

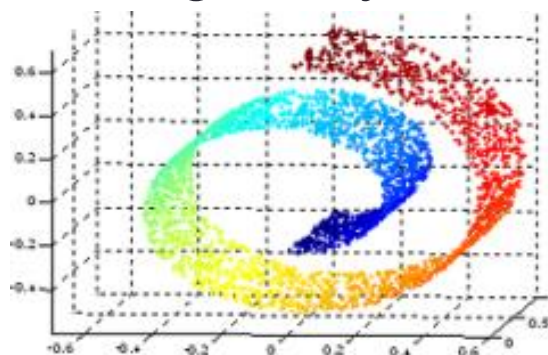
# Text Representation

---

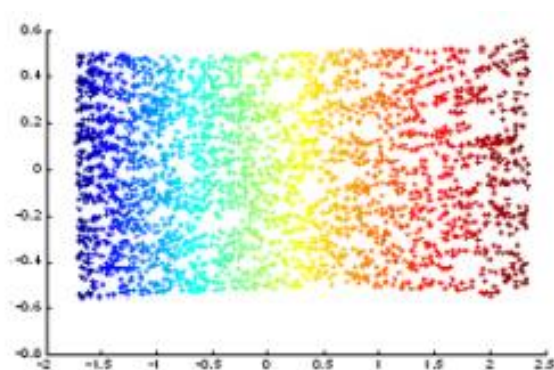
Representation learning

# What is an embedding?

- “Latent space” or “embedding space” refers to a low-dimensional representation of high-dimensional data
  - In neural network, the mapping from original data to the embedding space is often linear.
    - Ex of linear mapping/projection: PCA
- Mapping of these embeddings are one of the key tricks in deep learning today



High dimension



Low dimension

# Embeddings

- Can be trained by supervised or self-supervised techniques

## Self-Supervised Learning = Filling in the Blanks

Y. LeCun

- ▶ Predict any part of the input from any other part.

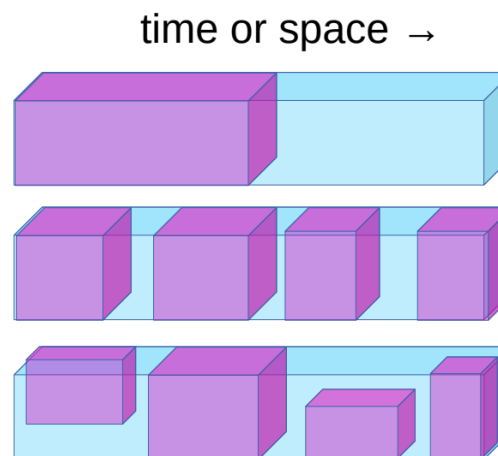
- ▶ Predict the **future** from the **past**.

- ▶ Predict the **masked** from the **visible**.

- ▶ Predict the **any occluded part** from **all available parts**.

- ▶ **Pretend there is a part of the input you don't know and predict that.**

- ▶ **Reconstruction = SSL when any part could be known or unknown**



# Outline

- Contrastive learning
- Sentence embeddings
  - MUSE
  - SimCSE
  - BGE
  - CLIP

# Contrastive learning

- An important technique for self-supervised training is contrastive learning
  - Similar things should have similar embeddings
  - Different things should have different embeddings
- Example: negative sampling loss in word2vec

$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

Context word  
(positive, +1)

Negative samples  
(negative, -1)

# Types of contrastive learning

- Triplet loss
- InfoNCE loss

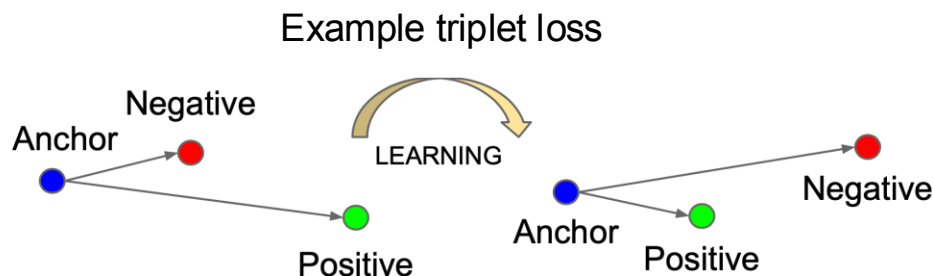
# Triplet loss

- Triplet loss considers an anchor, a positive, and a negative
- Requires mining of hard negative samples

$$\sum_i^N \left[ \|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]_+$$

margin

Take positive only max(0,x)



# Dealing with minibatches

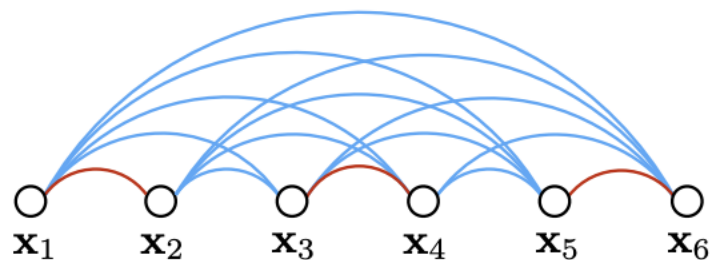
- Since we train in minibatches, most modern losses pair positive and negative samples within a minibatch for more efficient computation
  - Compute all pairwise distance within the minibatch



(a) Contrastive embedding



(b) Triplet embedding



(c) Lifted structured embedding



# NCE (Noise contrastive estimation) loss

- Maximize training data probability while reducing noise probability.
- Learn in a contrastive way to reduce overhead for normalization
  - $\text{Max Log}P(\text{data}) - \text{Log } P(\text{noise or negative samples})$
  - Ex: used to train word embeddings such as W2V, too many classes in the softmax output

# InfoNCE

- Similar to NCE but just for categorical cross entropy (instead of binary cross entropy)

<https://arxiv.org/pdf/1807.03748.pdf>

- Effectively maximize mutual information between **c** and positive **x**

$$L_{InfoNCE} = -E\left[\log \frac{f(x, c)}{\sum_{x'} f(x', c)}\right]$$

$$f(x, c) = \exp(\mathbf{z}^T W c)$$

$\mathbf{z}$  is encoded  $x$

- $f(\ )$  can be any function that describes similarity
- Can be extended to have multiple positive examples in a batch (soft nearest neighbor loss)

<https://arxiv.org/abs/1902.01889>

# Soft nearest neighbor loss

- Multiple positive and negative
- Adds temperature (either hyperparameter, or learned)
  - Weights the gradient size, helps model learn from hard negatives

**Definition.** The *soft nearest neighbor loss* at temperature  $T$ , for a batch of  $b$  samples  $(x, y)$ , is:

$$l_{sn}(x, y, T) = -\frac{1}{b} \sum_{i \in 1..b} \log \left( \frac{\sum_{\substack{j \in 1..b \\ j \neq i \\ y_i = y_j}} e^{-\frac{\|x_i - x_j\|^2}{T}}}{\sum_{\substack{k \in 1..b \\ k \neq i}} e^{-\frac{\|x_i - x_k\|^2}{T}}} \right) \quad (1)$$

# Contrastive summary

- The most common form you will see for contrastive learning is

$$\mathcal{L}^{\text{NT-Xent}} = -\frac{1}{n} \sum_{i,j \in \mathcal{MB}} \log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2n} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

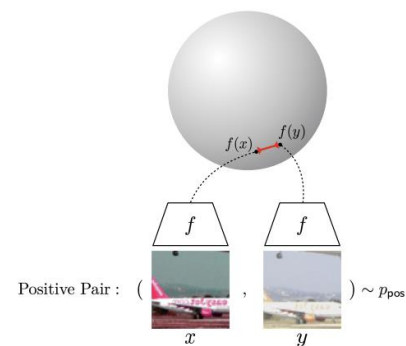
$$\tau \mathcal{L}^{\text{NT-Xent}} = \underbrace{-\frac{1}{n} \sum_{i,j} \text{sim}(\mathbf{z}_i, \mathbf{z}_j)}_{\mathcal{L}_{\text{alignment}}} + \underbrace{\frac{\tau}{n} \sum_i \log \sum_{k=1}^{2n} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}_{\mathcal{L}_{\text{distribution}}}$$

Encourage similar things to align

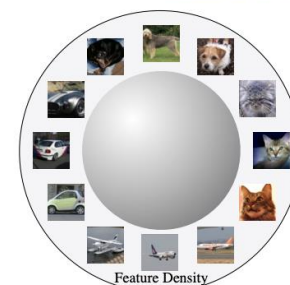
Encourage embeddings to spread uniformly in the hypersphere

- People often refer to this as contrastive loss, InfoNCE loss, normalized temperature scaled CE loss, ...

<https://arxiv.org/abs/2005.10242> <https://arxiv.org/abs/2011.02803> <https://arxiv.org/abs/2002.05709>



**Alignment:** Similar samples have similar features.  
(Figure inspired by Tian et al. (2019).)

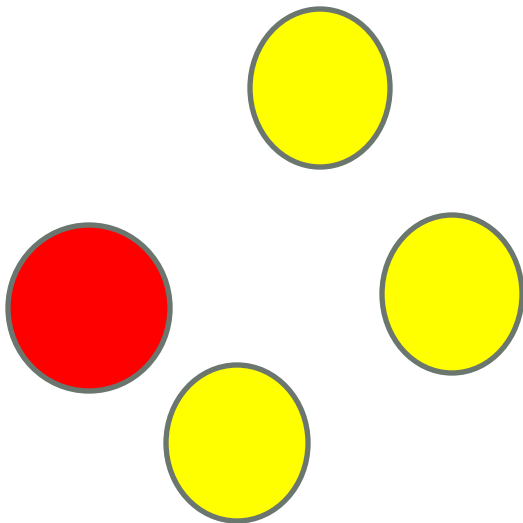


**Uniformity:** Preserve maximal information.

Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.

# Key details to contrastive loss works

- Large batch
- Hard/semi-hard negative mining
- Augmentation on the anchor and positive
- Other tricks includes - adding classification/supervised loss (CE/softmax loss)



# Outline

- Contrastive learning
- Sentence embeddings
  - MUSE
  - SimCSE
  - BGE
  - CLIP

# Sentence representation

- How would we create a sentence embedding?
- Compositionality from words/tokens!
  - Sum, max
  - Recurrence
  - Attention

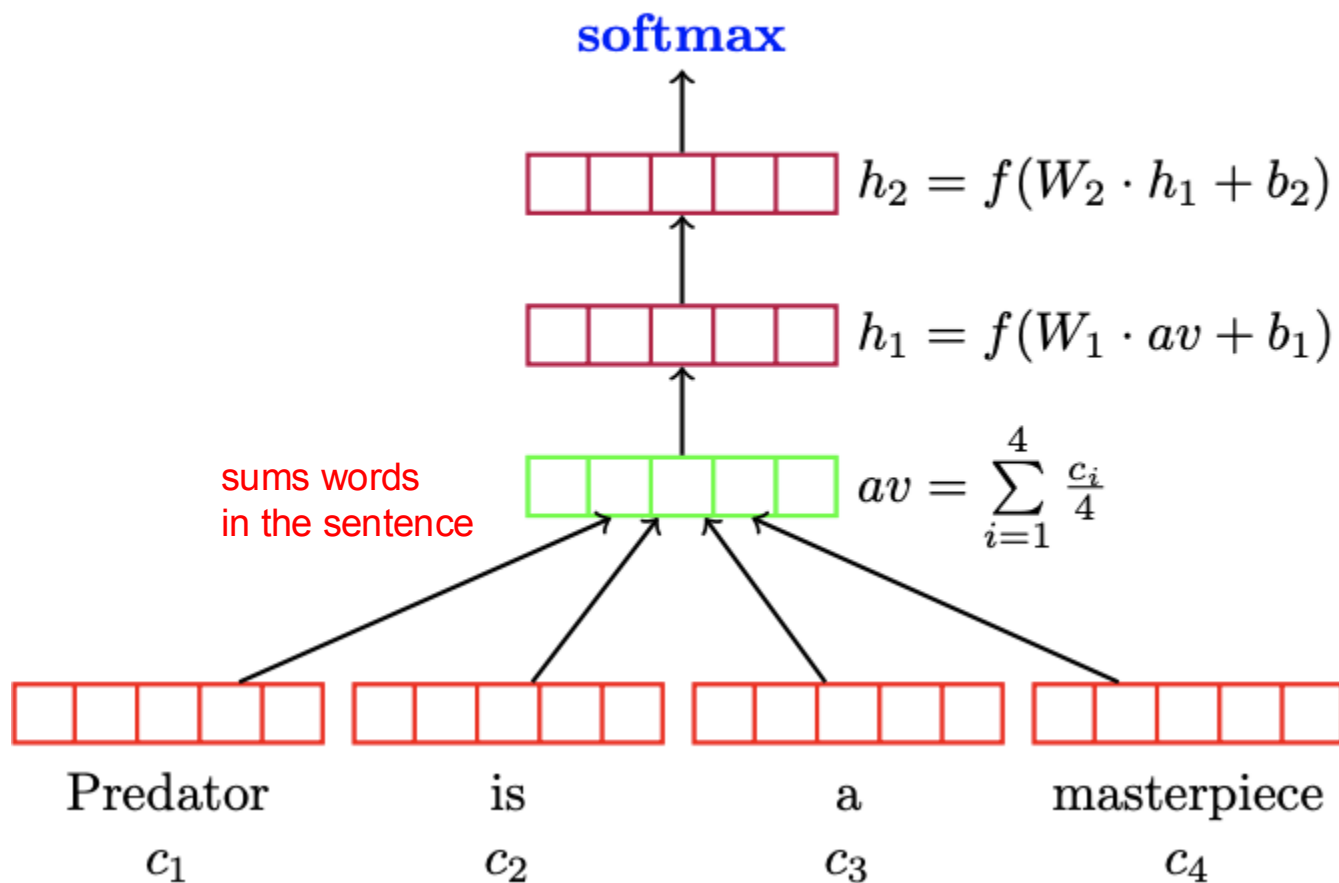
# MUSE

---



# Deep Averaging Networks (DAN)

## DAN



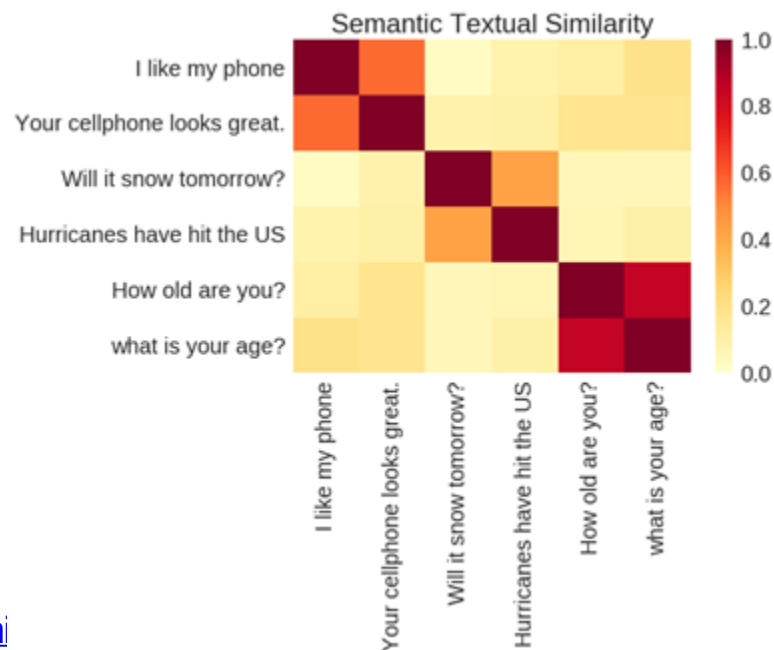
# Universal Sentence Encoder (USE)

A model focusing on sentence representation

Use sentencepiece tokenization

Pre-trained then used anywhere

Based on (1) DAN (lite version) or (2) Transformer



Official implementation with pretrained weights

<https://tfhub.dev/google/collections/universal-sentence-encoder/1>

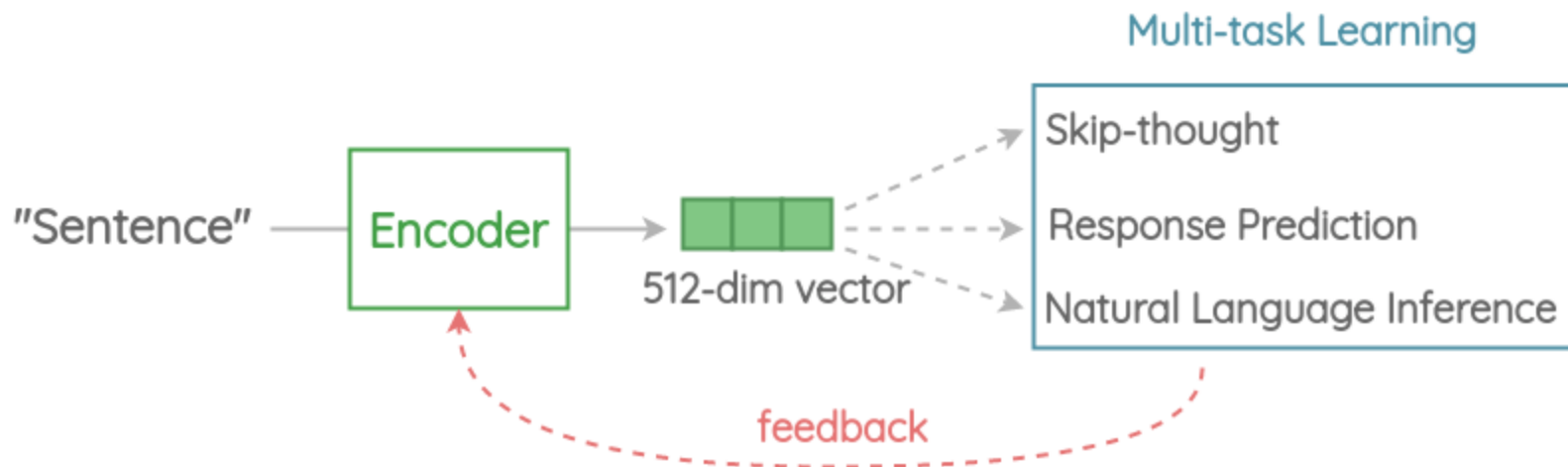
<https://ai.googleblog.com/2018/05/advances-in-semantic-textual-similarity.html>

<https://www.kaggle.com/google/universal-sentence-encoder>

# Pretraining USE

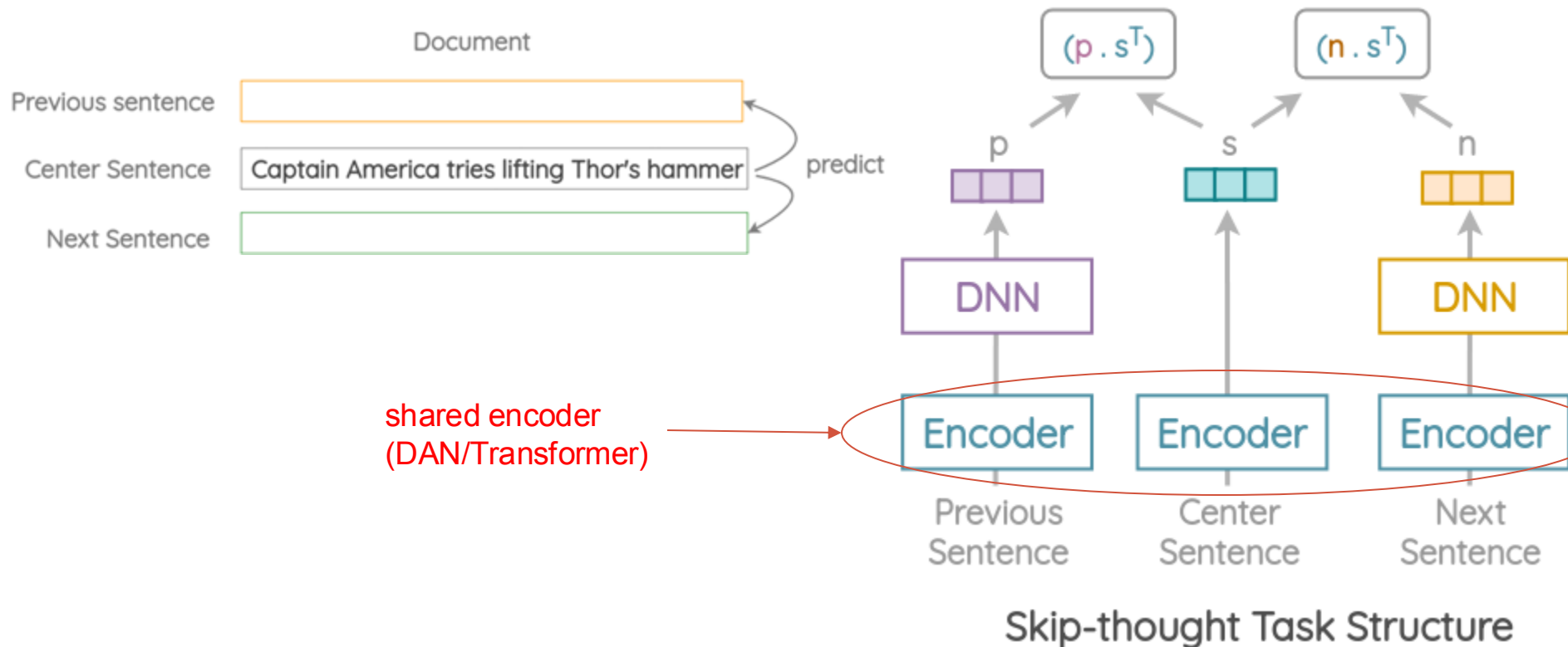
Training done using multi-task

- 1) Skip-thought
- 2) Response prediction
- 3) Natural language inference (NLI)



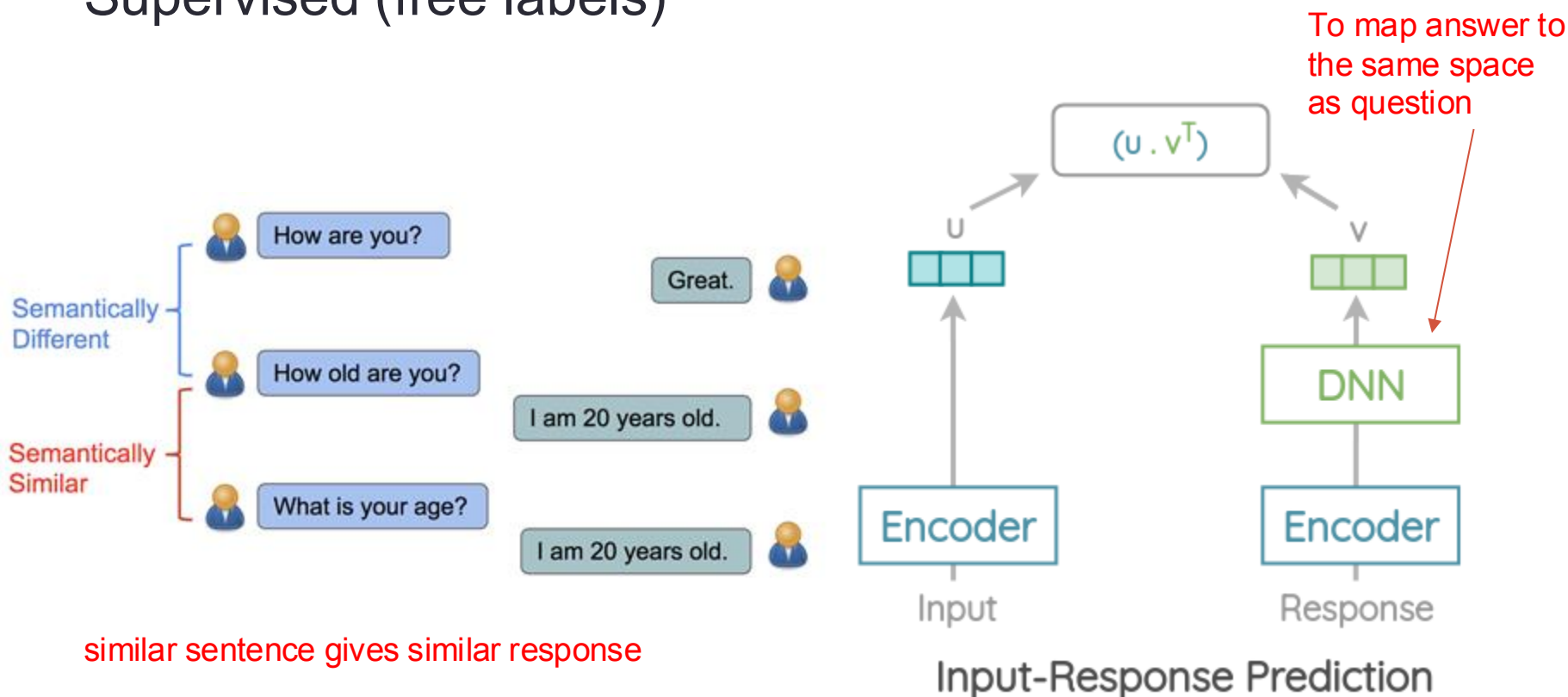
# Skip-thought task

Similar to skip-gram, use the middle to predict context  
Unsupervised



# Response prediction

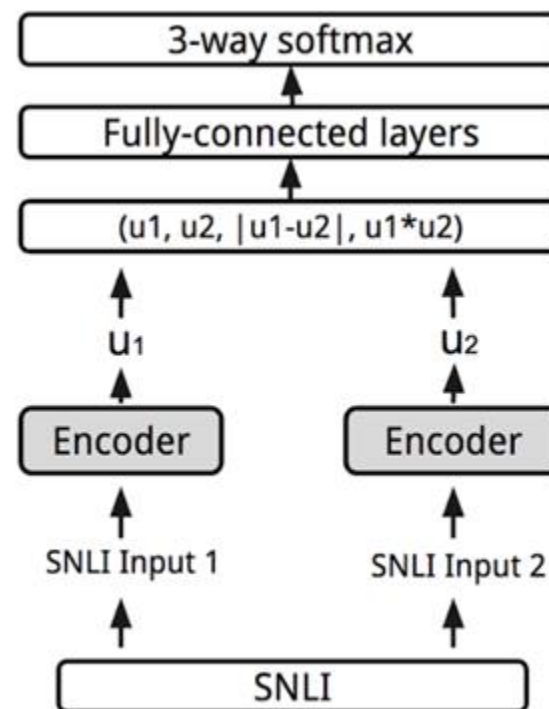
Match questions and answers in internet forum (scraped)  
Supervised (free labels)



# Natural Language Inference

Predict relationship between sentence  
Supervised

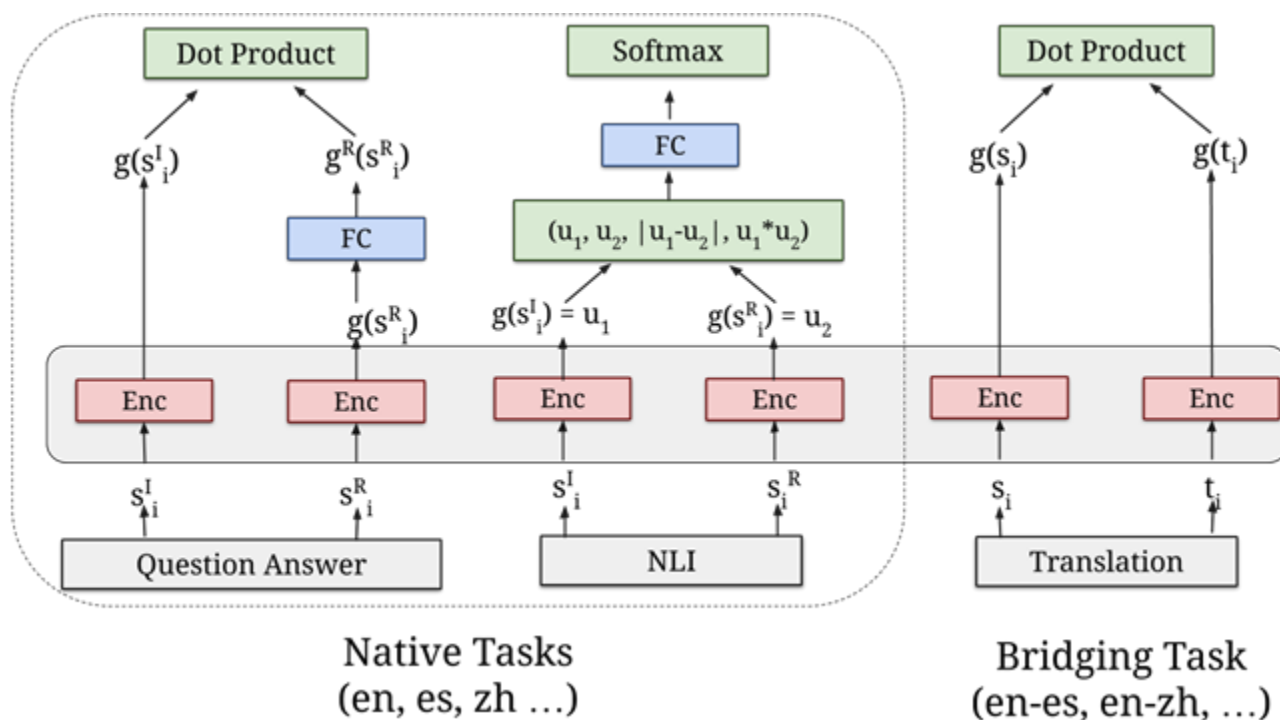
Premise	Hypothesis	Judgement
A soccer game with multiple males playing	Some men are playing a sport	entailment
I love Marvel movies	I hate Marvel movies	contradiction
I love Marvel movies	A ship arrived	neutral



# Multilingual USE

Can train to map multiple languages to the same presentation.

Can handle code switching, has Thai!



# Download-ables

cmlm-en-base

cmlm-en-large

cmlm-multilingual-base

cmlm-multilingual-base-br

cmlm-multilingual-preprocess

Trained with Masked Language Model loss (BERT-like)

large

multilingual

multilingual-large

multilingual-qa

qa

universal-sentence-encoder

Pytorch conversion

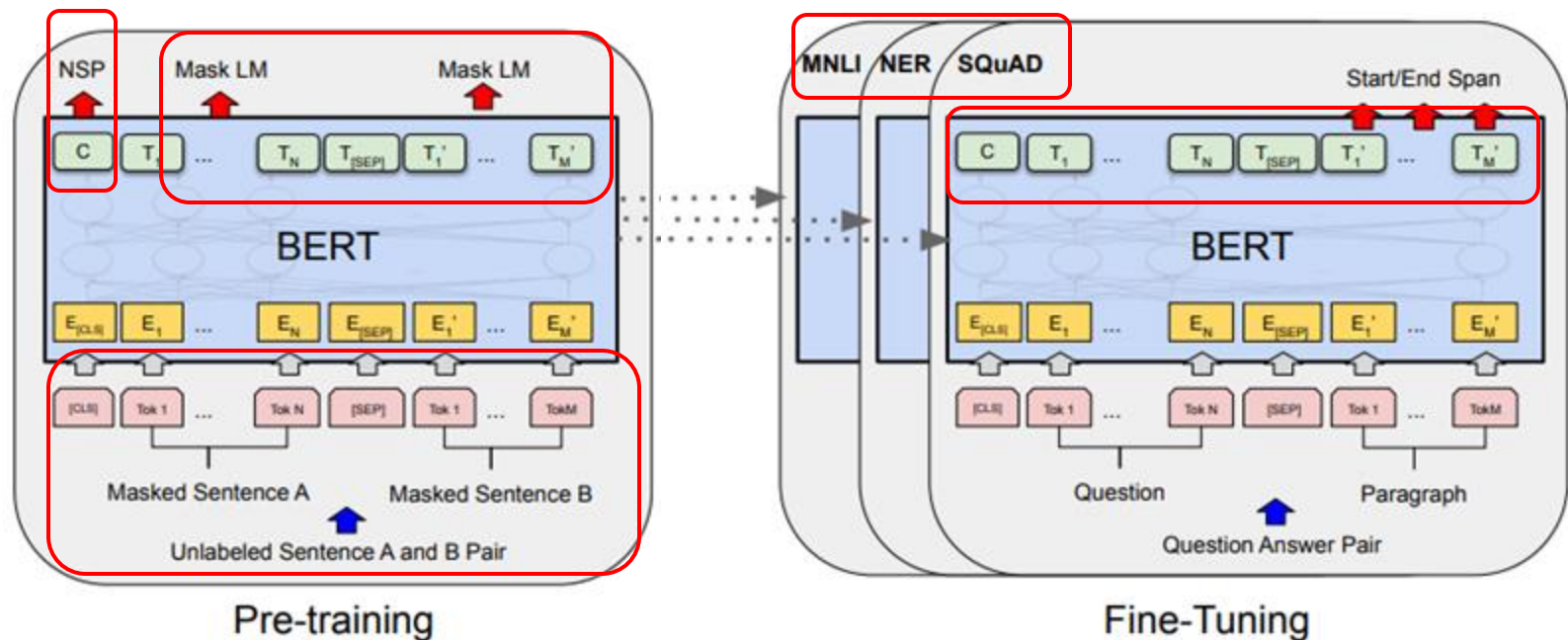
<https://huggingface.co/dayyass/universal-sentence-encoder-multilingual-large-3-pytorch>



# BERT-Based embeddings

---

# Sentence representation with BERT

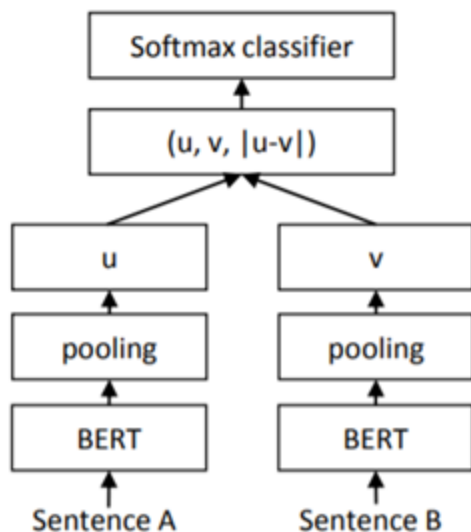


With BERT, we found that MLM training create good sentence representation too!

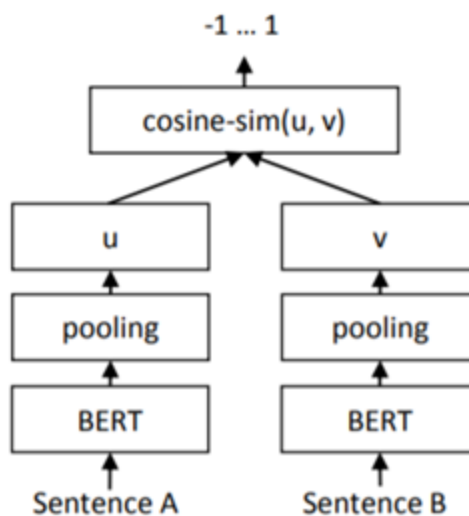
We can use NSP embedding or pool the token embeddings to create a sentence representation

# SBERT

Language Understanding



Semantic Understanding



Model	Spearman
<i>Not trained for STS</i>	
Avg. GloVe embeddings	58.02
Avg. BERT embeddings	46.35
InferSent - GloVe	68.03
Universal Sentence Encoder	74.92
SBERT-NLI-base	77.03
SBERT-NLI-large	79.23
<i>Trained on STS benchmark dataset</i>	
BERT-STSB-base	84.30 $\pm$ 0.76
SBERT-STSB-base	84.67 $\pm$ 0.19
SRoBERTa-STSB-base	<b>84.92 <math>\pm</math> 0.34</b>
BERT-STSB-large	<b>85.64 <math>\pm</math> 0.81</b>
SBERT-STSB-large	84.45 $\pm$ 0.43
SRoBERTa-STSB-large	85.02 $\pm$ 0.76
<i>Trained on NLI data + STS benchmark data</i>	
BERT-NLI-STSB-base	<b>88.33 <math>\pm</math> 0.19</b>
SBERT-NLI-STSB-base	85.35 $\pm$ 0.17
SRoBERTa-NLI-STSB-base	84.79 $\pm$ 0.38
BERT-NLI-STSB-large	<b>88.77 <math>\pm</math> 0.46</b>
SBERT-NLI-STSB-large	86.10 $\pm$ 0.13
SRoBERTa-NLI-STSB-large	86.15 $\pm$ 0.35

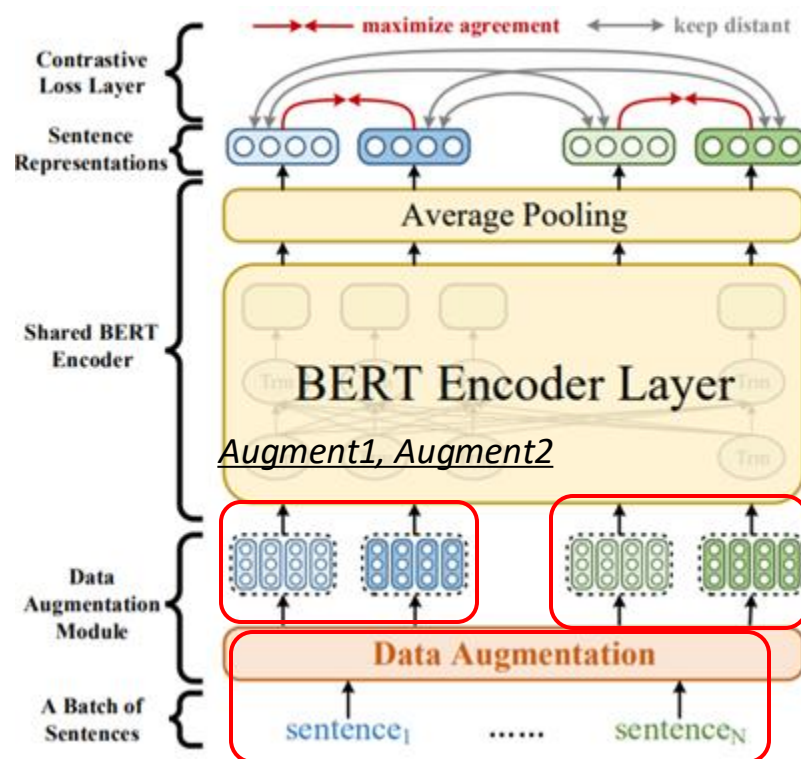
# Sentence level contrastive learning

- We can learn better sentence representation with some additional supervised (or unsupervised) sentence level contrastive learning

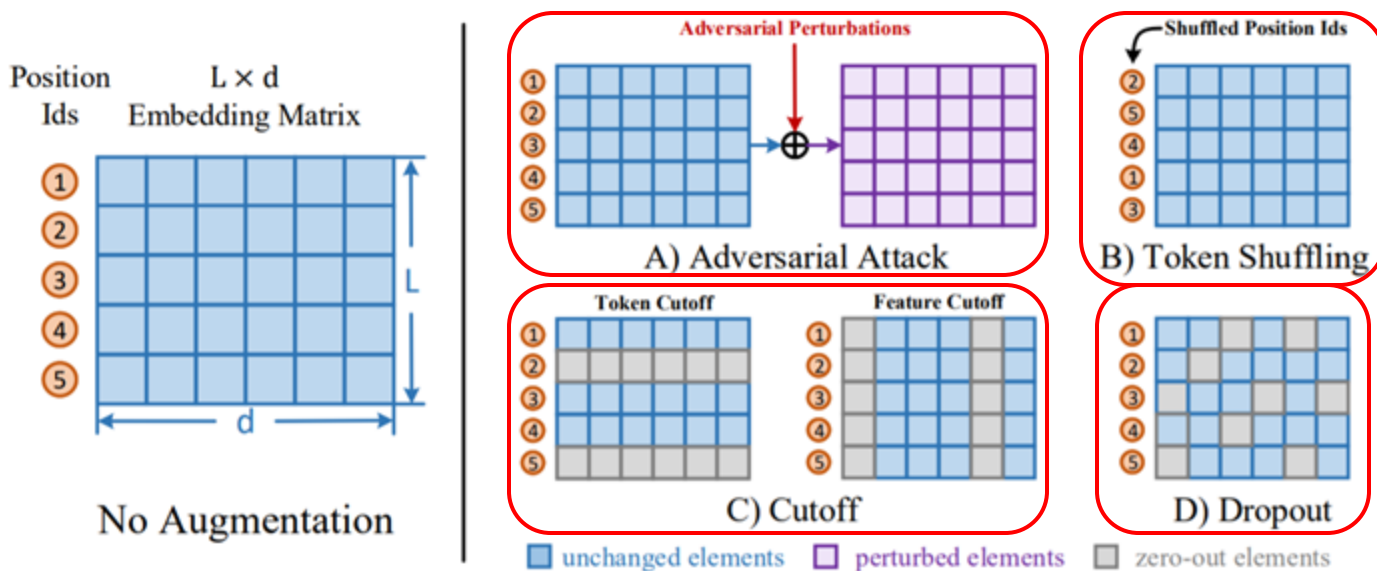
$$\mathcal{L}_{i,j} = -\log \frac{\exp(\text{sim}(r_i, r_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(r_i, r_k)/\tau)}$$

*Augment1, Augment2*

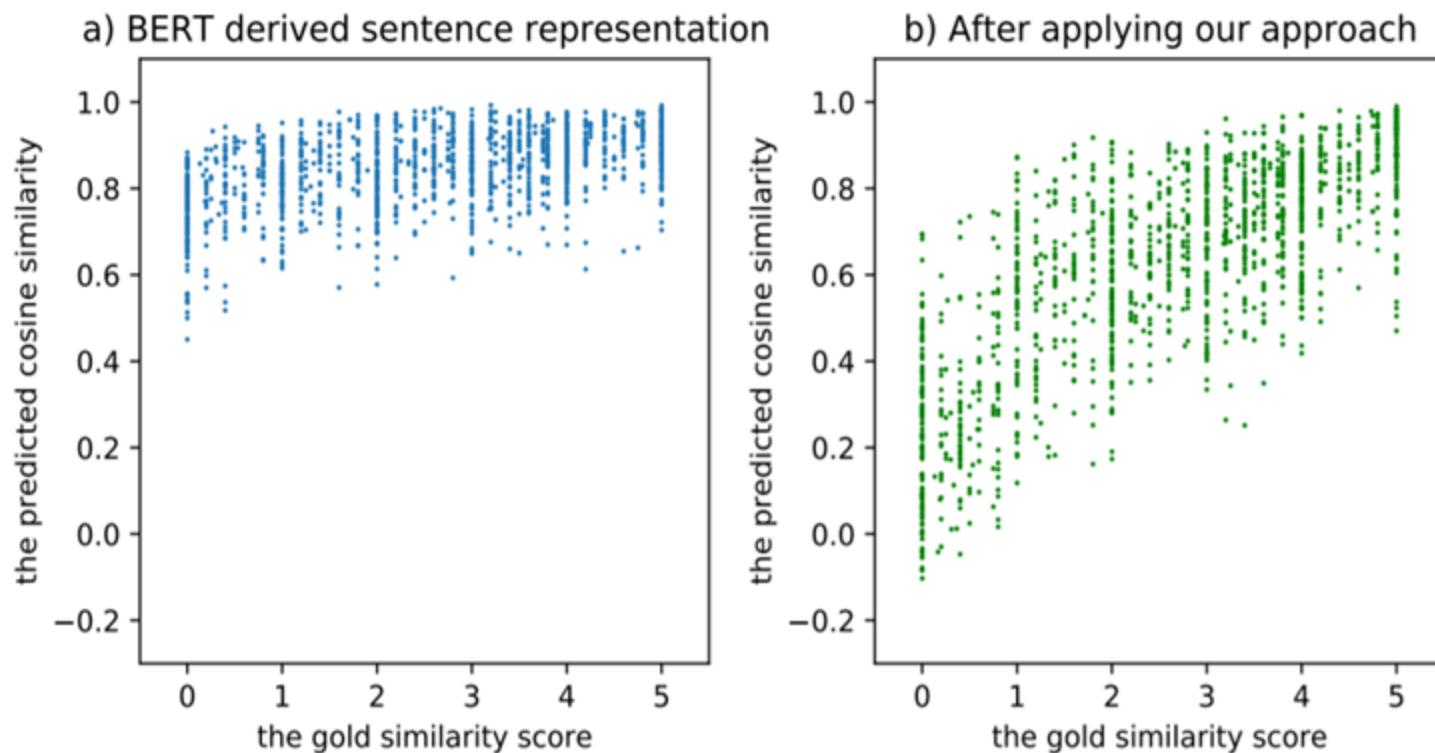
*Not augment2*



# ConSERT augmentations



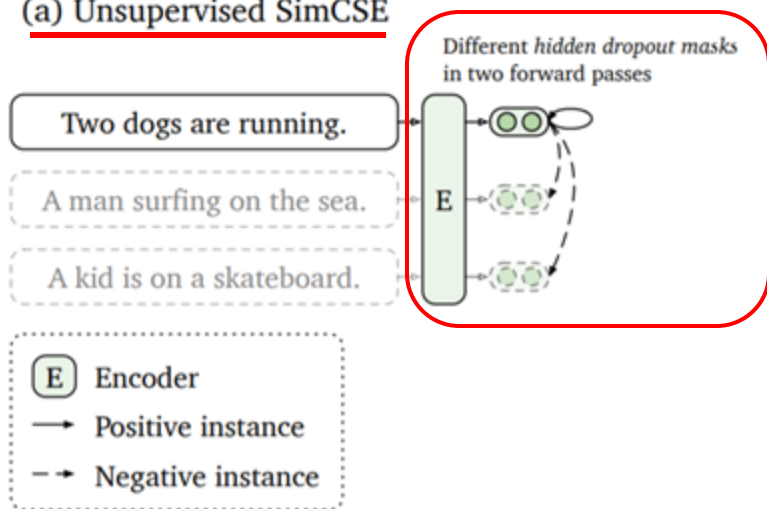
# ConSERT alignment



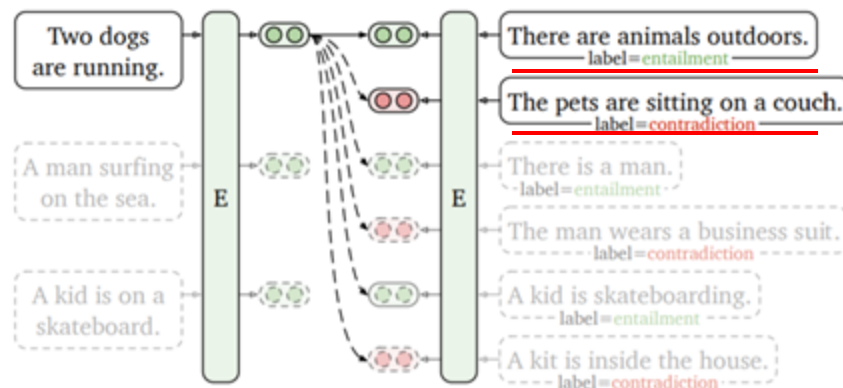
# SimCSE

- Use simple dropout in the model to create different versions of the same sentence

(a) Unsupervised SimCSE



(b) Supervised SimCSE



# SimCSE

Data augmentation			STS-B
None (unsup. SimCSE)			<b>82.5</b>
Crop	10%	20%	30%
	77.8	71.4	63.6
Word deletion	10%	20%	30%
	75.9	72.2	68.2
Delete one word			75.9
w/o dropout			74.2
Synonym replacement			77.4
MLM 15%			62.2

Other augmentations technique

Rather than contrastive, predict next sentence, 1 of 3 next sentences

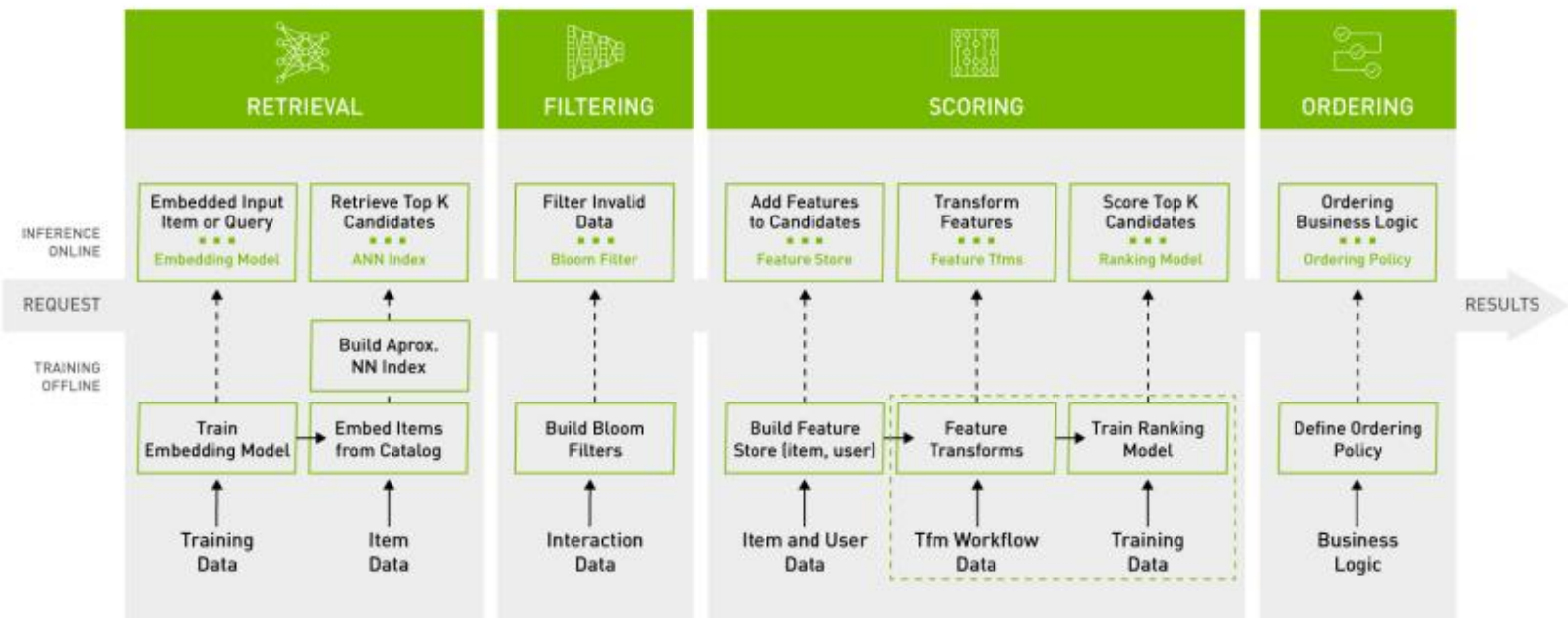
Training objective	$f_{\theta}$	$(f_{\theta_1}, f_{\theta_2})$
Next sentence	67.1	68.9
Next 3 sentences	67.4	68.8
Delete one word	75.9	73.1
Unsupervised SimCSE	<b>82.5</b>	80.7

$$\mathcal{L}_{i,j} = -\log \frac{\exp(\text{sim}(r_i, r_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(r_i, r_k)/\tau)}$$



# What's other use of embeddings?

- Retrieval and recommendation



# What's other use of embeddings?

- Learn joint embeddings between different modalities

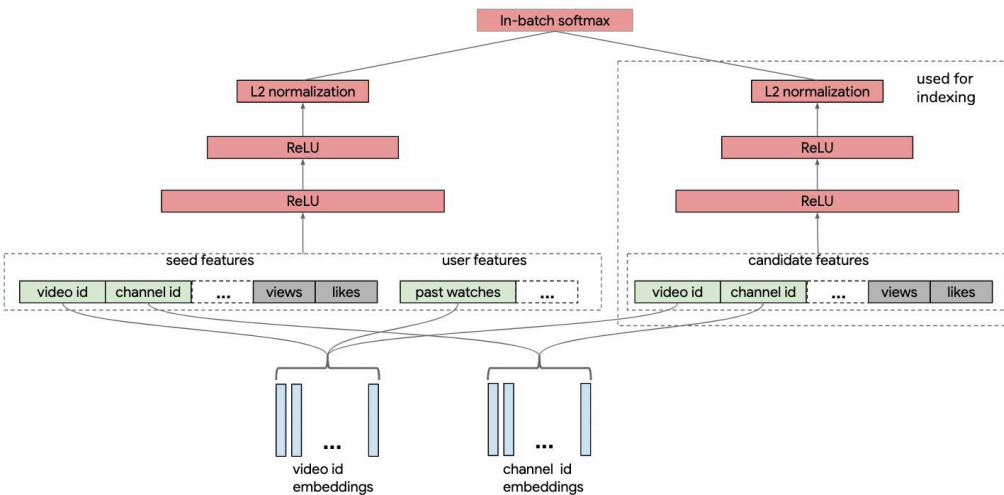
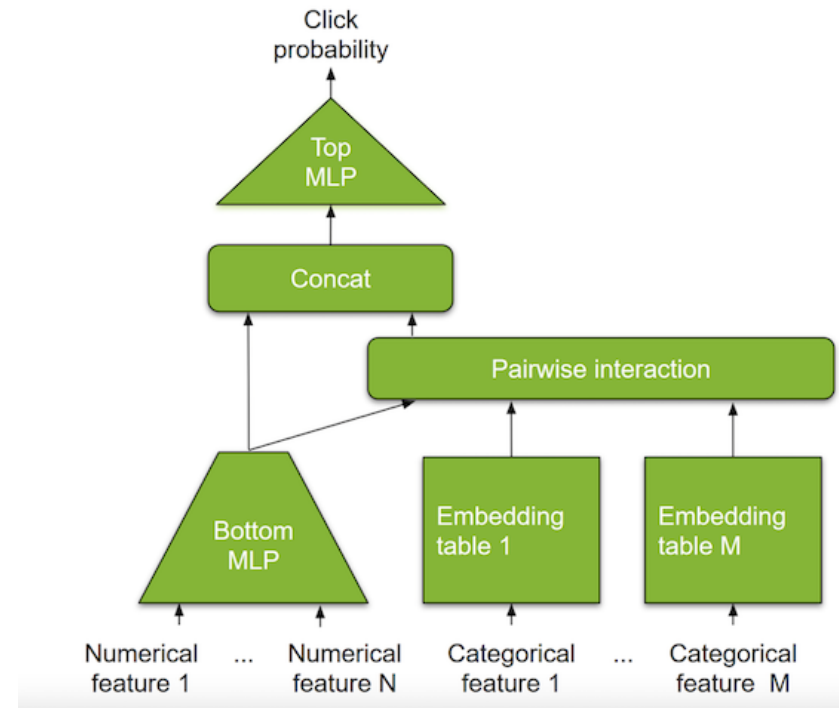


Figure 2: Illustration of the Neural Retrieval Model for YouTube.

Two tower model



Joint interaction model

# BGE-M3

- A retrieval model (Query -> Document)
- Built on top of BGE (Chinese embedding model)
  - BGE: Masked LM finetuned with contrastive and task specific losses
- BGE-M3 (multilingual, multifunction, multigranularity)
- Trained by multiple losses terms that utilizes different parts of the model embeddings

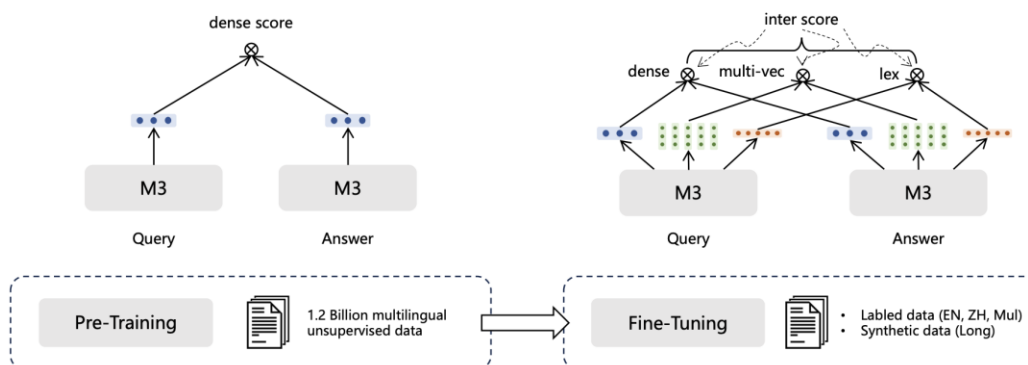
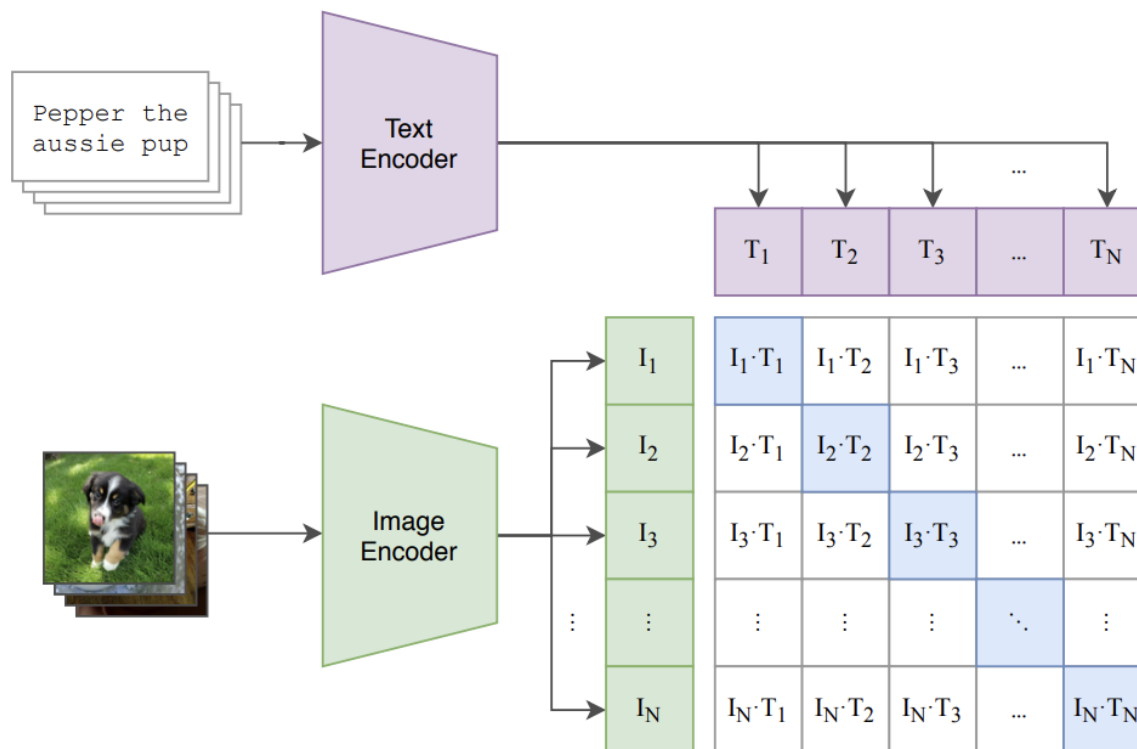


Figure 2: Multi-stage training process of M3-Embedding with self-knowledge distillation.

# CLIP

- Contrastive learning on image-text pairs

(1) Contrastive pre-training



# Outline

- Contrastive learning
- Sentence embeddings
  - MUSE
  - SimCSE
  - BGE
  - CLIP