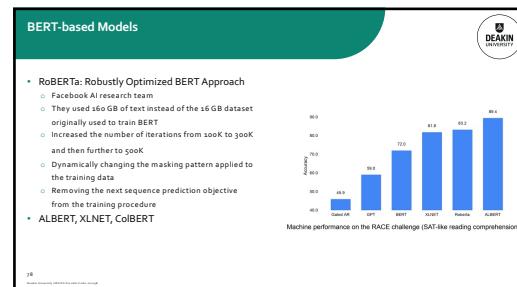
75



76



77



78

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

79

Zero-shot, One-shot and Few-shot, Contrasted with Traditional Fine-tuning

Traditional fine-tuning (not used for GPT-3)

Fine-tuning: The model is trained via repeated gradient updates using a large corpus of example texts.

example #1: sea otter == loutre de mer
example #2: peacock == paon bleu
example #3: penguin == pingouin
example #4: plain giraffe == girafe bancale
example #5: polar bear == ours polaire

choose or prompt

Language Models are Few-Shot Learners

Zero-shot: The model predicts the answer given only a natural language description of the task. No pre-train updates are performed.

Translate English to French
and description
choose or prompt

Few-shot: In addition to the task description, the model uses a single example of the task. No pre-train updates are performed.

Translate English to French
and description
choose or prompt

Few-shot: In addition to the task description, the model uses a few examples of the task. No pre-train updates are performed.

Translate English to French
and description
choose or prompt

80

Text-to-Text Transfer Transformer (T5)

translate English to German: That is good.
with sentence: The rhino grazed on the grass in the field.
summary: state authorities dispatched emergency crews Tuesday to help clear debris from a severe winter storm in Mississippi.
Dev test get...
not acceptable
3.0!
six people hospitalized after a storm in state county.

- A diagram of the text-to-text framework (T5)
- Every task, including translation, question answering, and classification, is cast as feeding T5 model text as input and training it to generate some target text

81

SIT330-770: Natural Language Processing

Week 10.12 – HuggingFace

Dr. Mohamed Reda Bouadjenek

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

82

HuggingFace

More than 5,000 organizations are using Hugging Face

Allegro Institute for AI	Meta AI	Grayphore	Google AI
Integrations - 100+ models	Commissary - 10+ models	Commodity - 10+ models	Cloudinary - 100+ models
Intel - 100+ models	Speechbrain - 100+ models	Microsoft - 200+ models	Grammy

Problems solvers

- Document classification and named entity recognition
- Text generation and language modeling

Future work

- HuggingFace has a course on NLP: <https://huggingface.co/course/chapter0/2?fw=pt>

83

Documentations and Tutorials

- 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/chapter0/2?fw=pt>

84

Getting Started

- 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
```

85

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

86

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, default to distilbert-base-uncased-finetuned-sst-2-english and revision a90f99b
# Using a pipeline without specifying a model name and revision in production is not recommended.
Downloading: 100% [██████████] 620620 [00:00<00:00, 26.9kB/s]
[{"label": "POSITIVE", "score": 0.9598849521446228}, {"label": "NEGATIVE", "score": 0.9994358693978459}]
```

```
classifier(
    ["I've been waiting for a HuggingFace course my whole life.", "I hate this so much!"])
[{"label": "POSITIVE", "score": 0.9598849521446228}, {"label": "NEGATIVE", "score": 0.9994358693978459}]
```

87

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'}
```

88



89

Summary

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

90