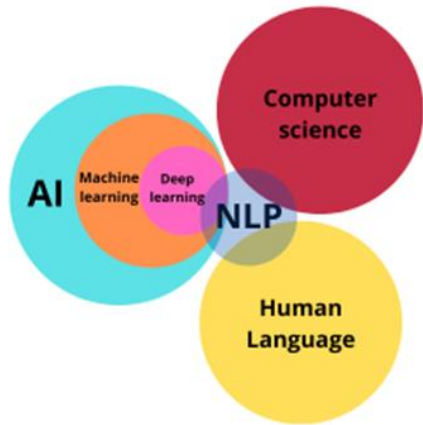# Grasping Deep Learning from Fundamentals to Applications
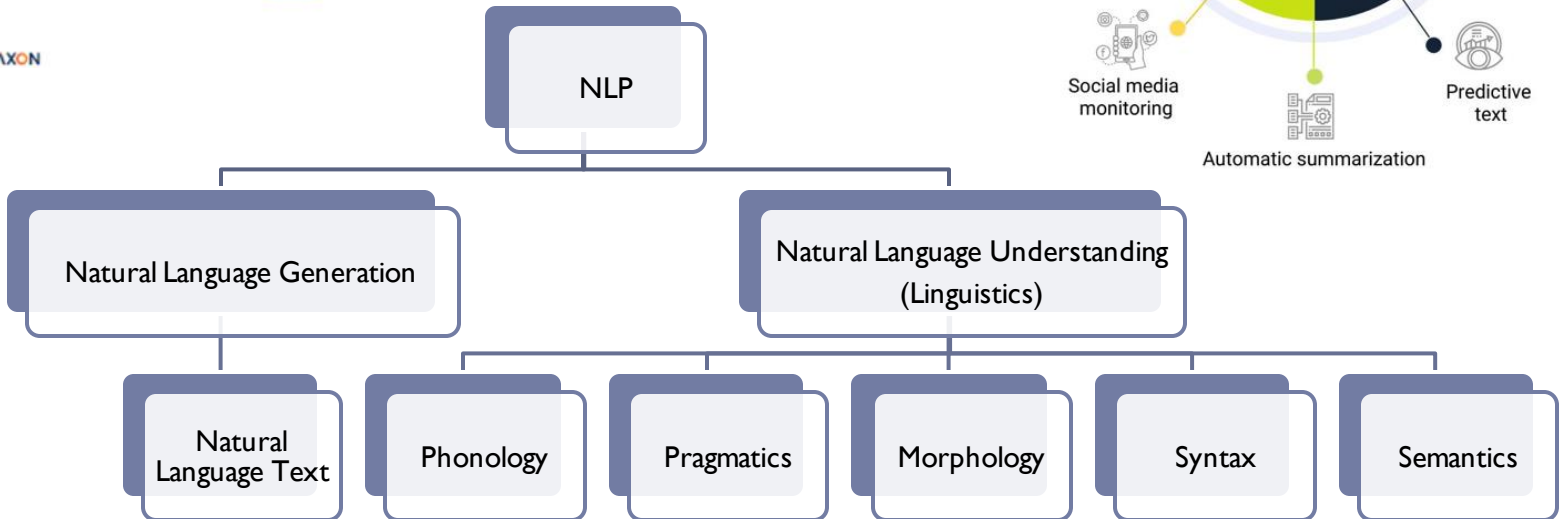
*June 15, 2023*

## Lecture 3 – Introduction to Natural Language Processing

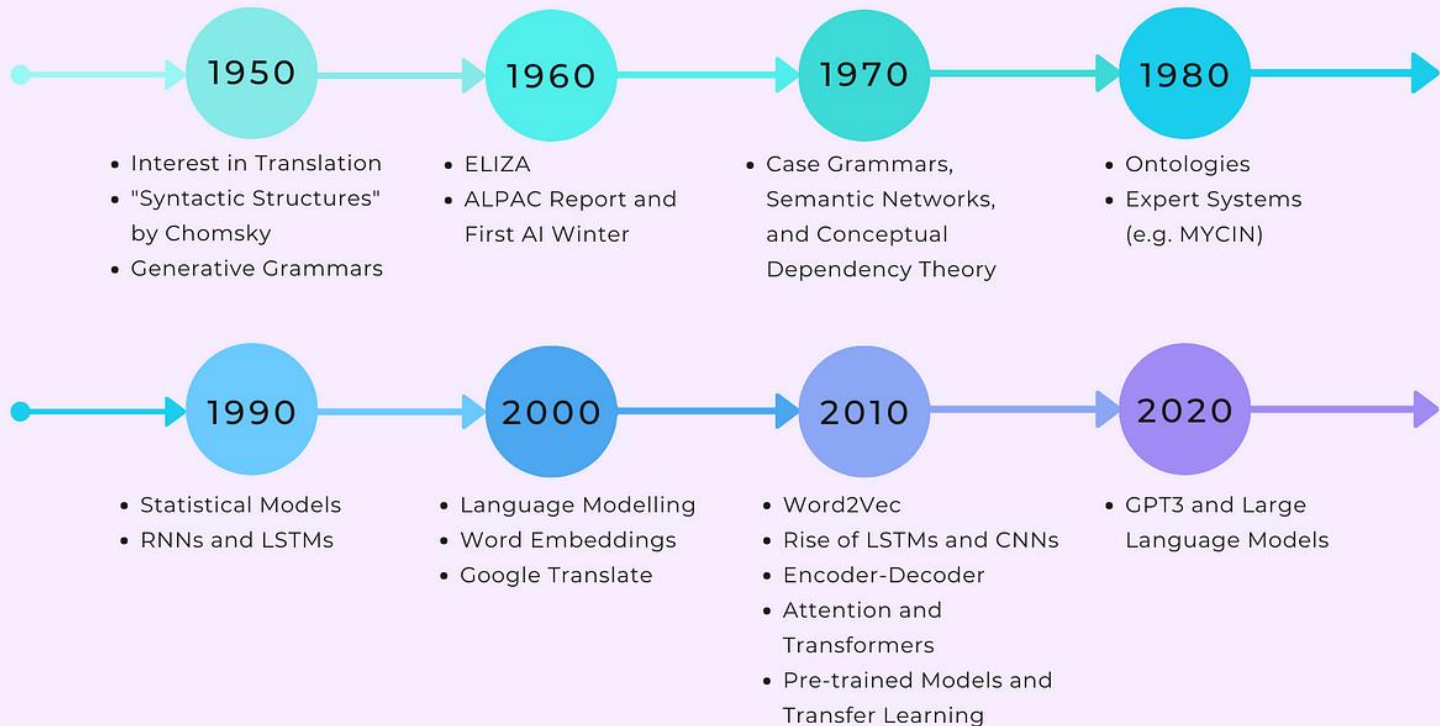Instructors: **Yufei Huang**, PhD; **Arun Das**, PhD

# Evolution of Natural Language Processing (NLP)



NLP

Natural Language Generation
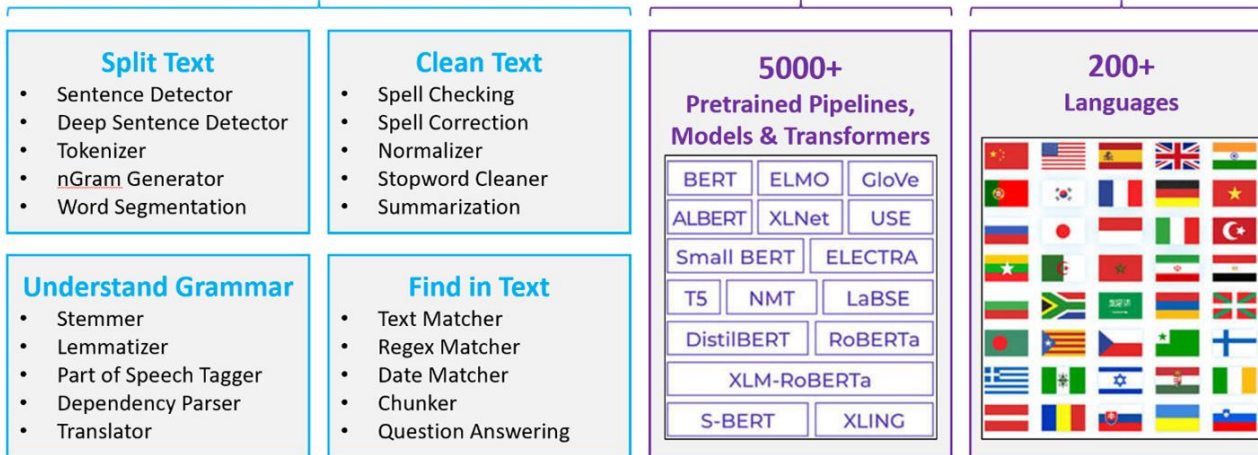
Natural Language Understanding (Linguistics)

Natural Language Text

Phonology

Pragmatics

Morphology

Syntax

Semantics

# A Brief Timeline of NLP

**1950**
- Interest in Translation
- "Syntactic Structures" by Chomsky
- Generative Grammars

**1960**
- ELIZA
- ALPAC Report and First AI Winter

**1970**
- Case Grammars, Semantic Networks, and Conceptual Dependency Theory

**1980**
- Ontologies
- Expert Systems (e.g. MYCIN)

**1990**
- Statistical Models
- RNNs and LSTMs

**2000**
- Language Modelling
- Word Embeddings
- Google Translate

**2010**
- Word2Vec
- Rise of LSTMs and CNNs
- Encoder-Decoder
- Attention and Transformers
- Pre-trained Models and Transfer Learning

**2020**
- GPT3 and Large Language Models

https://medium.com/nlplanet/a-brief-timeline-of-nlp-bc45b640f07d

University of **Pittsburgh**

# Heavy Investments!

| Entity Recognition | Information Extraction | Spelling & Grammar | Text Classification |
|---|---|---|---|
| Translation | Summarization | Question Answering | Emotion Detection |

**Annotators** | **Models** | **Languages**

**Split Text**
- Sentence Detector
- Deep Sentence Detector
- Tokenizer
- nGram Generator
- Word Segmentation

**Clean Text**
- Spell Checking
- Spell Correction
- Normalizer
- Stopword Cleaner
- Summarization

**5000+**
**Pretrained Pipelines, Models & Transformers**

| BERT | ELMO | GloVe |
| ALBERT | XLNet | USE |
| Small BERT | ELECTRA | |
| T5 | NMT | LaBSE |
| DistilBERT | RoBERTa | |
| XLM-RoBERTa | | |
| S-BERT | XLING | |

**200+**
**Languages**

**Understand Grammar**
- Stemmer
- Lemmatizer
- Part of Speech Tagger
- Dependency Parser
- Translator

**Find in Text**
- Text Matcher
- Regex Matcher
- Date Matcher
- Chunker
- Question Answering

| Trainable & Tunable | Scalable to a Cluster | Fast Inference | Hardware Optimized | Community |
|---|---|---|---|---|

Spark | intel | NVIDIA | NLP SUMMIT

ChatGPT

Microsoft Azure

University of Pittsburgh

Fintech Company A
→ Data, processes →
Chatbot SaaS Product
→ Chatbot for Company A

What is the fee for international money transfer?

*There is no extra fee for international transfers.*

Fintech Company B
→ Data, processes →
Chatbot SaaS Product
→ Chatbot for Company B

What is the fee for international money transfer?

*We charge a $5 flat fee for international transfers.*

University of Pittsburgh

# State-of-the-Art in Question Answering

QUESTION ANSWERING

ARITHMETIC

LANGUAGE UNDERSTANDING

**8 billion parameters**

https://ai.googleblog.com/2022/04/pathways-language-model-palm-scaling-to.html

# A Brief Timeline of NLP



**1950**
- Interest in Translation
- "Syntactic Structures" by Chomsky
- Generative Grammars

**1960**
- ELIZA
- ALPAC Report and First AI Winter

**1970**
- Case Grammars, Semantic Networks, and Conceptual Dependency Theory

**1980**
- Ontologies
- Expert Systems (e.g. MYCIN)

**1990**
- Statistical Models
- RNNs and LSTMs

**2000**
- Language Modelling
- Word Embeddings
- Google Translate

**2010**
- Word2Vec
- Rise of LSTMs and CNNs
- Encoder-Decoder
- Attention and Transformers
- Pre-trained Models and Transfer Learning

**2020**
- GPT3 and Large Language Models



Input = 0

Units=1

$h$

Prediction, $\tilde{y}$ = 12.1
Label, y = 32

$x$ → $w$ → $\tilde{y}$

Back Propagation

One Artificial Neuron

## Recurrent Neural Network (RNN)

output units

hidden units

input units

| time 1 output units | time 2 output units | time 3 output units |
| time 1 hidden units | time 2 hidden units | time 3 hidden units |
| time 1 input units | time 2 input units | time 3 input units |

*good for learning temporal associations in the input or input sequences

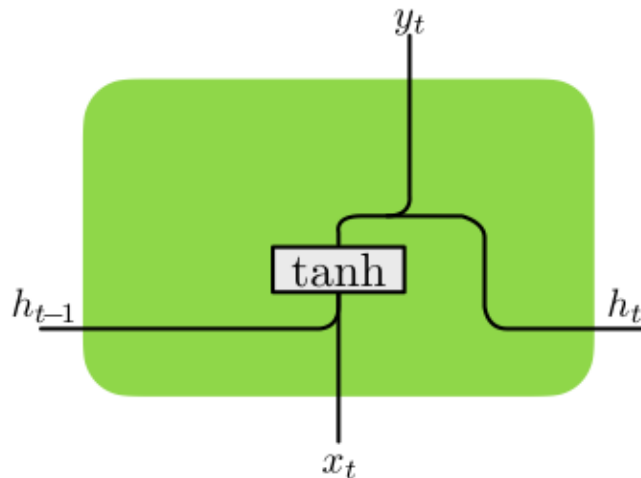University of Pittsburgh

# one to one    one to many    many to one    many to many    many to many

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_{t-1} + b_h)$$

$$y_t = W_{ho}h_t$$

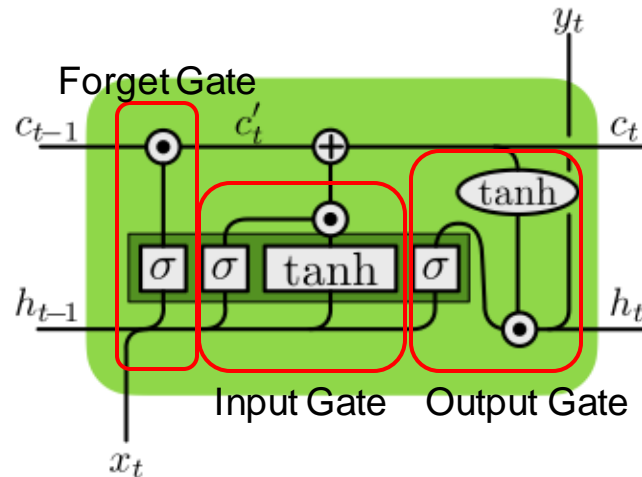University of Pittsburgh

# Long Short-Term Memory (LSTM)

- Adds cell state to the RNN.
- Adds four gates to control the flow of information.
- Carries computation sequentially in three steps.

# Long Short-Term Memory (LSTM)

- Forget Step: $c_t' = \sigma(W_{hf}h_{t-1} + W_{xf}x_t + b_f) \odot c_{t-1}$

- Input and Update Step: $g_t = \tanh(W_{hu}h_{t-1} + W_{xu}x_t + b_u)$
$$i_t = \sigma(W_{hi}h_{t-1} + W_{xi}x_t + b_i)$$
$$c_t = c_t' + i_t \odot g_t$$

- Output Step: $h_t = \sigma(W_{ho}h_{t-1} + W_{xo}x_t + b_o) \odot \tanh(c_t)$

\* Forgets specific information of the cell state.

\* Decide what and how much to add to the cell state.

\* Decide how much of the information stored in the cell state should be written to the new hidden state.



University of Pittsburgh

# Text Data Processing

•*Standardization*  refers to preprocessing the text, typically to remove punctuation or HTML elements to simplify the dataset.
•*Tokenization* refers to splitting strings into tokens (for example, splitting a sentence into individual words by splitting on whitespace).
•*Vectorization* refers to converting tokens into numbers so they can be fed into a neural network.

Textual data is usually preprocessed using the following 5 tasks:

1. Standardize each example (usually lowercasing + punctuation stripping)
2. Split each example into substrings (usually words)
3. Recombine substrings into tokens (usually ngrams)
4. Index tokens (associate a unique int value with each token)
5. Transform each example using this index into a vector of ints or a dense float vector.

University of
Pittsburgh

**Original Sentence:**

Yes it was a little low budget, but this movie shows love!

**Standardization:**

yes it was a little low budget but this movie shows love

**Tokenization:**

[yes, it, was, a, little, low, budget, but, this, movie, shows, love]

**Vectorization:**

[414    9  10  192   20    25   200    200   250    300  0 0 0 0 …]

# Subword Tokenizers

**Original Sentence:**
but they did n't test for curiosity .

Tokenization:
[b'[START]', b'but', b'they', b'did', b'n', b"'", b't', b'test', b'for', b'curiosity', b'.', b'[END]']

Vectorization:
[2, 87, 83, 149, 50, 9, 56, 664, 85, 2512, 15, 3]

\* Skipped standardization step in this example.

University of
Pittsburgh

# Typical DL Pipeline

| Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
|--------|--------|--------|--------|--------|
| Load and Preprocess Data | Data Visualization | Design the Deep Learning Model | Train the Model Define | Evaluate the Model |

**Deep Learning Training Process**



Test with your own text

This product was very bad!

Results

TAG            CONFIDENCE

Negative            99.7%

Classify Text

You can try this service from here: Semantic Analysis

- Weights are randomly initialized at the beginning.
- We know the actual labels (supervised).

1. Input data -> Batch of data.
2. We find the predictions.
3. We pass the predictions to the optimizer (the optimizer already knows the actual labels.)
4. We find the loss between predicted labels and actual labels.
5. We tune the "learnable" parameters according to the loss.
6. Go back to step 1.

University of Pittsburgh

Input sentence                                             label

**Yes it was a little low budget, but this movie shows love!**         1


Text Encoder

Vocabulary (** 1000 words ) – VOCAB SIZE

**[414   9  10  192   20    25    200    200    250    300  0 0 0 0 …]**

Vocabulary

```
array(['', '[UNK]', 'the', 'and', 'a', 'of', 'to', 'is', 'in', 'it', 'i',
       'this', 'that', 'br', 'was', 'as', 'for', 'with', 'movie', 'but'],
      dtype='<U14')
```

Encoded sentence

```
array([[ 10, 540,    4, ...,    0,    0,    0],
       [ 10, 120,   11, ...,    0,    0,    0],
       [414,    9,   14, ...,    0,    0,    0]])
```

* UNK token for unknown words that didn't fit in the set vocabulary size.
* Multiple words represented by same vector encoding.

University of Pittsburgh

Forget Gate

Input Gate   Output Gate

Bidirectional LSTM



University of Pittsburgh

Input

| That | was | a | great | movie |

TextVectorization

| 5 | 8 | 3 | 120 | 35 |

| E[5] | E[8] | E[3] | E[120] | E[35] |

Embedding          "Learnable" Parameter

                        or "trainable"

| 1 | 2 | 3 | 4 | 5 |

| 1 | 2 | 3 | 4 | 5 |

Bidirectional      Learnable Parameter

| 64 |

Dense              Learnable Parameter

| 0/1 |

Classification     Learnable Parameter

University of Pittsburgh

- Complex!
- May not have an exact match for phrases or words.
- Grammar/humor/context + cultural differences.

# Transformers

▸ Transformers are parallelizable.

▸ Transformers can capture distant or long-range contexts and dependencies.

▸ Transformers make no assumptions about the temporal/spatial relationships across the data.



**Attention Is All You Need**

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.
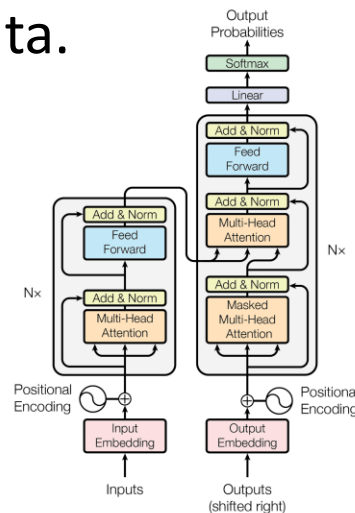
## 1 Introduction

Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [29, 2, 5]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [31, 21, 13].

## Portuguese

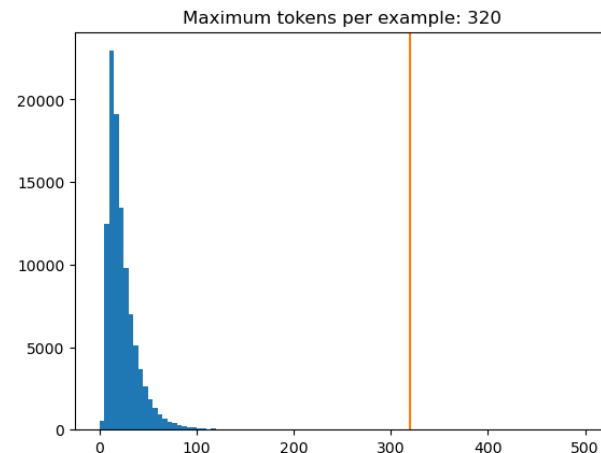*mas eles não tinham a curiosidade de me testar .*
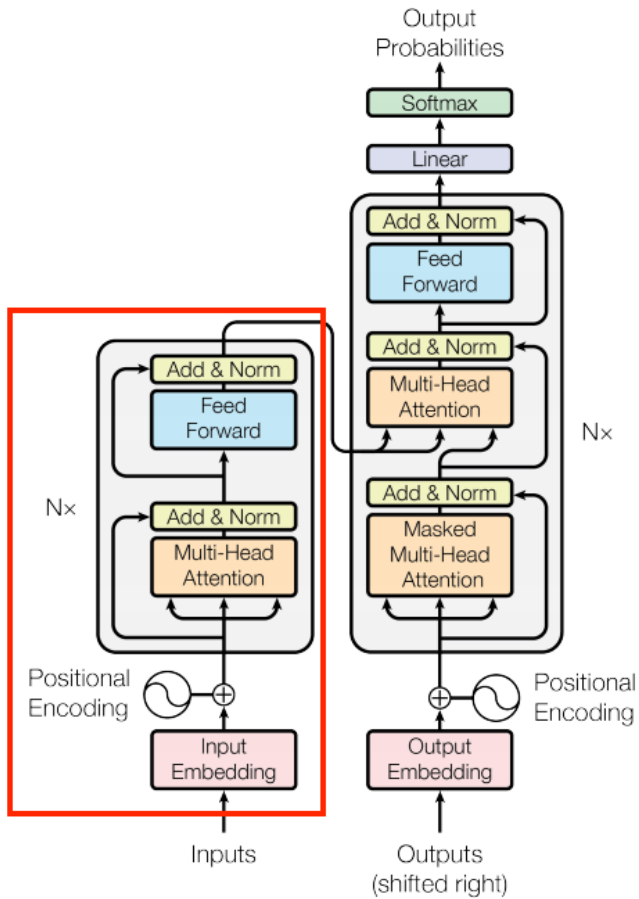
## English

*but they did n't test for curiosity .*

```
[2, 87, 83, 149, 50, 9, 56, 664, 85, 2512, 15, 3]
[b'[START]', b'but', b'they', b'did', b'n', b"'", b't', b'test', b'for', b'curiosity', b'.', b'[END]']
```
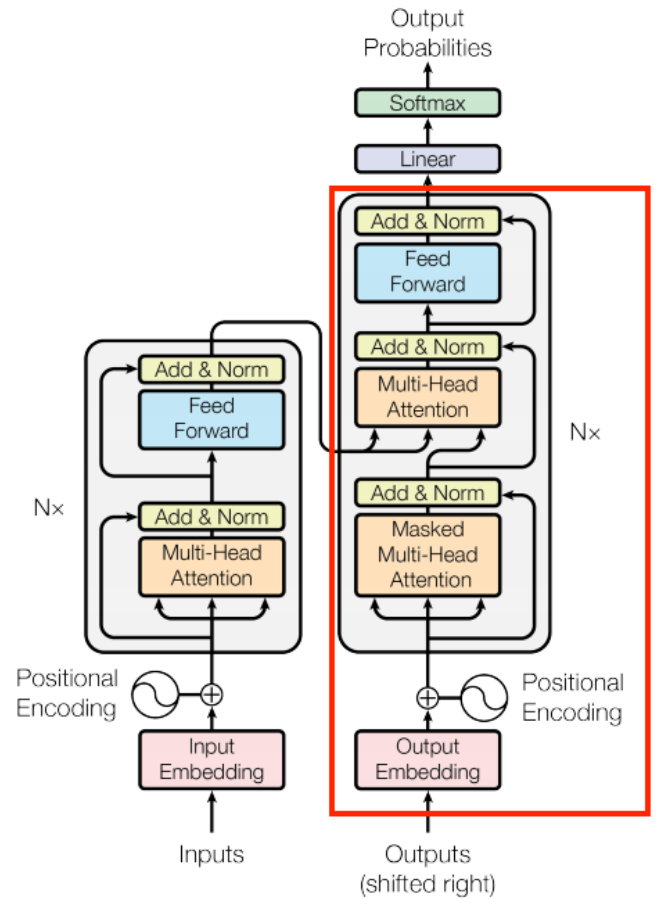
> Subwords: the word 'searchability' is decomposed into 'search' and '##ability', and the
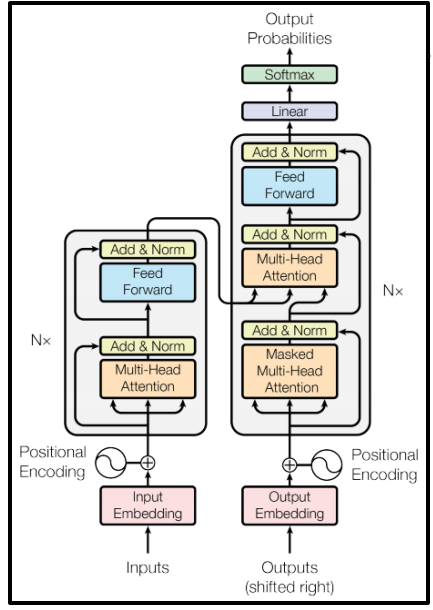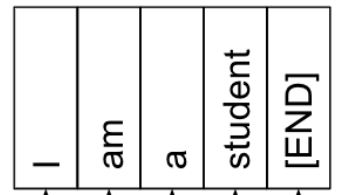word 'serendipity' into 's', '##ere', '##nd', '##ip' and '##ity'.



Maximum tokens per example: 320

# Encoder

# Decoder

* Generates the text one token at a time and feeds the output back to the input – Autoregressive model.

Labels

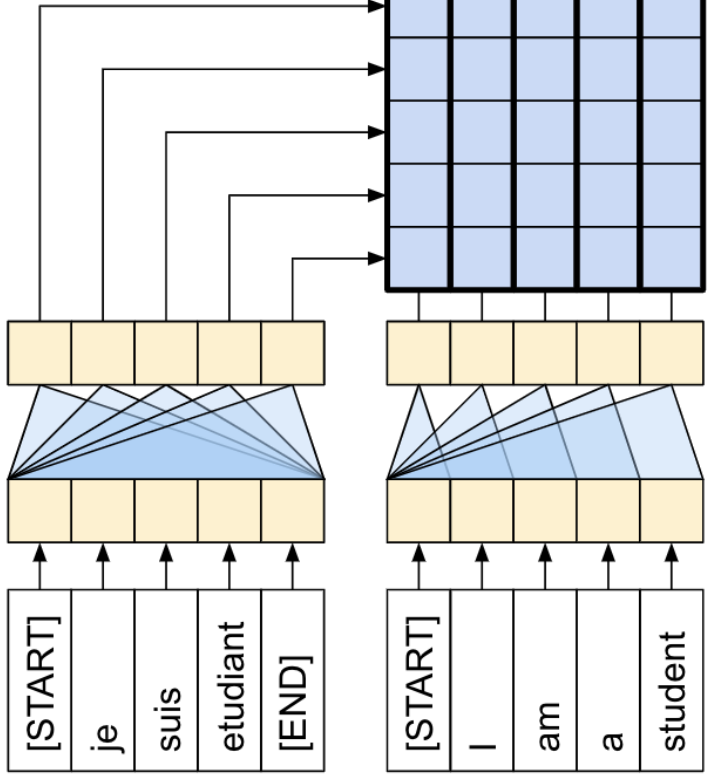* Inputs and Labels are shifted by 1.

**Cross Attention**

* Lets the decoder access the information extracted by the encoder.

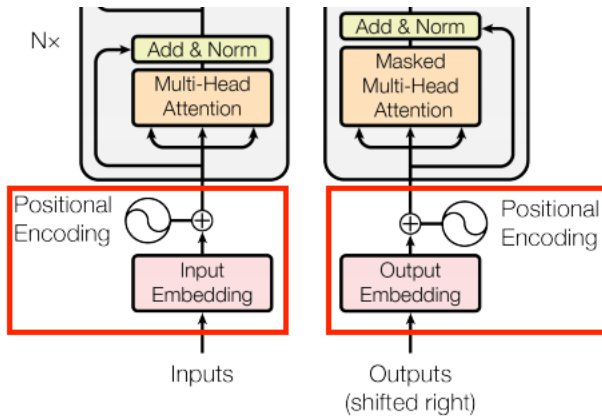* It computes a vector from the entire context sequence, and adds that to the decoder's output.

**Global Self-Attention**

* Responsible for processing the context sequence, and propagating information along its length.

**Causal Self-Attention**

* Makes sure output for each sequence element depends on the previous sequence elements.
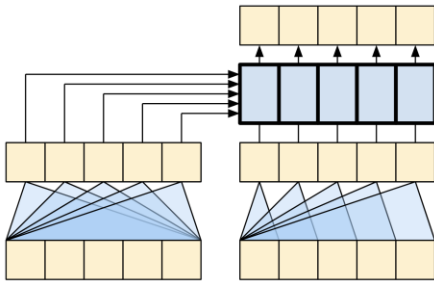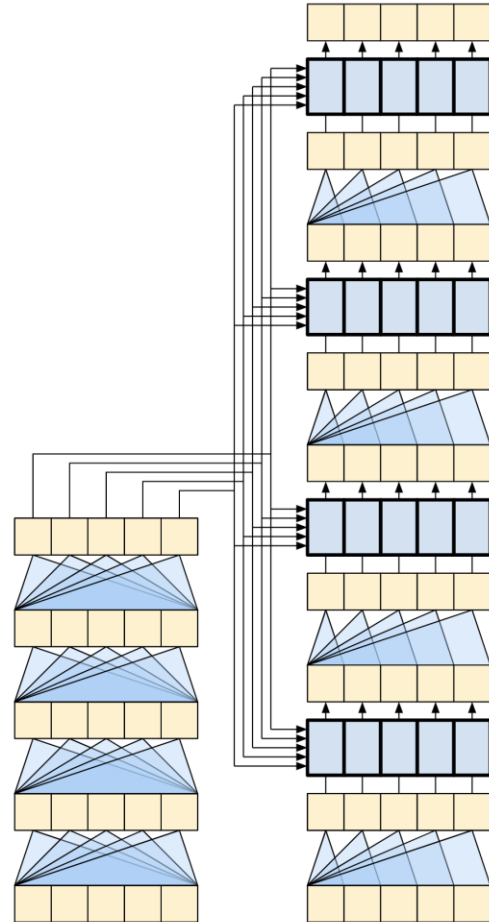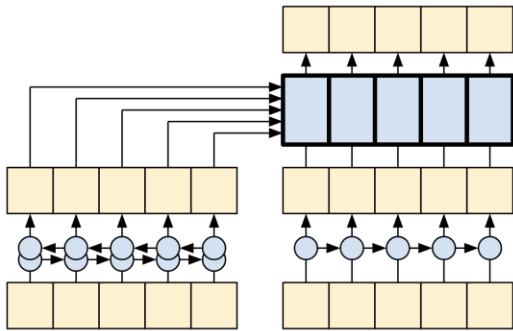
Inputs

University of Pittsburgh

## Positional Encoding



* A stack of sines and cosines that vibrate at different frequencies depending on their location along the depth of the embedding vector.

* The attention layers see their input as a set of vectors, with no order.

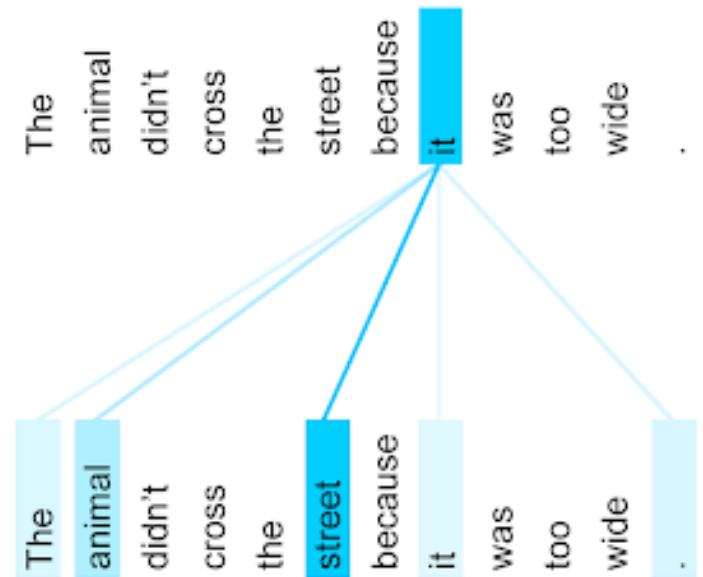* So, we add a positional encoding to the embeddings to force near-by elements to have similar positional encodings.

4-layer Transformer

1-layer Transformer

RNN+Attention

University of
Pittsburgh

The encoder self-attention distribution for the word "it" from the 5th to the 6th layer of a Transformer trained on English-to-French translation (one of eight attention heads).
Source: Google AI Blog.

# Code

# Step 1: Load and Preprocess Data

```python
VOCAB_SIZE = 1000
encoder = tf.keras.layers.TextVectorization(
    max_tokens=VOCAB_SIZE)
encoder.adapt(train_dataset.map(lambda text, label: text))
vocab = np.array(encoder.get_vocabulary())
vocab[:20]
```

```
array(['', '[UNK]', 'the', 'and', 'a', 'of', 'to', 'is', 'in', 'it', 'i',
       'this', 'that', 'br', 'was', 'as', 'for', 'with', 'movie', 'but'],
      dtype='<U14')
```

```python
encoded_example = encoder(example)[:3].numpy()
encoded_example
```

```
array([[ 10, 540,   4, ...,   0,   0,   0],
       [ 10, 120,  11, ...,   0,   0,   0],
       [414,   9,  14, ...,   0,   0,   0]])
```
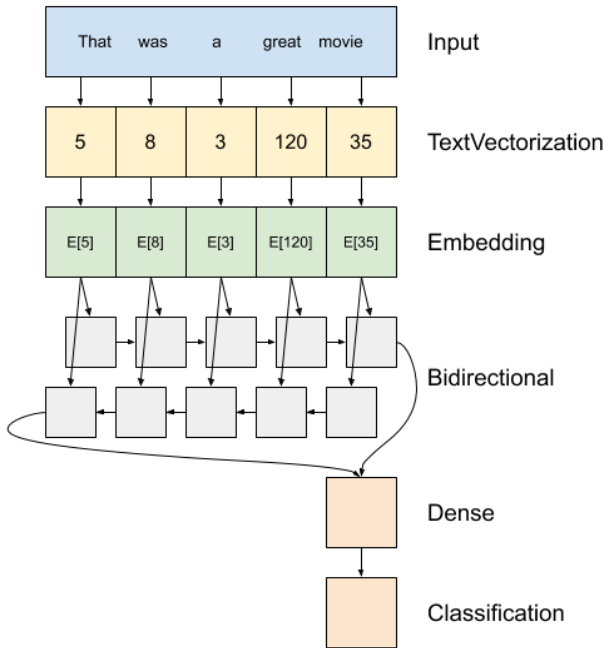
# Step 2: Visualize the dataset

```
[33]   1   num_words = 15
       2   words_in_the_sentence = str(example[n].numpy()).split(' ')[:num_words]
       3   encodeded_id_of_the_words = encoded_example[n][:num_words]
       4
       5   print("Encoding\tWord")
       6   for word, encoded_id in zip(words_in_the_sentence, encodeded_id_of_the_words):
       7     print(encoded_id, "\t\t", word)
```

```
Encoding        Word
10               b'I
86               first
1                encountered
11               this
120              show
51               when
10               I
14               was
1                staying
8                in
1                Japan
16               for
1                six
1                months
226              last
```

1000 VOCAB SIZE

Word Count or Bag of Words

University of
Pittsburgh

# Step 3: Design the NLP Model



```python
model = tf.keras.Sequential([
    encoder,
    tf.keras.layers.Embedding(
        input_dim=len(encoder.get_vocabulary()),
        output_dim=64,
        # Use masking to handle the variable sequence lengths
        mask_zero=True),
    tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(64)),
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(1)
])
```

The output of the Bidirectional LSTM is passed to a Dense layer with 64 nodes, and then further passed to the output layer for final binary classification.
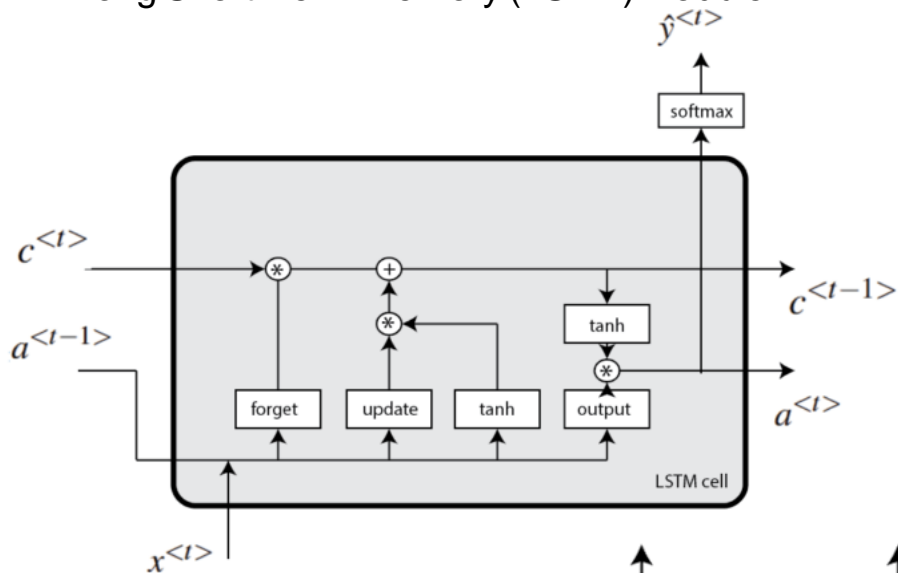
2nd word is processed based on the embedding at 2nd location in the sentence as well as the output of the first word.

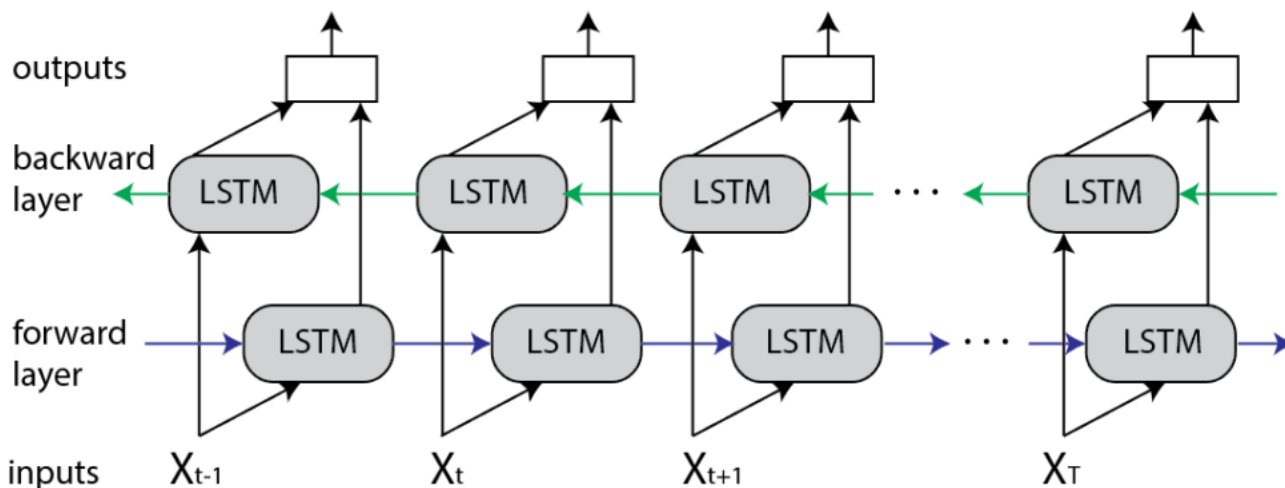3rd word is processed based on the embedding at the 3rd loc in the sentence as well as the output of the second word.

**Word2Vec** – Package by Google to create Embeddings.

University of Pittsburgh

# Long Short-Term Memory (LSTM) Module



## Bidirectional LSTM



University of Pittsburgh

# Step 4: Train the NLP Model

```python
model.compile(loss=tf.keras.losses.BinaryCrossentropy(from_logits=True),
              optimizer=tf.keras.optimizers.Adam(1e-4),
              metrics=['accuracy'])
```

Once we have defined the model, now we can compile the model with the loss and optimizer functions just like we did for the DNN and CNN examples last week. We can then fit the model on the train dataset to train the embedding layer, RNN, and dense layers. Note that the RNN layer has multiple layers inside which enables the temporal or sequential nature of learning. The overall parameters of the model is thus dependent on the embedding size, number and size of RNN layers, and the number and size of dense layers.

```python
history = model.fit(train_dataset, epochs=5)
```

# Step 5: Evaluate the Trained Model

```
test_loss, test_acc = model.evaluate(test_dataset)

print('Test Loss:', test_loss)
print('Test Accuracy:', test_acc)
```

We can run evaluate method on the model to find the test loss and accuracy.

Now, given a new input, we can understand if a movie review is positive or negative.

```
Question:  This is a fantastic movie.
Predicted label:  Positive
Question:  This is a bad movie.
Predicted label:  Negative
Question:  This movie was so bad that it was good.
Predicted label:  Negative
Question:  I will never say yes to watching this movie.
Predicted label:  Negative
Question:  Skip this movie.
Predicted label:  Negative
Question:  Don't waste your time.
Predicted label:  Negative
```

We can now experiment by adding multiple RNN layers to the network and trying out different types of RNN layers.

University of Pittsburgh