# Multimodal Machine Learning
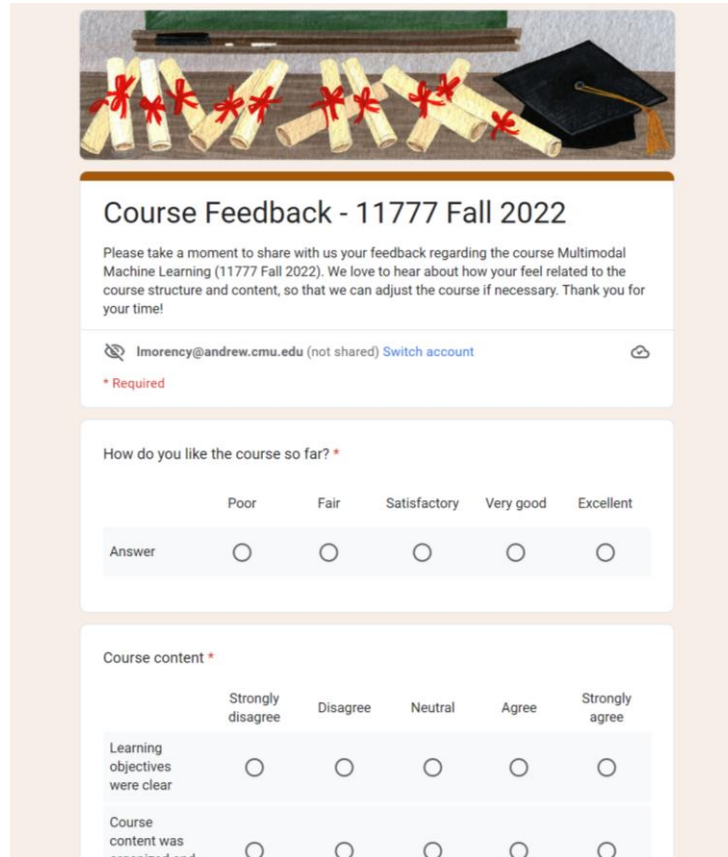
## Lecture 6.2: Multimodal Aligned Representations

**Louis-Philippe Morency**

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*

# Administrative Stuff

Carnegie Mellon University

# Share Your Thoughts!

## Deadline
Please submit your feedback about this course before this Wednesday 10/5

Optional,
but greatly appreciated! ☺

Anonymous, by default.
- You can optionally share your email address if you want us to follow-up with you directly.

# Second Project Assignment (Due Monday 10/10)

Main goals:

1. Help clarify and expand your research ideas
   - Build qualitative intuitions by directly studying the original data
   - Perform analyses on your dataset, relevant to your research ideas
2. Understand the structure in your data and modalities
   - Perform analyses and visualizations to understand each modality
   - Study representations from CNNs, word2vec, BERT, …

Two types of analyses:
   - Idea-oriented analyses
   - Modality-oriented analyses

# Multimodal Machine Learning

## Lecture 6.2: Alignment and Representation

Louis-Philippe Morency

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*

# Objectives of today's class

- Transformer pre-training
  - BERT: Bidirectional Encoder Representations from Transformers
- Multimodal transformers (Image and language)
  - Concatenated transformers (VisualBERT, Uniter)
  - Crossmodal transformers (ViLBERT, LXMERT
  - Modality-shift transformer (MAG-BERT)
  - Video and language transformers (VideoBERT, ActBERT)
- Visual transformers
  - Vision transformer, Masked Auto-Encoder
  - Visual-and-language transformer (ViLT, ALBEF)

# BERT: Transformer Pre-training

# Transformer Self-Attention

# Transformer – Residual Connection

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

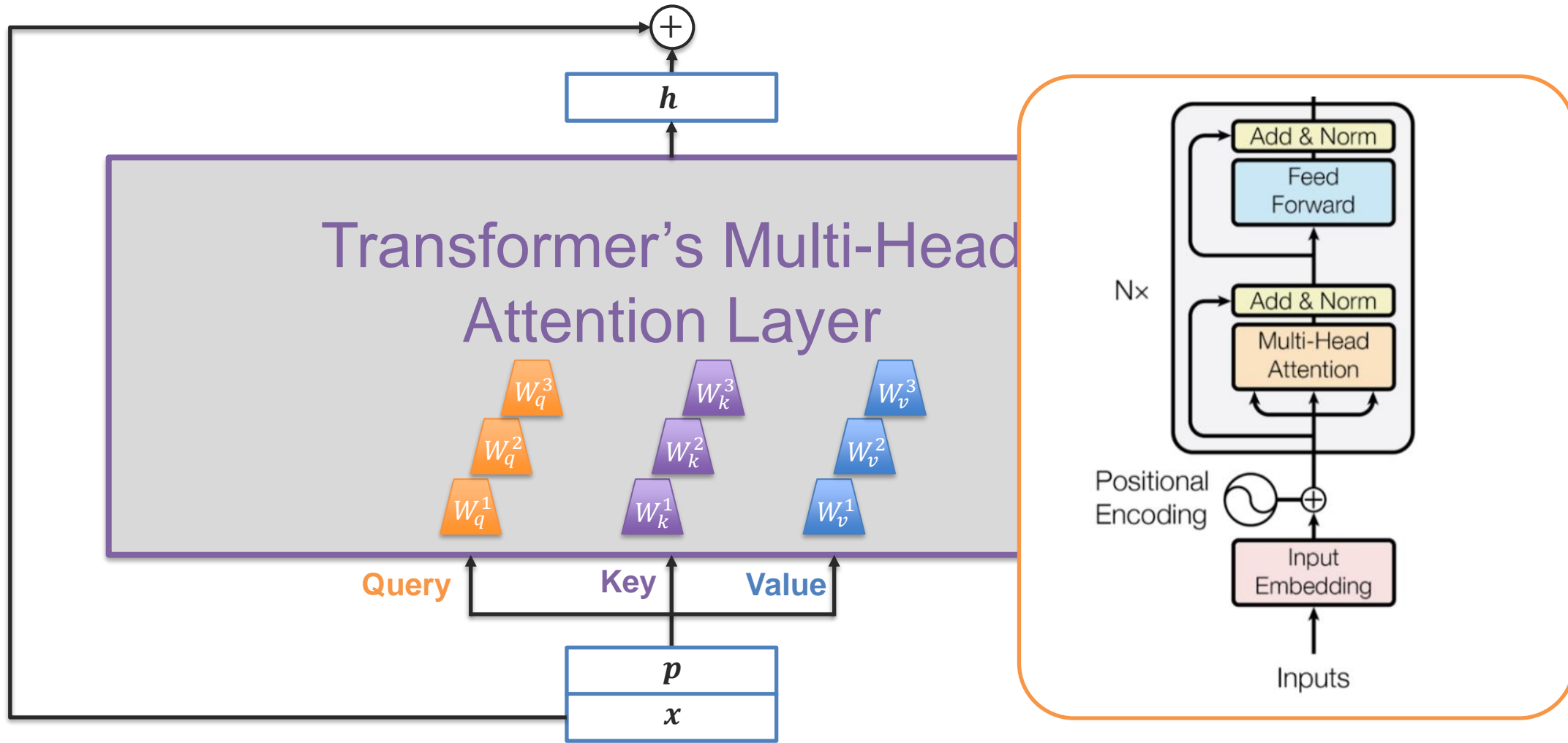① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

③ Can be used learn relationship between sentences

④ Models bidirectional and long-range interactions between tokens

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**How can we do all this?**

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

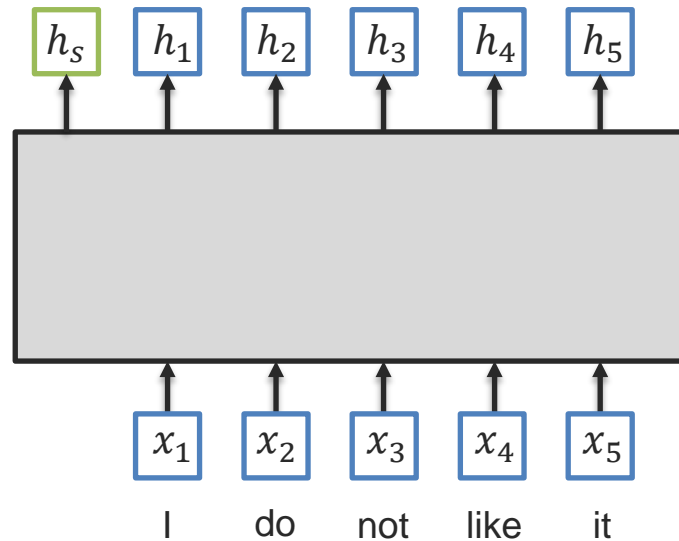I    do    not    like    it        I    enjoy    my    time    here

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

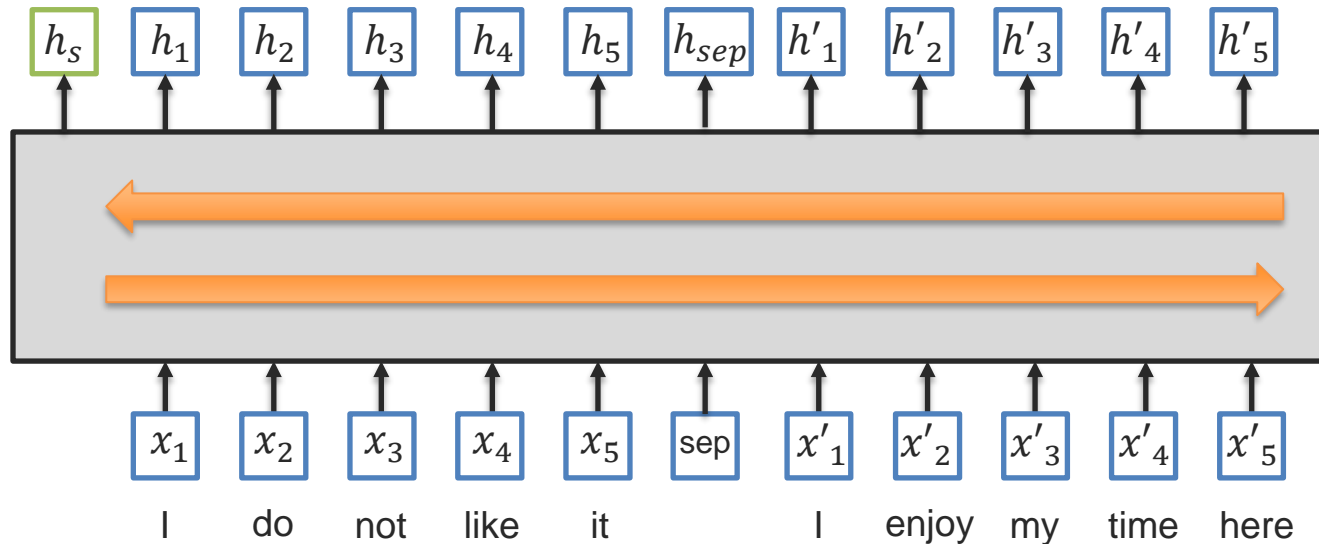③ Can be used learn relationship between sentences

④ Models bidirectional interactions between tokens



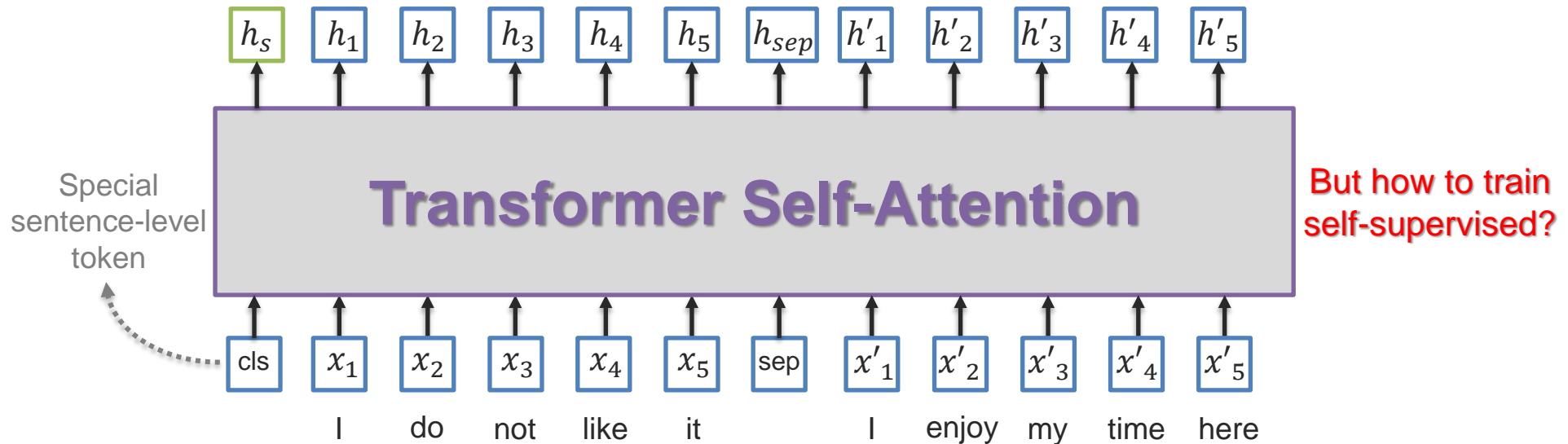Special sentence-level token

But how to train self-supervised?

# Pre-training BERT Model

**1** **Masked Language Model**

Randomly mask input tokens and then try to predict them

**What is the loss function?**

# Pre-training BERT Model

② **Next Sentence Prediction**

Given two sentences, predict if this is the next one or not

IsNext
**- or -**
NotNext

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I    do    not    like    it            I    enjoy    my    time    here

# Fine-Tuning BERT

**1** Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**How?**



| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

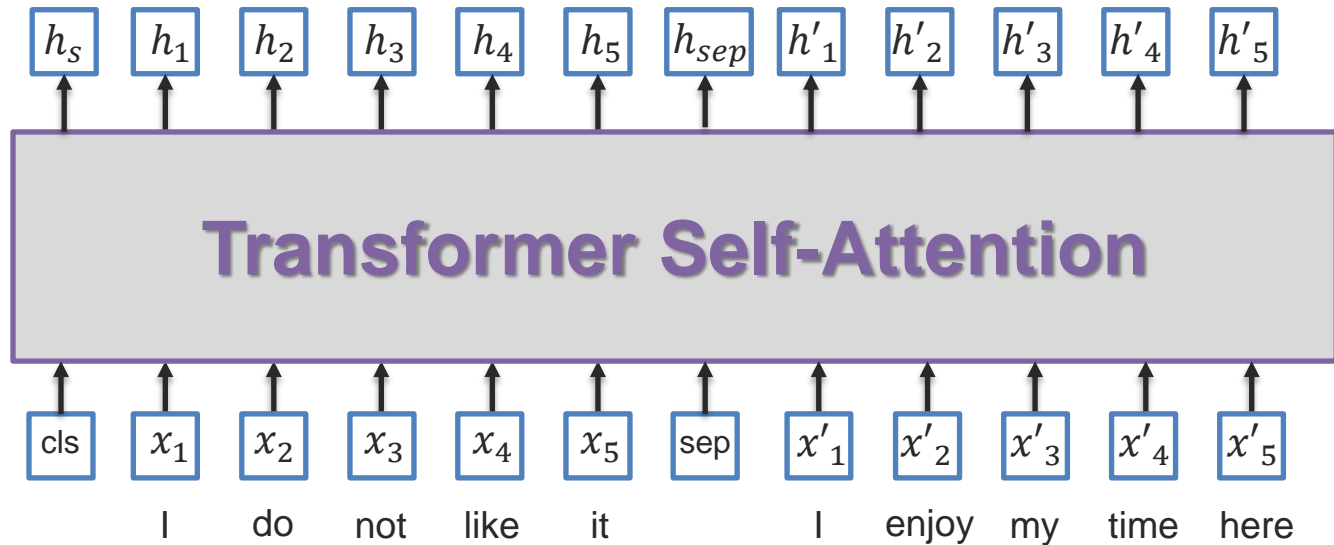I    do    not    like    it         I    enjoy    my    time    here

# Fine-Tuning BERT

**1** Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification



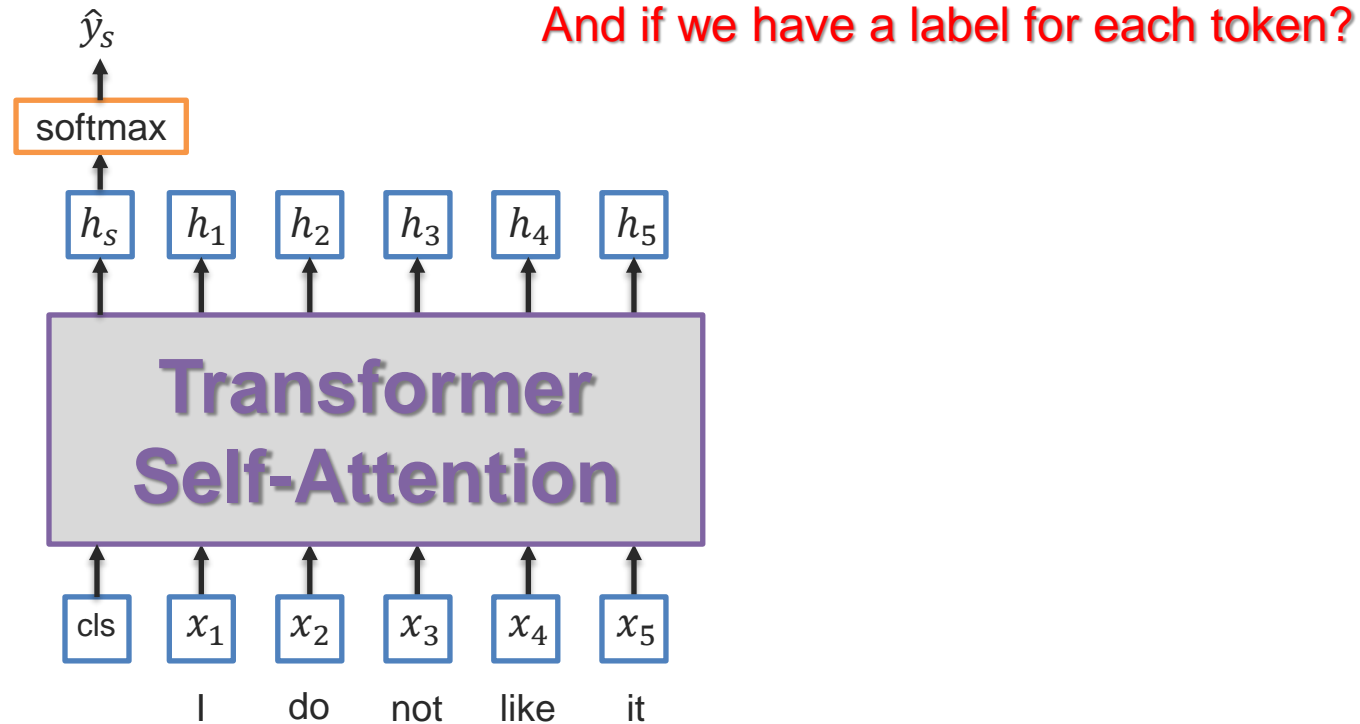And if we have a label for each token?

$\hat{y}_s$

softmax

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |

## Transformer Self-Attention

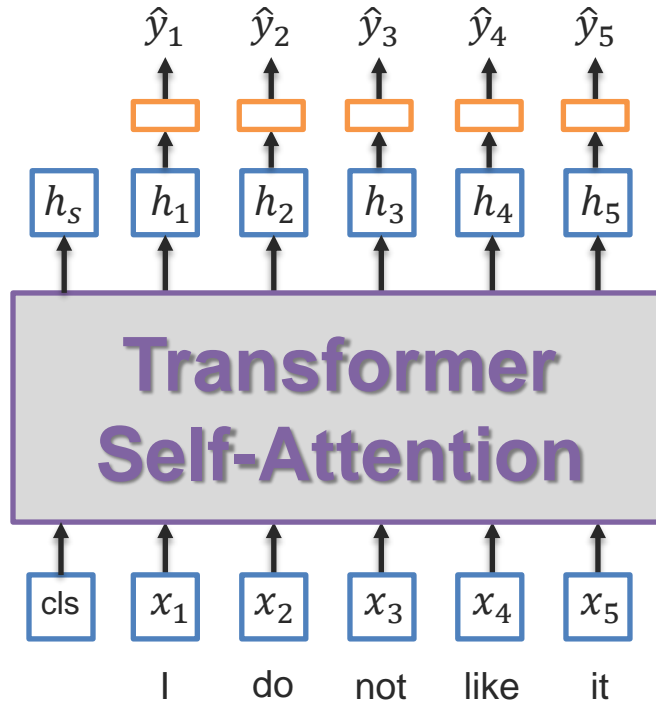| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

I    do   not   like   it

# Fine-Tuning BERT

② Token-level classification for only one sentence

Examples: part-of-speech tagging, slot filling



How to compare two sentences?

# Fine-Tuning BERT

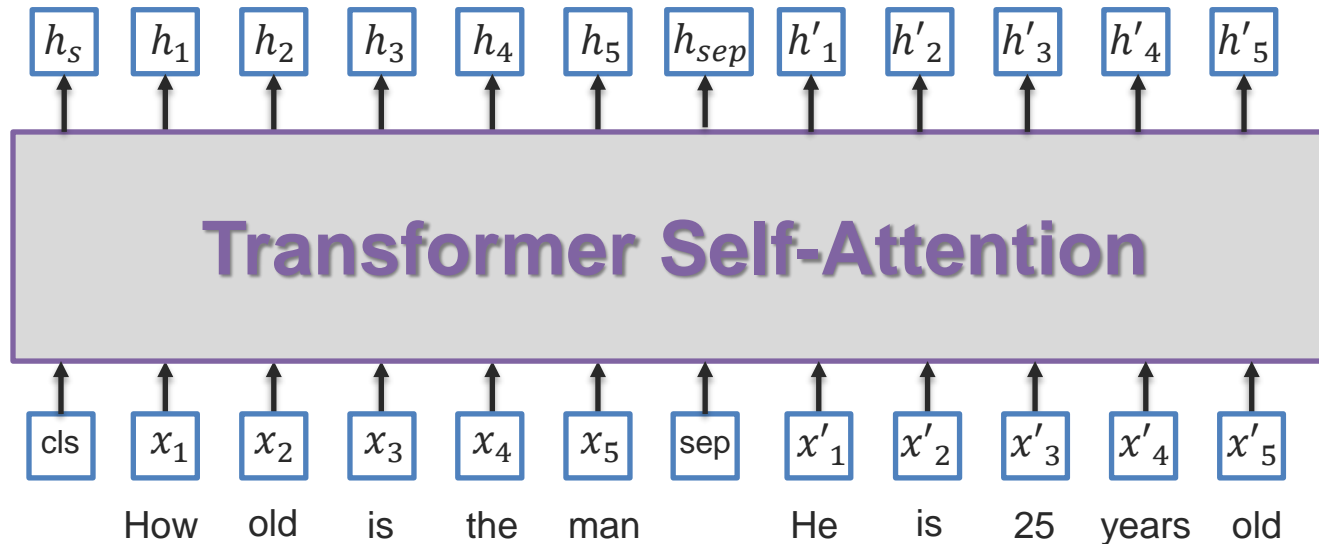**Question-answering: find start/end of the answer in the document**

**Paragraph:** "...Other legislation followed, including the Migratory Bird Conservation Act of 1929, a *1937 treaty* prohibiting the hunting of right and gray whales, and the *Bald Eagle Protection Act of 1940*. These *later laws* had a low cost to society—the species were relatively rare—and little *opposition* was raised."

**Question 1:** "*Which laws faced significant opposition?*"
**Plausible Answer:** *later laws*
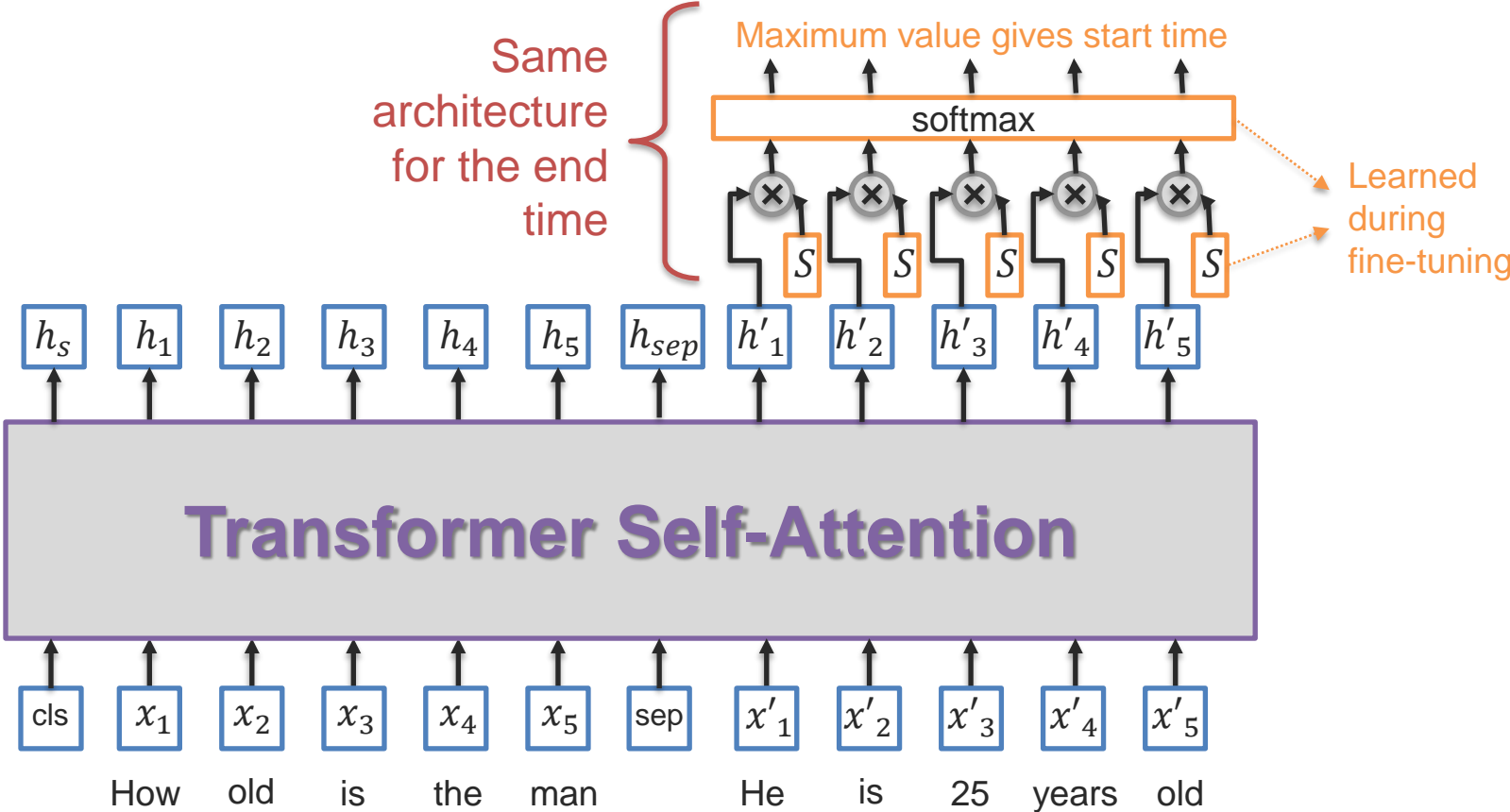
**Question 2:** "*What was the name of the 1937 treaty?*"
**Plausible Answer:** *Bald Eagle Protection Act*

$h_s$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_{sep}$ $h'_1$ $h'_2$ $h'_3$ $h'_4$ $h'_5$

**Transformer Self-Attention**

How?

cls $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ sep $x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$
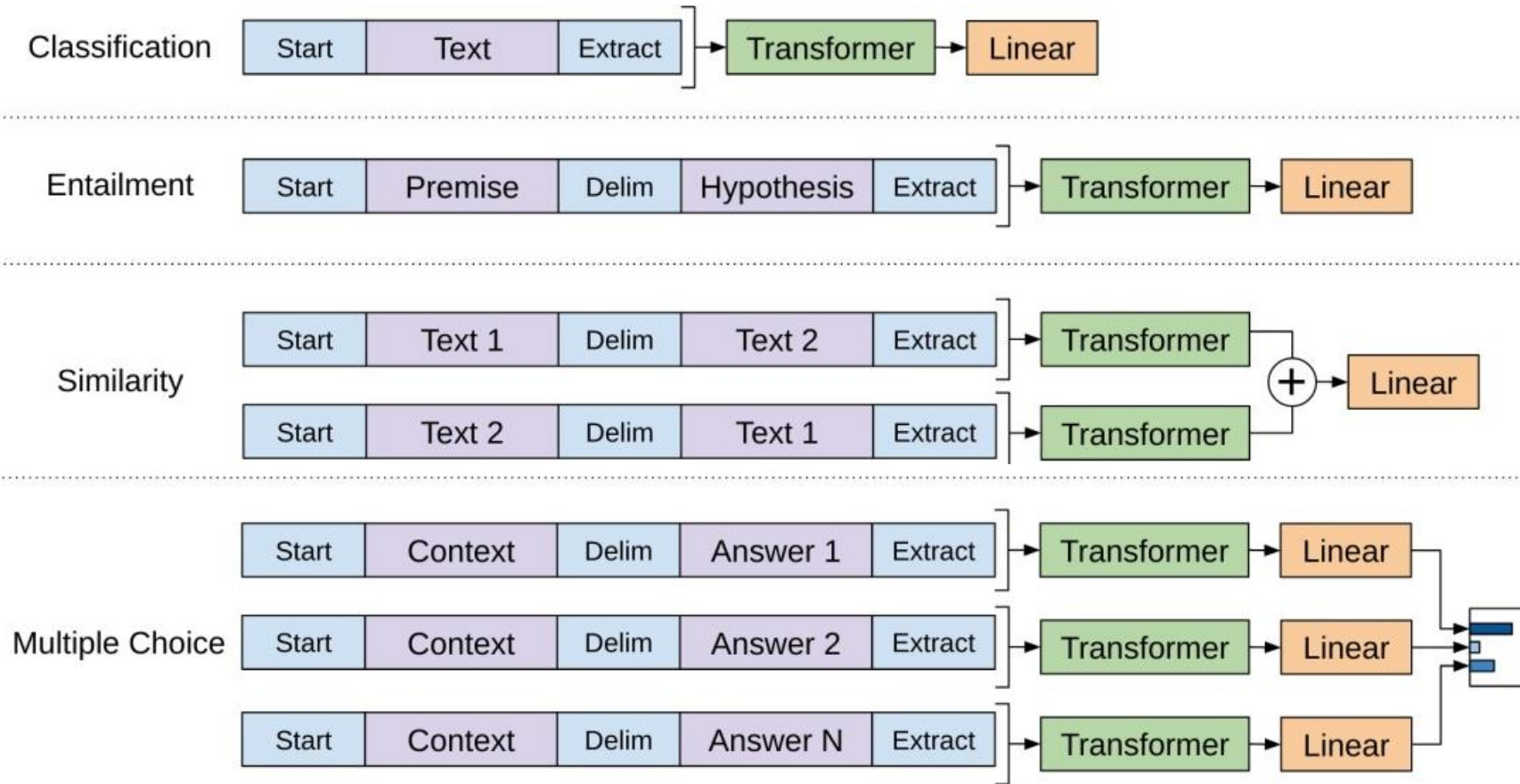
How old is the man     He is 25 years old

# Fine-Tuning BERT

④ Question-answering: find start/end of the answer in the document

# Other Fine-tuning Approaches



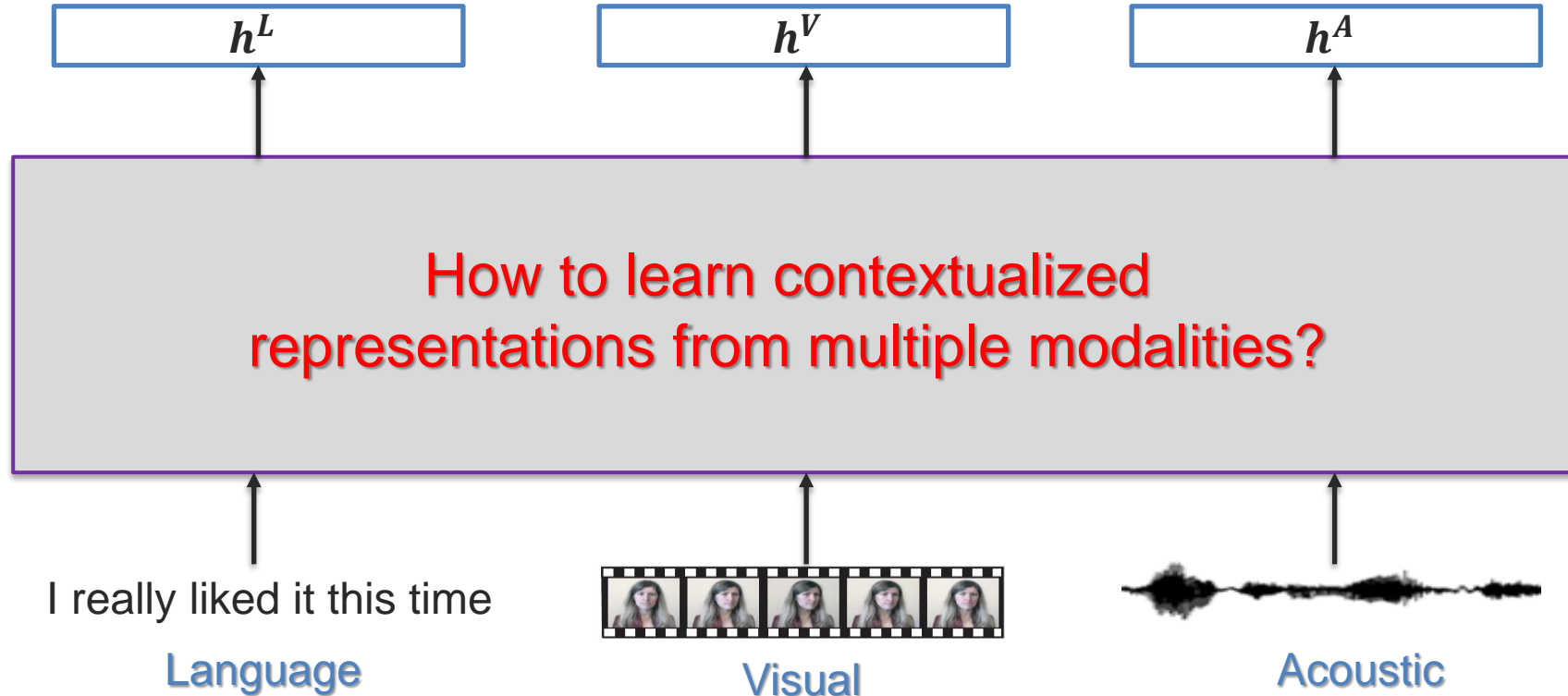https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

# Multimodal Transformers

Language Technologies Institute

Carnegie Mellon University

# Multimodal Embeddings
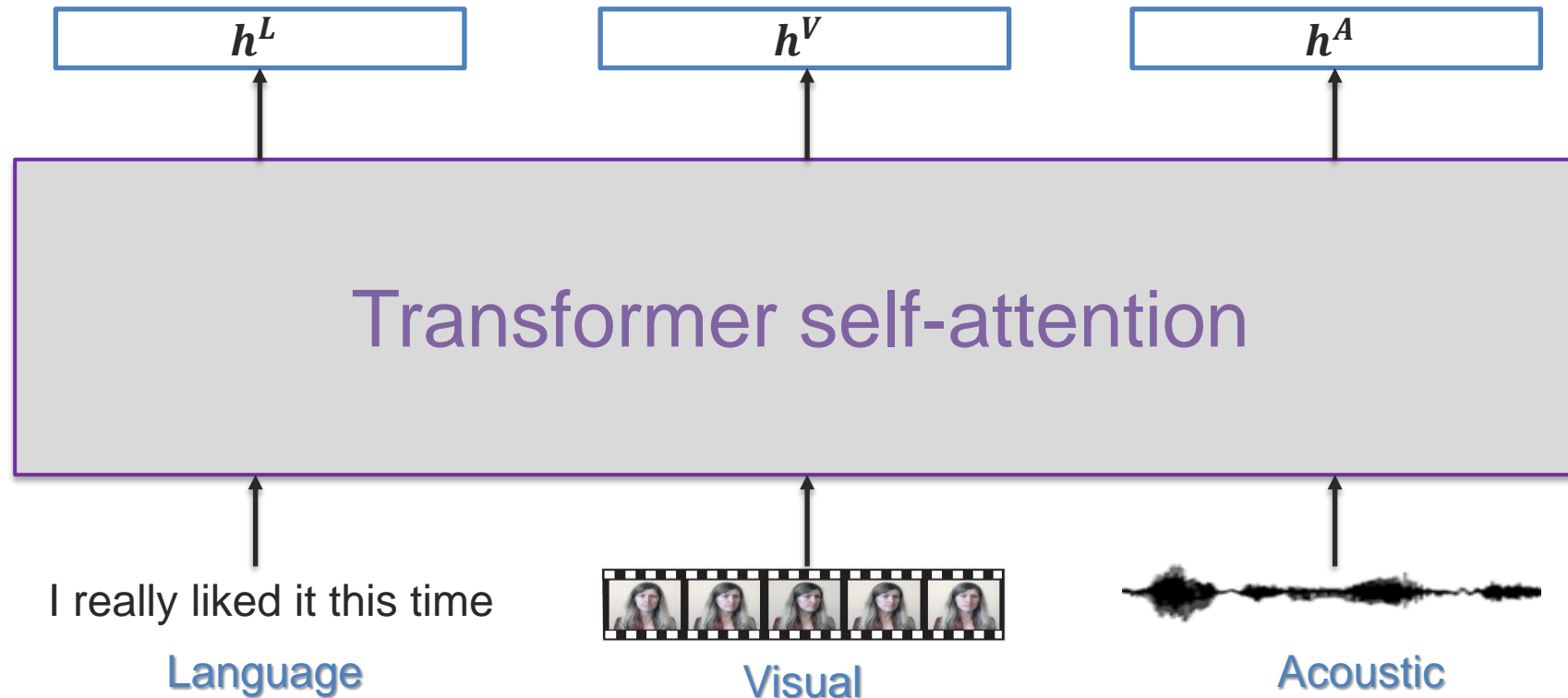


Option 1: Concatenate modalities and learn BERT transformer

# Simple Solution: Contextualized Multimodal Embeddings

| $h^L$ | $h^V$ | $h^A$ |
|---|---|---|

Transformer self-attention

I really liked it this time

Language                    Visual                    Acoustic
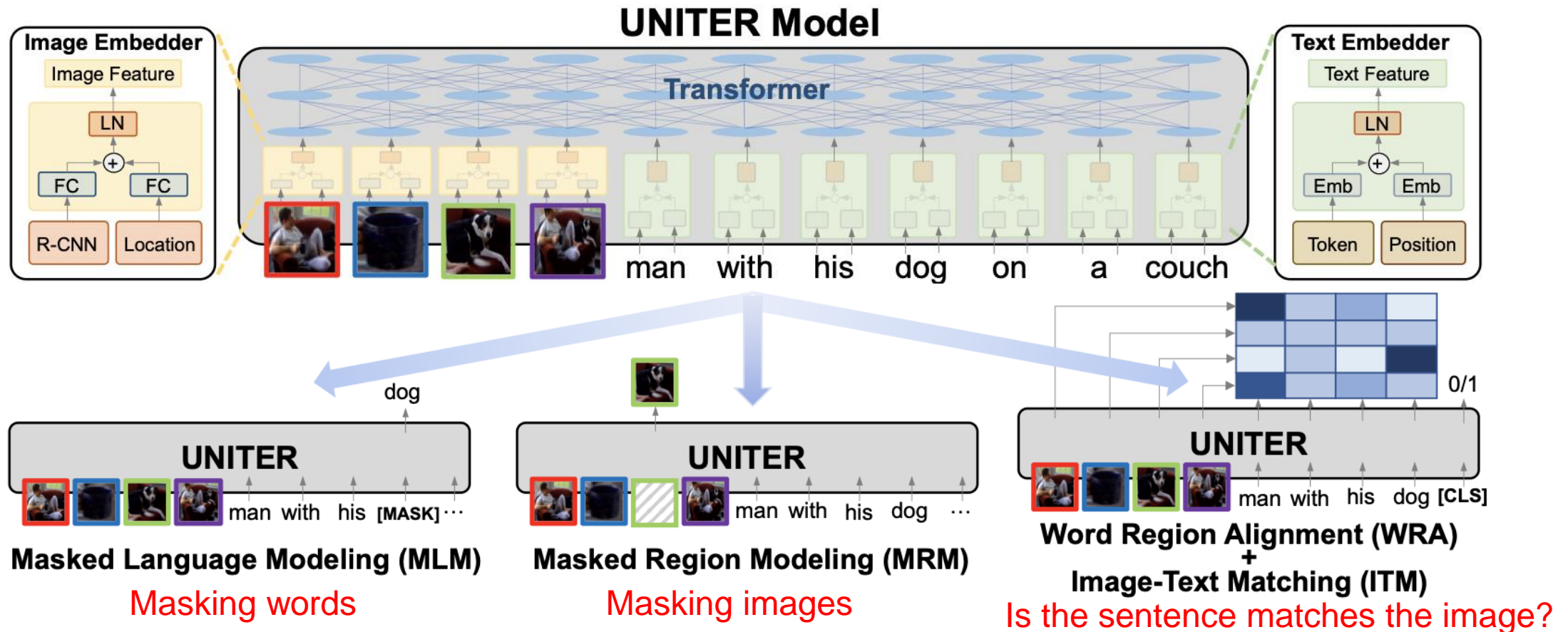
# VisualBERT



Li, Liunian Harold, et al. "Visualbert: A simple and performant baseline for vision and language." *arXiv* (2019).

# UNITER

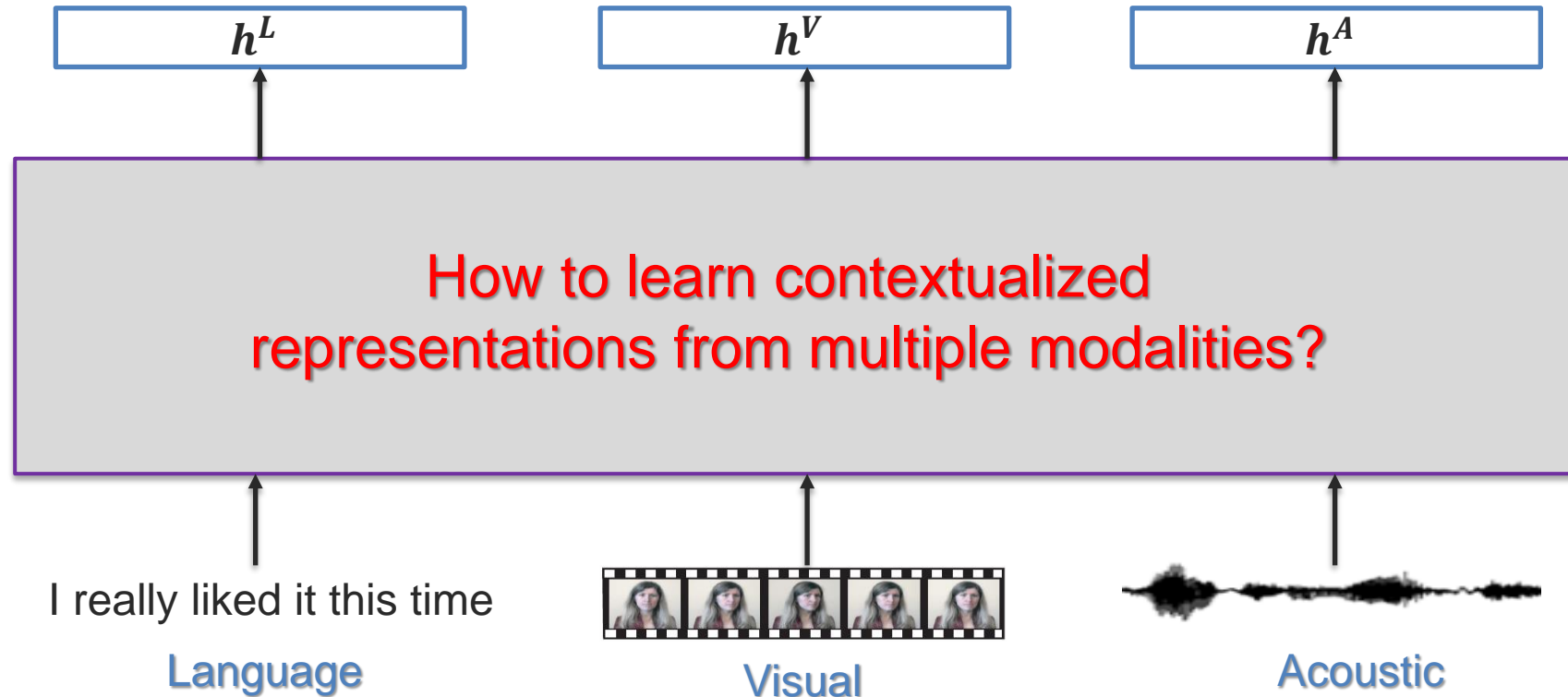Similar Transformer architecture to BERT and VisualBERT… but with slightly different optimization



**UNITER Model**

Image Embedder — Image Feature — LN — FC + FC — R-CNN, Location

Text Embedder — Text Feature — LN — Emb + Emb — Token, Position

Transformer

man with his dog on a couch

**Masked Language Modeling (MLM)**
Masking words

**Masked Region Modeling (MRM)**
Masking images

**Word Region Alignment (WRA) + Image-Text Matching (ITM)**
Is the sentence matches the image?

Chen, Yen-Chun, et al. "Uniter: Universal image-text representation learning." *European conference on computer vision*. 2020.

# Multimodal Embeddings

$$h^L \qquad h^V \qquad h^A$$

How to learn contextualized
representations from multiple modalities?

I really liked it this time

**Language**  **Visual**  **Acoustic**

Option 2: Look at pairwise interactions between modalities

# Multimodal Transformer – Pairwise Cross-Modal

# Cross-Modal Transformer Module ($V \to L$)
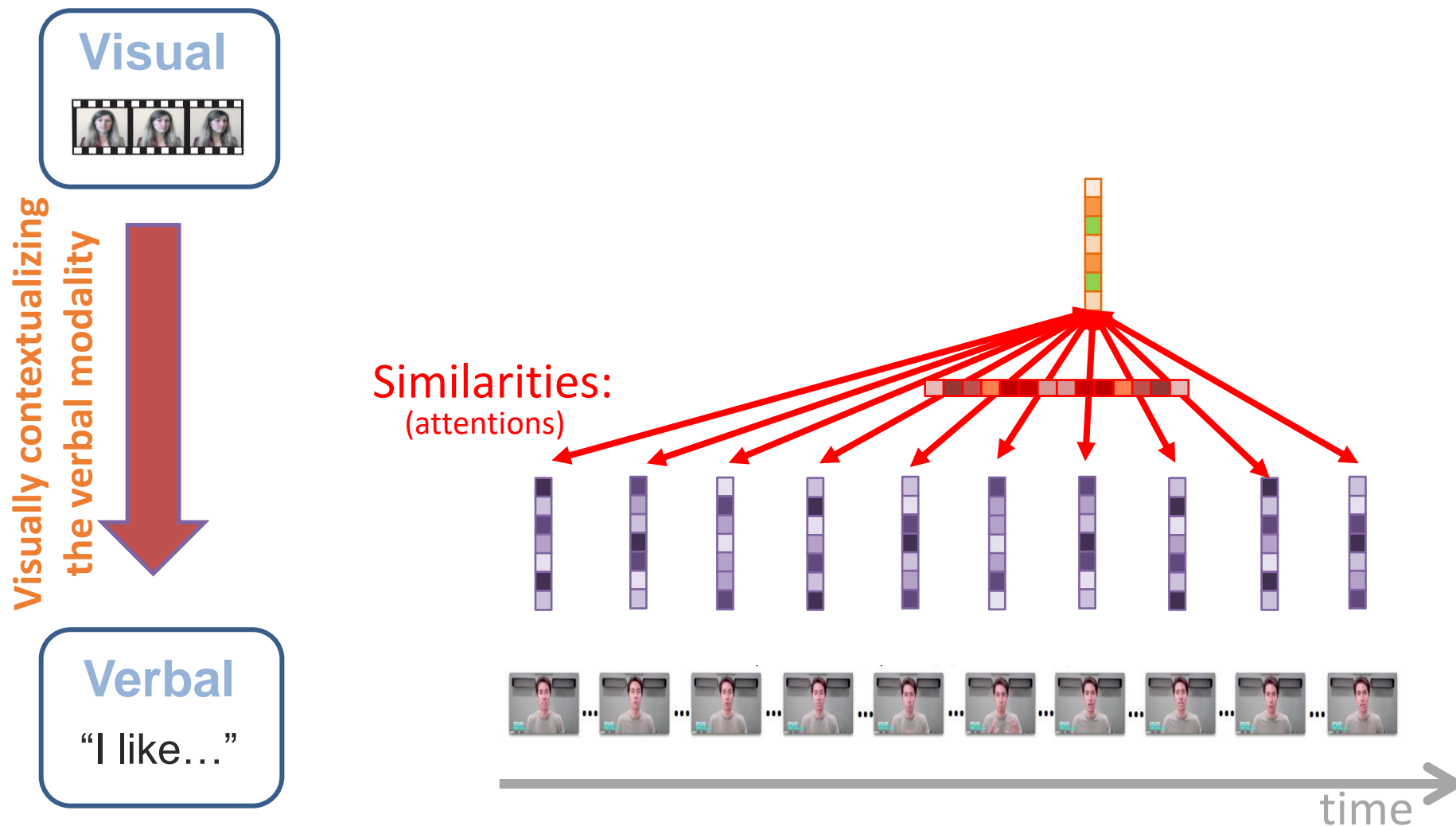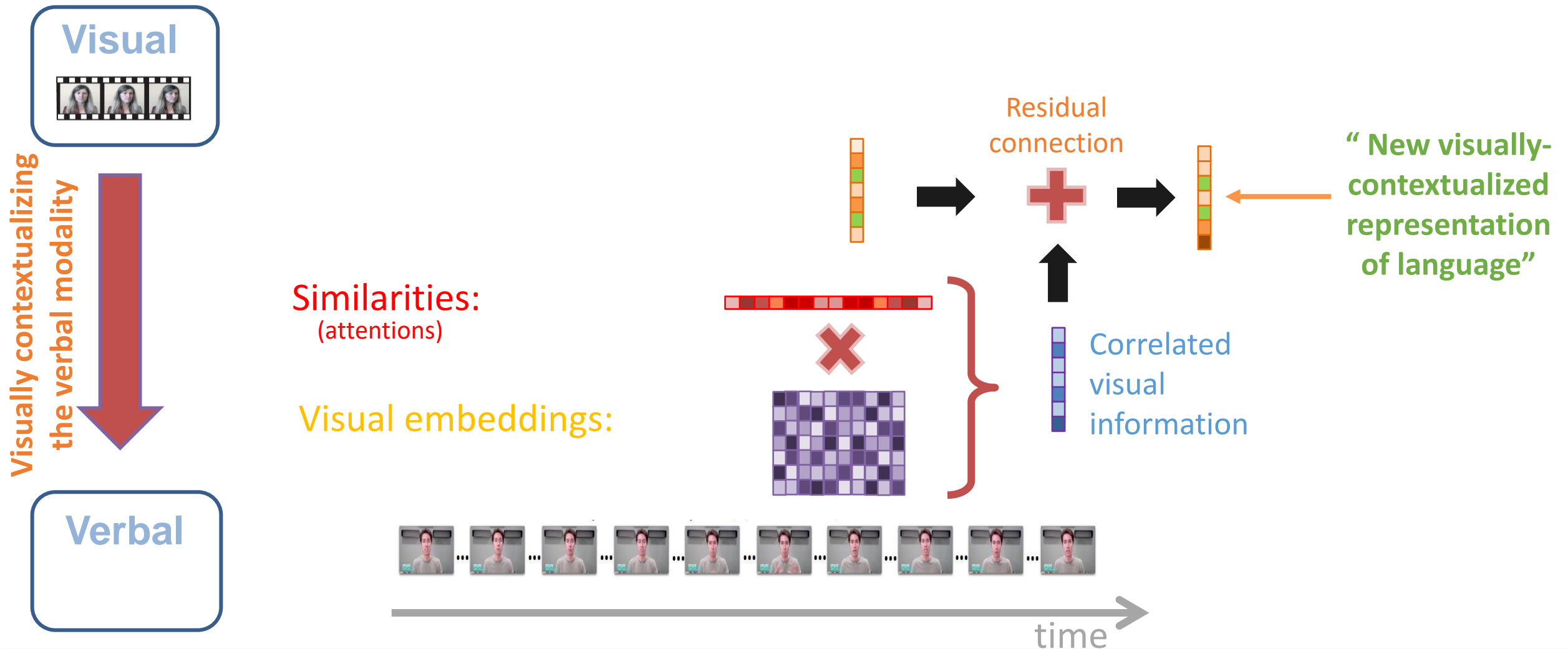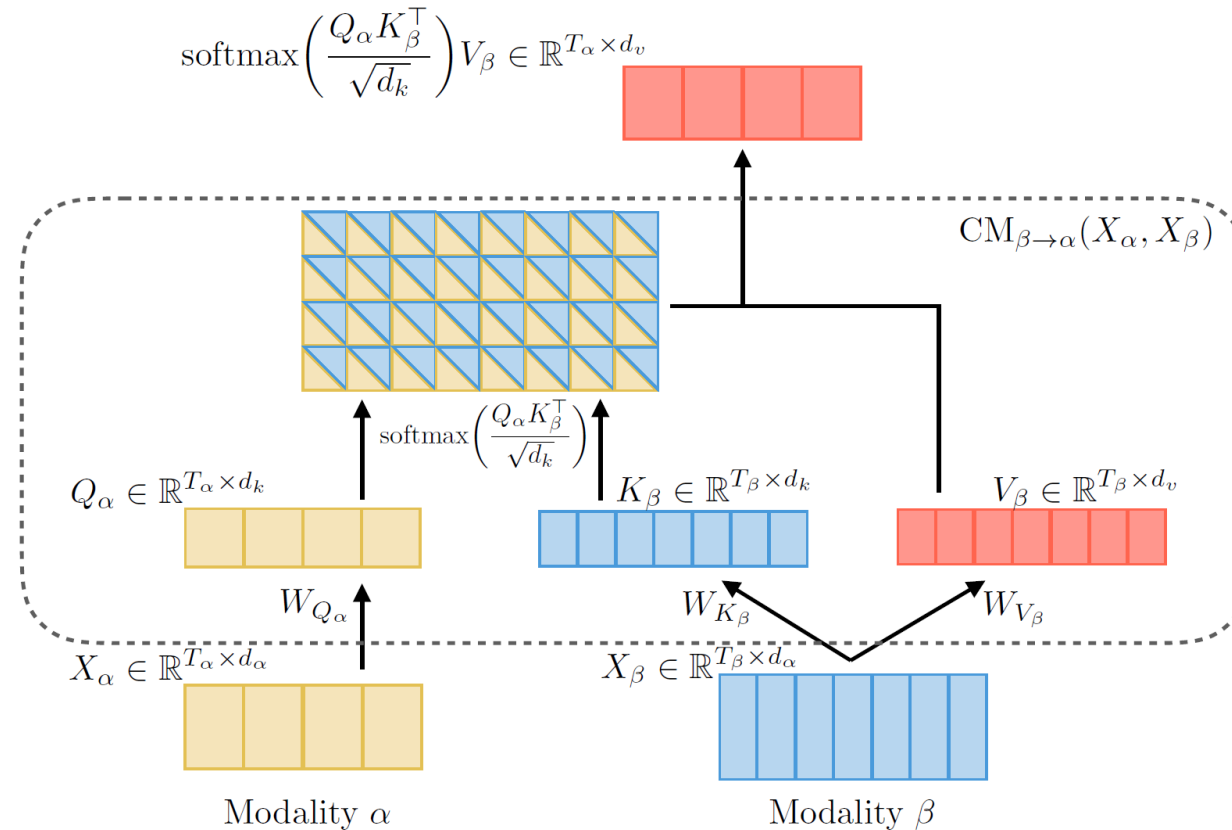
# Cross-Modal Transformer Module ($V \rightarrow L$)

# Cross-Modal Transformer Module ($\beta \to \alpha$)



$$\text{softmax}\left(\frac{Q_\alpha K_\beta^\top}{\sqrt{d_k}}\right) V_\beta \in \mathbb{R}^{T_\alpha \times d_v}$$

$\text{CM}_{\beta \to \alpha}(X_\alpha, X_\beta)$

$$\text{softmax}\left(\frac{Q_\alpha K_\beta^\top}{\sqrt{d_k}}\right)$$

$Q_\alpha \in \mathbb{R}^{T_\alpha \times d_k}$

$K_\beta \in \mathbb{R}^{T_\beta \times d_k}$

$V_\beta \in \mathbb{R}^{T_\beta \times d_v}$

$W_{Q_\alpha}$

$W_{K_\beta}$

$W_{V_\beta}$

$X_\alpha \in \mathbb{R}^{T_\alpha \times d_\alpha}$

$X_\beta \in \mathbb{R}^{T_\beta \times d_\alpha}$

Modality $\alpha$

Modality $\beta$

Tsai et al., Multimodal Transformer for Unaligned Multimodal Language Sequences, ACL 2019

# ViLBERT



**Cross-Modal Transformer Modules**

**Unimodal Transformer**

Lu, Jiasen, et al. "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks." *arXiv* (August 6, 2019).
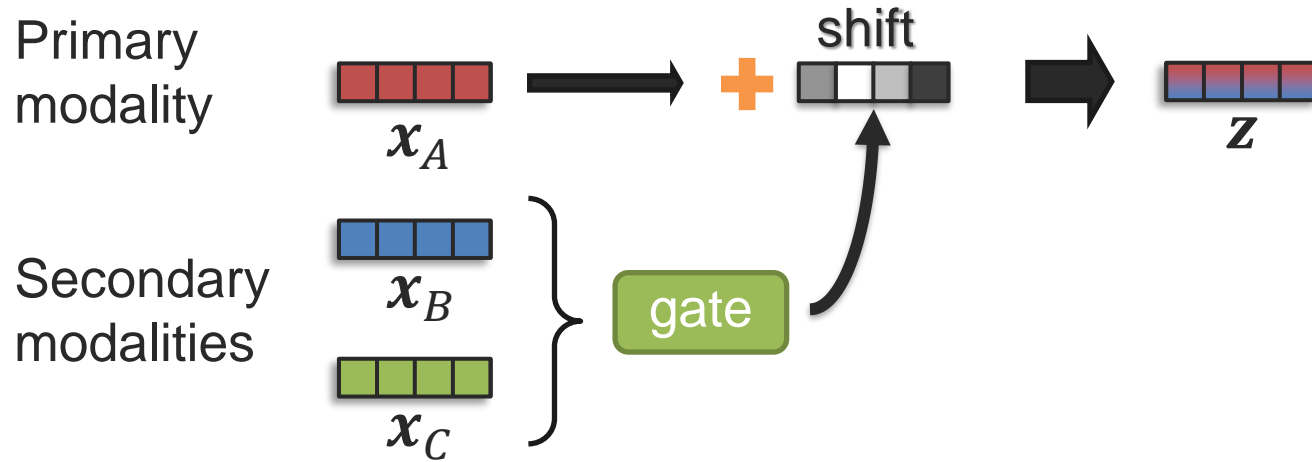
# LXMERT



Tan, Hao, and Mohit Bansal. "Lxmert: Learning cross-modality encoder representations from transformers." *arXiv* (August 20, 2019).
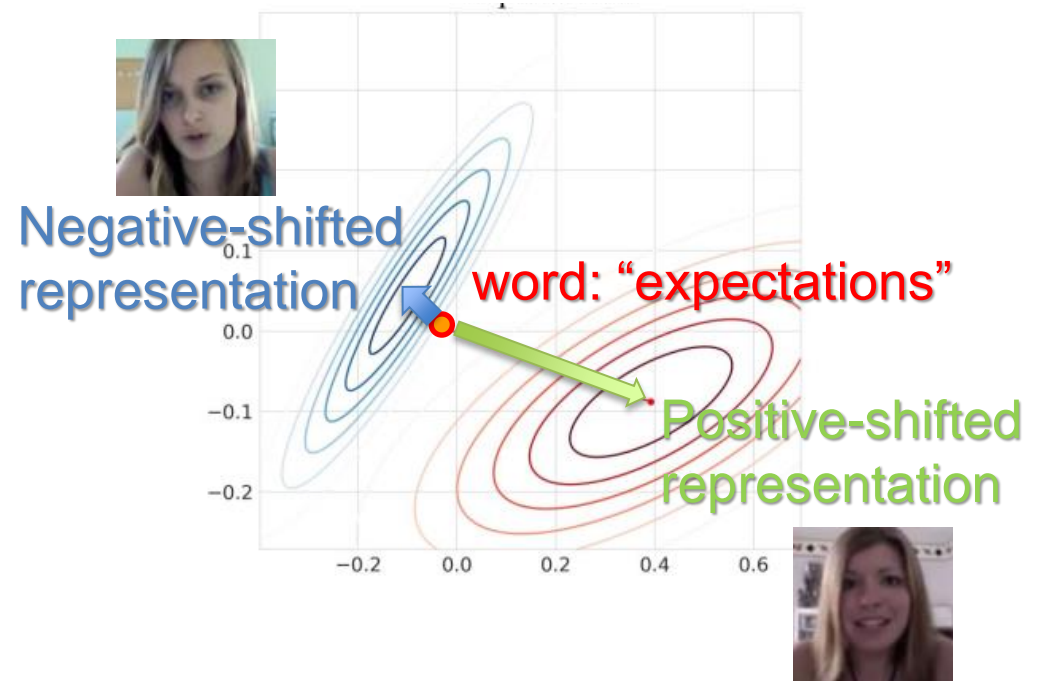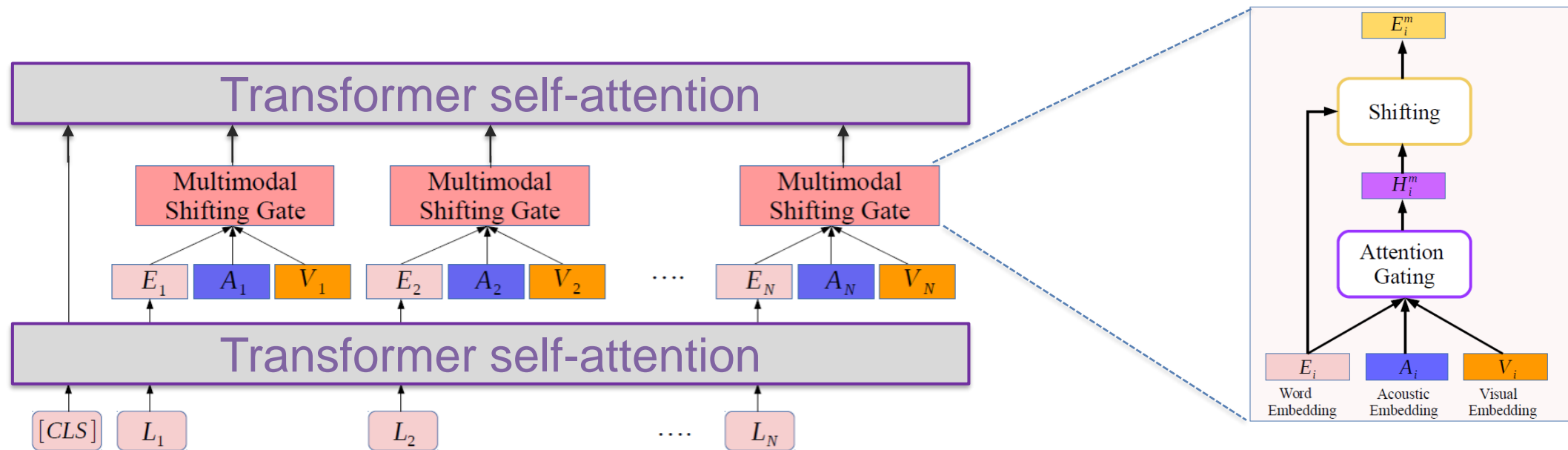
# Reminder: Modality-Shifting Fusion

Primary modality

Secondary modalities

shift

$x_A$

$x_B$

$x_C$

gate

$z$

## Example with language modality:

Primary modality: language

Secondary modalities: acoustic and visual

Negative-shifted representation

word: "expectations"

Positive-shifted representation

Wang et al., Words Can Shift: Dynamically Adjusting Word Representations Using Nonverbal Behaviors, AAAI 2019

# Modality-Shifting with Transformers

Multimodal Adaptation Gate (MAG) + BERT



Rahman et al., Integrating Multimodal Information in Large Pretrained Transformers, ACL 2020

# Video-based Representation and Alignment

## HowTo100M benchmark dataset



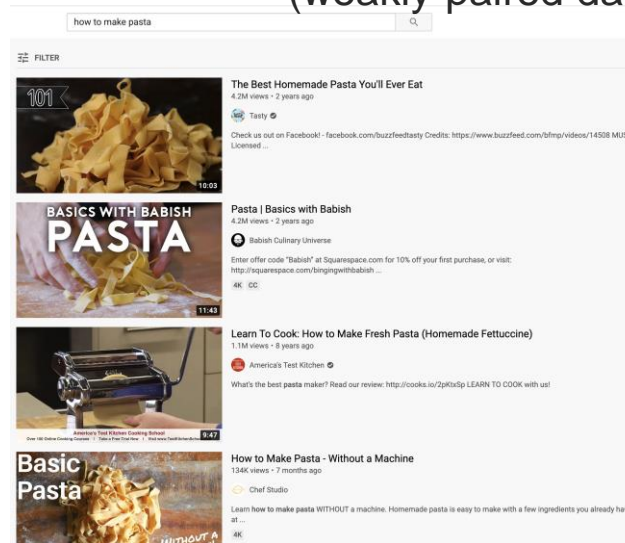| Category | Tasks | Videos | Clips |
|---|---|---|---|
| Food and Entertaining | 11504 | 497k | 54.4M |
| Home and Garden | 5068 | 270k | 29.5M |
| Hobbies and Crafts | 4273 | 251k | 29.8M |
| Cars & Other Vehicles | 810 | 68k | 7.8M |
| Pets and Animals | 552 | 31k | 3.5M |
| Holidays and Traditions | 411 | 27k | 3.0M |
| Personal Care and Style | 181 | 16k | 1.6M |
| Sports and Fitness | 205 | 16k | 2.0M |
| Health | 172 | 15k | 1.7M |
| Education and Communications | 239 | 15k | 1.6M |
| Arts and Entertainment | 138 | 10k | 1.2M |
| Computers and Electronics | 58 | 5k | 0.6M |
| Total | 23.6k | 1.22M | 136.6M |

https://www.di.ens.fr/willow/research/howto100m/

# Visual Representations from Uncurated Instructional Videos

**Goal:** Learn better visual representations…

… by taking advantage of large-scale video+language resources

### Instructional videos
#### (weakly-paired data)



*it's turning into a much thicker mixture*



*The biggest mistake is not kneading it enough*
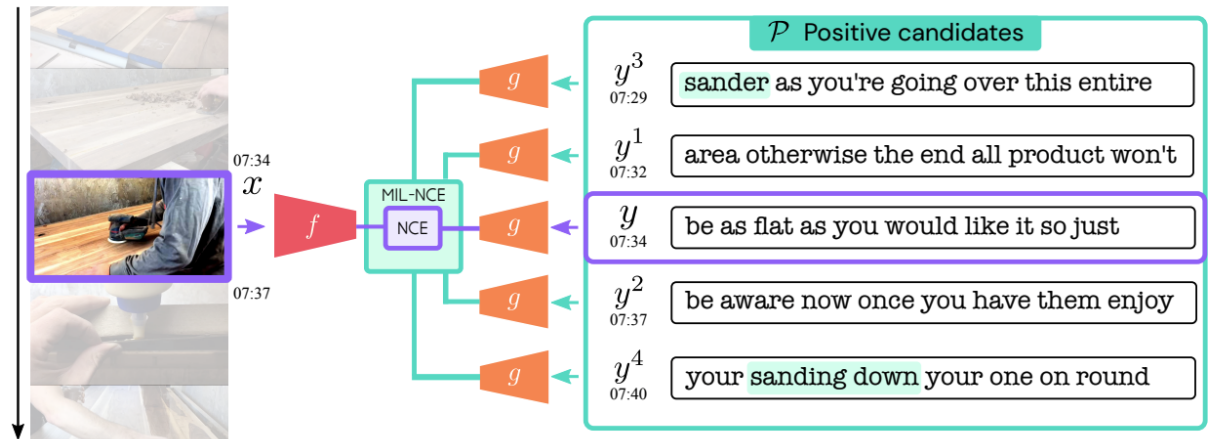


…

End-to-End Learning of Visual Representations from Uncurated Instructional Videos
Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman – CVPR 2020

# Weakly Paired Data

**Data point:** "a short 3.2 seconds video clip (32 frames at 10 FPS) together with a small number of words (not exceeding 16)"
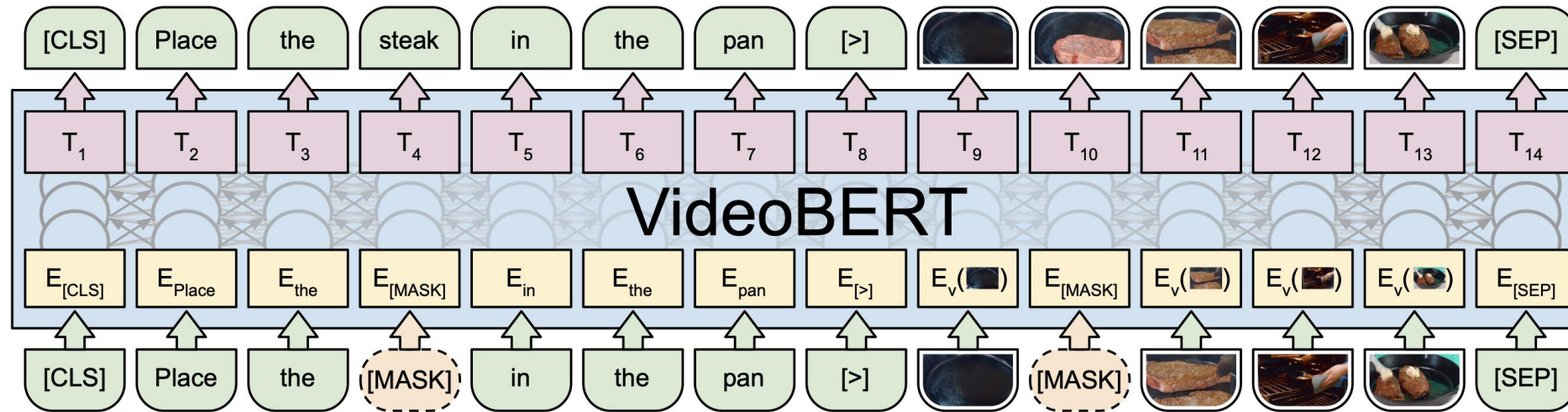


How to handle this misalignment?    Multi-instance learning!

How to do it self-supervised?    Contrastive learning!

End-to-End Learning of Visual Representations from Uncurated Instructional Videos
Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman – CVPR 2020

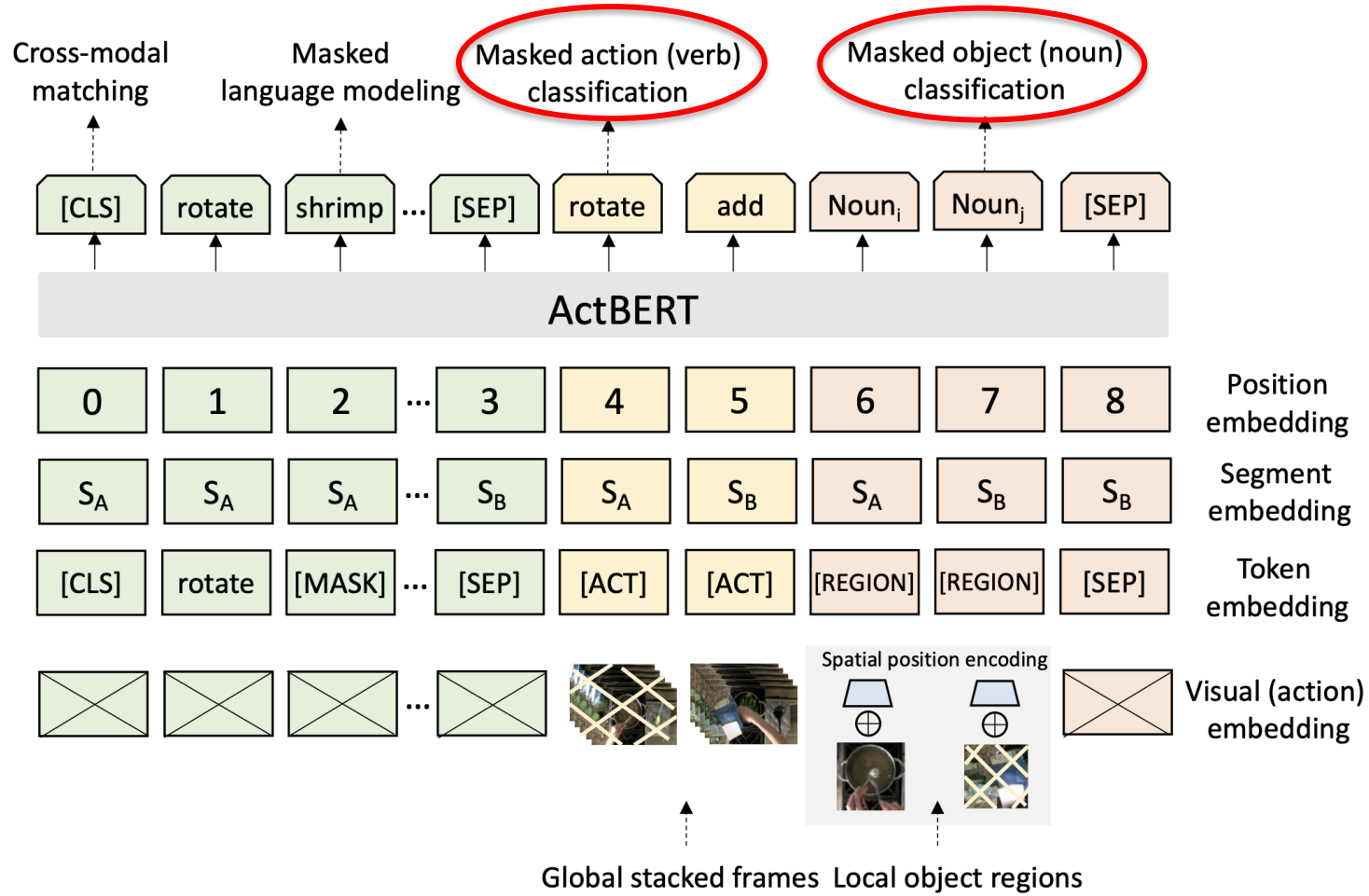# Another Approach for Weakly-Paired Video Data



How do we get visual words now?

**K-mean clustering + centroid**

Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, Cordelia Schmid; VideoBERT: A Joint Model for Video and Language Representation Learning ICCV, 2019
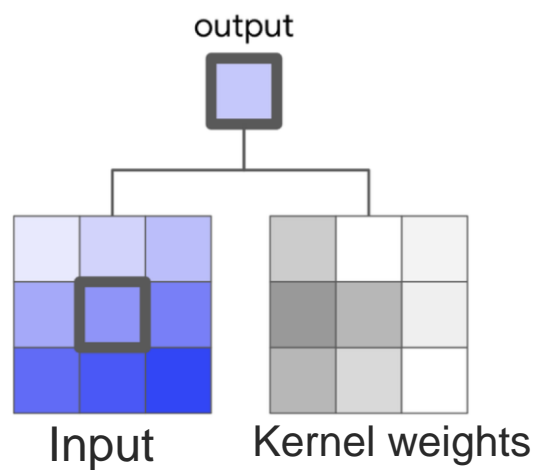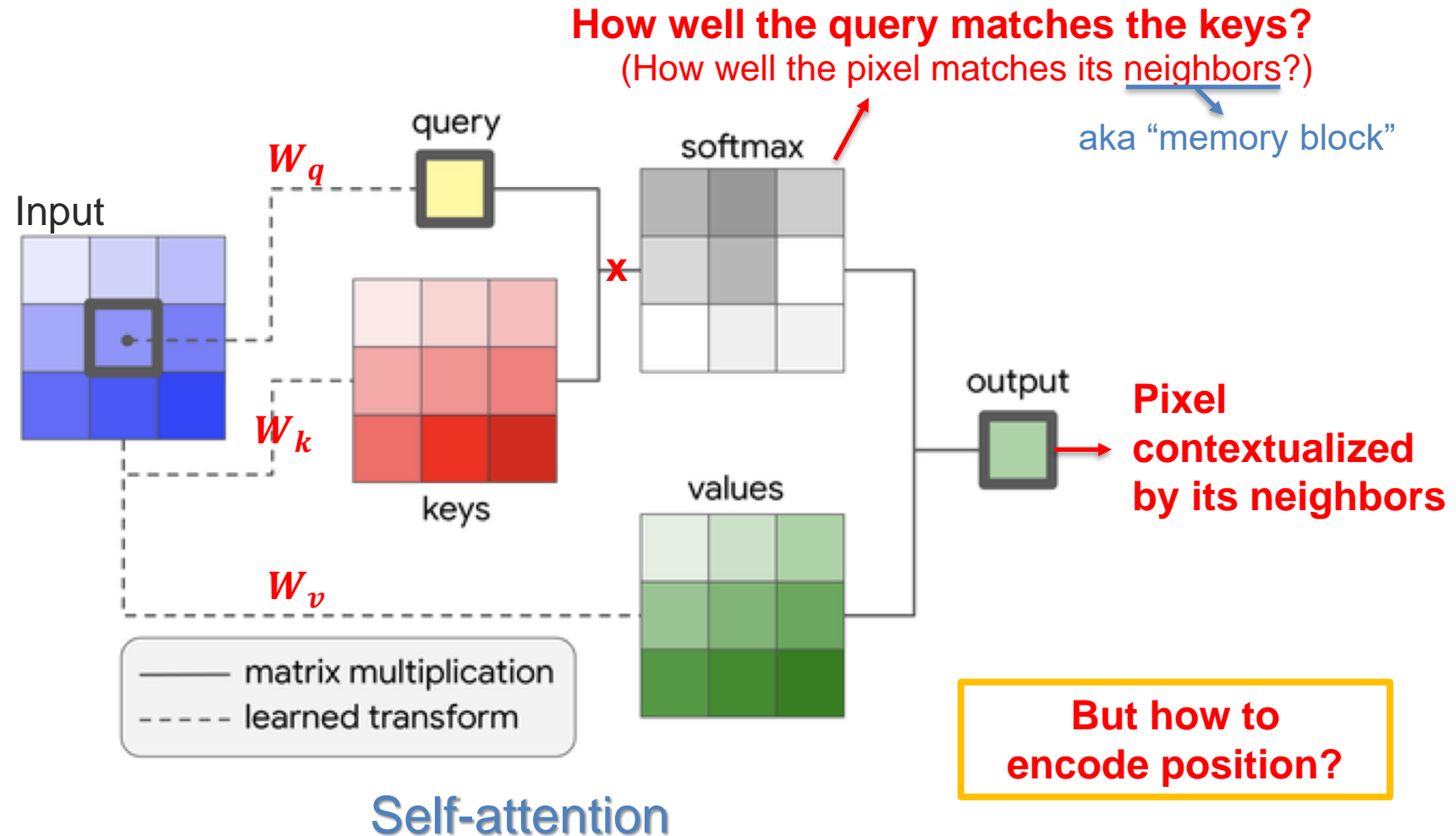
# ActBERT



Cross-modal matching  Masked language modeling  **Masked action (verb) classification**  **Masked object (noun) classification**

| [CLS] | rotate | shrimp | ... | [SEP] | rotate | add | Noun$_i$ | Noun$_j$ | [SEP] |

**ActBERT**

| 0 | 1 | 2 | ... | 3 | 4 | 5 | 6 | 7 | 8 | Position embedding |
| $S_A$ | $S_A$ | $S_A$ | ... | $S_B$ | $S_A$ | $S_B$ | $S_A$ | $S_B$ | $S_B$ | Segment embedding |
| [CLS] | rotate | [MASK] | ... | [SEP] | [ACT] | [ACT] | [REGION] | [REGION] | [SEP] | Token embedding |

Spatial position encoding

Visual (action) embedding

Global stacked frames   Local object regions

Zhu and Yang, ActBERT: Learning Global-Local Video-Text Representations, CVPR 2020

# Going Beyond CNNs… Vision Transformers (and more!)

# Replacing a CNN w/ Self-Attention



How well the query matches the keys?
(How well the pixel matches its neighbors?)

aka "memory block"

Convolution

Self-attention

$W_q$   $W_k$   $W_v$

Pixel contextualized by its neighbors

But how to encode position?

https://arxiv.org/abs/1906.05909

# Replacing a CNN w/ Self-Attention

Image patch



2D relative position embedding



Position embedding is added to the key:

$$y_{ij} = \sum_{a,b \in \mathcal{N}_k(i,j)} \texttt{softmax}_{ab} \left( q_{ij}^{\top} k_{ab} + q_{ij}^{\top} r_{a-i,b-j} \right) v_{ab}$$
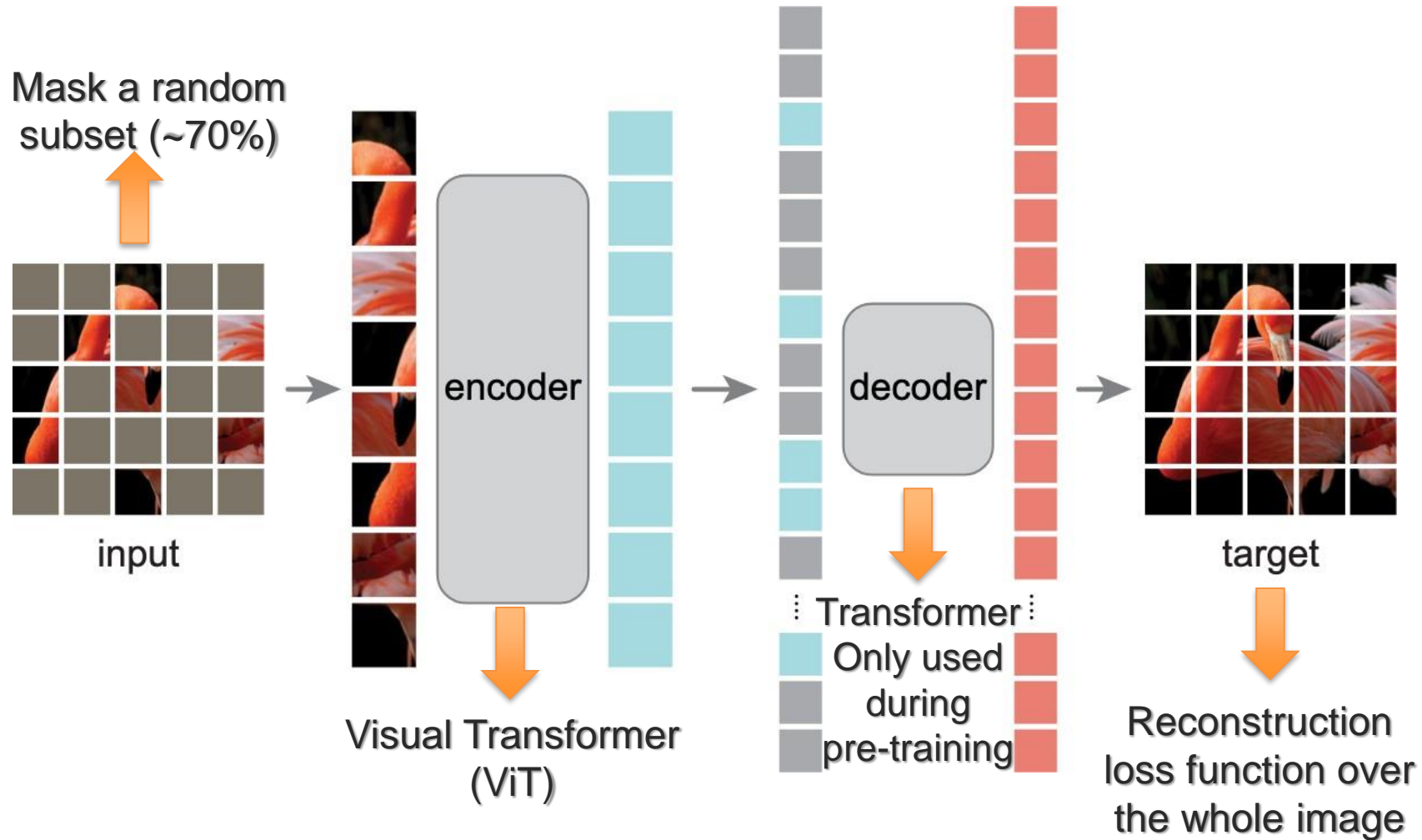
Carnegie Mellon University

# Vision Transformer (ViT)

Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv* (2020).
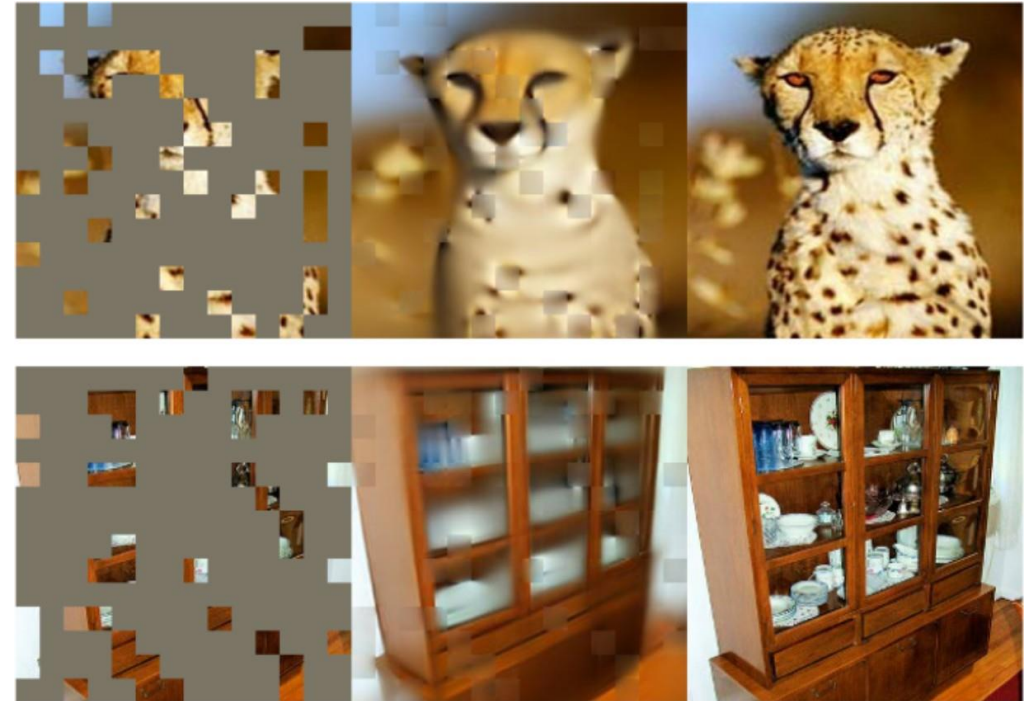
# Vision Transformer (ViT)



Embedding for the whole image
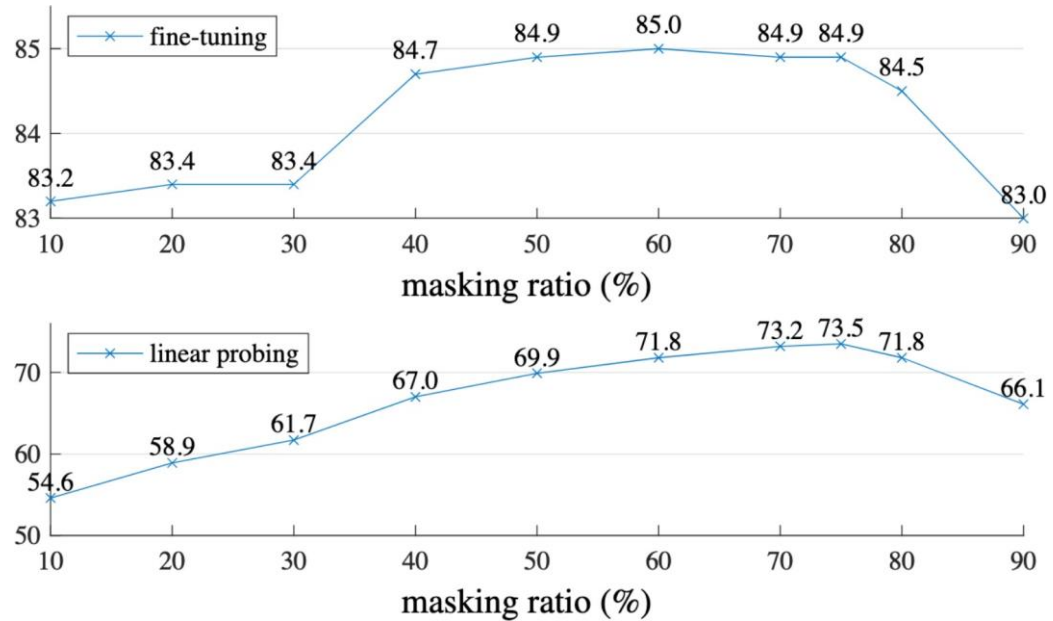
16x16 image patches

Flattening the image patches

Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv* (2020).
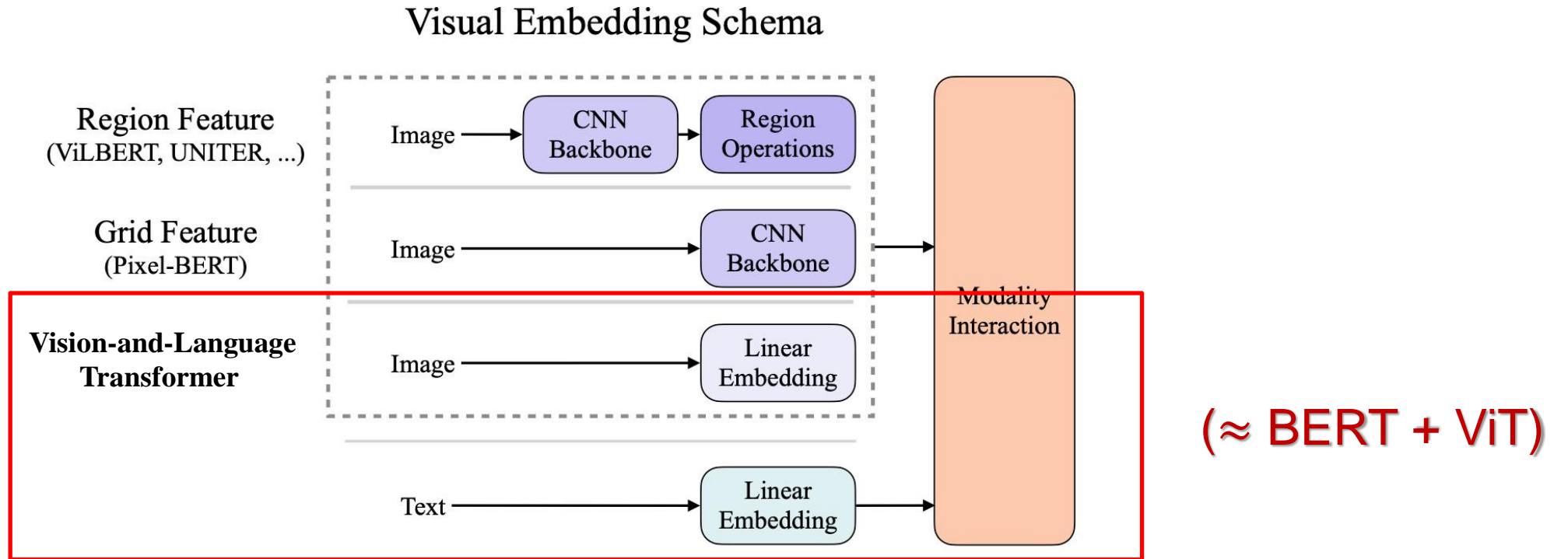
# Masked Auto-Encoder (MAE)



He et al., Masked Autoencoders Are Scalable Vision Learners, CVPR 2022
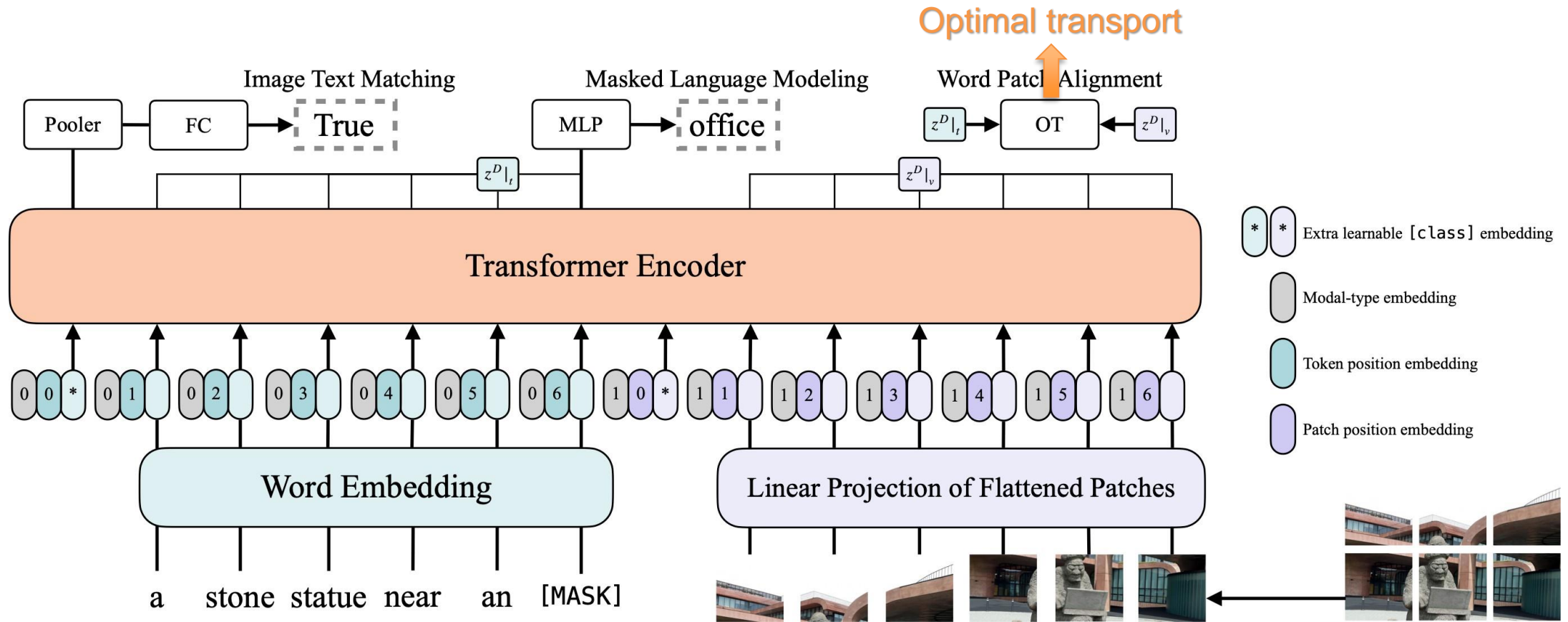
# Masked Auto-Encoder (MAE)



He et al., Masked Autoencoders Are Scalable Vision Learners, CVPR 2022

# Visual Transformers for Multimodal Learning



Visual Embedding Schema

(≈ BERT + ViT)

https://arxiv.org/abs/2102.03334

# Visual-and-Language Transformer (ViLT) (≈ BERT + ViT)

# Visual-and-Language Transformer (ViLT)

Example of alignment between modalities:

https://arxiv.org/abs/2102.03334

# ALBEF: Align Before Fusion (≈ BERT + ViT + CLIP-ish)



Li et al., Align before Fuse: Vision and Language Representation Learning with Momentum Distillation, Neurips 2021