



多模态技术在学术界的研究现状以及 在工业界的落地应用

付杰

jem.fu@jina.ai

AI Engineer@Jina AI



About Jina AI



Jina AI 成立于 2020 年 2 月，是一家业内领先的开源企业，致力于通过人工智能和深度学习技术，打造下一代云原生神经搜索框架，帮助开发者和企业简化非结构化数据的搜索难题。

Jina AI 总部位于德国柏林，在北京、深圳、巴塞罗那等地均设有办公室，海外员工人数比例超过三分之二。

2021 年，Jina AI 完成了由 Canaan Partners 领投的 3000 万美元 A 轮融资，截至目前融资总额达 3900 万美元。





DocArray

将非结构化数据，统一成同一种数据结构



CLIP-as-service

基于 CLIP 的图像和文本跨模态编码服务



NOW

图像搜索的无代码解决方案



Jina

适用于所有数据类型的云原生神经搜索框架



Finetuner

微调深度神经网络，优化神经搜索任务中的向量



Hub

分享和发现神经搜索应用的可复用组件

**Creator of
Neural Search**

**Contributor to
Open Source**

1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

3

多模态技术在工业界的应用

4

Jina在多模态场景下的应用

多模态数据VS单模态数据

- 什么是单模态 / 多模态数据？

我们生活的世界是一个多模态的世界，包含丰富的文本、视觉、听觉、嗅觉信息，通常我们所接触的物体都不只有一个模态，例如短视频、电商物品等



- 多模态数据处理起来有哪些难点？

1. **表征**: 利用多模态的互补性和冗余性的特质来表示和总结多模态数据
2. **翻译**: 从一种模式转换为另一种模式
3. **对齐**: 确定来自两种或两种以上不同模式的(分)要素之间的直接关系
4. **融合**: 将来自两个或多个模式的信息连接起来
5. **协同学习**: 从一种模式的知识学习可以帮助一个计算模型训练另一个不同的模态

多模态数据VS单模态数据

- 为什么要研究多模态数据？

图表3：中国短视频行业产业链全景图



哪些行业涉及到多模态数据？

- 短视频数据：
 - 如何理解用户上传的视频
 - 如何实现内容的精准/多样性分发/推荐
- 电商搜索：
 - 如何更好的理解商品
 - 针对用户搜索历史如何更好地推荐商品

1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

- 多模态表征

- 多模态预训练

3

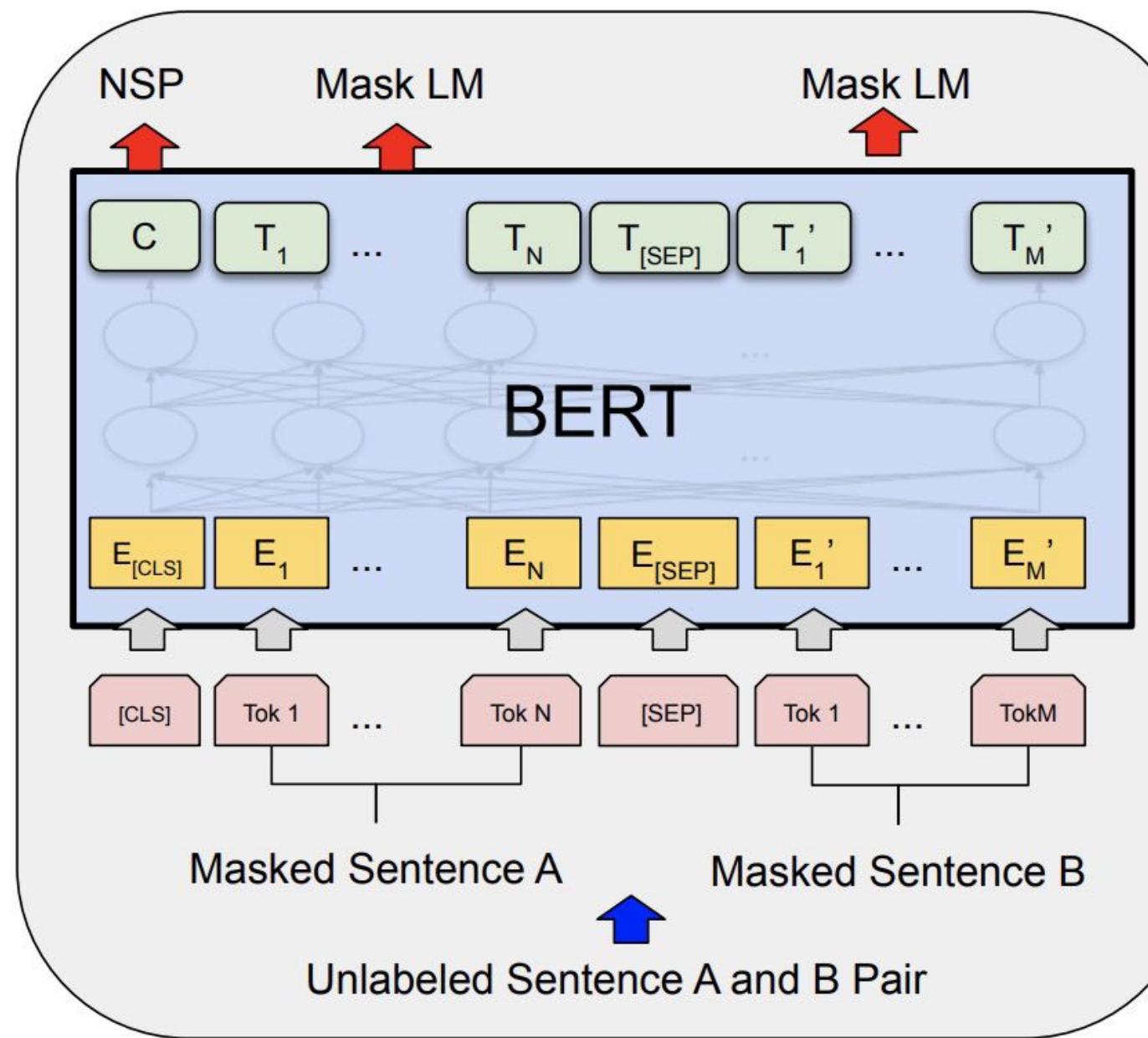
多模态技术在工业界的应用

4

Jina在多模态场景下的应用

多模态表征

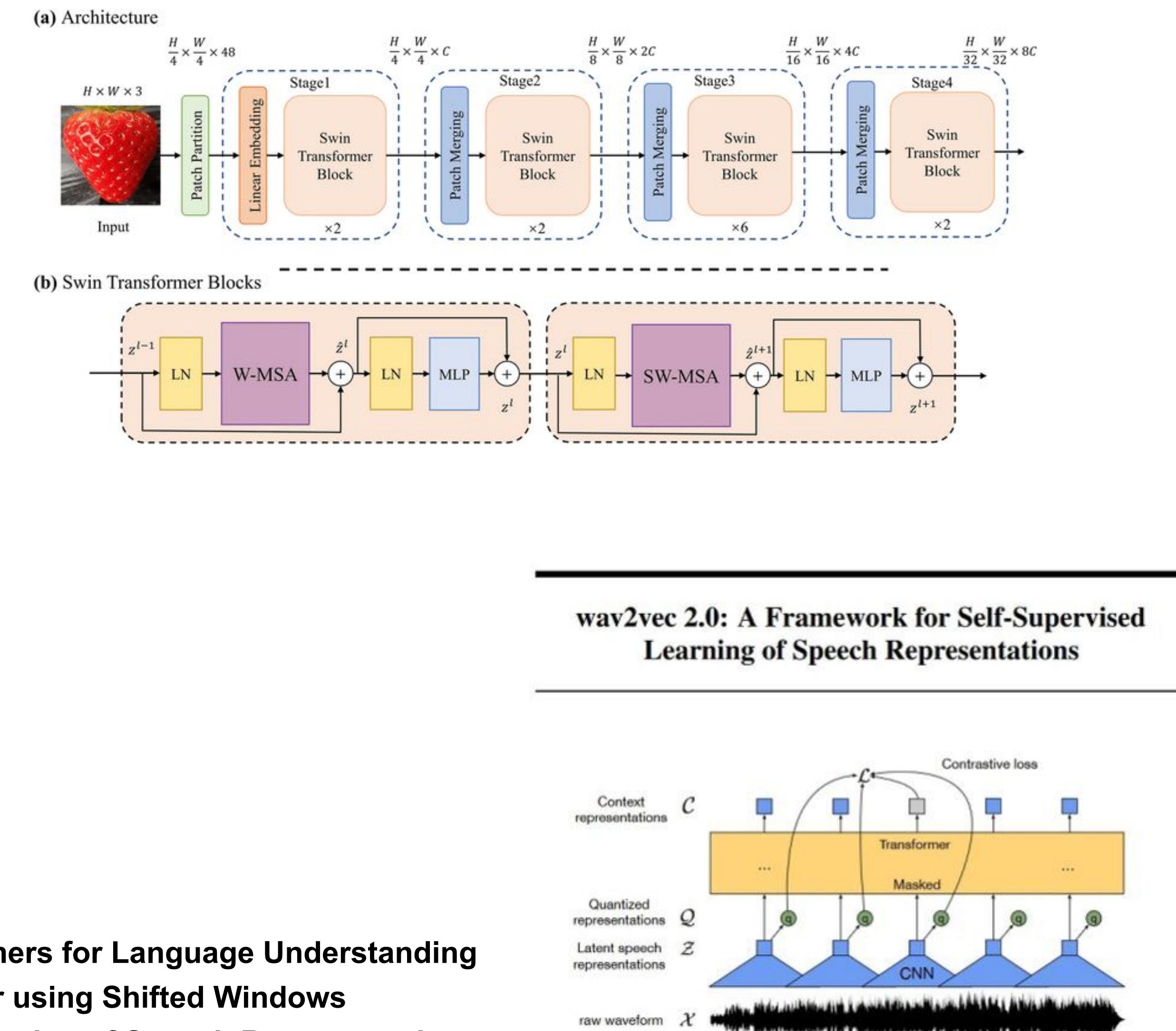
- 单模态表征



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations



多模态表征

- 多模态的特征对齐:引入传统统计学习方法进行模态对齐

Canonical Correlation Analysis: 将两组数据降维到一维后, 其相关系数最大

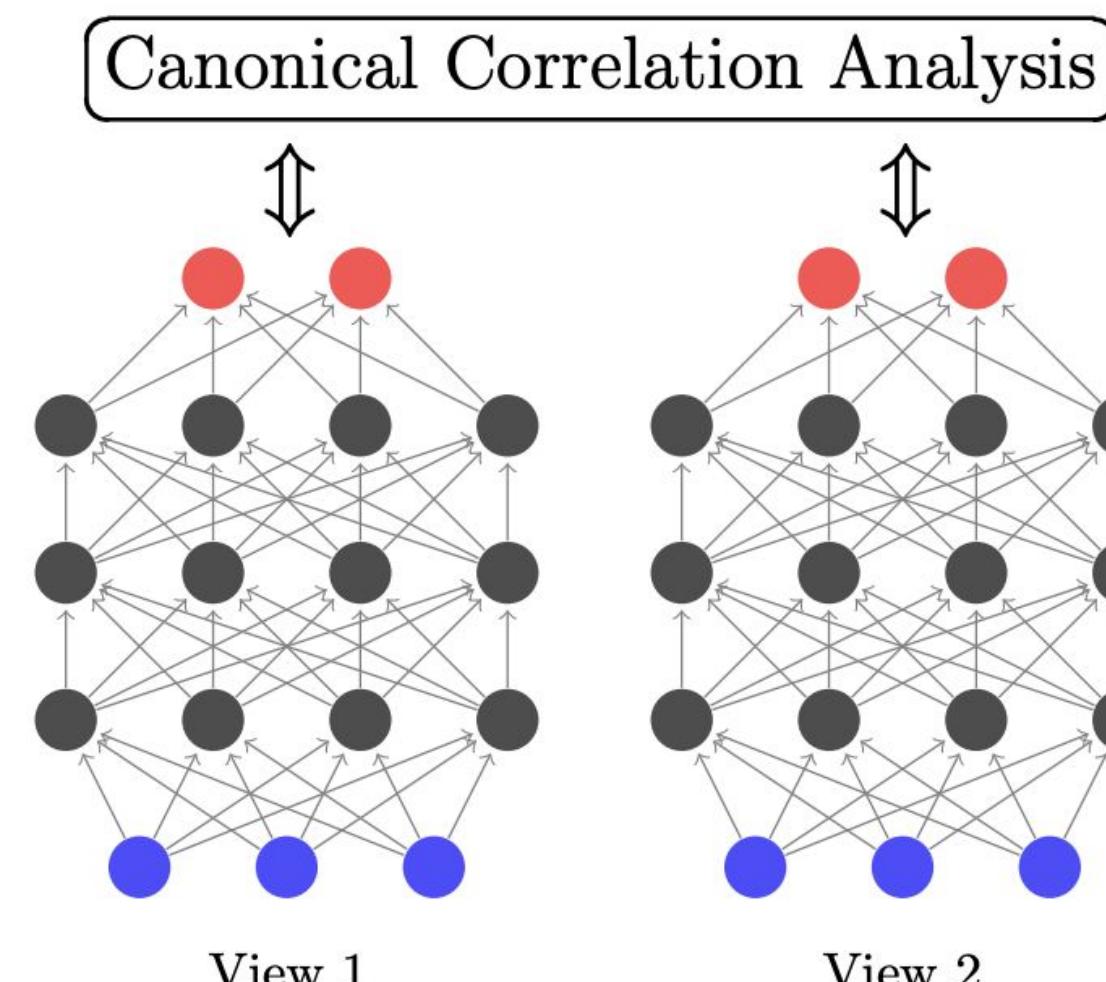
$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{D(X)} \sqrt{D(Y)}}$$

$$X' = a^T X, Y' = b^T Y$$

$$\xrightarrow{\quad} \underbrace{\arg \max_{a,b}}_{a^T S_{XX} a = 1, b^T S_{YY} b = 1} \frac{a^T S_{XY} b}{\sqrt{a^T S_{XX} a} \sqrt{b^T S_{YY} b}}$$

$$\xrightarrow{\quad} \underbrace{\arg \max_{a,b}}_{a^T S_{XX} a = 1, b^T S_{YY} b = 1} a^T S_{XY} b$$

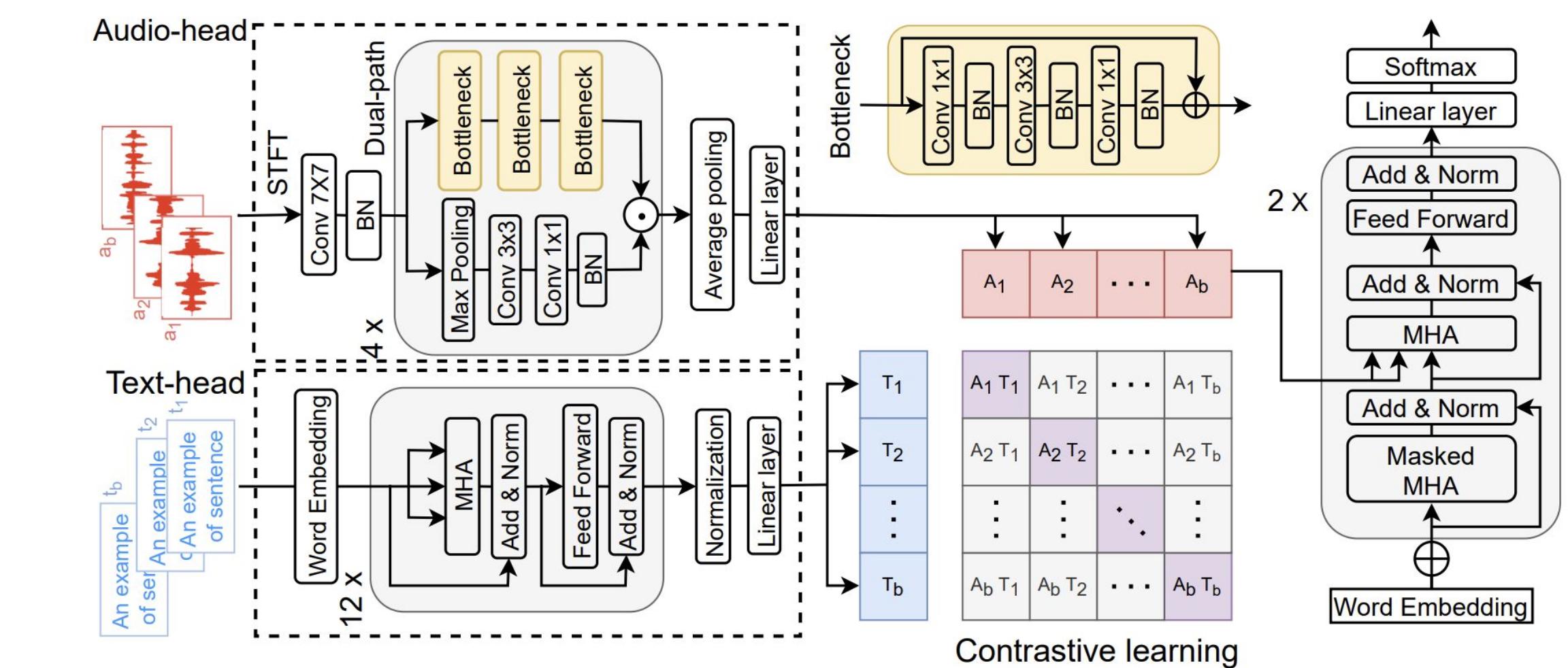
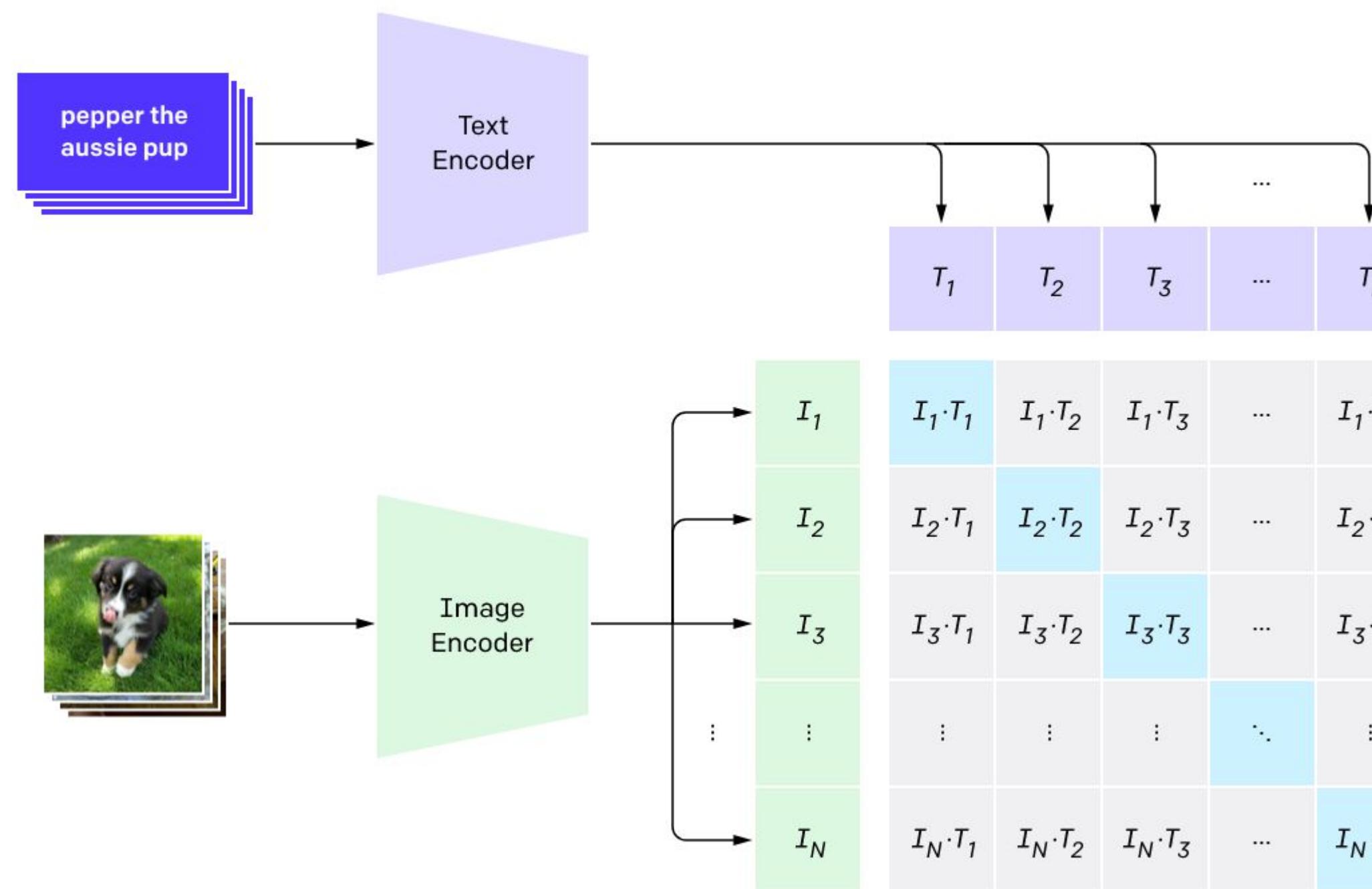
Deep Canonical Correlation Analysis: 在深度网络后加入CCA,
使得网络参数能够自动对相同类别但不同view的数据学习到对齐特征



多模态表征

- 多模态的特征对齐: 基于对比学习进行模态对齐

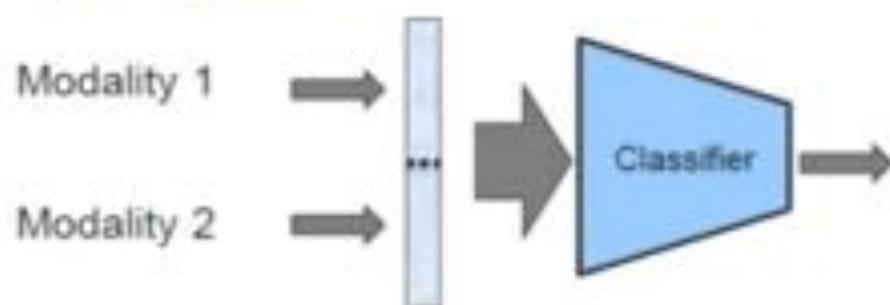
1. Contrastive pre-training



多模态表征

- 多模态的特征融合: 基于矩阵/注意力的方法

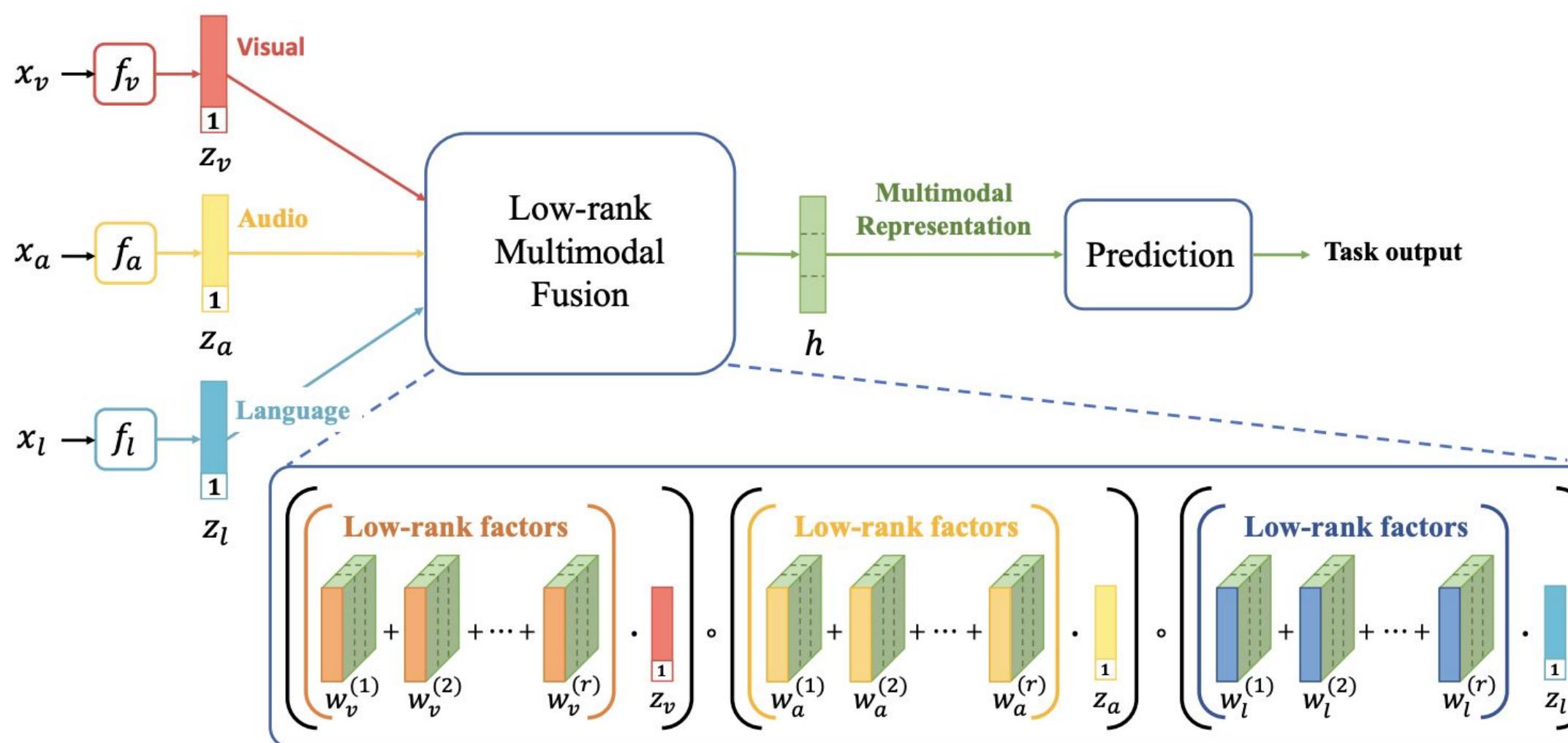
1) Early Fusion



2) Late Fusion



基于矩阵方法



Efficient Low-rank Multimodal Fusion with Modality-Specific Factors

基于注意力的方法

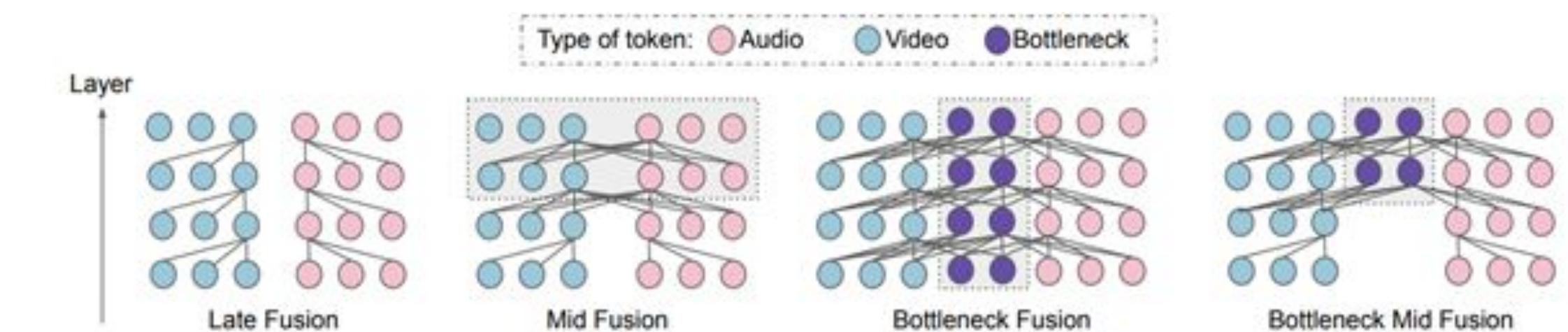


Figure 1: **Cross-modal Fusion.** Unlike late fusion (left), where no cross-modal information is exchanged in the model until after the classifier, we investigate two pathways for the exchange of cross-modal information. The first is via standard pairwise self attention across all hidden units in a layer, but applied only to later layers in the model – mid fusion (middle, left). We also propose the use of ‘fusion bottlenecks’ (middle, right) that restrict attention flow within a layer through tight latent units. Both forms of restriction can be applied in conjunction (Bottleneck Mid Fusion) for optimal performance (right). We show $B = 2$ bottleneck units and 3 hidden units per modality. Grey boxes indicate tokens that receive attention flow from both audio and video tokens.

Attention Bottlenecks for Multimodal Fusion

1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

- 多模态表征
- 多模态预训练

3

多模态技术在工业界的应用

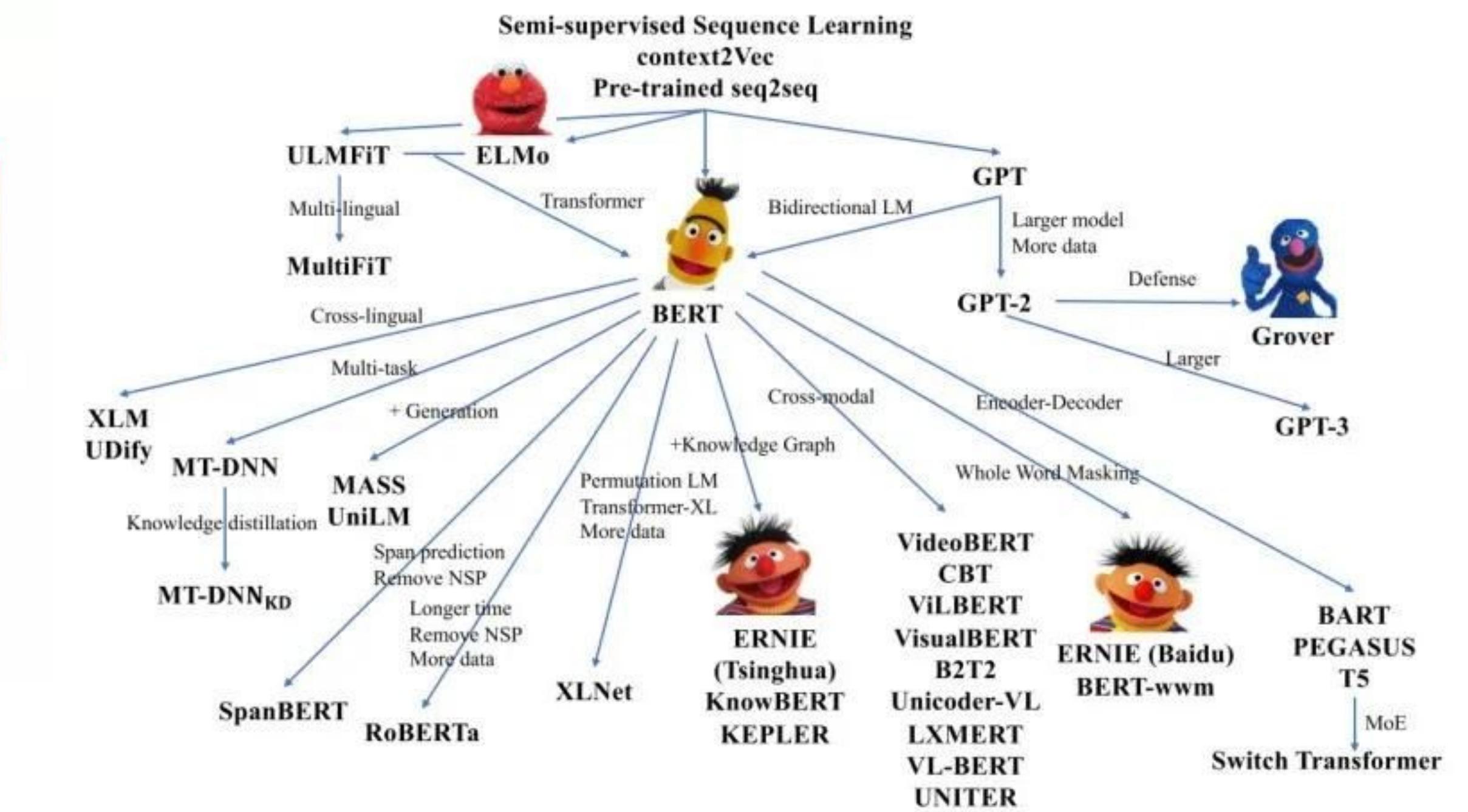
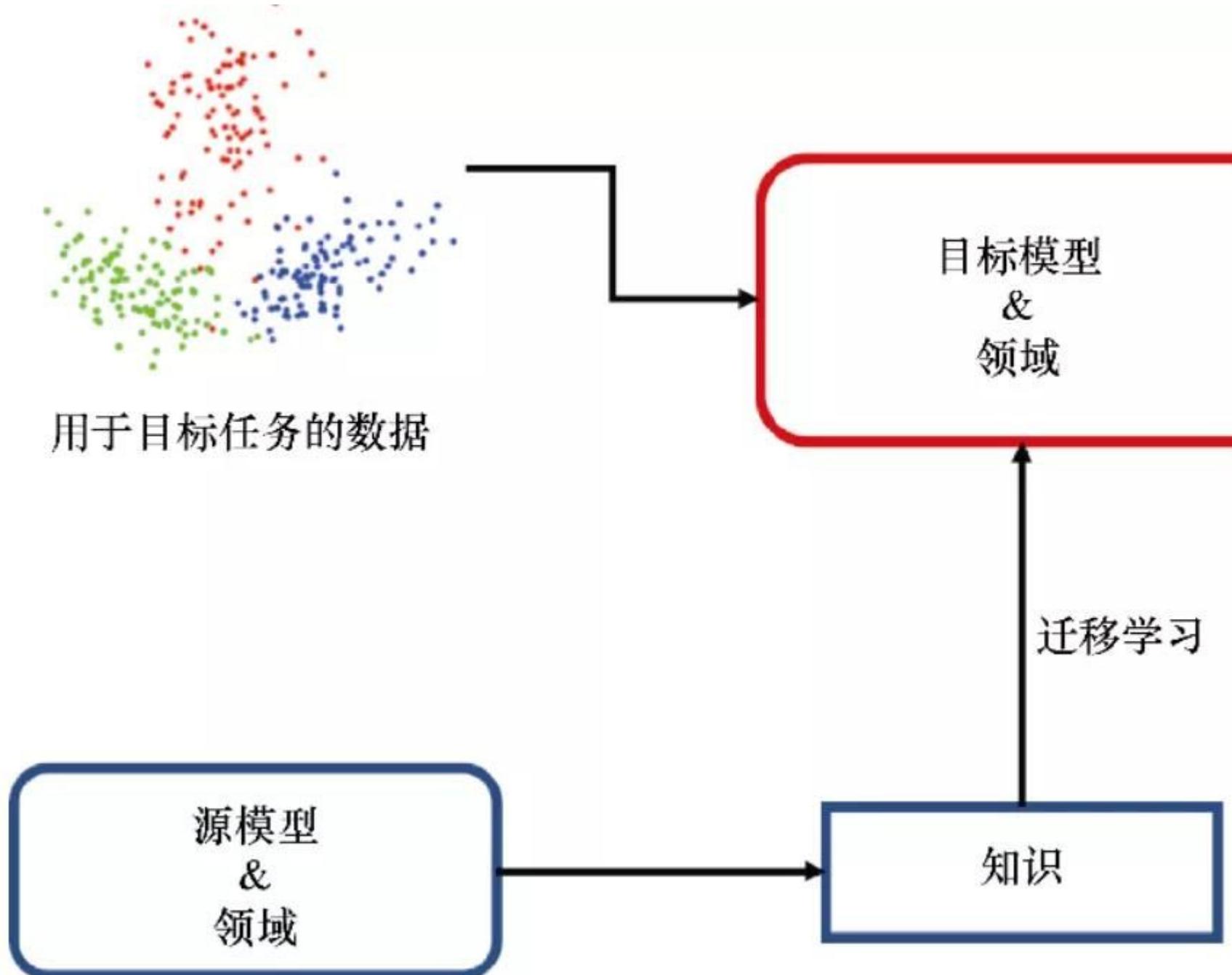
4

Jina在多模态场景下的应用

多模态预训练

- 什么是预训练

把人类的语言知识，先学了一个东西，然后再代入到某个具体任务，就顺手了



多模态预训练

- 常用预训练任务

Masked Language Modeling (MLM)

$$\mathcal{L}_{\text{MLM}} = -\mathbb{E}_{(\mathcal{W}, \mathcal{V}) \in \mathcal{D}} \log P_{\theta} (w_m | w_{\setminus m}, \mathcal{V}), \quad (1)$$

where $w_m, w_{\setminus m}$ represent the masked tokens and unmasked tokens, respectively, and $(\mathcal{W}, \mathcal{V}) \in \mathcal{D}$ represents a text \mathcal{W} and an image \mathcal{V} sampled from dataset \mathcal{D} .

Masked Region Prediction (MRP)

- Masked Region Classification (MRC)
- Masked Region Feature Regression (MRFR)

Image-Text Matching (ITM)

Cross-Modal Contrastive Learning (CMCL)

$$\mathcal{L}_{\text{MRC}} = \mathbb{E}_{(\mathcal{W}, \mathcal{V}) \in \mathcal{D}} \sum_{i=1}^m \text{CE} (\text{softmax}(\text{FC}(h_{v_i}), c(v_i))), \quad (2)$$

where m is the length of image features; $c(v_i)$ represents the true label of the masked region, such as the object detection output or the (pre-defined) visual tokens.

$$\mathcal{L}_{\text{MRFR}} = \mathbb{E}_{(\mathcal{W}, \mathcal{V}) \in \mathcal{D}} \sum_{i=1}^m \|\text{FC}(h_{v_i}) - \hat{E}_v(v_i)\|^2. \quad (3)$$

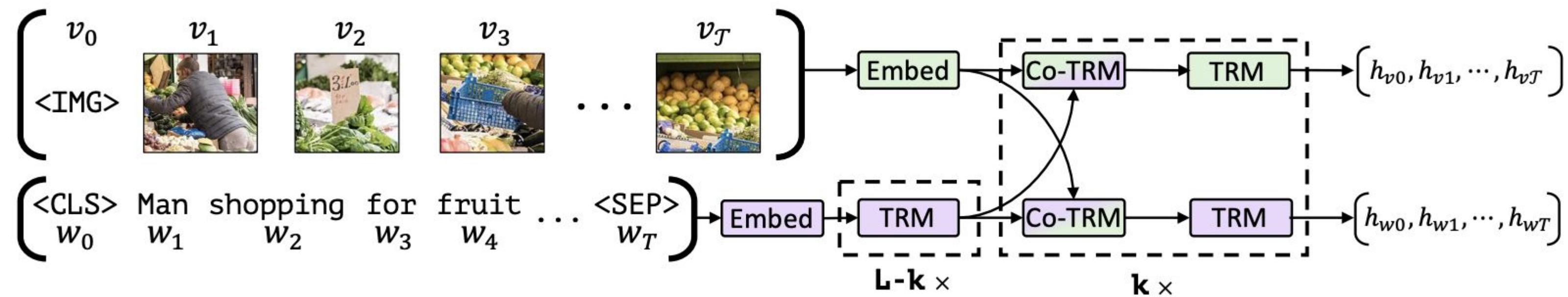
$$\begin{aligned} \mathcal{L}_{\text{VLM}} = & -\mathbb{E}_{(\mathcal{W}, \mathcal{V}) \in \mathcal{D}} [y \log s_{\theta}(\mathcal{W}, \mathcal{V}) \\ & + (1 - y) \log (1 - s_{\theta}(\mathcal{W}, \mathcal{V}))]. \end{aligned} \quad (4)$$

$$\mathcal{L}_{\text{i2t}} = -\frac{1}{n} \sum_{i=1}^N \log \frac{\exp(x_i^\top y_i / \sigma)}{\sum_{j=1}^N \exp(x_i^\top y_j / \sigma)}, \quad (5)$$

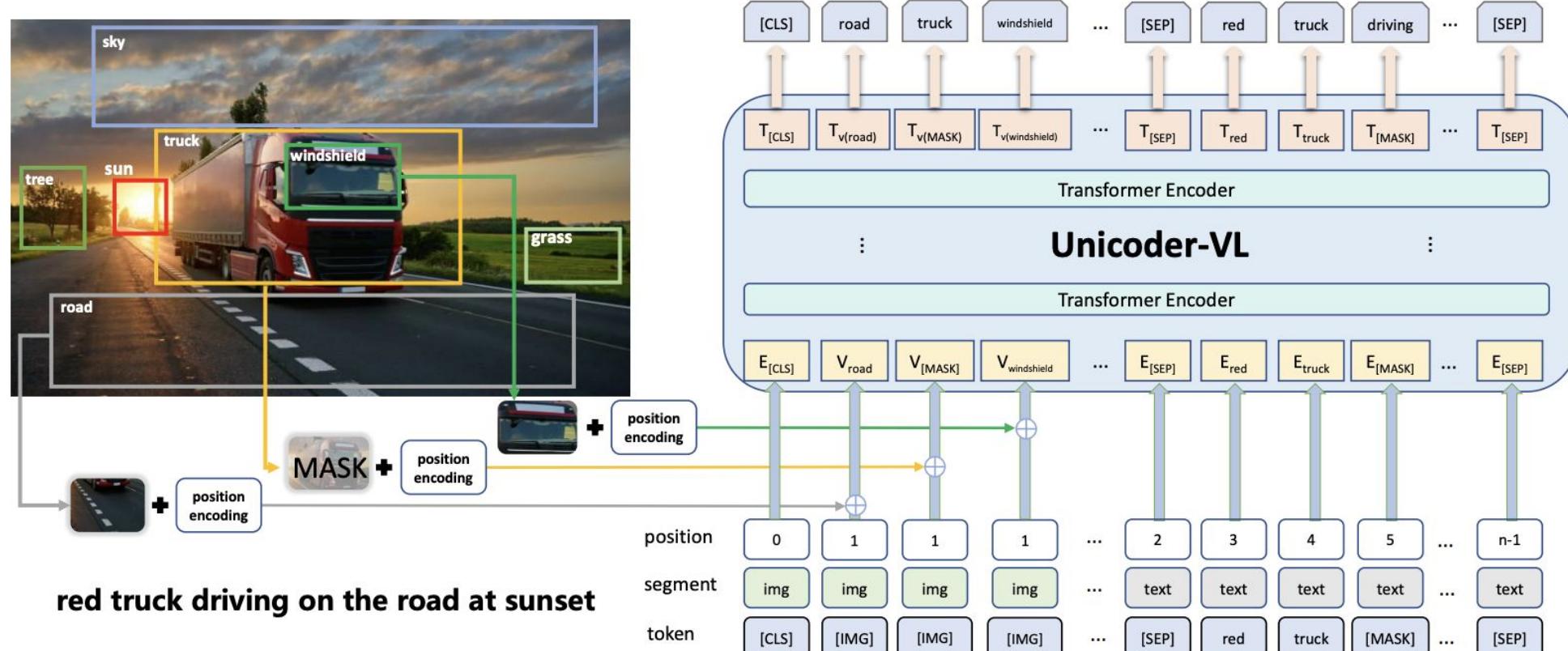
多模态预训练

- 经典的双塔/单流模型

Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks



A Universal Encoder for Vision and Language by Cross-modal Pre-training

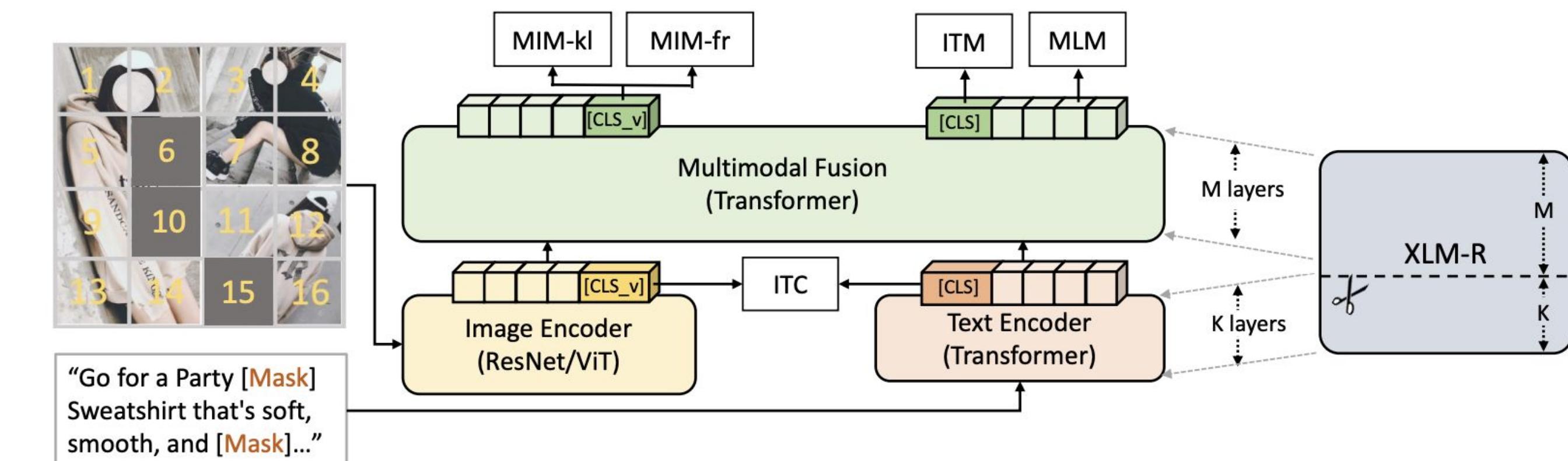
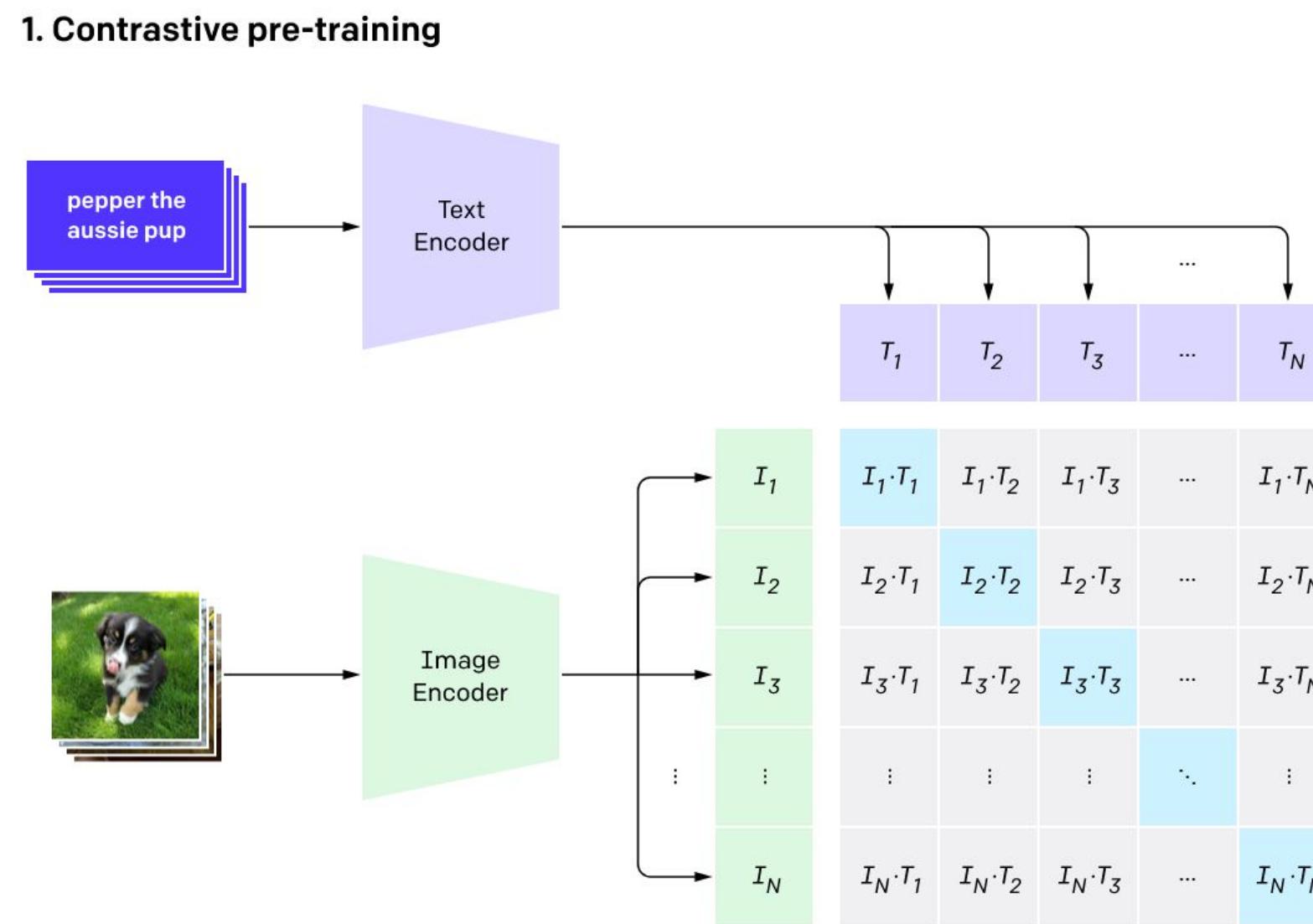


双塔模型无需计算cross-attention,
速度更快, 能满足线上需求, 但性能相
对较单流模型更低

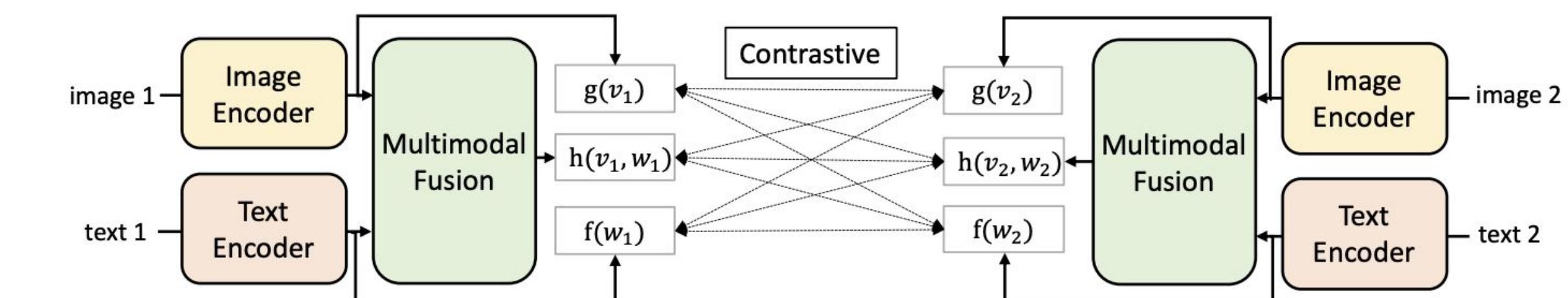
多模态预训练

- 基于对比学习的预训练模型

CommerceMM: Large-Scale Commerce MultiModal Representation Learning with Omni Retrieval



(a) CommerceMM on the input (image, text) pair and 5 image-text pre-training tasks.



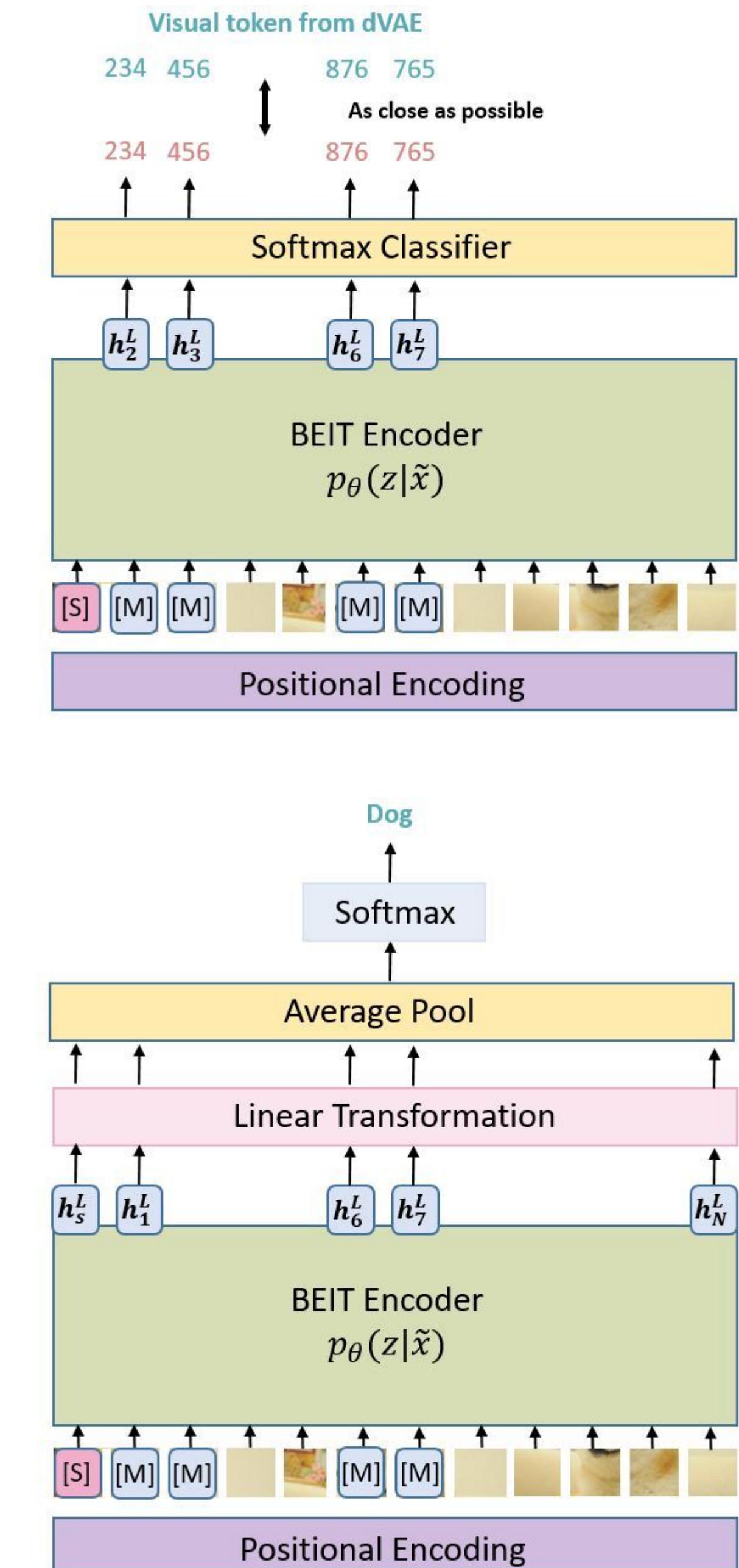
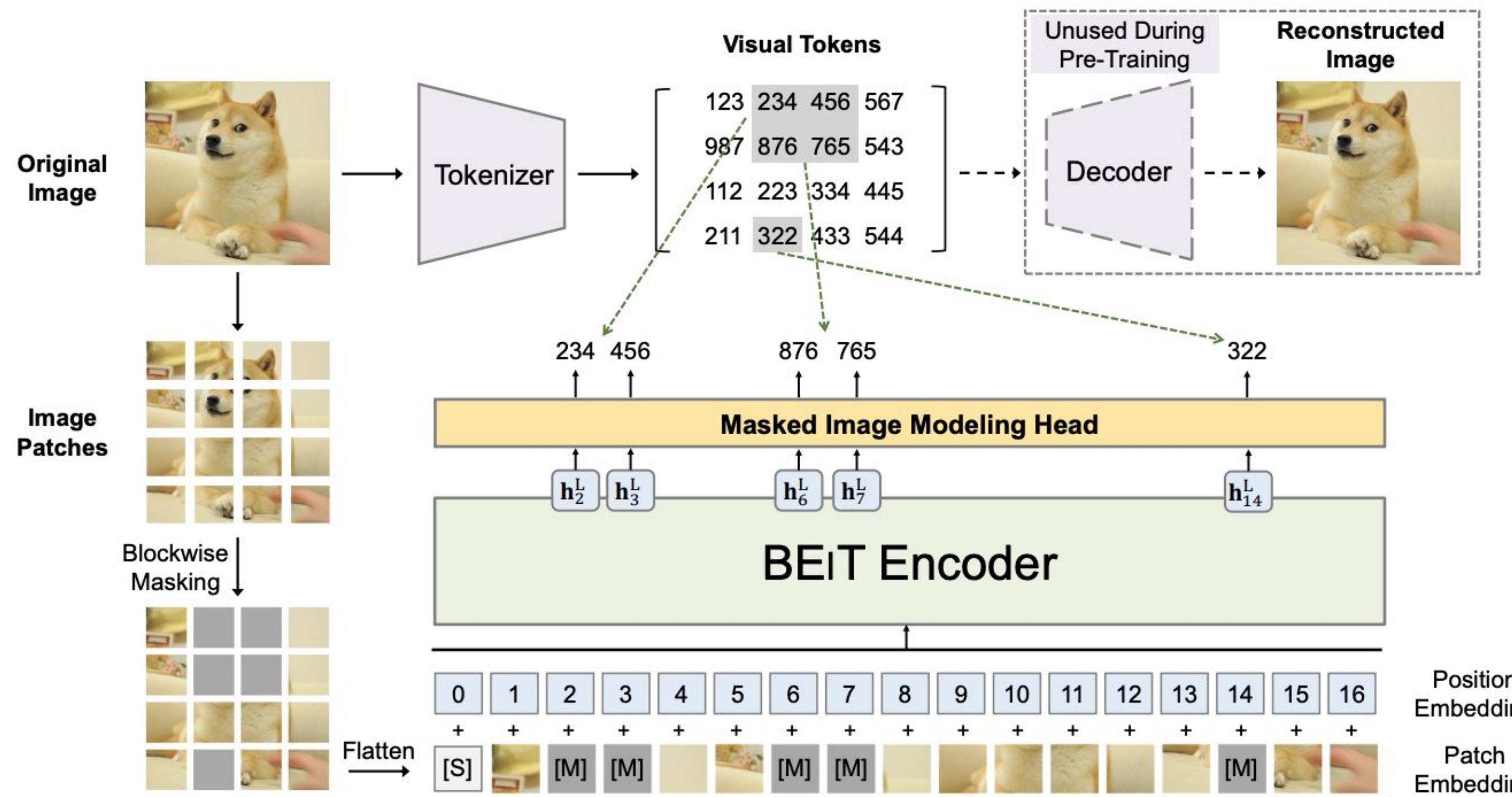
(b) Omni-Retrieval Pre-training on the cross-modal and cross-pair data.

Figure 2: CommerceMM Model Architecture with the image-text pre-training and omni-retrieval tasks.

多模态预训练

- 基于predict masked patch的预训练模型

BEIT: BERT Pre-Training of Image Transformers



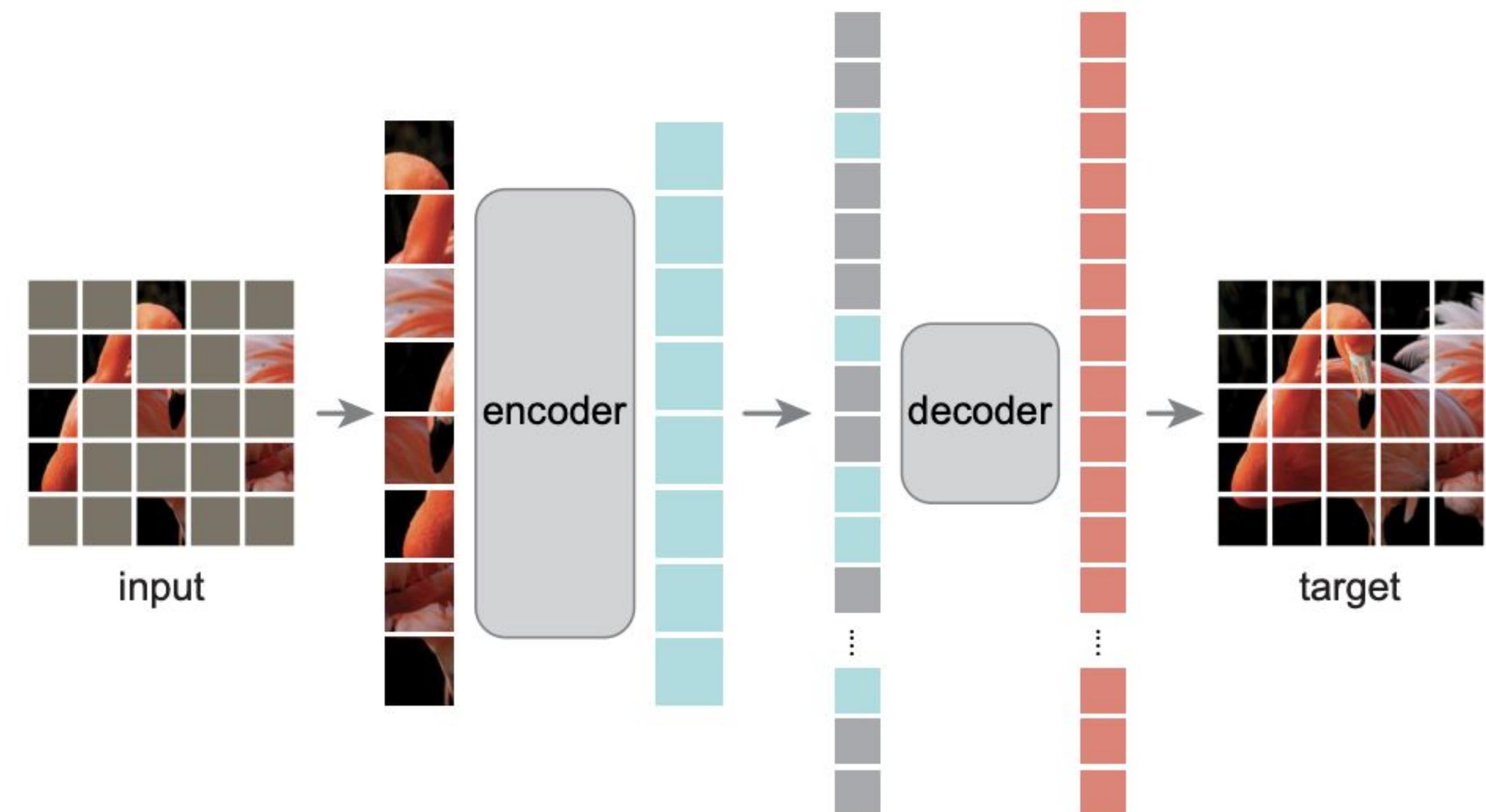
多模态预训练

- 基于**predict masked patch**的预训练模型

Masked Autoencoders Are Scalable Vision Learners

掩码自编码在视觉和文本上有什么不同？

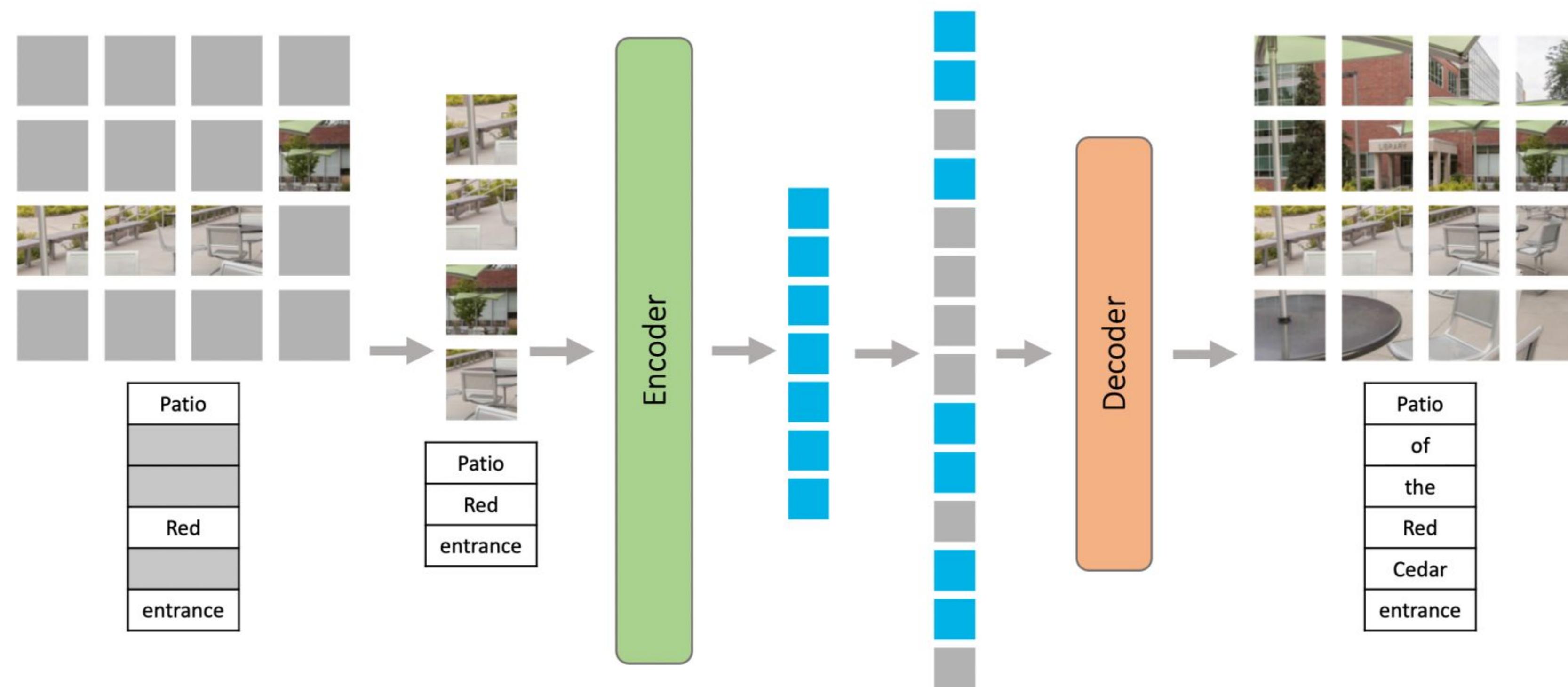
1. 模型结构有差异:transformer vs. conv
2. 信息密度差异较大:15% vs. 75%
3. 解码目标不一致:token vs. pixel



多模态预训练

- 基于predict masked patch的预训练模型

Multimodal Masked Autoencoders Learn Transferable Representations



1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

3

多模态技术在工业界的应用

4

Jina在多模态场景下的应用

多模态技术在工业界的应用

- 爱奇艺/优酷——视频搜索，视频内容理解

基于视频内容召回



基于标题以出现召回不足



基于视频图像以下信息召回

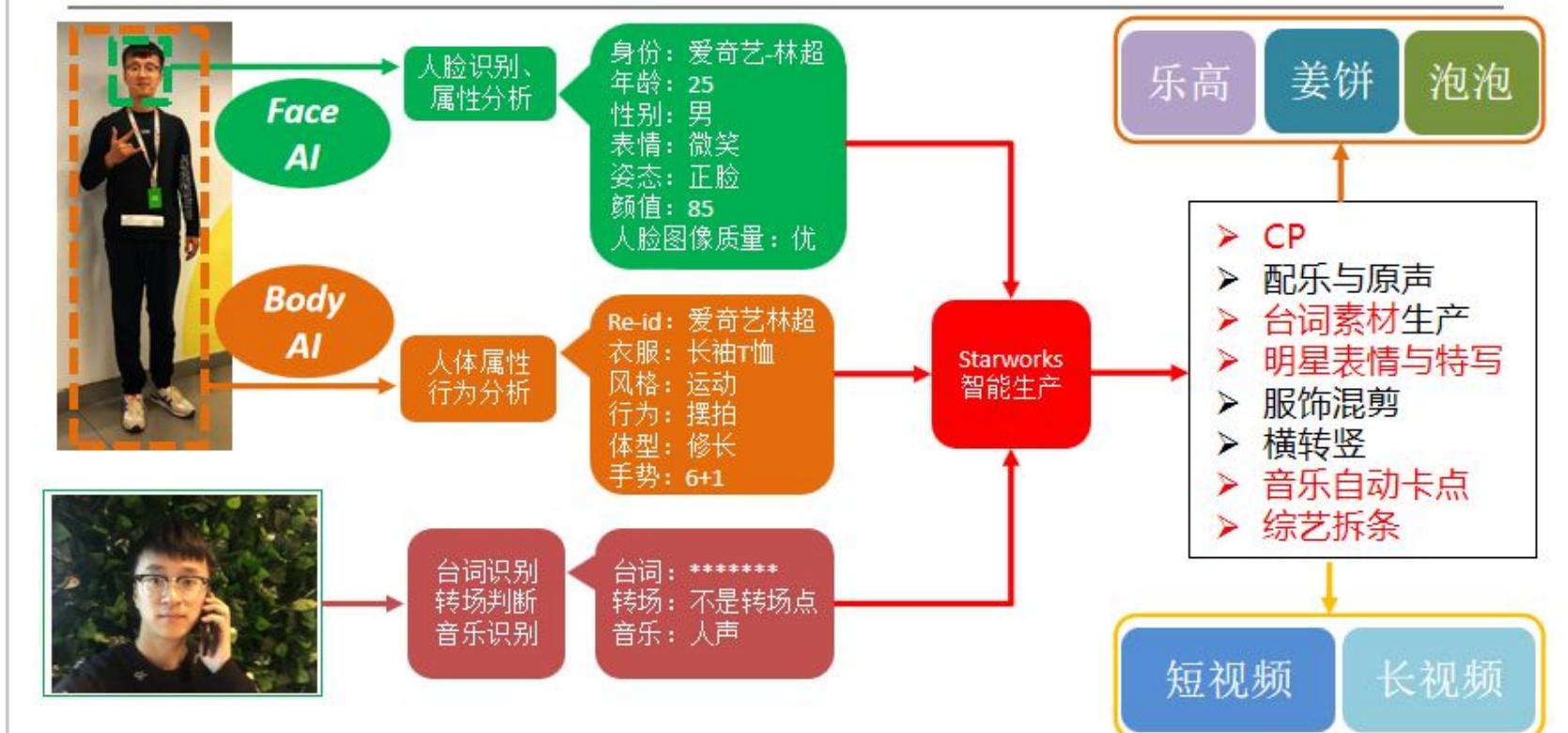
1. 类目标签
2. 事件/场景标签
3. 物体/人物标签
4. 字幕/OCR/ASR形成NLP标签

- 多模态召回辅助文本召回
- 只看TA
- 明星视频自动化生成

...

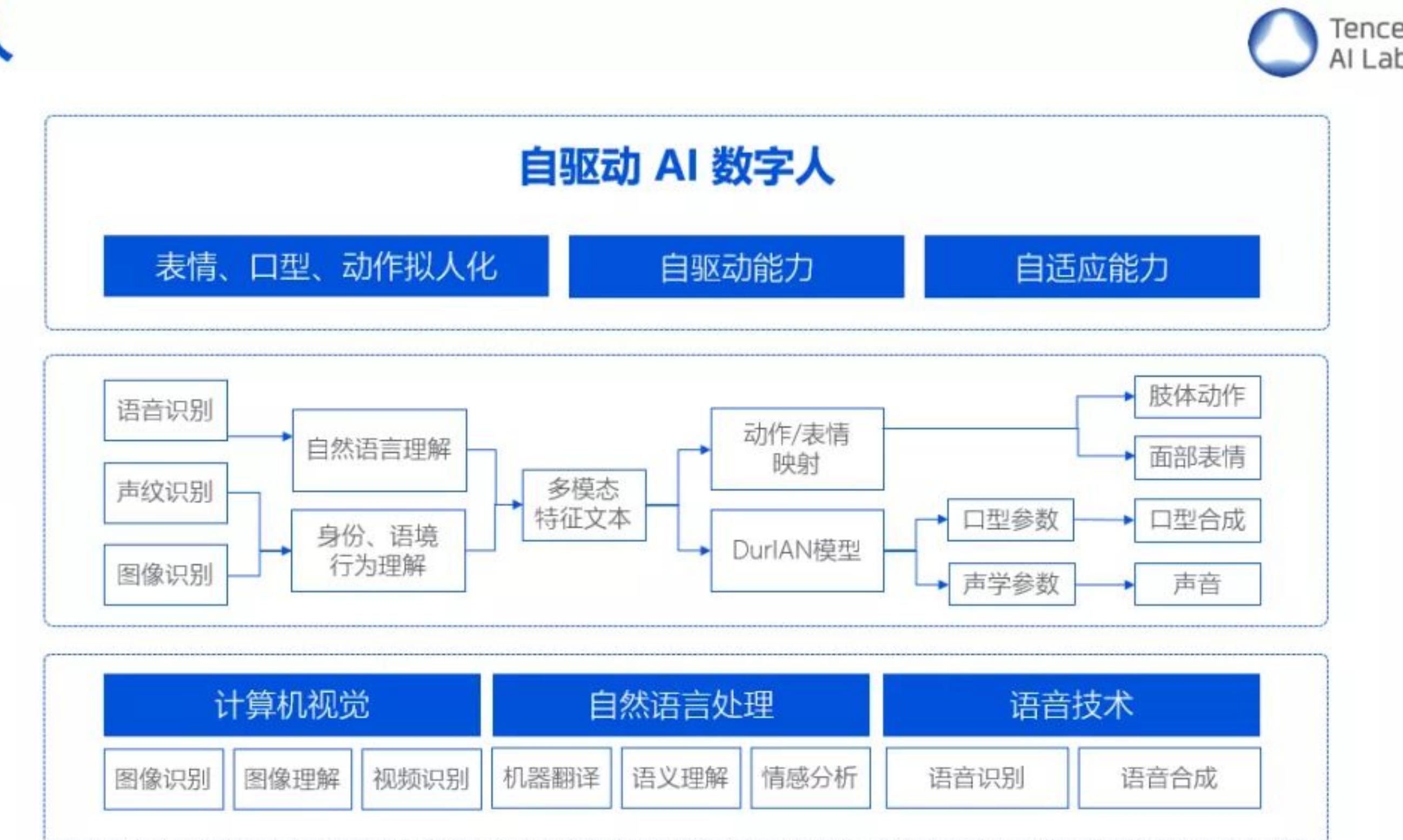
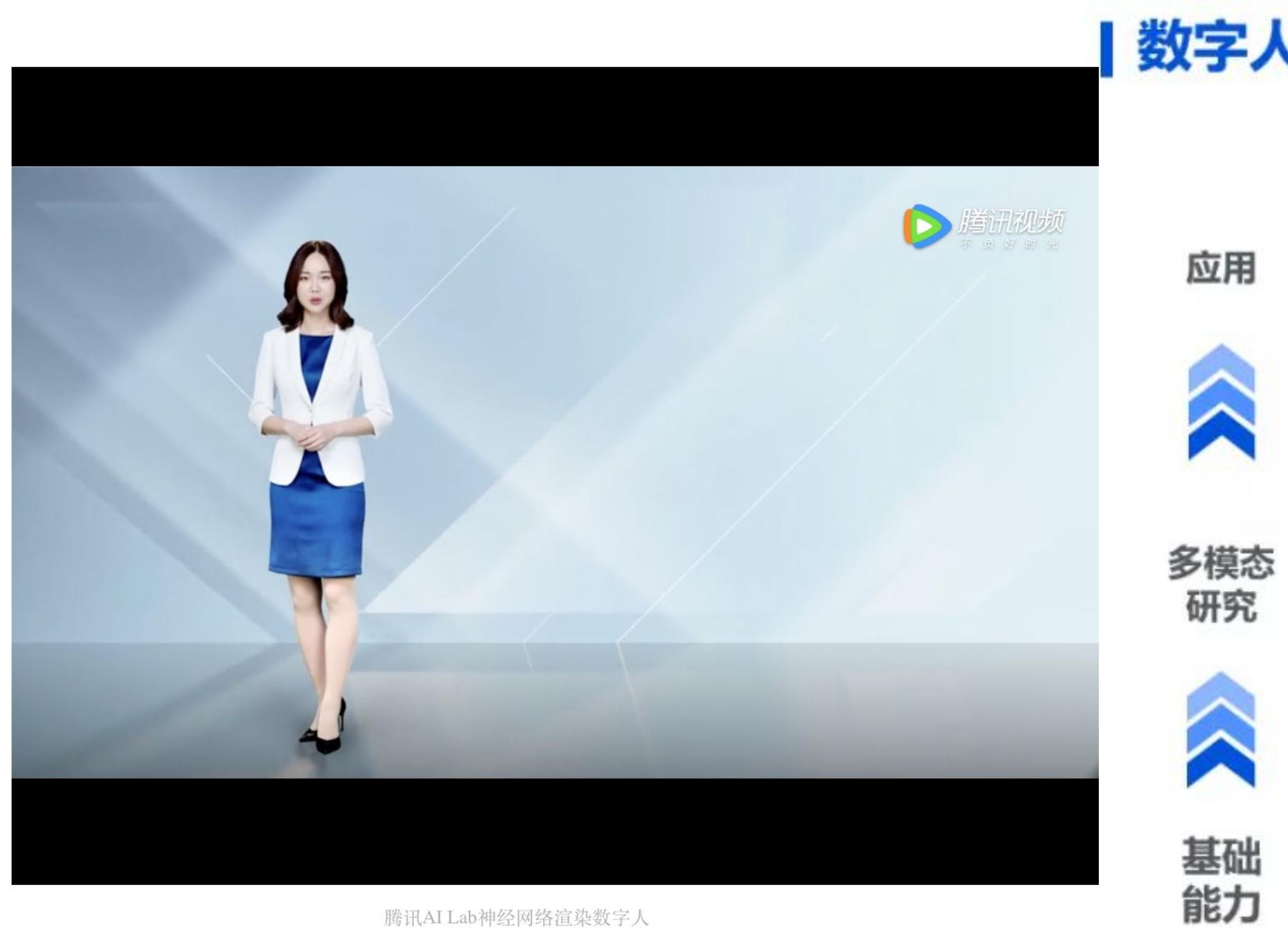


PersonAI - Starworks



多模态技术在工业界的应用

● 腾讯AILab——数字人



1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

3

多模态技术在工业界的应用

4

Jina在多模态场景下的应用

- 搭建视频搜索**demo**
- DalleFlow

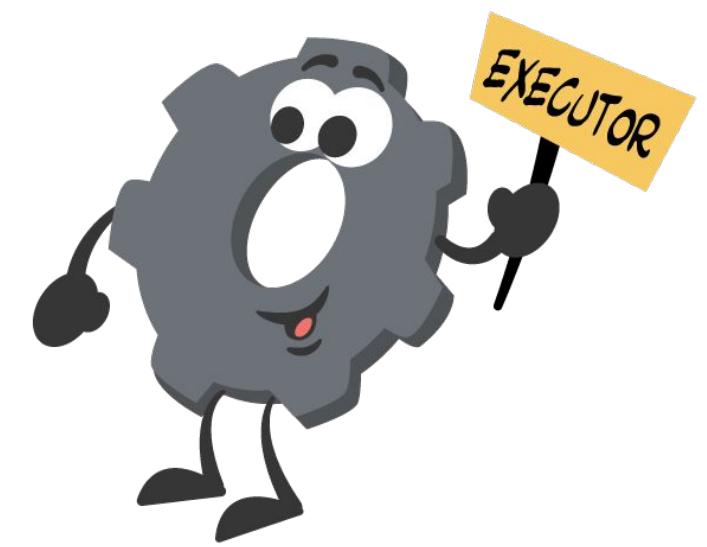
用Jina快速搭建视频搜索demo

- Jina的几个核心概念



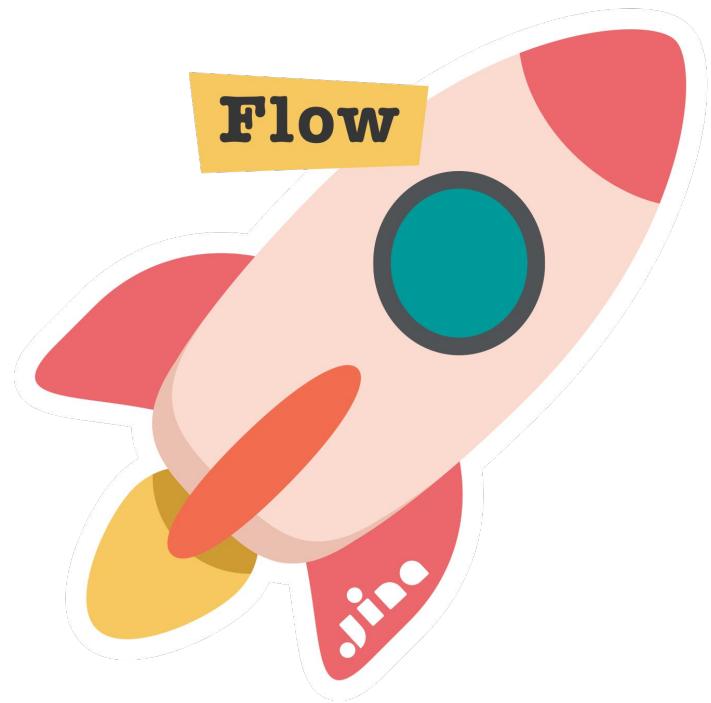
DocArray:

Basic data type to represent a piece of information(text, image, etc)



Executor:

How Jina process Documents.
(Segment, encode, index, etc)



Flow:

How Jina streamlines and distributes Executors

用Jina快速搭建视频搜索demo

- 场景:up主想要快速在成千上万的视频素材库里搜索到和自己脚本内容相匹配的片段



- [Bilibili:用 AI 搜索梦华录刘亦菲口红同款](#)
- [线上测试:<https://lipstick-db.senses.chat>](#)
- [Github:lipstick-db](#)



用Jina快速搭建视频搜索demo

- 场景:up主想要快速在成千上万的视频素材库里搜索到和自己脚本内容相匹配的片段

Query: eating hot pot and singing songs



[Bilibili:输入关键词就能自动剪视频？我写了一个AI视频搜剪神器？](#)

Query: mahjong mask

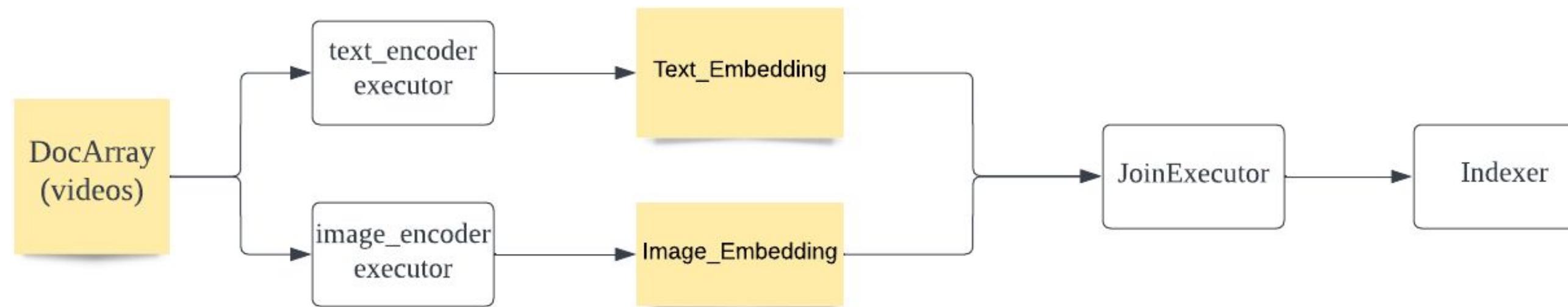


[Github:Jina-Clip](#)

用Jina快速搭建视频搜索demo

- How to do it?

INDEX



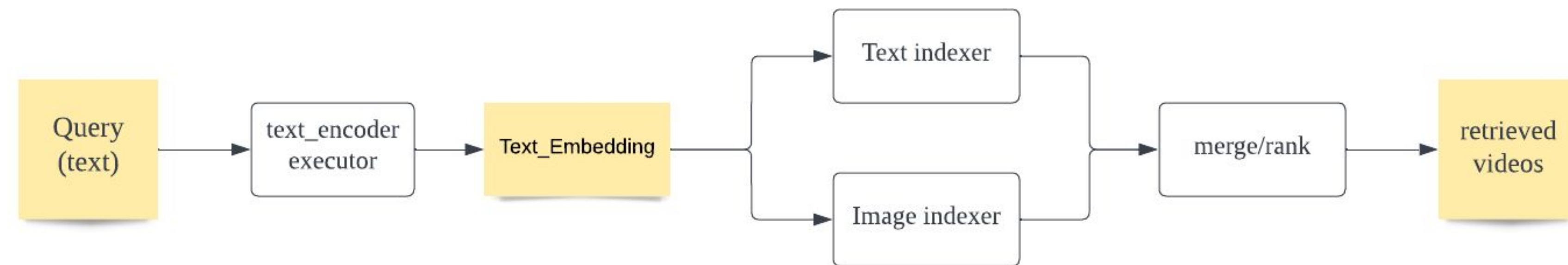
```
@requests
def encode(self, docs: DocumentArray, **kwargs):
    """
    Encode all documents with the `text` attribute and store the embeddings in the
    `embedding` attribute.
    :param docs: DocumentArray containing the Documents to be encoded
    """
    for docs_batch in docs.batch(batch_size=128):
        text_batch = docs_batch.texts
        with torch.inference_mode():
            embeddings = self.model(text_batch)
            for doc, e in zip(docs_batch, embeddings):
                doc.embedding = e
```

```
@requests(on='/index')
def index(self, bdocs: 'DocumentArray', **kwargs):
    """
    Add All Documents to the DocumentArray
    :param docs: the docs to add
    """
    if docs:
        self._index.extend(docs)
```

用Jina快速搭建视频搜索demo

- How to do it?

QUERY



The screenshot shows the Jina Hub interface with the following sections:

- Data Types:** Text, Image, Video, Audio, 3D Mesh, PDF, PowerPoint.
- Libraries:** MongoDB, PostgreSQL, Numpy, PaddlePaddle, PyTorch, spaCy, TensorFlow, Transformers.
- Keywords:** encoder (23), indexer (23), text (18), image (16), PyTorch (13), en (9).
- Executors (84):** Trending section listing executors:
 - TransformerTorchEncoder**: Editors' Pick. Description: "TransformerTorchEncoder" encodes text using transformer models. Updated a month ago, 4231 downloads.
 - CLIPImageEncoder**: Description: Image encoder that wraps the image embedding functionality using the CLIP ... Updated a month ago.
 - CLIPTextEncoder**: Editors' Pick. Description: Text encoder that wraps the text embedding functionality using the CLIP ... Updated 3 months ago, 1471 downloads.
 - HNSWPostgresIndexer**: Description: A complete indexer based on HNSW and Postgres ... Updated 9 days ago.
 - ImageTorchEncoder**: Editors' Pick. Description: Image encoder that wraps the image embedding functionality using the PyTorch ... Updated 3 months ago.
 - SimpleIndexer**: Description: A simple indexer that stores documents in memory and provides a basic search interface. Updated 3 months ago.

1

多模态数据**vs**单模态数据

2

学术界如何处理多模态数据

3

多模态技术在工业界的应用

4

Jina在多模态场景下的应用

- 搭建视频搜索**demo**
- **DALLE Flow**

DALLE Flow

Quick Facts

- Client-server architecture
- Combine DALLE & Diffusion using Jina
- 288 Lines of Python
- 42 Lines of YAML
- a weekend-work

README.md

DALL·E FLOW

A Human-in-the-Loop? workflow for creating HD images from text

Slack 3.1k Open in Colab docker build failing image size 17.9 GB

DALL·E Flow is an interactive workflow for generating high-definition images from text prompt. First, it leverages DALL·E-Mega to generate image candidates, and then calls CLIP-as-service to rank the candidates w.r.t. the prompt. The preferred candidate is fed to GLID-3 XL for diffusion, which often enriches the texture and background. Finally, the candidate is upscaled to 1024x1024 via SwinIR.

DALL·E Flow is built with Jina in a client-server architecture, which gives it high scalability, non-blocking streaming, and a modern Pythonic interface. Client can interact with the server via gRPC/Websocket/HTTP with TLS.

Why Human-in-the-Loop? Generative art is a creative process. While recent advances of DALL·E unleash people's creativity, having a single-prompt-single-output UX/UI locks the imagination to a *single* possibility, which is bad no matter how fine this single result is. DALL·E Flow is an alternative to the one-liner, by formalizing the generative art as an iterative procedure.

```
20     - executors/glid3/executor.py
21
22     env:
23         CUDA_VISIBLE_DEVICES: 0 # change this if you have multiple GPU
24         XLA_PYTHON_CLIENT_ALLOCATOR: platform #
25             https://jax.readthedocs.io/en/latest/gpu_memory_allocation.html
26
27     replicas: 1
28     needs: [gateway]
29
30     - name: rerank
31     uses: ReRank
32     uses_with:
33         clip_server: grpcs://demo-cas.jina.ai:2096
34
35     py_modules:
36         - executors/rerank/executor.py
37
38     needs: [dalle, diffusion]
39
40     - name: upscaler
41     uses: SwinIRUpscaler
42
43     py_modules:
44         - executors/swinir/executor.py
45
46     uses_with:
47         swinir_path: ../SwinIR
48         store_path: dalle.db
49
50     env:
51         CUDA_VISIBLE_DEVICES: 0 # change this if you have multiple GPU
52
53     - name: store
54     uses: MyStore
55
56     py_modules:
57         - executors/store/executor.py
58
59     uses_with:
60         store_path: dalle-flow.db
```

upscaler

SwinIRUpscaler

SwinIRUpscaler

SwinIRUpscaler

SwinIRUpscaler

SwinIRUpscaler

store

MyStore

gateway



```
1 jtype: Flow
2 with:
3   protocol: grpc
4   port: 51005
5   env:
6     JINA_LOG_LEVEL: debug
7 executors:
8   - name: dalle
9     uses: executors/dalle/config.yml
10    timeout_ready: -1 # slow download speed often leads to timeout
11    env:
12      CUDA_VISIBLE_DEVICES: 0 # change this if you have multiple GPU
13      XLA_PYTHON_CLIENT_ALLOCATOR: platform #
14      https://jax.readthedocs.io/en/latest/gpu_memory_allocation.html
15    - name: diffusion
16      uses: GLID3Diffusion
17      uses_with:
18        glid3_path: ../glid-3-xl
19        steps: 100
20      py_modules:
21        - executors/glid3/executor.py
22      env:
23        CUDA_VISIBLE_DEVICES: 0 # change this if you have multiple GPU
24        XLA_PYTHON_CLIENT_ALLOCATOR: platform #
25        https://jax.readthedocs.io/en/latest/gpu_memory_allocation.html
26      replicas: 1
27      needs: [gateway]
28    - name: rerank
29      uses: ReRank
30      uses_with:
31        clip_server: grpcs://demo-cas.jina.ai:2096
32      py_modules:
33        - executors/rerank/executor.py
34      needs: [dalle, diffusion]
35    - name: upscaler
36      uses: SwinIRUpscaler
37      py_modules:
38        - executors/swinir/executor.py
39      uses_with:
40        swinir_path: ../SwinIR
41        store_path: dalle.db
42      env:
43        CUDA_VISIBLE_DEVICES: 0 # change this if you have multiple GPU
44    - name: store
45      uses: MyStore
46      py_modules:
47        - executors/store/executor.py
48      uses_with:
49        store_path: dalle-flow.db
```

Flow configuration

- port
- protocol
- environment variables

Flow spec



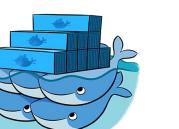
DALLE

- init kwargs
- GPU config
- executor config



Diffusion

- init kwargs
- GPU config
- executor config



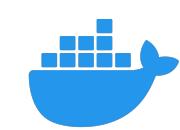
Rerank

- init kwargs
- executor config



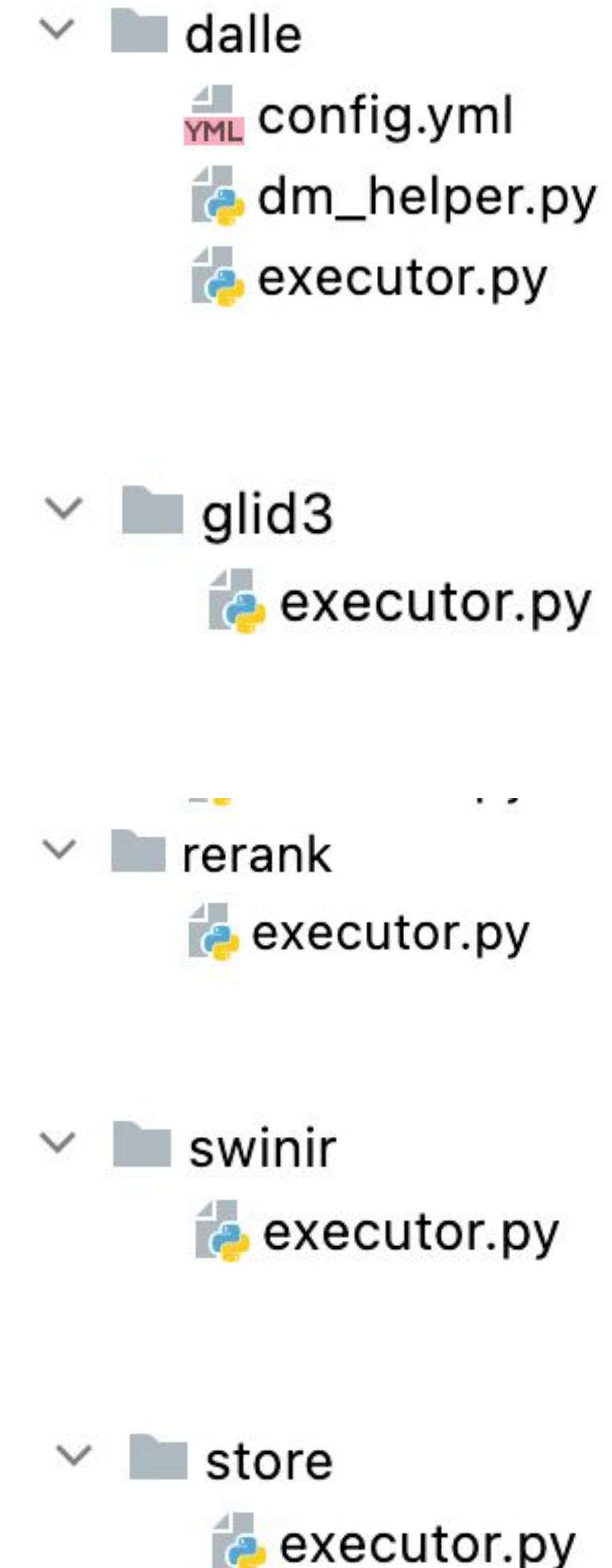
Upscaler

- init kwargs
- executor config



Store

- init kwargs
- executor config



"ocean beach front view in Van Gogh style"



.jina

***"the statue of liberty wearing a
vr headset"***





欢迎加入 Jina AI 开源社区 解锁数据应用的可能性

官网 : <https://jina.ai/>

GitHub : <https://github.com/jina-ai>

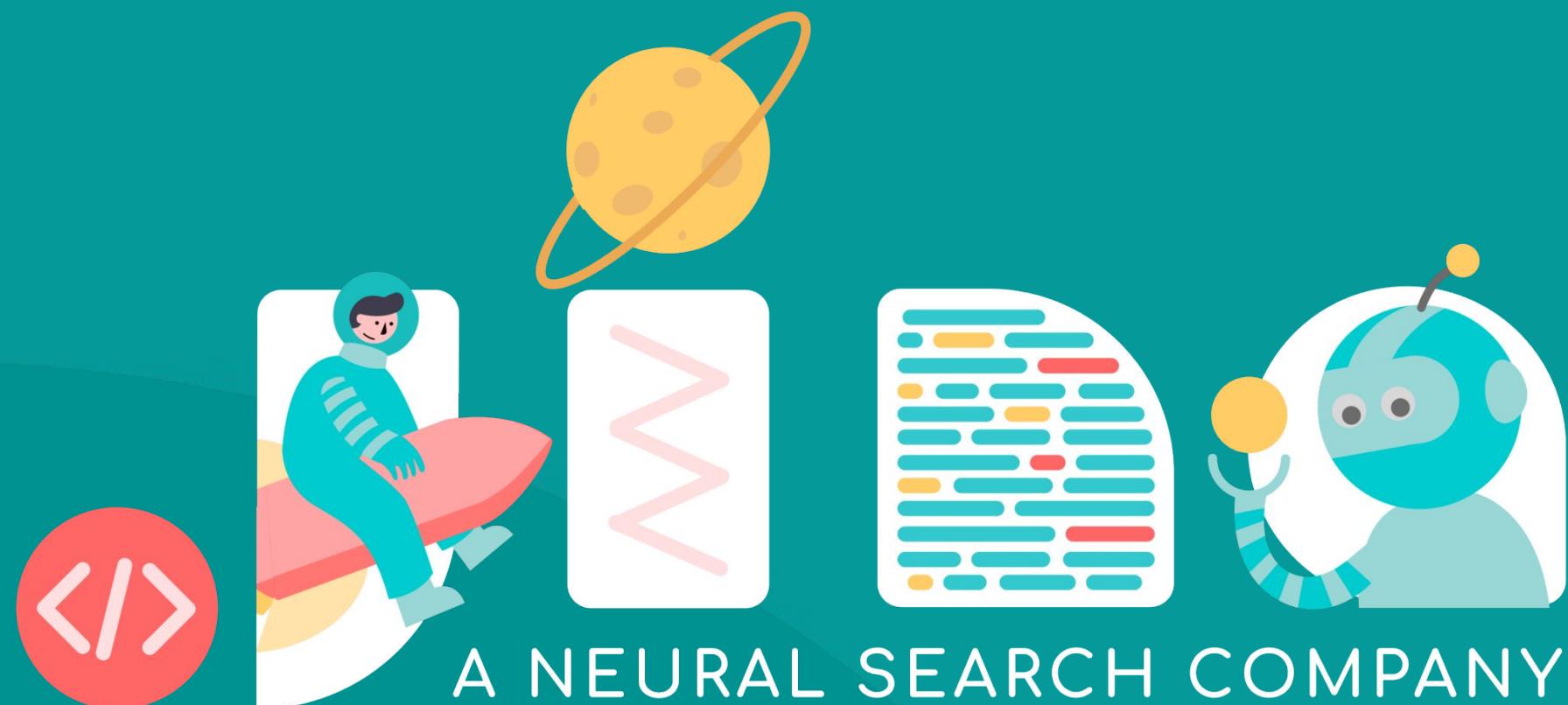
加入全球开发者社区 : <https://slack.jina.ai/>



扫码关注 Jina AI 公众号



添加小助手加入技术交流群



Thank you!