

# EEP 596: LLMs: From Transformers to GPT || Lecture 6

Dr. Karthik Mohan

Univ. of Washington, Seattle

January 21, 2026

# Logistics

- Assignment 2 due Friday
- How's Kaggle going?
- Mini Project 1 assigned Friday, due 2 weeks later
- MP1 focused on embeddings, SBERT, semantic search, etc
- Paper Presentations begining today (last 15 mins of class)
- You will be graded on clarity of presentation and your responses on Q&A
- Anything else?

# Last lecture

- Recap of Embeddings and Cosine Similarity
- Glove Embeddings
- Sentence Embeddings with Glove
- Semantic Search Demo

# Today's Lecture

- Word2Vec/Glove Embeddings
- Product2Vec
- Embedding Theory
- Transformers
- BERT
- Sentence Transformers and Sentence BERT

# Deep Learning References

## Deep Learning

Great reference for the theory and fundamentals of deep learning: Book by Goodfellow and Bengio et al [Bengio et al](#)

## [Deep Learning History](#)

## Embeddings

[SBERT and its usefulness](#) [SBert Details](#)

# Word2Vec

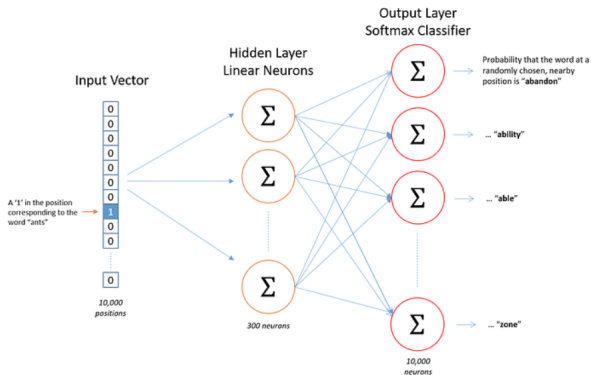
## Skip Gram Model

Is based on the skip-gram model! How is training done? It's semi-supervised!!

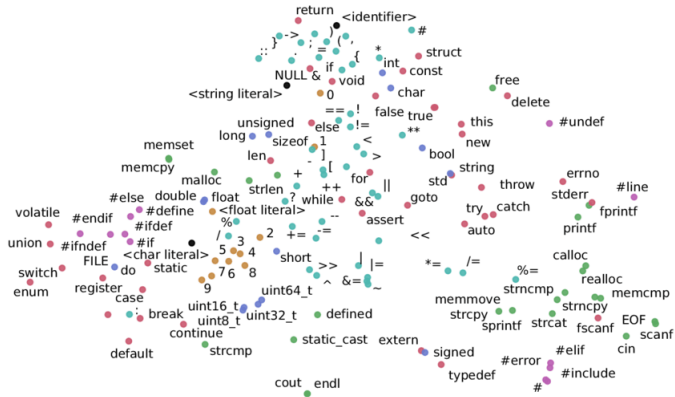
Source Text	Training Samples					
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)		
The	quick	brown				
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	The	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)	
The	quick	brown	fox			
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	The	quick	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The	quick	brown	fox	jumps		
The <table><tr><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
quick	brown	fox	jumps	over		

# Word2Vec

## Architecture



## Word2Vec representation





# ICE #1


What do the embedding dimensions of word2vec represent?

- ① Fixed words decided by word2vec
- ② Topics that are common among the words
- ③ Parts of speech of the words (nouns, adjectives, etc)
- ④ Book titles that these words came from

# Product2Vec

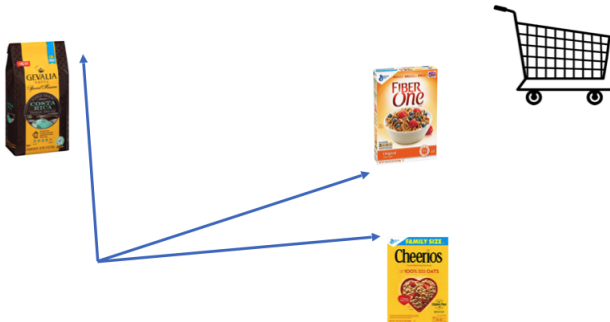


Represent products in product space with a large matrix of embedding coordinate vectors " $L$ "

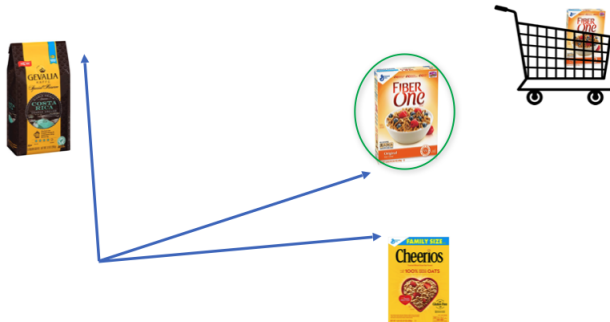

$$L = \begin{pmatrix} 1.5 & 1.9 & 1.8 & 1.4 & \cdots & 0.4 \\ 0.6 & 0.1 & 1.0 & 1.6 & \cdots & 1.9 \\ 0.6 & 1.6 & 1.6 & 1.6 & \cdots & 1.8 \\ 0.6 & 1.0 & 0.1 & 1.6 & \cdots & 0.6 \\ 0.8 & 1.4 & 1.9 & 0.8 & \cdots & 0.7 \end{pmatrix}$$

We obtain these embedding vectors from the Product2Vec service [London et al, 2017]

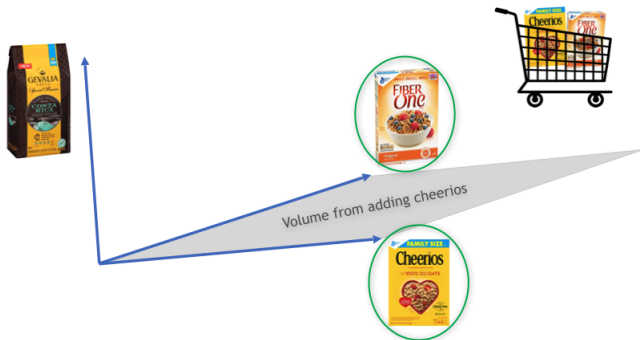
# Product2Vec application



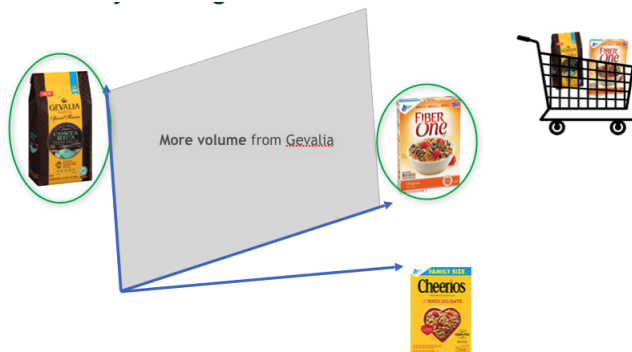
# Product2Vec application



# Product2Vec application



# Product2Vec application



# Breakout 1: Discuss your favorite X2Vec!

## X2Vec

In your group - Discuss an application that requires machine learning. Be specific about it - Example, data, features, the type of problem (classification, clustering, etc). Can you see how X2Vec would benefit your application. What would be your X in this case? How would you learn X2vec for your application? And how would you use it?

Let's list out some X's in X2Vec!




# ICE #2

## View Similarity or Purchase Similarity

Consider a company that sells products online. As we know, embedding representations for words or products in this case are learned from data. The question is which data to use? These are referred to as signals sometimes. So the question is, does **view similarity** of products represent a better signal for learning embeddings or **purchase similarity**? Remember: Good embeddings embed similar products close to each other dis-similar products away from each other.

# ICE #2: View or Purchase Similarity?




SPORTS RESEARCH  
**BIOTIN**  
Extra Strength<sup>®</sup>  
5000MG PER SERVING  
SUPPORTS HEALTHY SKIN & HAIR<sup>®</sup>  
VEGAN | NON-GMO | WITH COCONUT OIL  
120 Veggie Softgels | Dietary Supplement

Roll over image to zoom in


- PREMIUM, POTENT FORMULAS:** Take your pick between our High Potency Biotin 2,500 mcg, Extra Strength Biotin 5,000 mcg, and Max Strength Biotin 10,000 mcg—all gluten-free, soy-free, and made with cold-pressed 100% organic coconut oil.
- THE SPORTS RESEARCH DIFFERENCE:** Founded in Southern California in 1980, Sports Research is a family-owned business born from a passion for fitness and wellness. Our goal is to embrace the sport of life through research-backed products created for every body—inside and out.

[Report an issue with this product or seller](#)

**Consider a similar item**






Amazon Elements Vegan Biotin 5000 mcg - Hair, Skin, Nails, 130 Capsules (4 month supply) (Packaging may vary)  
130 Count (Pack of 1)  
★★★★★ (32081)  
\$5.46 (\$0.07/Count) [prime](#)  
[Veg Climate Pledge Friendly](#)



Sports Research Vitamin K2 as MK-7 100mcg with Coconut MCT Oil - 120 Veggie Softgels (4 Month Supply)..  
★★★★★ 11  
\$33.95 [prime](#)

## Frequently bought together


**Total price: \$66.16**  
Add all 3 to Cart

**This item:** Sports Research Vegan Biotin 5000mcg with Organic Coconut Oil - Extra Strength...  
117<sup>th</sup> (\$0.15/Count) [prime](#)


Sports Research Triple Strength Omega-3 Fish Oil - Burgeless Fish Oil Supplement w/EPA & DHA...  
124<sup>th</sup> (\$0.23/Count) [prime](#)

Sports Research Vitamin D3 K2 with 5000iu of Vegan D3 & 100mcg of Vitamin K2 as MK7...  
125<sup>th</sup> (\$0.40/Count) [prime](#)


## Products related to this item




Biotin | Collagen | Keratin | Hyaluronic Acid - Hair Growth Support Pills, 25000 mcg...  
★★★★★ 5,705  
**Amazon's Choice** [prime](#)  
**At Best Seller**  
\$27.89 (\$0.40/Count)  
[prime](#)  
[Veg Climate Pledge Friendly](#)




PURE RESEARCH Liquid Biotin & Collagen Hair Growth Drops 60,000mcg - Biotin...  
★★★★★ 28,024  
\$21.89 (\$0.36/FL OZ)  
[prime](#)




Sports Research Vitamin D3 K2 with 5000iu of Vegan D3 & 100mcg of Vitamin K2 as MK7...  
★★★★★ 35,526  
\$23.95 (\$0.43/Count)  
[prime](#)




Gorvus Biotin & Collagen Hair Growth Support Drops - Hair Supplement - Healthy Skin...  
★★★★★ 3,445  
\$19.99 (\$10.00/FL OZ)  
[prime](#)



Biotin & Collagen | Biotin 5000mcg, Keratin & Collagen - Hair...  
★★★★★ 957  
**Amazon's Choice** [prime](#)  
\$22.97 (\$0.38/Count)  
[prime](#)



Superhair Hair Vitamins Extra Strength Biotin 6000mcg, Vitamin C, E, Coconut Oil, 2...  
★★★★★ 47,314  
\$29.99 (\$0.30/Count)  
[prime](#)



Biotin Capsules with Collagen and Keratin - 25000MCGS Per Serving - Biotin Vitamins...  
★★★★★ 896  
\$19.89 (\$0.33/Count)  
[prime](#)

Page 1 of 58

# Generating Sentence Embeddings from Glove

**Averaging embeddings of words:** If we have a word embedding, how do we generate the sentence embedding?

# Generating Sentence Embeddings from Glove

**Averaging embeddings of words:** If we have a word embedding, how do we generate the sentence embedding?

**Simple Solution: Just average the word embeddings**

# How do we improve Sentence Embeddings?

## Sentence Embeddings

As you are probably observing in your Mini-Project 1 assignment - Averaging word embeddings doesn't "perform" as well. So we need sentence embeddings that do better than just averaging word embeddings

- Perhaps, capture the sequence of information flow in a sentence.

# How do we improve Sentence Embeddings?

## Sentence Embeddings

As you are probably observing in your Mini-Project 1 assignment - Averaging word embeddings doesn't "perform" as well. So we need sentence embeddings that do better than just averaging word embeddings - Perhaps, capture the sequence of information flow in a sentence.

### Example 1

Sentence 1: "Me loves my friend"

Sentence 2: "My friend loves me"

Should they have the exact same sentence embeddings?

# How do we improve Sentence Embeddings?

## Sentence Embeddings

As you are probably observing in your Mini-Project 1 assignment - Averaging word embeddings doesn't "perform" as well. So we need sentence embeddings that do better than just averaging word embeddings - Perhaps, capture the sequence of information flow in a sentence.

### Example 1

Sentence 1: "Me loves my friend"

Sentence 2: "My friend loves me"

Should they have the exact same sentence embeddings?

### Example 2

Sentence 1: "I like chocolate milk"

Sentence 2: "I like milk chocolate"

Should they have the same sentence embeddings?

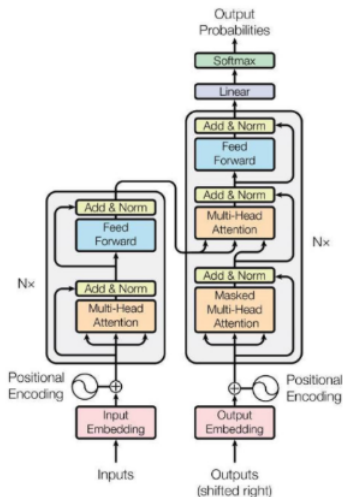
# Next Topic: Transformers, BERT and connections to Embeddings

## Capturing Sequence of information

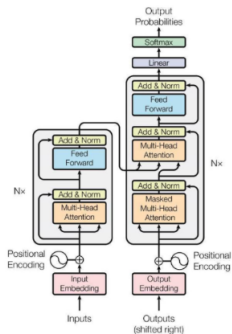
As we discussed in the history of Deep Learning - RNNs and LSTMs are DL archs that are able to capture sequence information in a sentence to some extent (the **chocolate milk** vs **milk chocolate** example). On the other hand, they weren't robust to larger context or multiple sentences and could only operate with smaller sentence lengths. This is where the advent of Transformers was a breakthrough for ML/DL and AI in general - They could do much better in capturing context, sequential information, supported multiple sentences and paragraphs, etc.



# Transformers - Encoder and Decoder Architecture

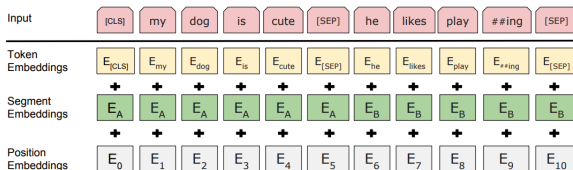


# Encoder and Encoder Embeddings



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	#ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{\#ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

# Understanding Encoder/BERT at high-level



# BERT - Bi-directional Encoders from Transformers

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	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 <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

# Parsing the Embeddings and Encoder/Decoder Terminology

- 1 **Encoder:** The **architecture component** of the transformer that transforms inputs through a series of Neural layers into a vector (embedding). This vector can then be useful for downstream tasks: Emotion detection, Classification, etc

# Parsing the Embeddings and Encoder/Decoder Terminology

- ① **Encoder:** The **architecture component** of the transformer that transforms inputs through a series of Neural layers into a vector (embedding). This vector can then be useful for downstream tasks: Emotion detection, Classification, etc
- ② **Decoder:** The **architecture component** of the transformer that transforms inputs one at a time to generate words in sequence. Example: You ask a question of ChatGPT and it starts generating words, one at a time - Just like it were typing. That's decoder in action.

# Parsing the Embeddings and Encoder/Decoder Terminology

- 1 **Encoder:** The **architecture component** of the transformer that transforms inputs through a series of Neural layers into a vector (embedding). This vector can then be useful for downstream tasks: Emotion detection, Classification, etc
- 2 **Decoder:** The **architecture component** of the transformer that transforms inputs one at a time to generate words in sequence. Example: You ask a question of ChatGPT and it starts generating words, one at a time - Just like it were typing. That's decoder in action.
- 3 **Input Embedding:** How an input gets embedded as a vector. What would be the starting point to embed "chocolate milk" as a vector?

# Parsing the Embeddings and Encoder/Decoder Terminology

- ❶ **Encoder:** The **architecture component** of the transformer that transforms inputs through a series of Neural layers into a vector (embedding). This vector can then be useful for downstream tasks: Emotion detection, Classification, etc
- ❷ **Decoder:** The **architecture component** of the transformer that transforms inputs one at a time to generate words in sequence. Example: You ask a question of ChatGPT and it starts generating words, one at a time - Just like it were typing. That's decoder in action.
- ❸ **Input Embedding:** How an input gets embedded as a vector. What would be the starting point to embed “chocolate milk” as a vector?
- ❹ **Token Embedding:** This refers to the individual token embeddings or word embeddings (or sub-word embeddings)



# Parsing the Embeddings and Encoder/Decoder Terminology

- ① **Segment Embedding:** This refers to a generic embedding that says this was segment 1 or segment 2 of the input

# Parsing the Embeddings and Encoder/Decoder Terminology

- ① **Segment Embedding:** This refers to a generic embedding that says this was segment 1 or segment 2 of the input
- ② **Position Embedding:** This adds in the position information into the embedding. Did “chocolate” come in at the beginning of the sentence or middle or the end?

# Parsing the Embeddings and Encoder/Decoder Terminology

- ① **Segment Embedding:** This refers to a generic embedding that says this was segment 1 or segment 2 of the input
- ② **Position Embedding:** This adds in the position information into the embedding. Did “chocolate” come in at the beginning of the sentence or middle or the end?
- ③ **BERT Embedding:** This is the embedding or the vector used after having gone through the Encoder Architecture

# Parsing the Embeddings and Encoder/Decoder Terminology

- 1 **Segment Embedding:** This refers to a generic embedding that says this was segment 1 or segment 2 of the input
- 2 **Position Embedding:** This adds in the position information into the embedding. Did “chocolate” come in at the beginning of the sentence or middle or the end?
- 3 **BERT Embedding:** This is the embedding or the vector used after having gone through the Encoder Architecture
- 4 **Sentence BERT (sBERT) Embedding:** This is the embedding that you are using in Mini-Project 1, an **encoding** into a vector that's optimized for sentence similarity! (More on this in a bit)

# Path of a sentence through transformer layers

---

## 1 Raw Input Sentence

Example input:

"I love machine learning"

At this stage, it's just a string—no math yet.

---

# Path of a sentence through transformer layers

## 2 Text Normalization (Basic Preprocessing)

Before tokenization, BERT applies **basic text cleaning**:

- Lowercasing (*for uncased models*)
- Unicode normalization
- Whitespace cleanup

Result:

```
arduino
```

```
"i love machine learning"
```

# Path of a sentence through transformer layers

## 3 WordPiece Tokenization

BERT uses **WordPiece tokenization**, which breaks words into subword units.

### Steps:

1. Split on whitespace
2. Decompose unknown or rare words into subwords
3. Prefix continuation pieces with `##`

Example:

CSS

```
["i", "love", "machine", "learning"]
```

If a word were rare:

arduino

```
"unhappiness" → ["un", "##happi", "##ness"]
```

# Path of a sentence through transformer layers

## 4 Add Special Tokens

BERT **requires special tokens** for context and classification.

Token	Purpose
[CLS]	Sentence-level representation
[SEP]	Sentence separator / end marker

Final token sequence:

CSS

[CLS] i love machine learning [SEP]



# Path of a sentence through transformer layers

## 5 Convert Tokens → Token IDs

Each token is mapped to an integer via BERT's vocabulary.

Example (IDs are illustrative):

```
yaml
```

```
[101, 1045, 2293, 3698, 4083, 102]
```

Now the sentence is numerical.

---

# Path of a sentence through transformer layers

## 6 Create Input Embeddings (Critical Step)

Each token gets 3 embeddings, which are summed element-wise:

### a) Token Embeddings

Learned vector for each word/subword.

### b) Position Embeddings

Encode token order:

```
csharp
```

```
[CLS]=pos0, i=pos1, love=pos2, ...
```

### c) Segment (Token Type) Embeddings

Used to distinguish sentence A vs B.

- Single sentence → all zeros

Final embedding per token:

```
ini
```

```
Embedding = Token + Position + Segment
```

# Path of a sentence through transformer layers

## 7 Input Matrix to Transformer

At this point, your sentence is a matrix:

```
rust  
  
(sequence_length × hidden_size)  
= (6 × 768) ← for BERT-Base
```

Each row corresponds to one token.

---

# Path of a sentence through transformer layers

## 8 Pass Through Transformer Encoder Layers

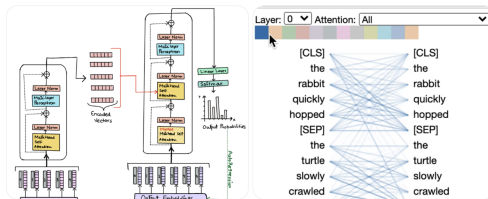
BERT-Base has 12 identical Transformer layers.

Each layer contains:

### ◆ Multi-Head Self-Attention

Every token attends to **every other token** (including itself).

This lets [CLS] gather information from the entire sentence.



# Path of a sentence through transformer layers

---

## ◆ Feed-Forward Network (FFN)

Applies a non-linear transformation **independently** to each token vector.

---

## ◆ Residual Connections + LayerNorm

Helps with:

- Gradient flow
  - Stable training
-

# Path of a sentence through transformer layers

## 10 Extract the [CLS] Token

The first vector corresponds to [CLS] :

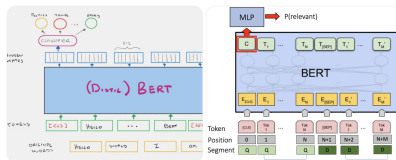
```
ini  
  
CLS_embedding = last_hidden_state[:, 0, :]
```

Shape:

```
SCSS  
  
(768,)
```

This is the **sentence embedding** used for:

- Classification
- Regression
- Sentence-level tasks

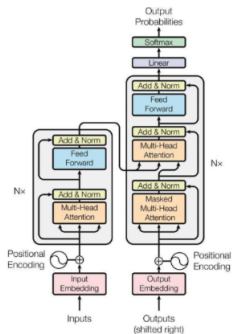


# Llama3 High-level (Decoder Model)

## Llama 3 70B Model

- Parameters: Approximately 70 billion.
- Vocabulary Size: 128,000 tokens.
- Total tokens seen: 15 trillion
- Context Window: Supports sequences up to 8,192 tokens.
- Architecture:
- Layers: 80 transformer blocks.
- Attention Heads: 64 heads per layer.
- Hidden Dimension: 10,649.

# Encoder and Encoder Embeddings



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	#ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{\#ing}$	$E_{[SEP]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$



# In-class Coding exercise - Carries Bonus points for grade

## Compute Sentence Embedding (15 mins)

**Instructions:** Write a python module to compute sentence embeddings given words and their embeddings.

Encapsulate your code in a class structure. With one method for each of the sentence embedding computations asked below. Each method takes in as input a sentence and outputs a sentence embedding. Also have class variables to store vocab embeddings, filler words and attention weights (part 3).

**Submit:** Your code on canvas link that will be opened up shortly for bonus points that can be used towards your grade.

**Vocabulary:** Assume the vocabulary consists of 5 words:  
*I, lot, love, chocolate, milk.*

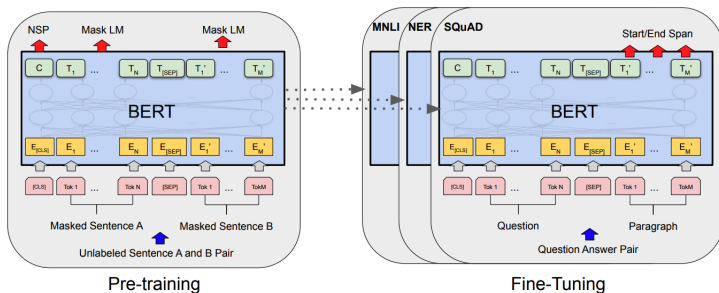
**Sentence:** Consider the following sentence: I love chocolate milk as well! Assume the embedding(*i*th word) =  $[i - 1, i + 1]$  where *i* is 0-indexed.

# In-class Coding exercise - Carries Bonus points for grade

## Compute Sentence Embedding (15 mins)

- 1. Simple Average Embedding:** Compute the sentence embedding by simple averaging. For any word not in vocabulary, use the embedding as  $[-1, -1]$ . What's the resulting sentence embedding with this approach?
- 2. Skip filler words embedding:** Let filler words be *[I, and, as, for, it, or, maybe]*. If you encounter any filler words in a sentence, skip the filler word for the embedding computation. What's the sentence embedding with this approach?
- 3. Learned Sentence embedding:** Assume you learned a set of weights that in word embeddings as input and give out sentence embedding as output. This is done by taking a weighted average of the embeddings. For the above sentence, let the weights learned could be from a self-attention layer. Assume the weights for words in the vocabulary are:  $[0.5, 1, 0.7, 0.9, 0.3, 0.2, 0.1]$ . What would be the resulting embedding?

# BERT - Bi-directional Encoders from Transformers



# BERT pre-training

## Two Tasks

- ① **Masked LM Model:** Mask a word in the middle of a sentence and have BERT predict the masked word
- ② **Next-sentence prediction:** Predict the next sentence - Use both positive and negative labels. How are these generated?

# BERT pre-training

## Two Tasks

- 1 **Masked LM Model:** Mask a word in the middle of a sentence and have BERT predict the masked word
- 2 **Next-sentence prediction:** Predict the next sentence - Use both positive and negative labels. How are these generated?

## ICE: Supervised or Un-supervised?

- 1 Are the above two tasks supervised or un-supervised?

# BERT pre-training

## Two Tasks

- 1 **Masked LM Model:** Mask a word in the middle of a sentence and have BERT predict the masked word
- 2 **Next-sentence prediction:** Predict the next sentence - Use both positive and negative labels. How are these generated?

## ICE: Supervised or Un-supervised?

- 1 Are the above two tasks supervised or un-supervised?

## Data set!

English Wikipedia and book corpus documents!

# Loss Function for Masked Language Model (MLM)

Loss Function for MLM mimicks which type of classic ML model?

---

# Loss Function for Masked Language Model (MLM)

Loss Function for MLM mimicks which type of classic ML model?

Cross-Entropy

$$L(p, \hat{p}) = - \sum_i [p_i \log(\hat{p}_i) + (1 - p_i) \log(1 - \hat{p}_i)]$$



# Loss Function for Masked Language Model (MLM)

Loss Function for MLM mimicks which type of classic ML model?

Cross-Entropy

$$L(p, \hat{p}) = - \sum_i [p_i \log(\hat{p}_i) + (1 - p_i) \log(1 - \hat{p}_i)]$$

ICE: What is the loss function for Binary Classification?

# BERT - Bi-directional Encoders from Transformers

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	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 <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

# Sentence BERT a.k.a sBERT

Uses Siamese Twins architecture

---

# Sentence BERT a.k.a sBERT

Uses Siamese Twins architecture

Advantages of sBERT

More optimized for Sentence Similarity Search.

# SBERT - Siamese BERT architecture

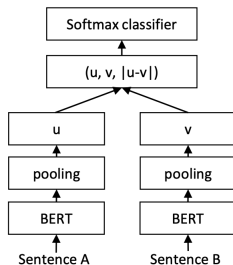


Figure 1: SBERT architecture with classification objective function, e.g., for fine-tuning on SNLI dataset. The two BERT networks have tied weights (siamese network structure).

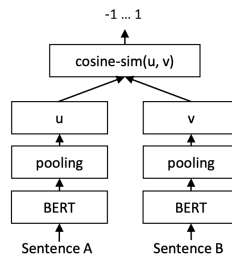


Figure 2: SBERT architecture at inference, for example, to compute similarity scores. This architecture is also used with the regression objective function.

# Breakouts Time #2

## Retrieving Tables with Chat bots — 7 mins

You are building a chat-bot product at your company where queries come in from customers that own data in your company's cloud service. Your chat-bot responds and retrieves the right table or combination of tables (through merge/filter operations) that contains this information or returns back with follow up questions to get more precise information or get back with a "Sorry, I don't have that information" response. How would you go about building a chat-bot like this? What data would you use? What data stores/data bases would be appropriate? What Deep Learning models would you use, would it be supervised or un-supervised learning? What would be your evaluation metric? How would you test if your chat bot is accurate in its responses?