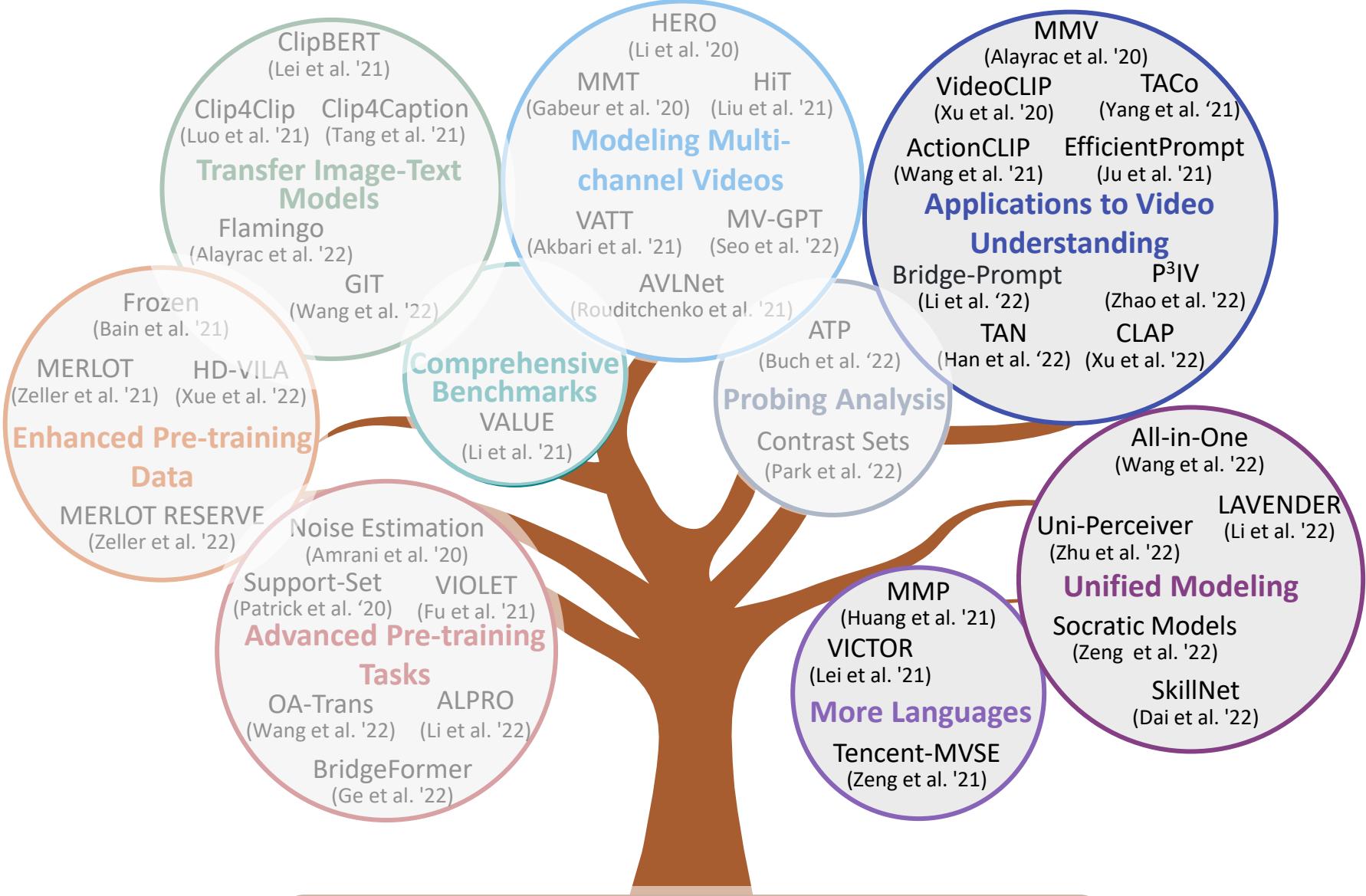


Advanced Topics in Video-Text Pre-training

Chung-Ching Lin



Pioneering work in Video-Text Pre-training

VideoBERT
(Sun et al. '19)

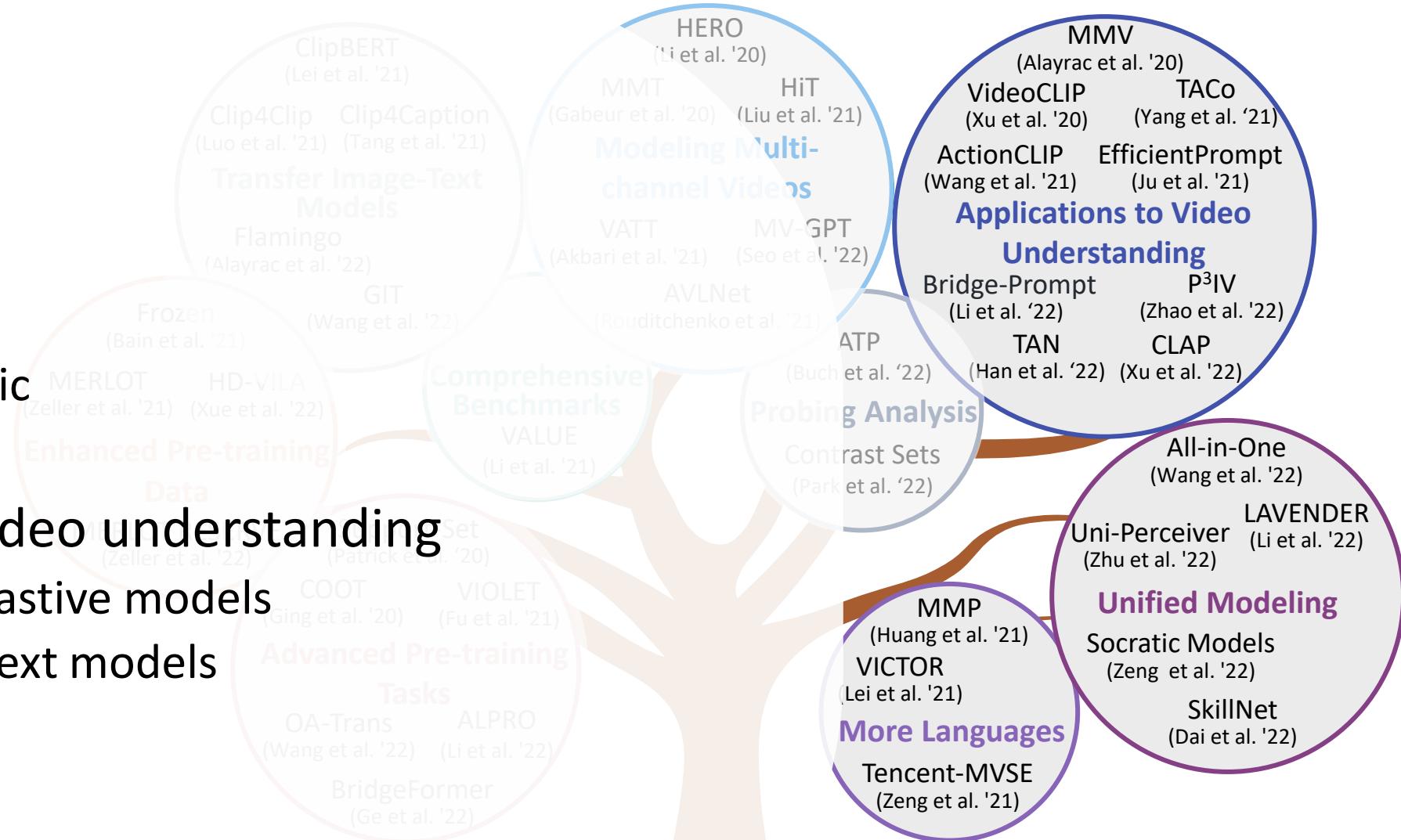
UniVL
(Luo et al. '20)

HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)

Agenda

- Unified modeling
 - Modality-agnostic
 - Task-agnostic
- Applications to video understanding
 - Video-text contrastive models
 - Transfer image-text models
 - Prompt tuning
- More languages



Pioneering work in Video-Text Pre-training

VideoBERT
(Sun et al. '19)

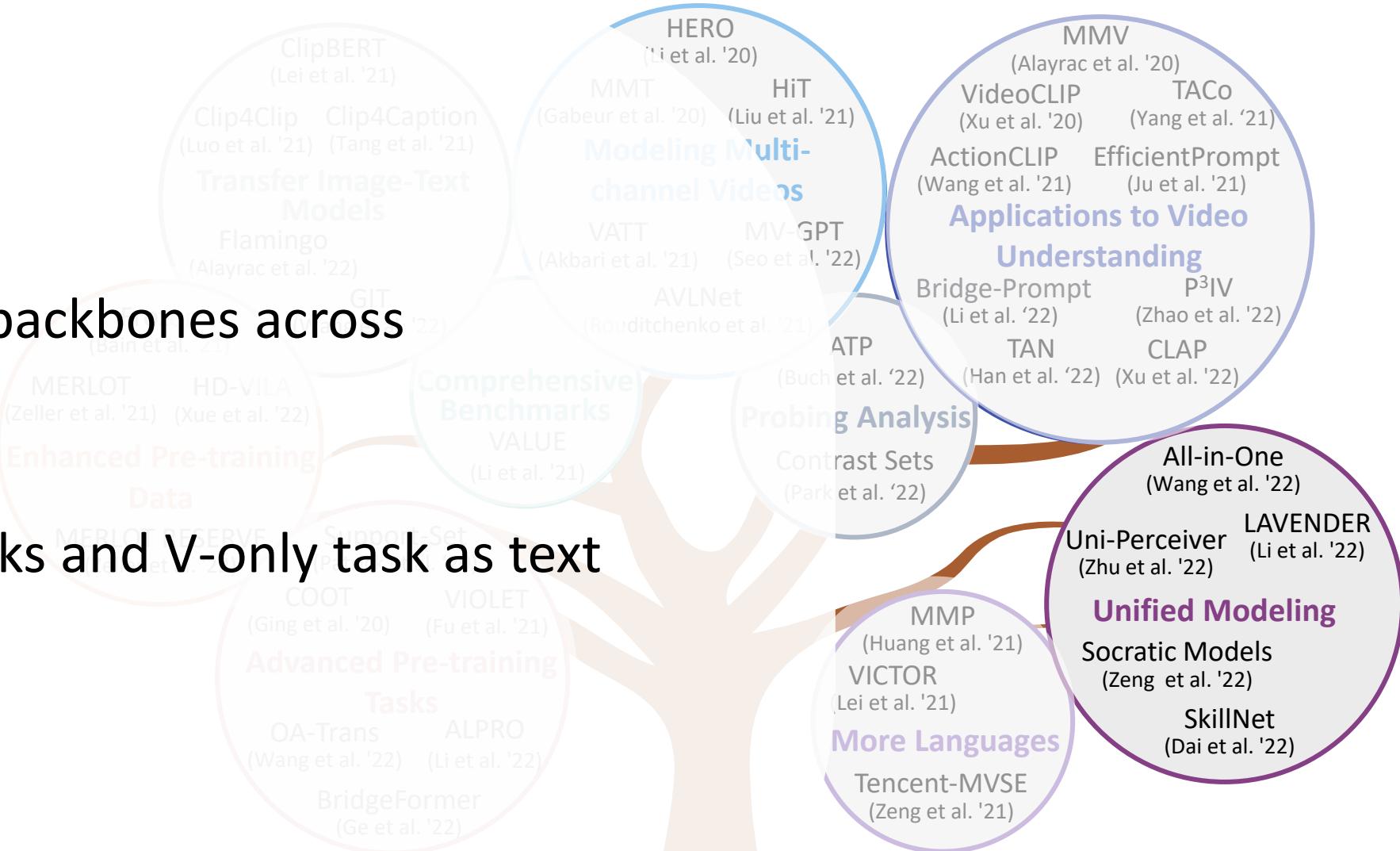
ActBERT
(Zhu and Yang '20)

HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)

Unification

- Modalities: unify backbones across modalities
- Tasks: unify VL tasks and V-only task as text generation



Pioneering work in Video-Text Pre-training

VideoBERT
(Sun et al. '19)

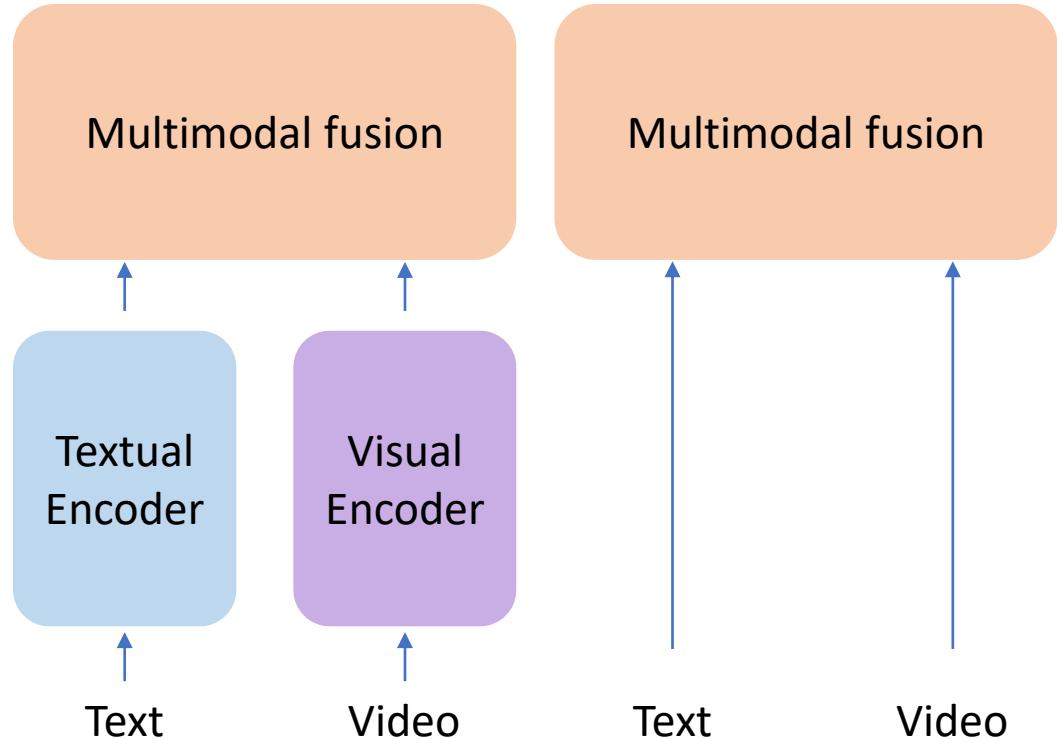
ActBERT
(Zhu and Yang '20)

HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)

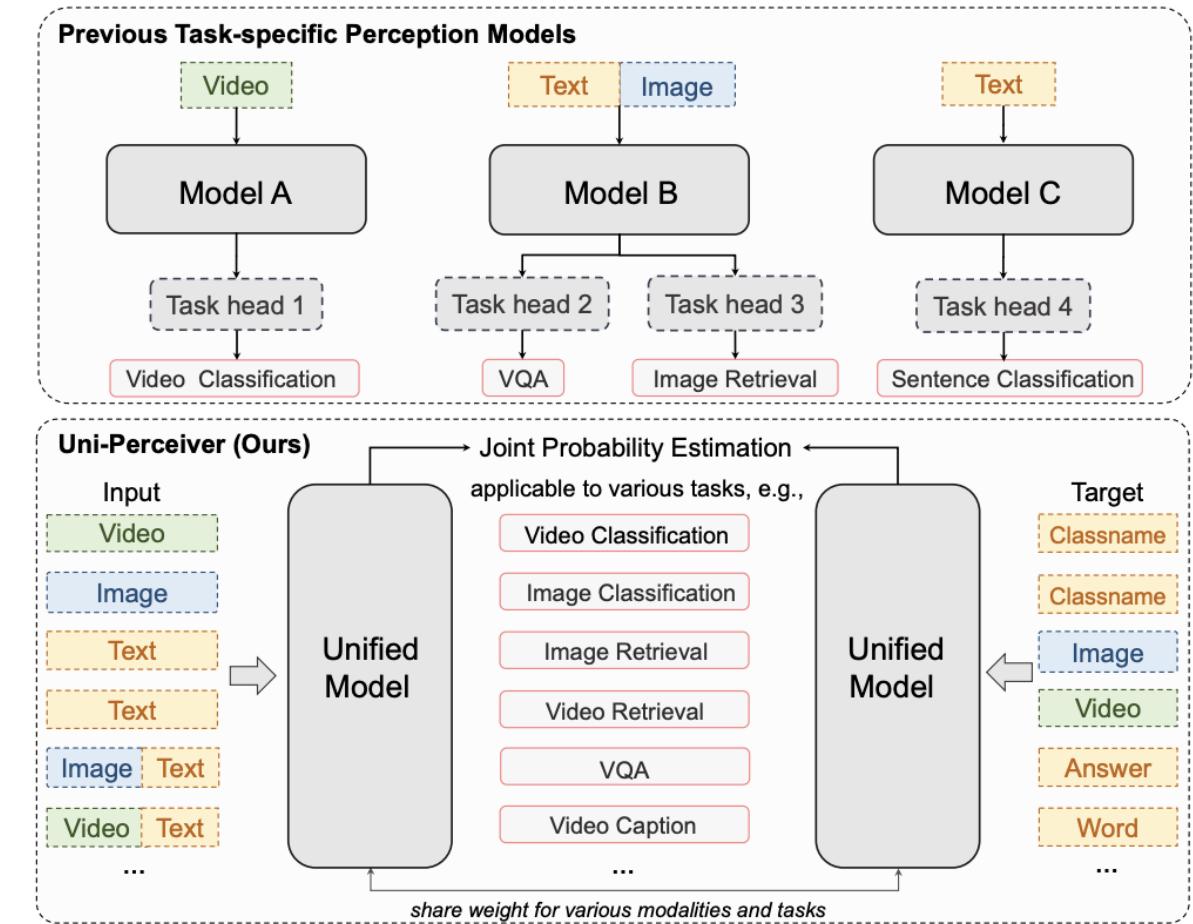
Unification

- Modalities: unify backbones (text, image, video, audio)
 - Uni-Perceiver, All in One, SkillNet, VATT
- Tasks: unify VL tasks V-only task as text generation
 - LAVENDER, Socratic Models, GIT, Flamingo



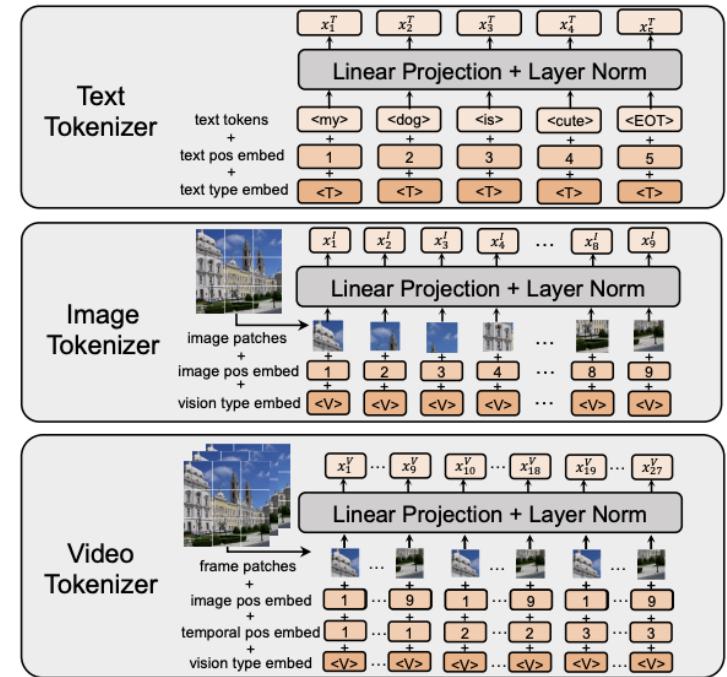
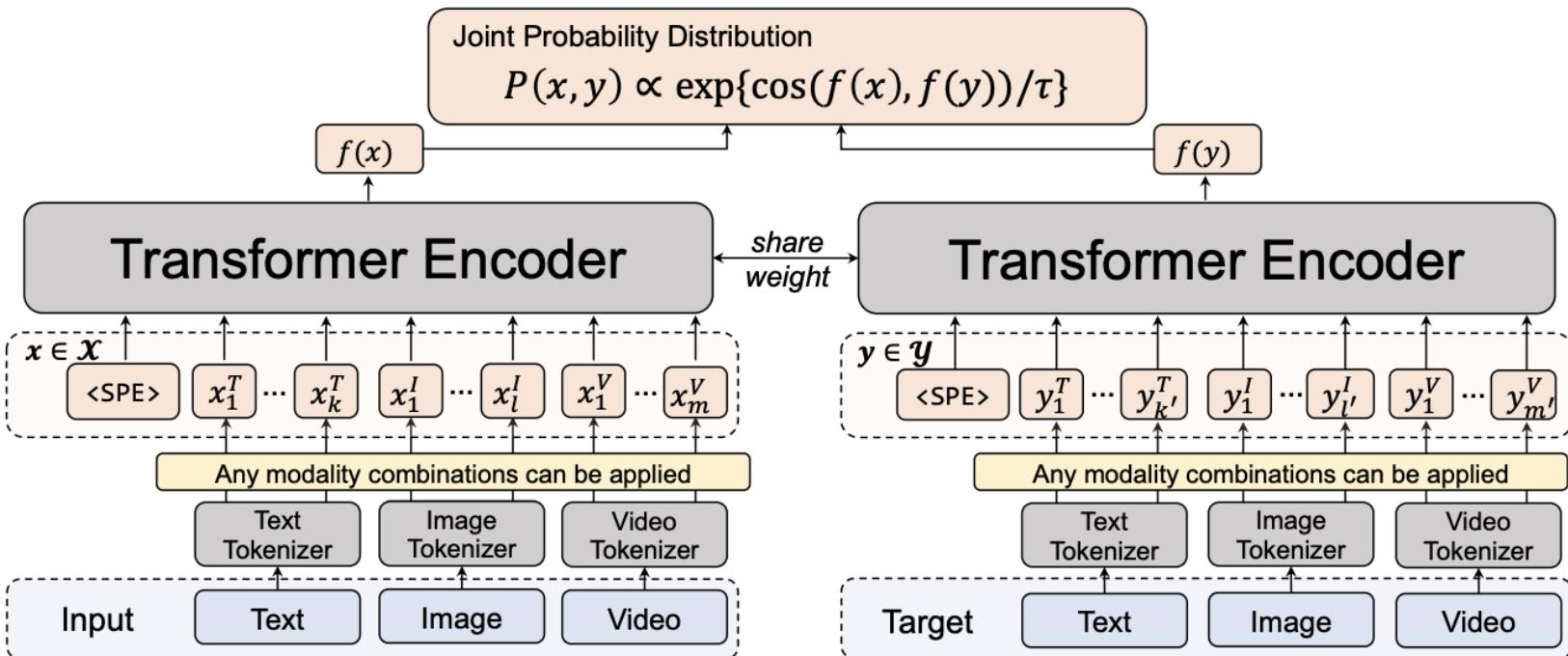
Modality-specific to modality-agnostic

- A modality-agnostic Transformer encoder and lightweight modality-specific tokenizer
- Encodes different task inputs and targets from arbitrary modalities into a unified representation



Modality-agnostic model

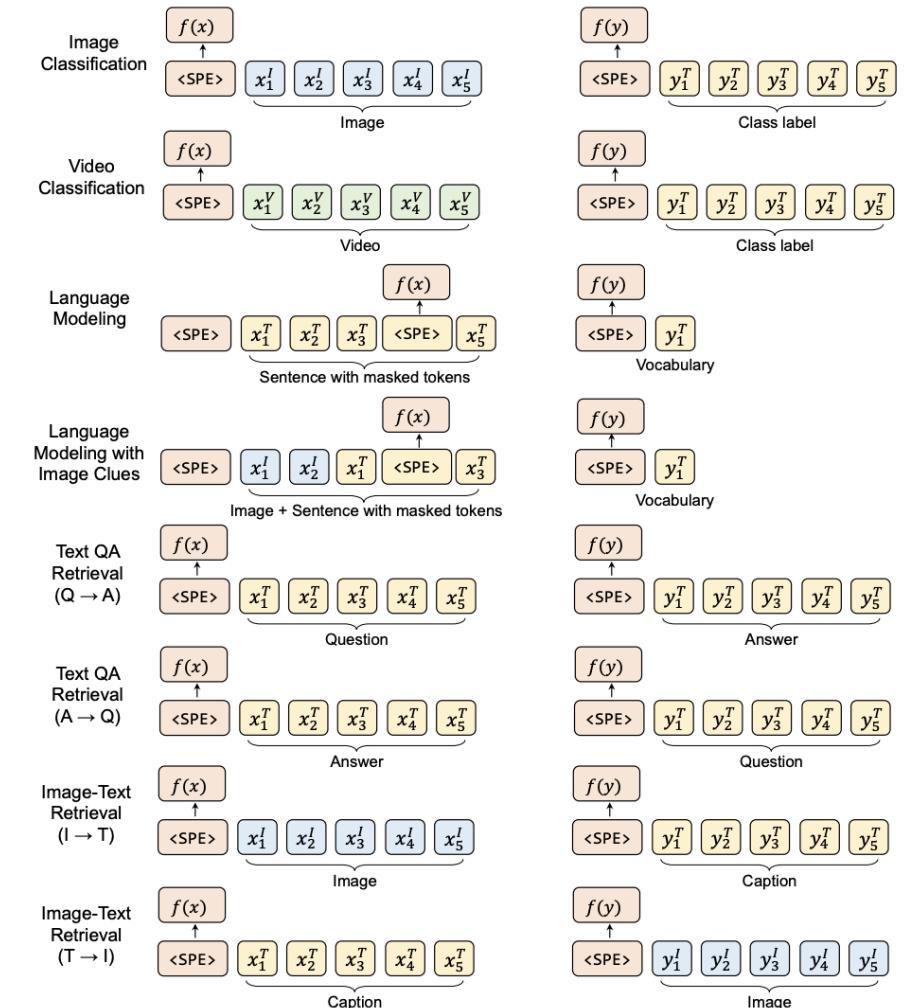
- Encodes different task inputs and targets from arbitrary modalities into a unified representation



Pre-training on Multi-Modal Tasks

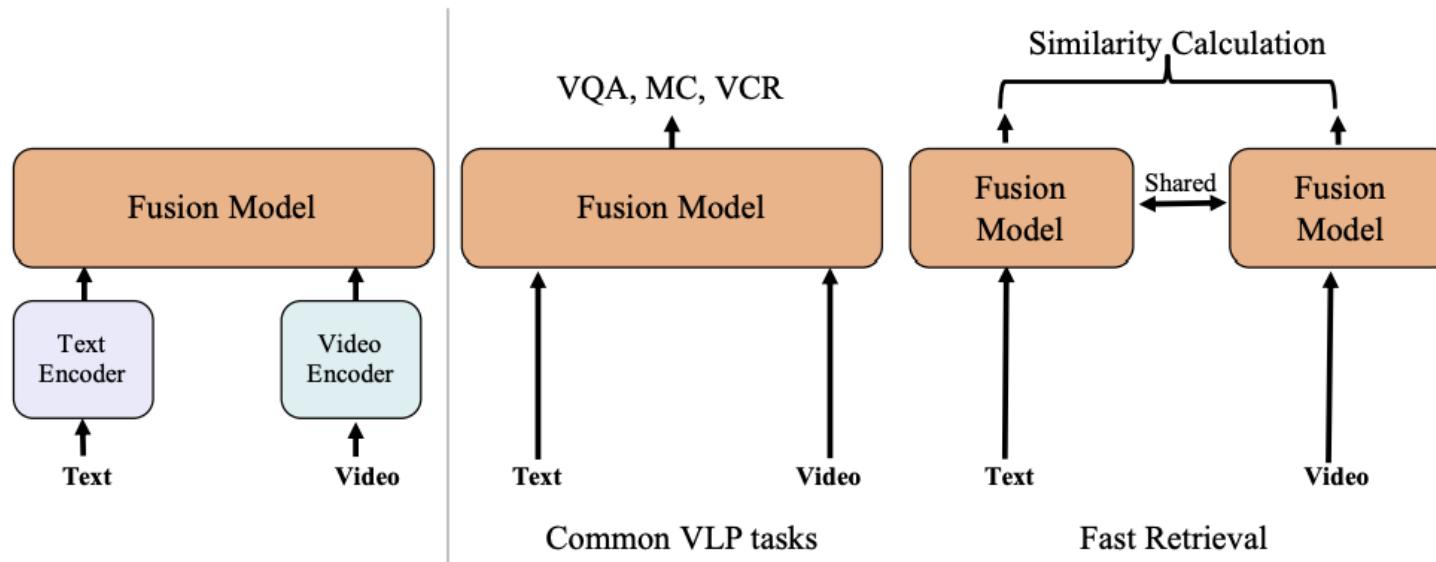
Pre-training dataset

Dataset	#Images	#Videos	#Text
ImageNet-21k [17]	14.2M	0	21K
Kinetics-700 [32]	0	542K	700
Moments in Time [57]	0	792K	339
Books&Wiki [93]	0	0	101M
PAQ [41]	0	0	65M
CC3M [68]	3.0M	0	3.0M
CC12M [9]	11.1M	0	11.1M
COCO Caption [12]	113K	0	567K
Visual Genome [36]	108K	0	5.41M
SBU [58]	830K	0	830K
YFCC [31]	14.8M	0	14.8M



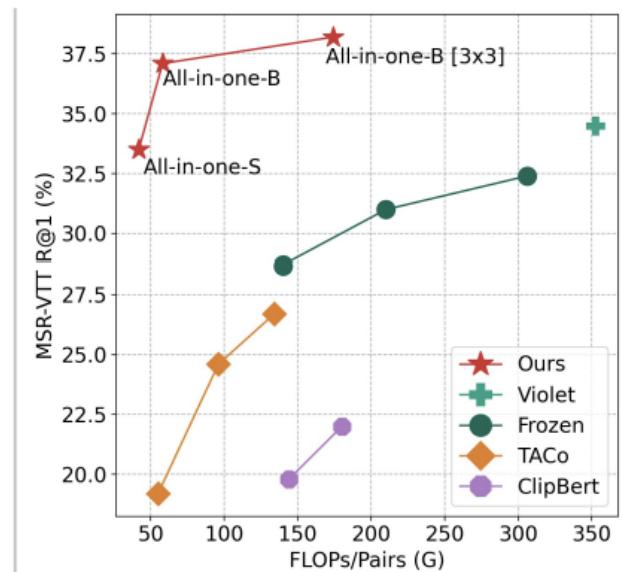
Modality-Agnostic Video-Language Pre-training models

- An unified backbone architecture



(a). Conventional Pre-training Framework

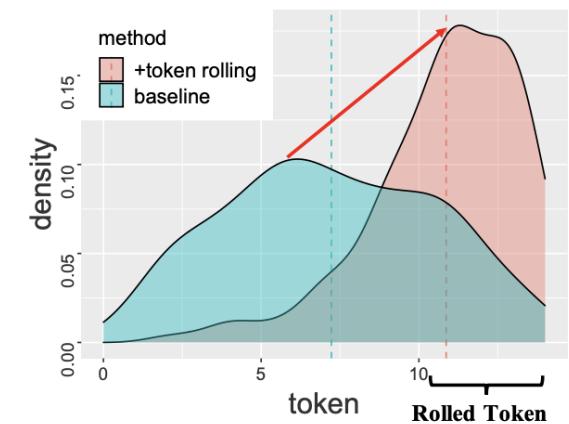
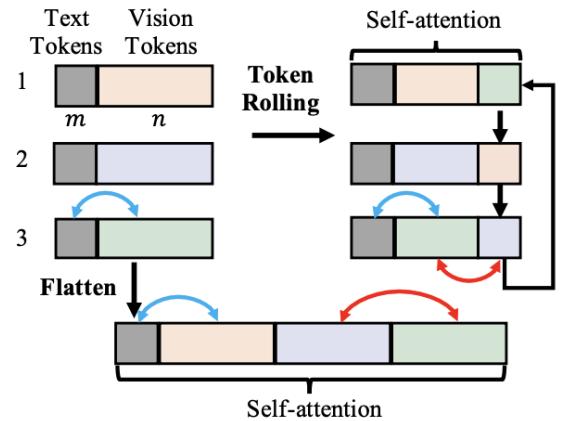
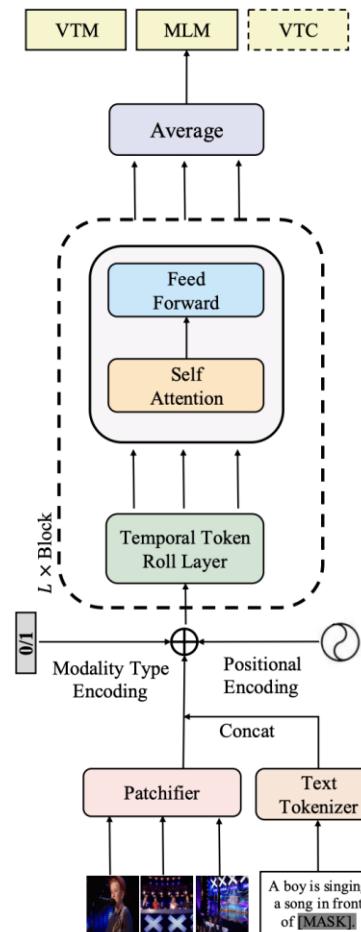
(b). Ours All-in-one Transformer



(c). Flops & Performance Comparison

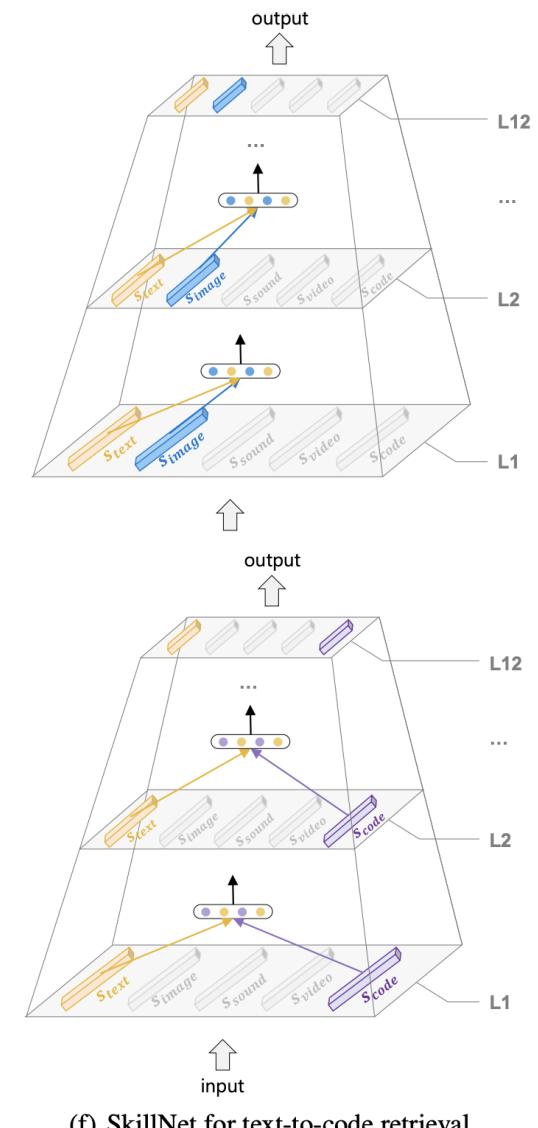
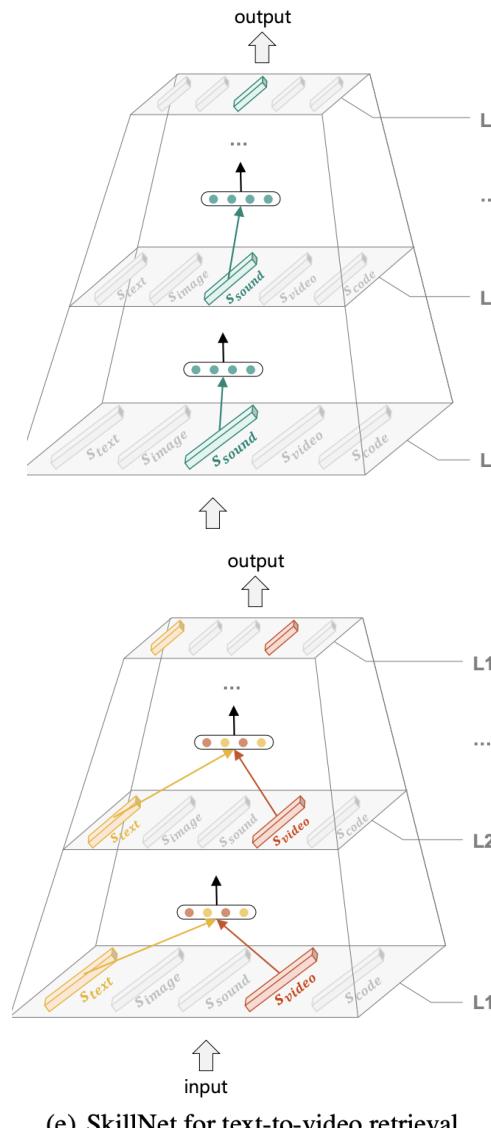
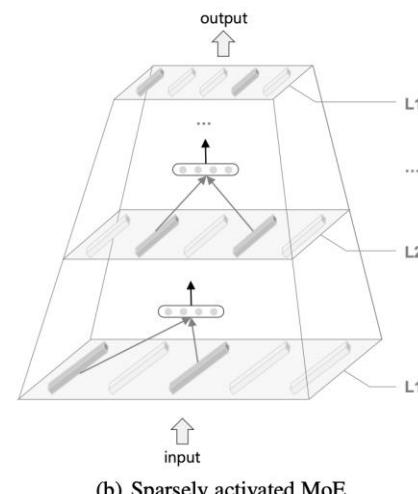
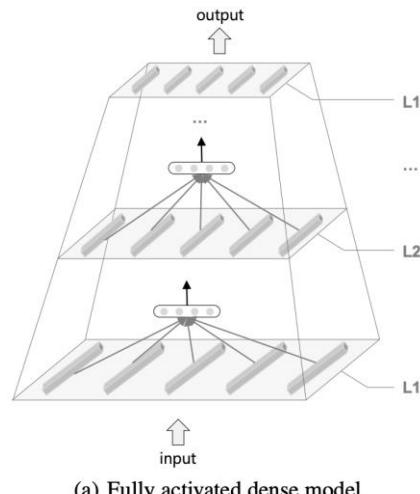
Modality-Agnostic Video-Language Pre-training models

- Token rolling operation to encode temporal representations from video clips in a non-parametric manner

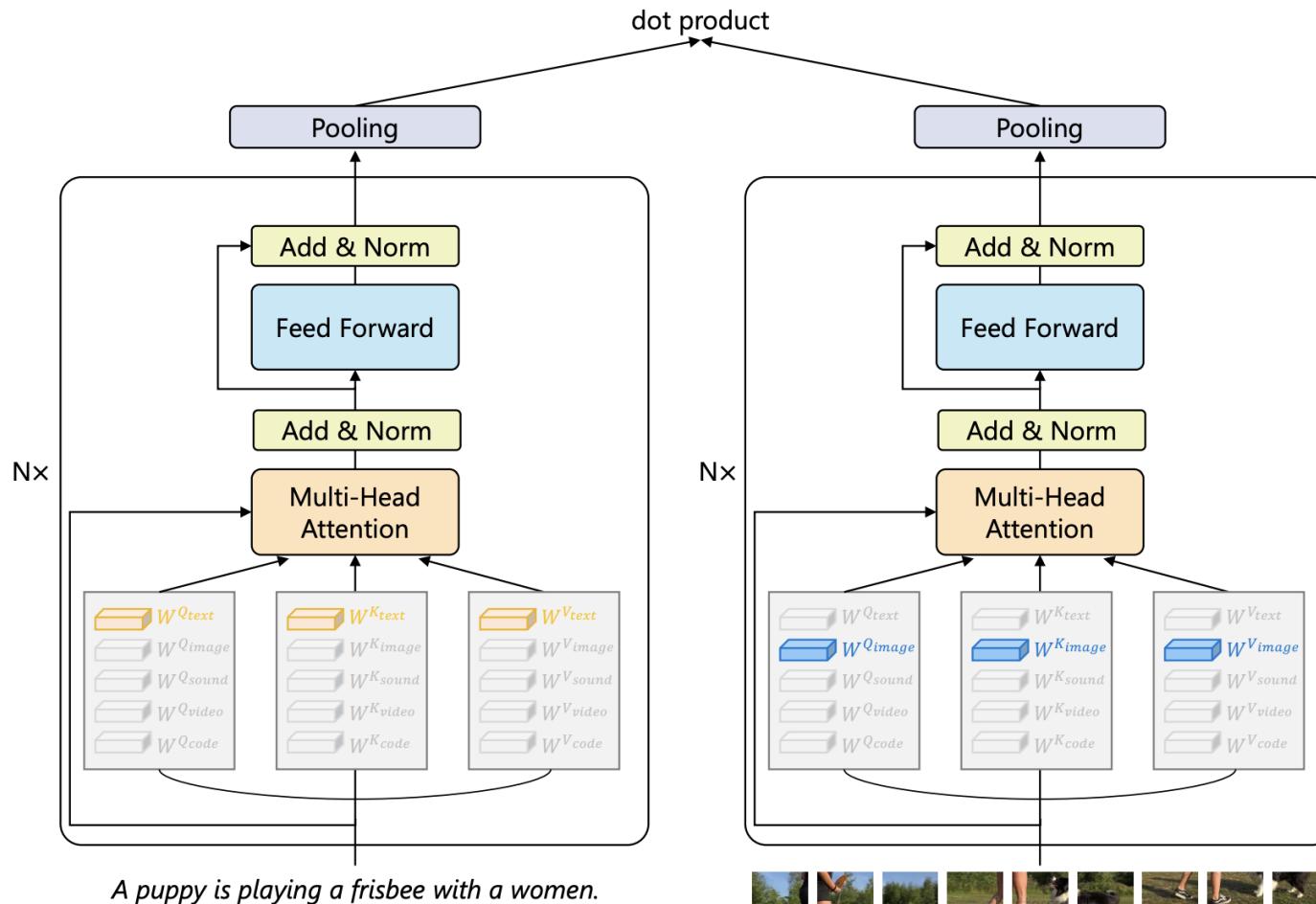


Sparse multimodal Mixture-of-Expert (MoE)

- Handling multiple modalities of information with a single model
- Different parts of the parameters are specialized for processing different modalities

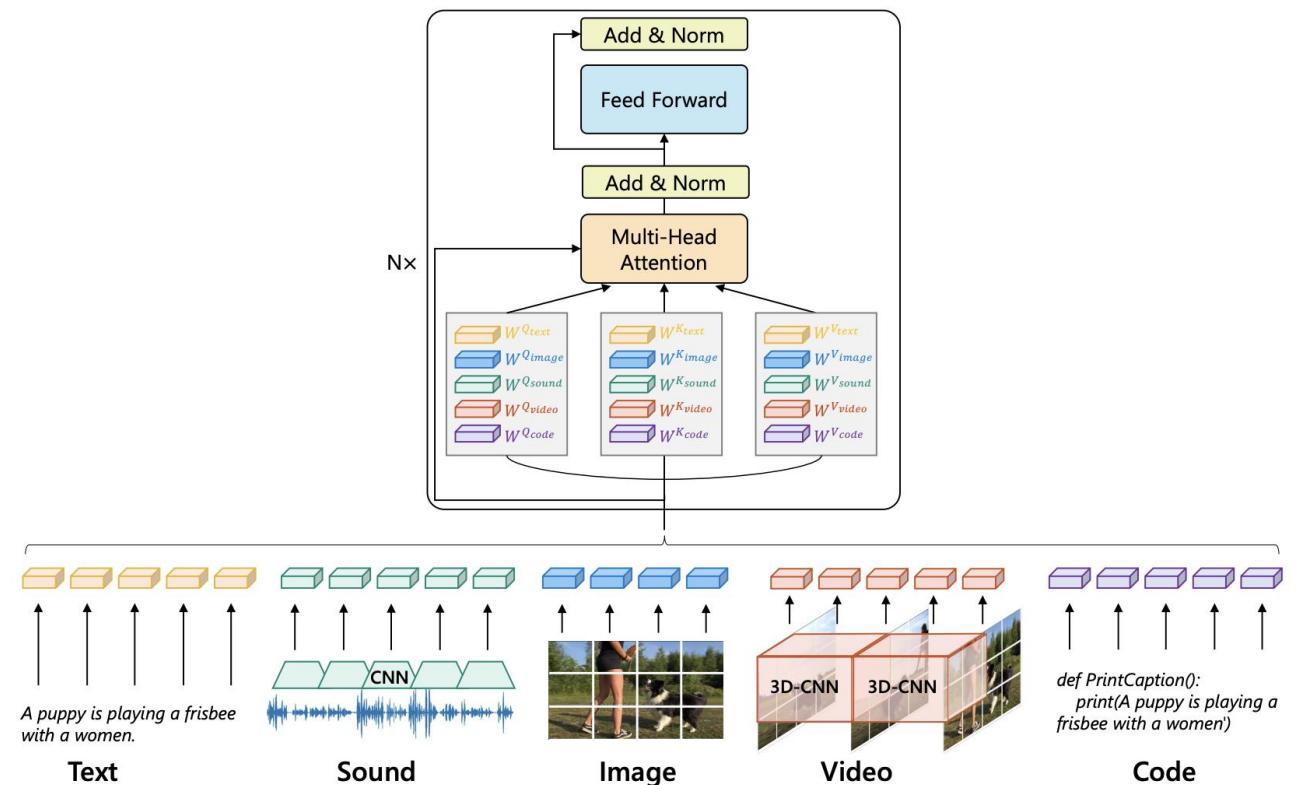


Sparse multimodal Mixture-of-Expert (MoE)



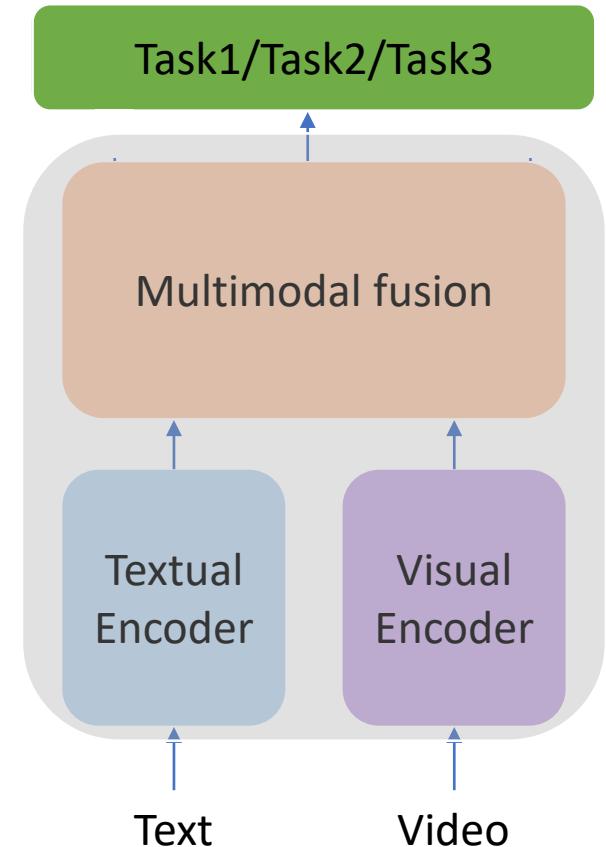
Sparse multimodal Mixture-of-Expert (MoE)

- The embeddings of different modalities
 - Text
 - Sound
 - Image
 - Video
 - Code



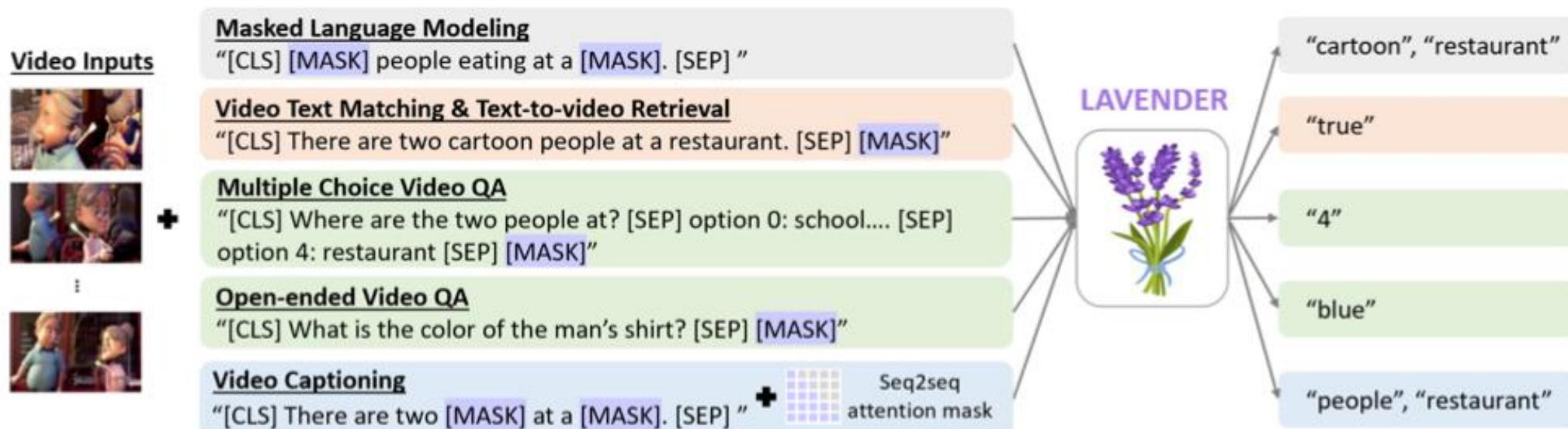
Unification

- Modalities: architecture (across audio, code)
 - Uni-Perceiver, All in One, SkillNet
- Tasks: unify VL tasks/V-only tasks as text generation
 - LAVENDER, Socratic Models, GIT, Flamingo



Unify VidL tasks as text generation

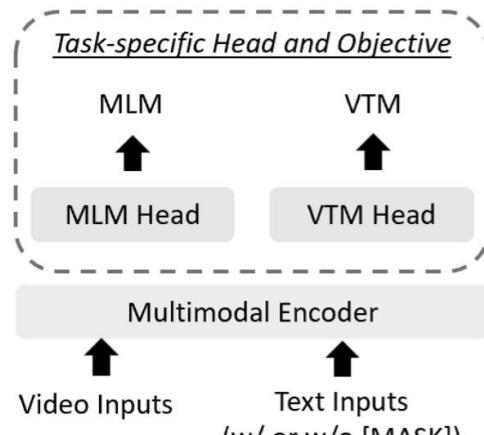
- Unify both pre-training and downstream finetuning as Masked Language Modeling



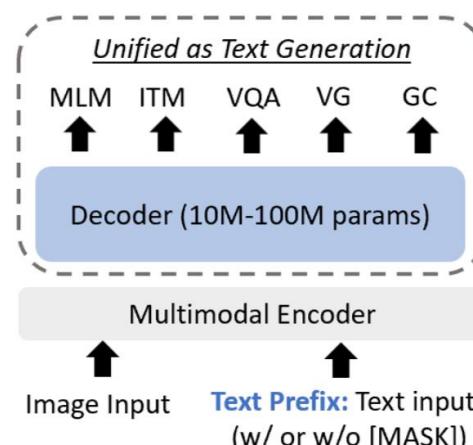
Unify VidL tasks as text generation

- Adopt **an encoder-only architecture**, with **a lightweight MLM head**, instead of the heavy decoder in unified image-text models

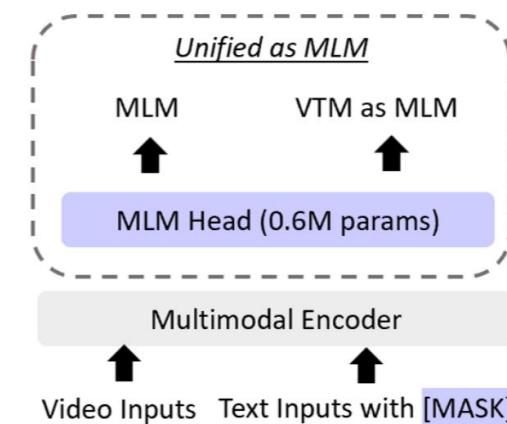
Comparison to existing methods during **pre-training**



(a) Existing VidL Methods



(b) Existing Unified Image-text Models

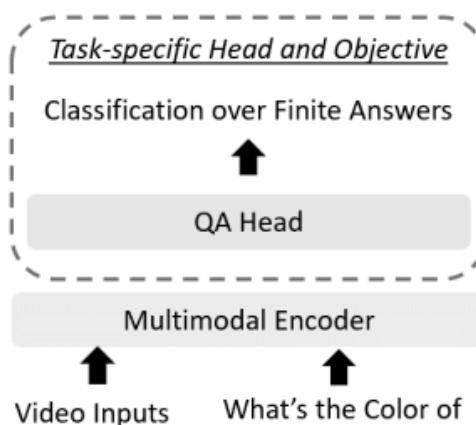


(c) LAVENDER

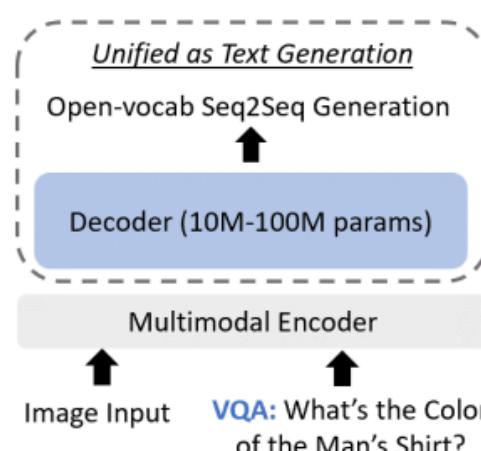
Unify VidL tasks as text generation

- Adopt **an encoder-only architecture**, with **a lightweight MLM head**, instead of the heavy decoder in unified image-text models

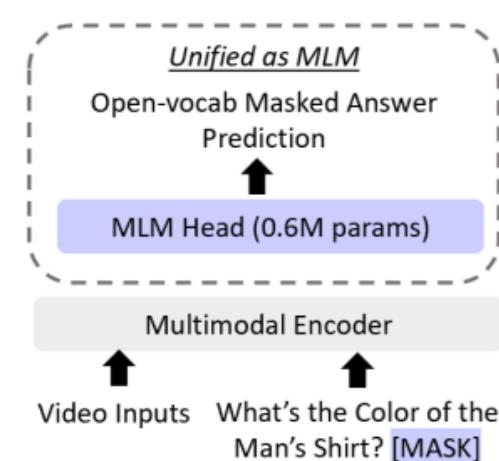
Comparison to existing methods **on downstream image/video question answering task**



(a) Existing VidL Methods



(b) Existing Unified Image-text Models



(c) LAVENDER

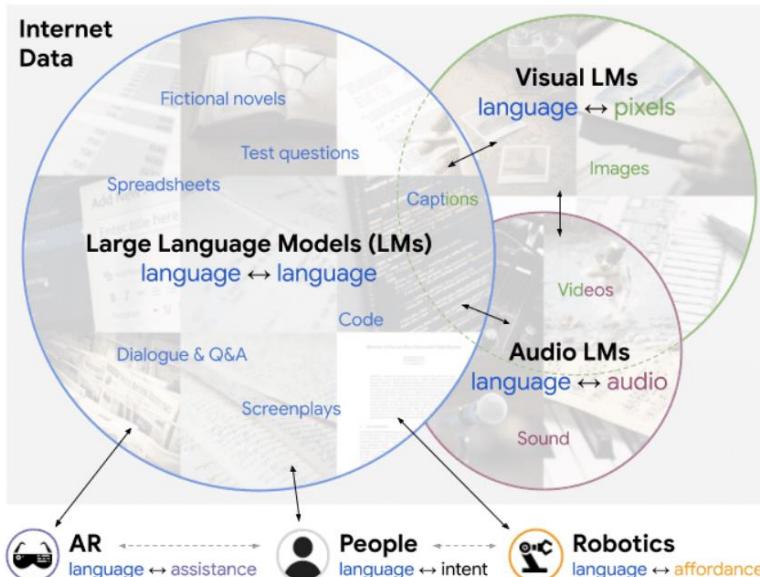
Unify VidL tasks as text generation

- LAVENDER surpasses previous methods with task-specific designs on 12 video-language tasks

	TGIF			MSRVTT		LSMDC		MSVD		Captioning		Retrieval	
	Action	Transition	Frame	MC	QA	MC	FiB	QA	MSRVTT	MSVD	DiDeMo	LSMDC	
Published	94.0	96.2	69.5	90.9	43.1	81.7	52.9	46.3	60.0	120.6	65.1	41.9	
SOTA	[84]	[84]	[84]	[84]	[84]	[84]	[84]	[78]	[57]	[35]	[6]	[12]	
LAVENDER	96.6	99.1	73.5	97.4	45.0	87.0	57.1	56.6	60.1	150.7	72.4	43.3	
Δ	2.6↑	2.9↑	4.0↑	6.5↑	1.9↑	5.3↑	4.2↑	10.3↑	0.1↑	30.1↑	7.3↑	1.4↑	

Composing Zero-Shot Multimodal Reasoning with Language

- Different pretrained models store different forms of commonsense knowledge across different domains
 - Visual-language models (VLMs) are trained on Internet-scale image captions
 - Large language models (LMs) are further trained on Internet-scale text with no images (e.g., spreadsheets, SAT questions, code)
- The diversity is symbiotic



I am an intelligent image captioning bot. This image is a {img_type}. There {num_people}. I think this photo was taken at a {place1}, {place2}, or {place3}. I think there might be a {object1}, {object2}, {object3},... in this {img_type}. A creative short caption I can generate to describe this image is:



SM (ours): This image shows an inviting dining space with plenty of natural light.



SM (ours): People gather under a blossoming cherry tree, enjoying the beauty of nature together.



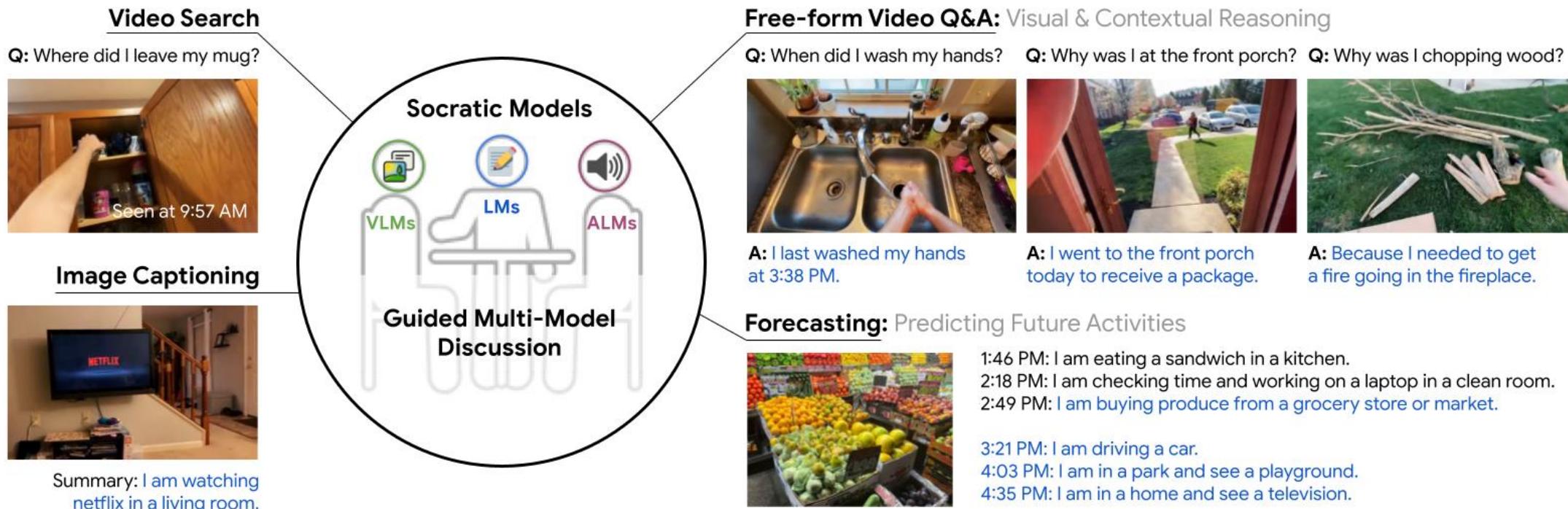
SM (ours): At the outdoor market, you can find everything from plantains to Japanese bananas.

ClipCap: Students enjoying the cherry blossoms.

ClipCap: A bunch of bananas sitting on top of a table.

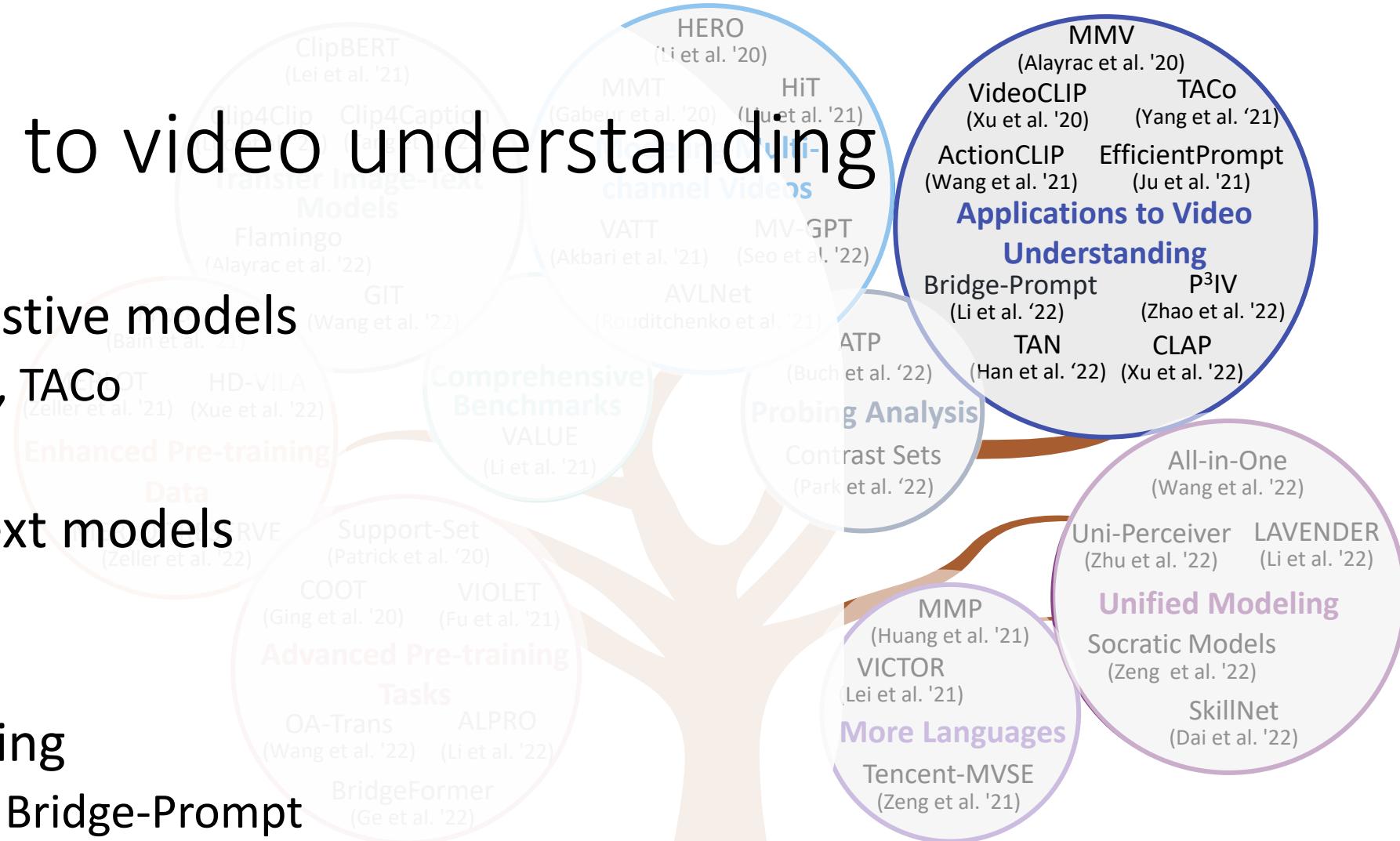
Composing Zero-Shot Multimodal Reasoning with Language

- A modular framework in which multiple pretrained models may be composed zero-shot
 - via multimodal-informed prompting, to exchange information with each other and capture new multimodal capabilities, without requiring finetuning.



Applications to video understanding

- Video-text contrastive models
 - MMV, VideoCLIP, TACo
- Transfer image-text models
 - ActionCLIP
- With Prompt Tuning
 - EfficientPrompt, Bridge-Prompt



Pioneering work in Video-Text Pre-training

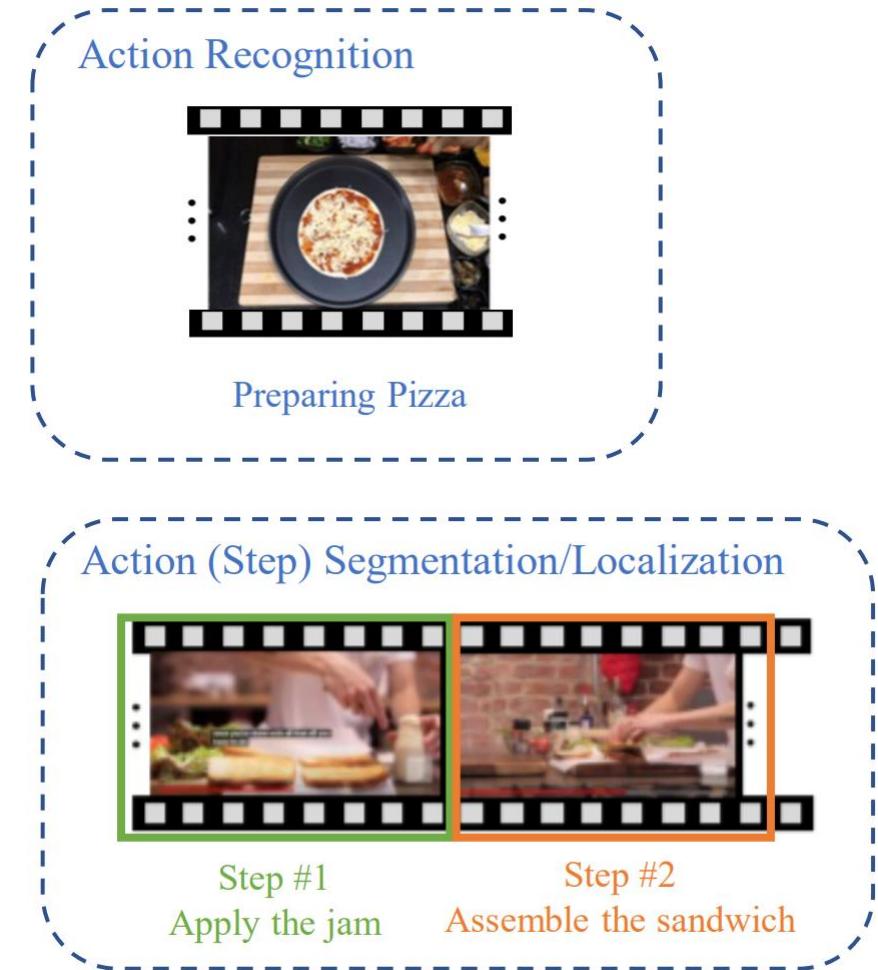
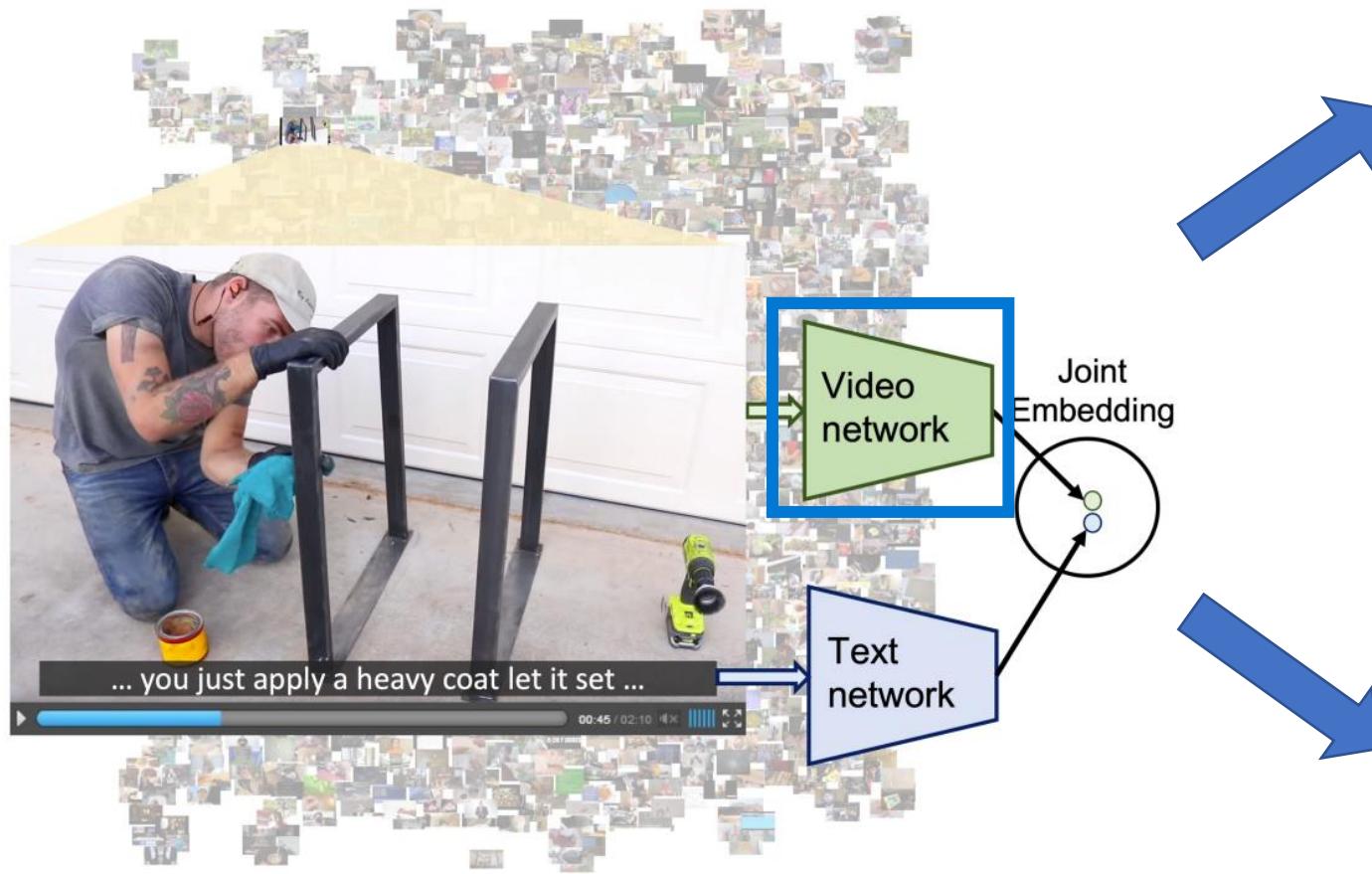
VideoBERT
(Sun et al. '19)

ActBERT
(Zhu and Yang '20)

HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)

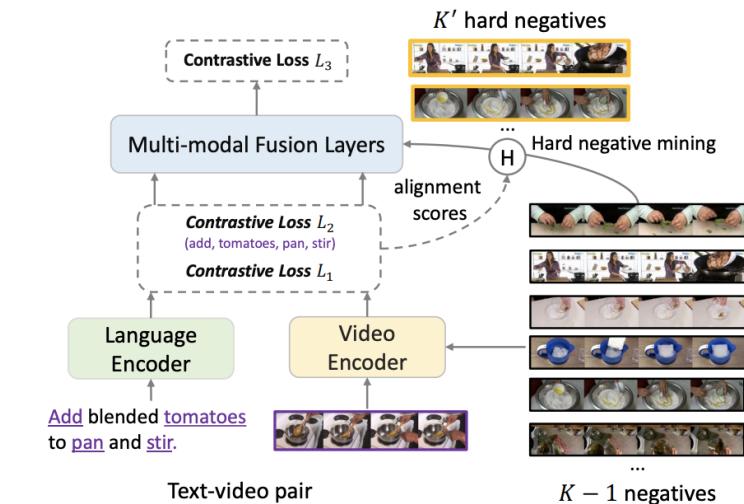
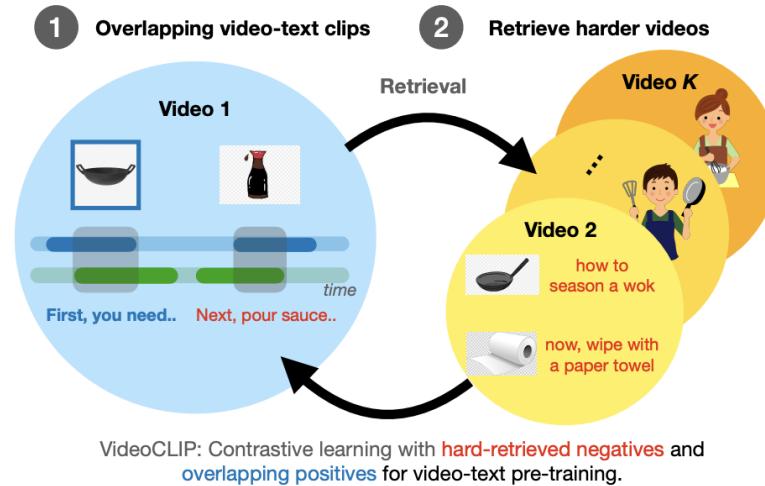
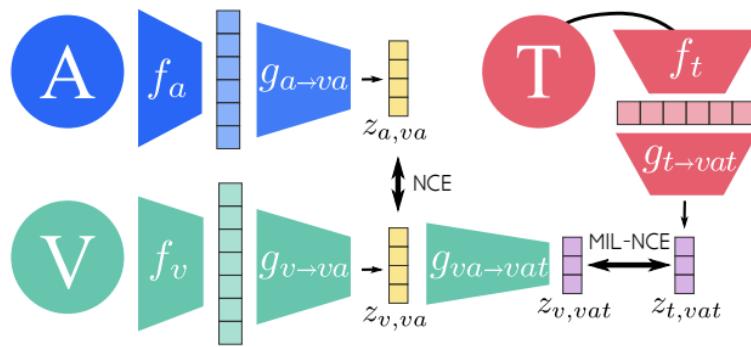
Video-text contrastive models



[HTM, Miech et al., 2019], [MIL-NCE, Miech et al., 2020]

Figure credit: Howto100m: Learning a text-video embedding by watching hundred million narrated video clips, ICCV 2019

Video-text contrastive models

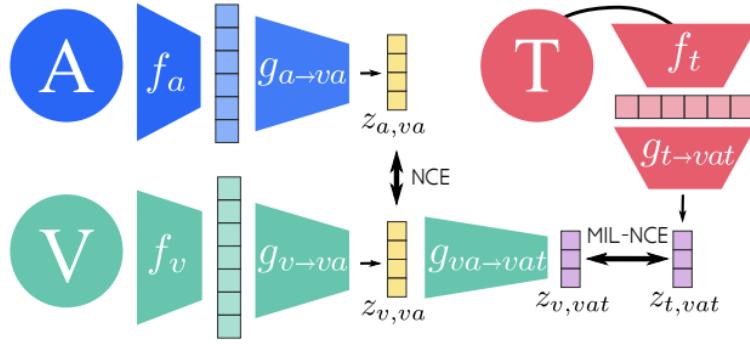


MMV: Self-Supervised MultiModal Versatile Networks,
NeurIPS 2020

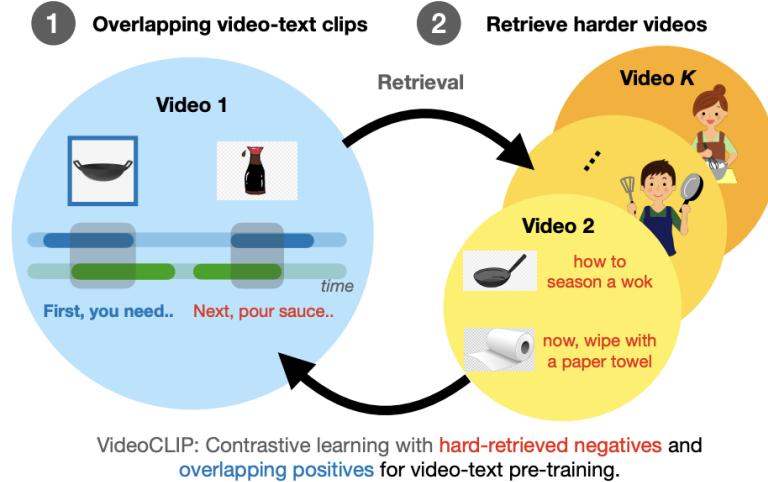
VideoCLIP: Contrastive Pre-training for Zero-shot Video-Text Understanding, EMNLP 2021

TACo: Token-aware Cascade Contrastive Learning for Video-Text Alignment, ICCV 2021

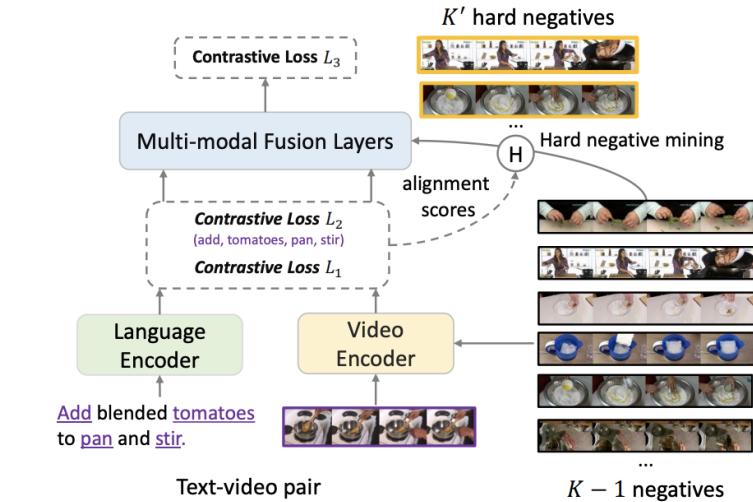
Video-text contrastive models



MMV: Self-Supervised MultiModal Versatile Networks,
NeurIPS 2020



VideoCLIP: Contrastive Pre-training for Zero-shot Video-Text Understanding, EMNLP 2021



$$\mathcal{L}(x) = \lambda_{va} \text{NCE}(x_v, x_a) + \lambda_{vt} \text{MIL-NCE}(x_v, x_t)$$

$$\mathcal{L} = - \sum_{(v,t) \in B} \left(\log \text{NCE}(z_v, z_t) + \log \text{NCE}(z_t, z_v) \right)$$

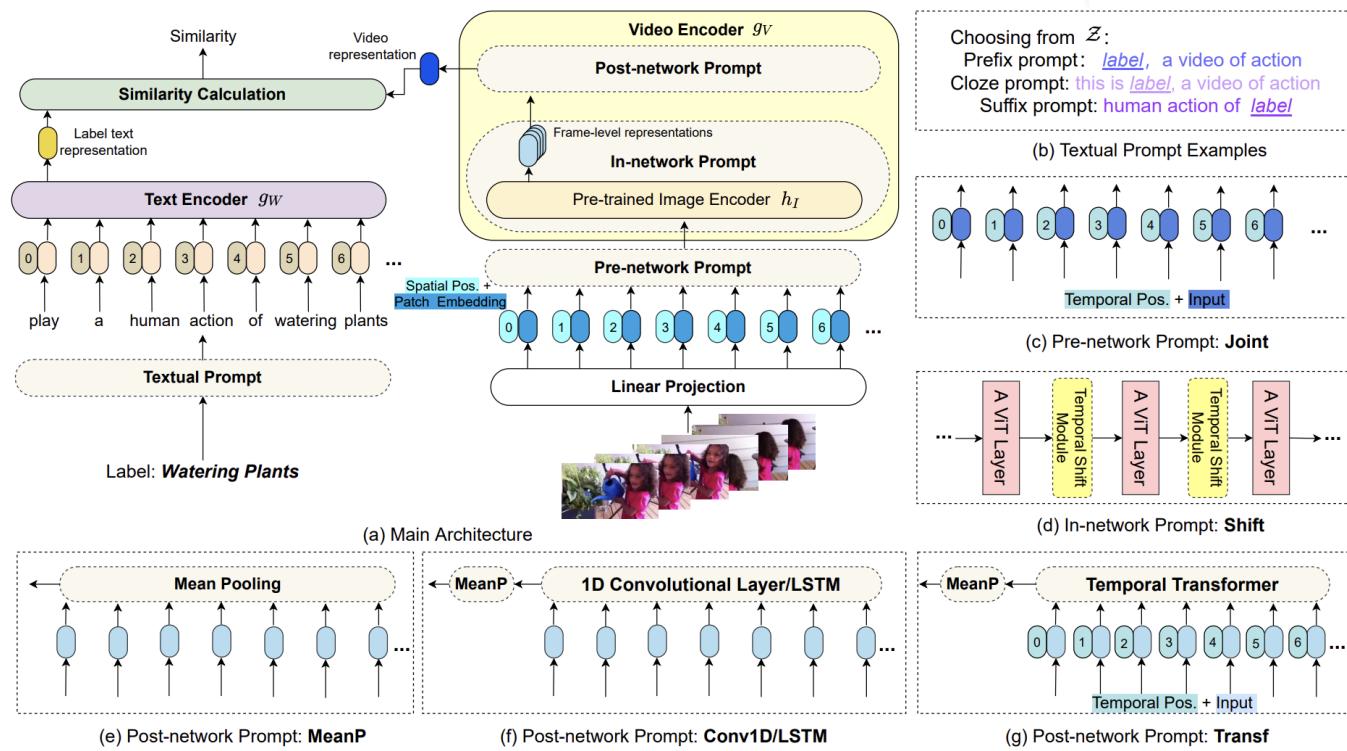
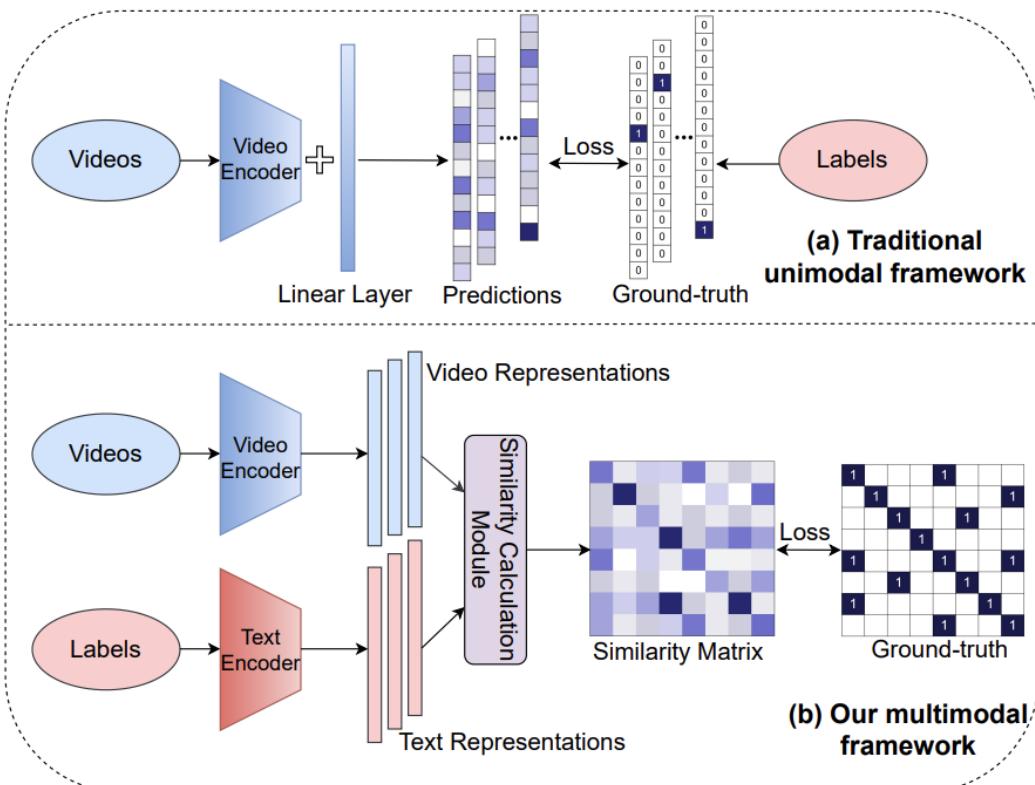
$$L_1 = - \sum_{i=1}^K \log \left(\frac{\exp^{\bar{x}_i \cdot \bar{y}_i / \tau_1}}{\exp^{\bar{x}_i \cdot \bar{y}_i / \tau_1} + \sum_{j \neq i} \exp^{\bar{x}_j \cdot \bar{y}_i / \tau_1}} \right)$$

$$L_2 = - \sum_{i=1}^K \sum_{p \in \mathcal{P}_i} \log \left(\frac{\exp^{s(x_i, y_i^p) / \tau_2}}{\exp^{s(x_i, y_i^p) / \tau_2} + \sum_{j \neq i} \exp^{s(x_j, y_i^p) / \tau_2}} \right)$$

$$L_3 = - \sum_{i=1}^K \log \left(\frac{\exp^{w \cdot z_{i,i}^{cls}}}{\exp^{w \cdot z_{i,i}^{cls}} + \sum_{j \neq i} \exp^{w \cdot z_{j,i}^{cls}}} \right)$$

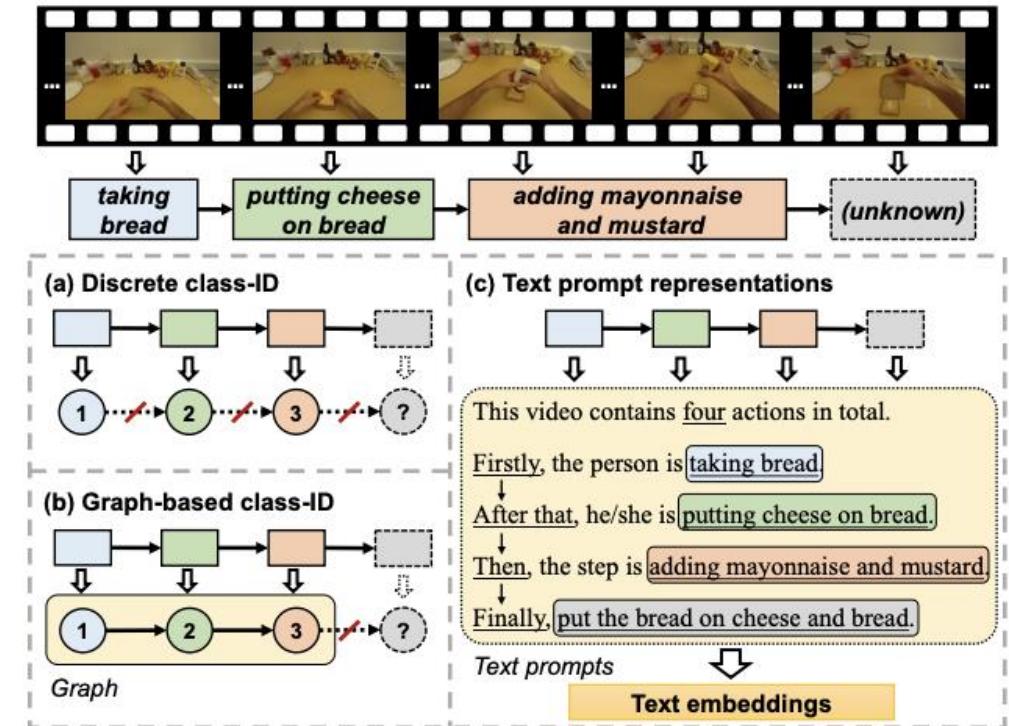
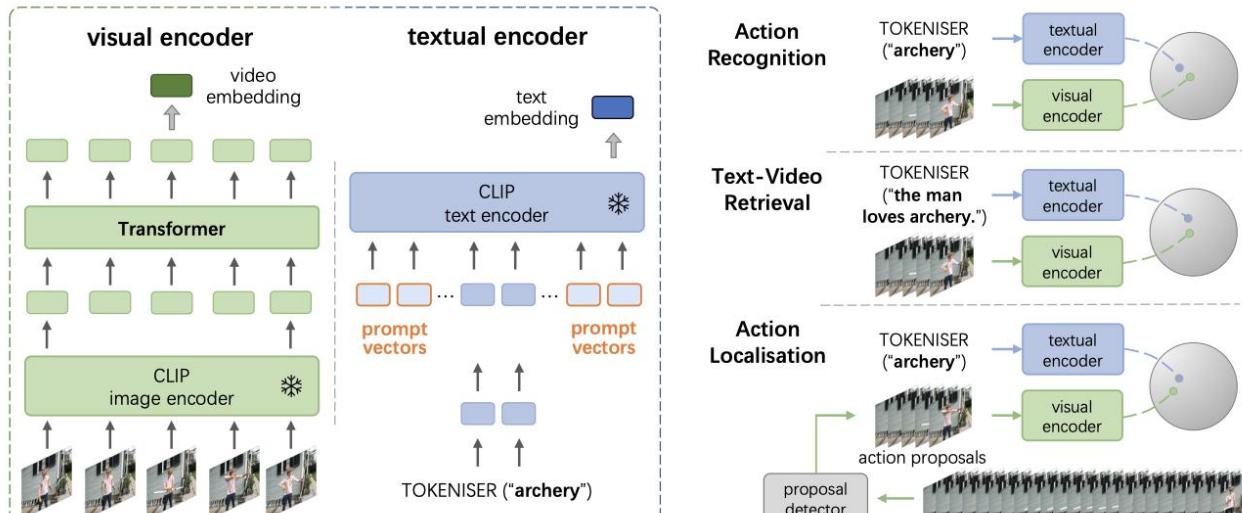
Transfer image-text models

- Attaching importance to the semantic information of label texts rather than simply mapping them into numbers.



Prompting

- Efficiently adapt one pre-trained visual-language model to novel tasks with minimal training

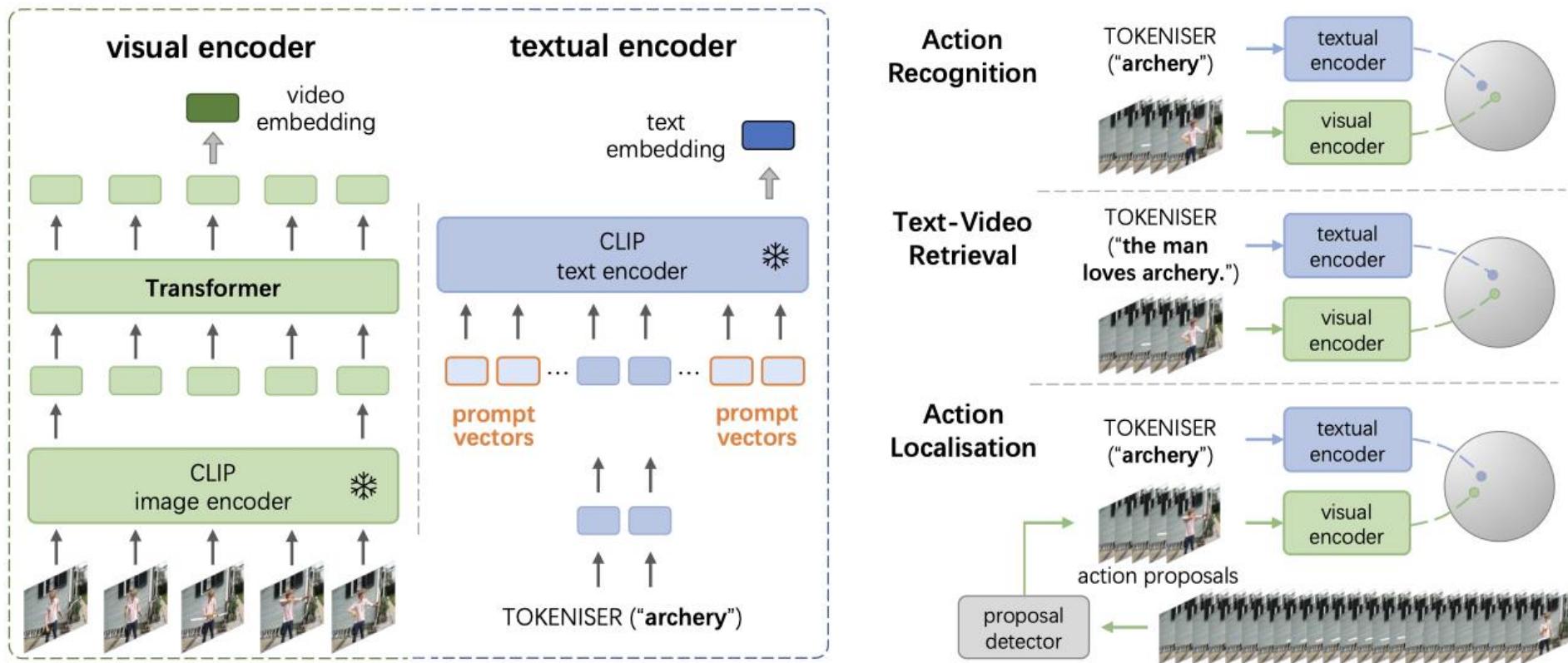


Efficient Prompting: Prompting Visual-Language Models for Efficient Video Understanding, arXiv 2021

Bridge-Prompt: Towards Ordinal Action Understanding in Instructional Videos, CVPR 2022

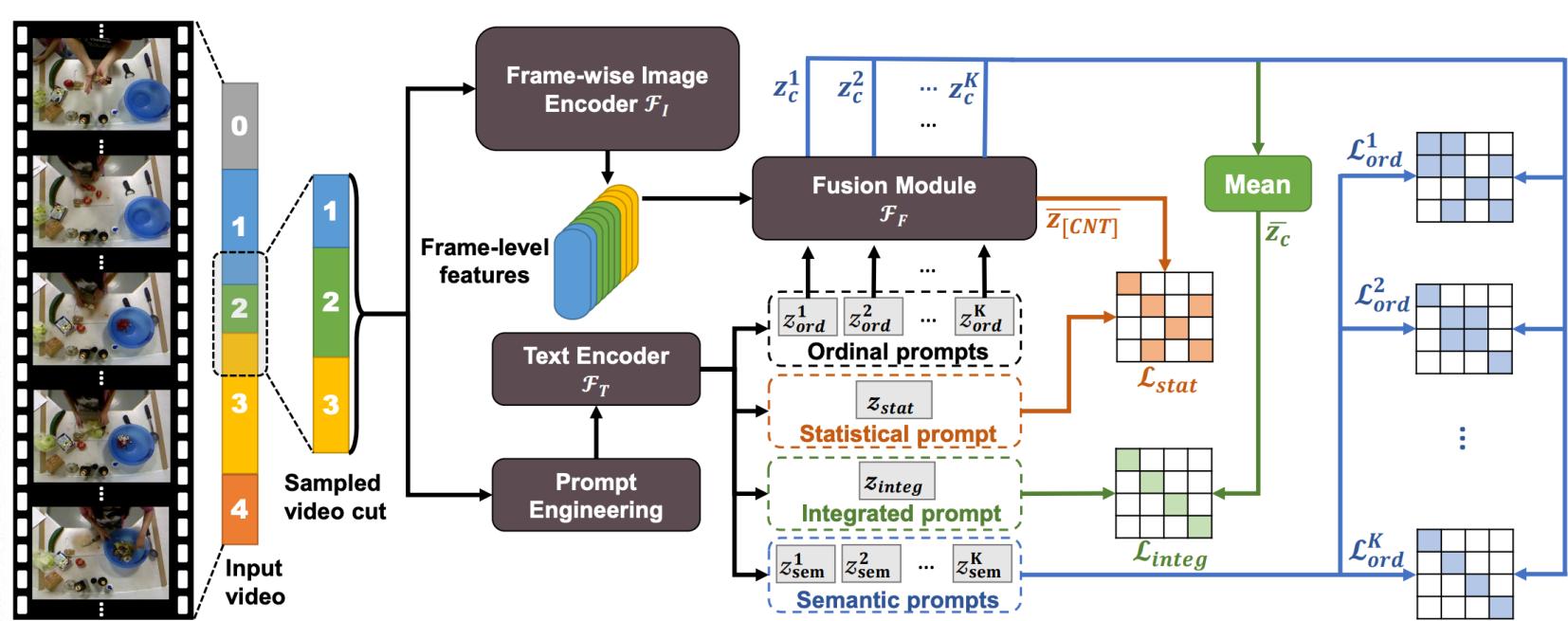
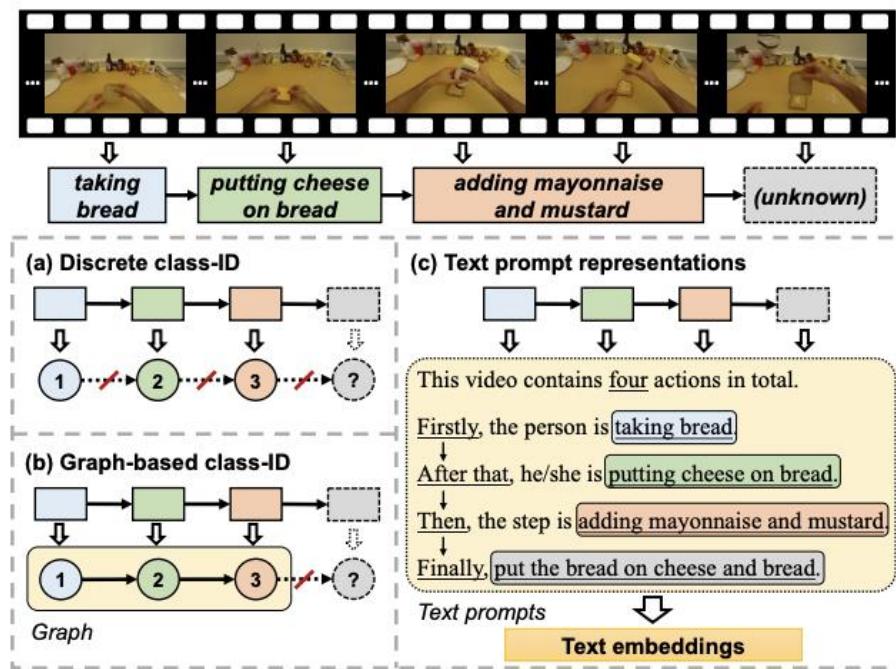
Continuous prompt vectors

- Optimize a few random vectors that convert the novel tasks into the same format as the pre-training objectives



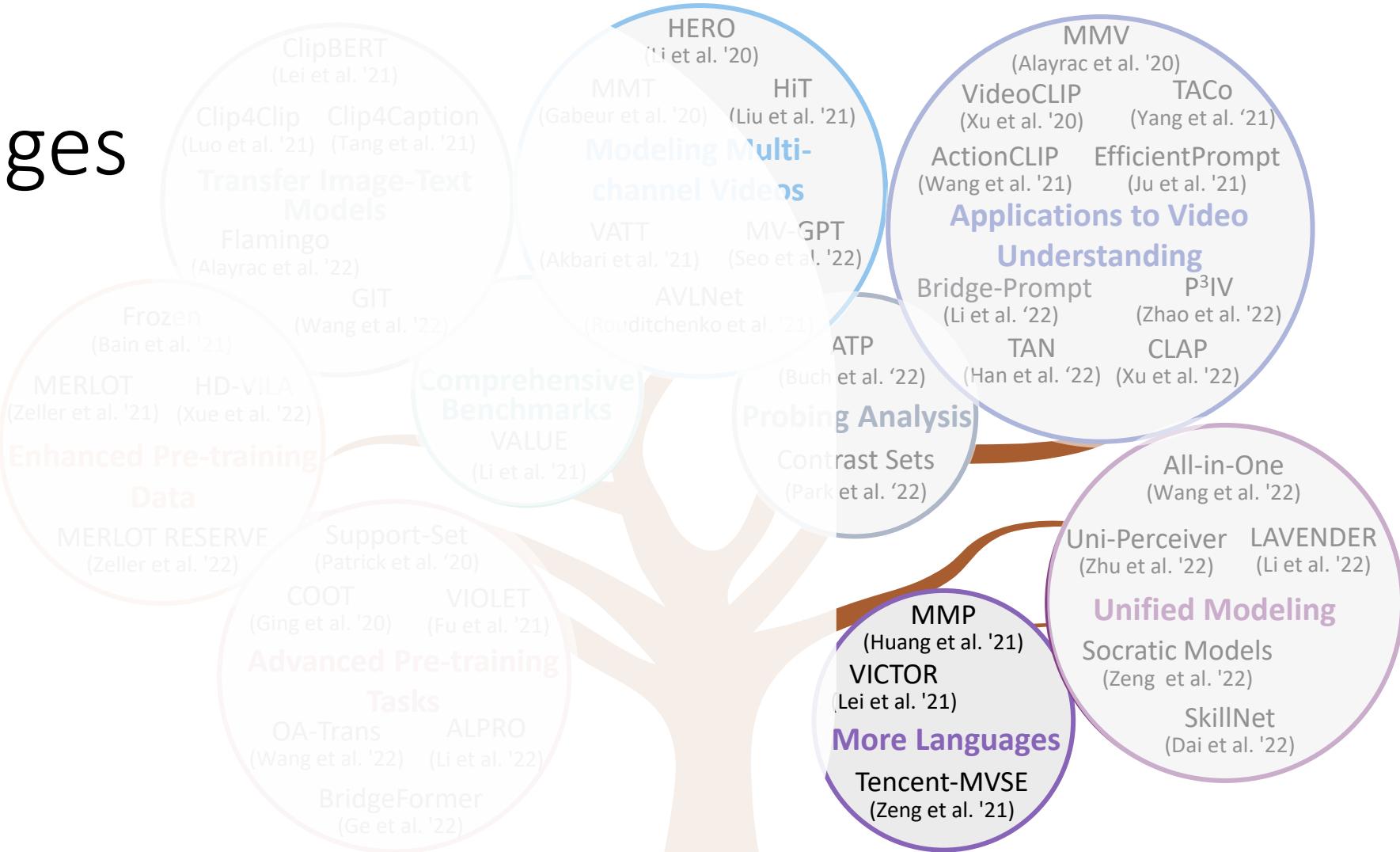
Model the semantics across adjacent actions

- Simultaneously exploits both out-of-context and contextual information from a series of ordinal actions in instructional videos



More languages

- MMP
- VICTOR
- Tencent-MVSE



Pioneering work in Video-Text Pre-training

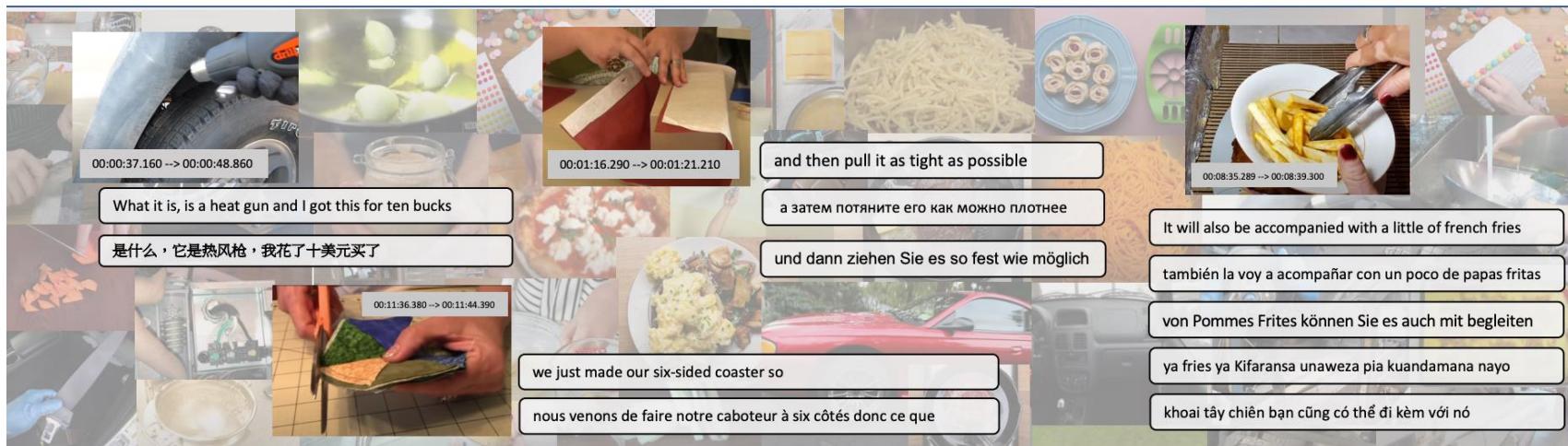
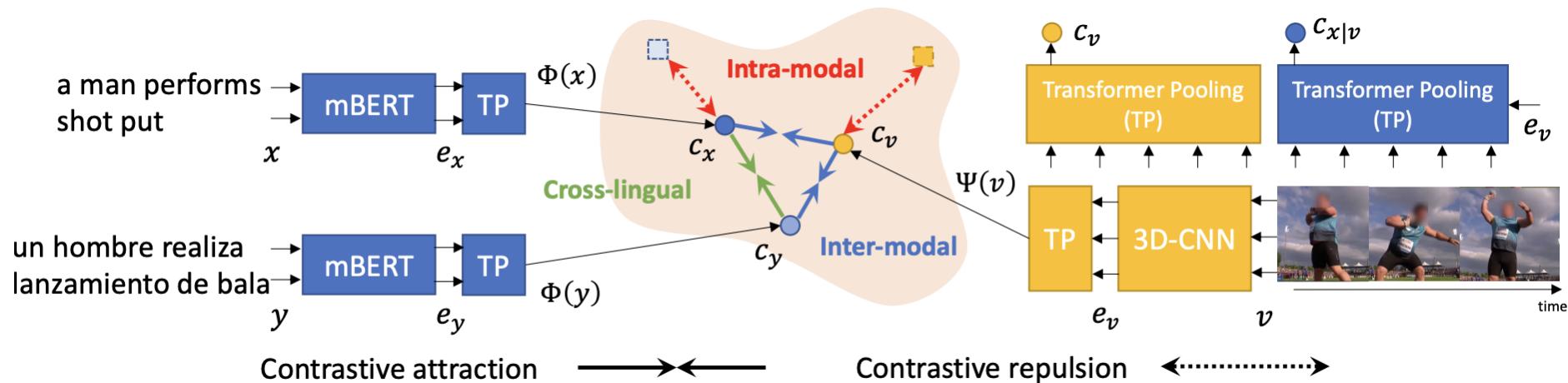
VideoBERT
(Sun et al. '19)

ActBERT
(Zhu and Yang '20)

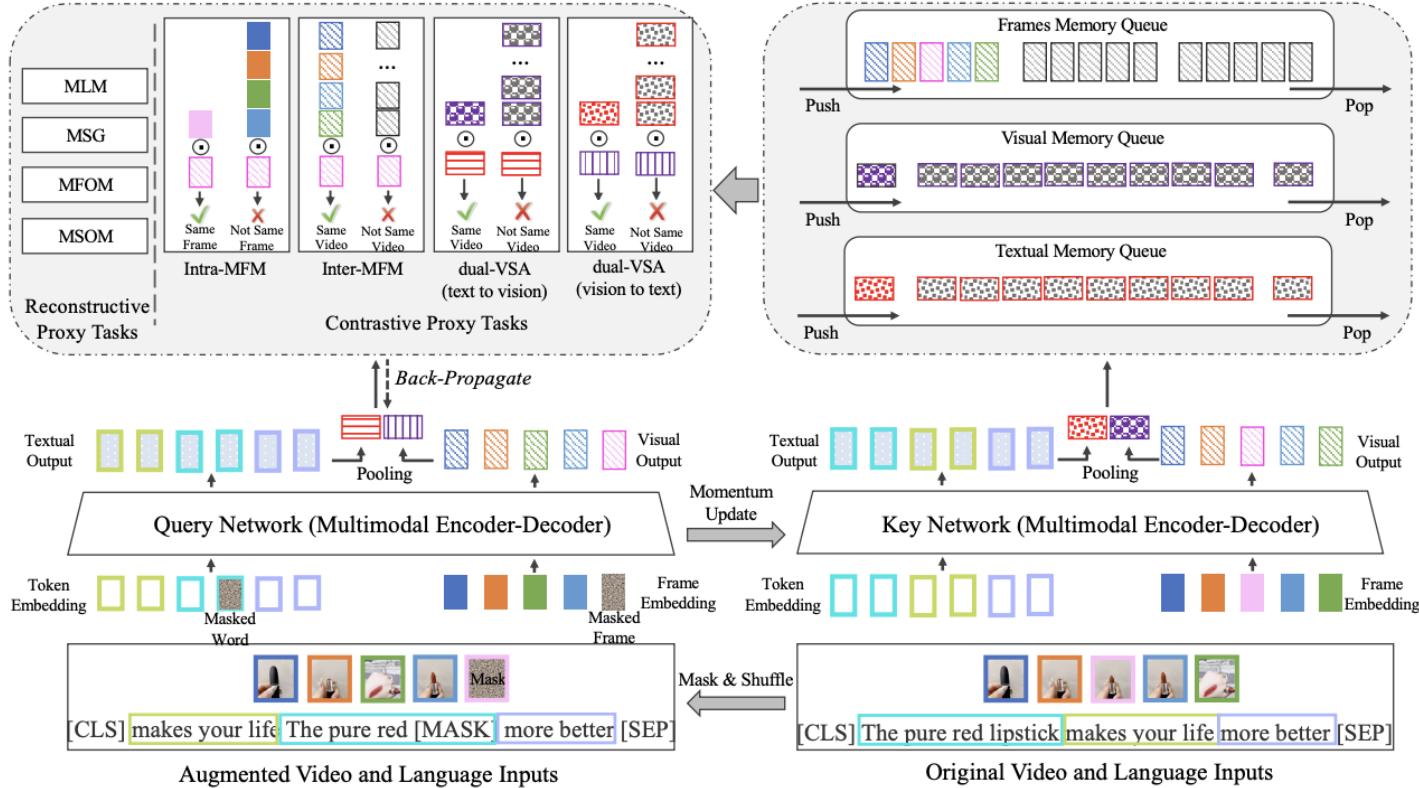
HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)

Multilingual Multimodal Pre-training



Understanding Chinese Video and Language via Contrastive Multimodal Pre-Training



A Large-Scale Benchmark Dataset for Multi-Modal Video Similarity Evaluation

梅西告诉你，点球不是非要射门。成人之美也是一种美德！(Messi tells you, we don't have to shoot when playing stop kick, helping others is also a virtue!)



面对这样的门将，连穆勒都无可奈何，足球场上运气也很重要！(Facing such a goalkeeper, even Muller has no alternative, luck is also important on the football field!)



德云社“大小姐”郭麒麟，你这是要出师了么？这么狠呢(Deyunshe “young madam” Qilin Guo, are you going to finish your apprenticeship? So cruel!)



用了这个灵敏度，不仅能压平底锅，还能压信号枪！(Using such sensitivity, not only can press pan, but also can press flare gun!)



德云社“大小姐”郭麒麟，你这是要出师了么？这么狠呢(Deyunshe “young madam” Qilin Guo, are you going to finish your apprenticeship? So cruel!)



光子，你卖的什么破伞，我要退钱(Guangzi, what broken parachute you sell, I want a refund)



Title: 和平精英：如何对待胆小的敌人？没有什么是一颗雷解决不了的
(PLUG: How to treat coward? Nothing can be solved by a grenade.)

ASR: 但他好像并没有来找我的勇气，面对这种胆小鬼只能主动出击...
(But he seems do not to have the courage to find me. Faced with this situation, I have to be proactive...)

Category: 游戏-手游 (game-mobile game)

Tag: 海岛地图；生存游戏；和平精英；军事题材；射击游戏；达人解说
(island map; survival game; PLUG; military subject; shooting game; master explanation)



Title: 大货车体验了一把VIP的待遇！一条船就装一辆车！

ASR: (Truck experiences VIP treatment! A ship only holds a truck!)

Category: 生活-生活记录 (life-life recording)

Tag: 生活随手拍；记录生活；货车；船；随记
(life snapshot; record life; truck; shop; snapshot)



Title: 烧烤小妙招，怕肉质太硬，那就用这一物，让肉质鲜嫩！

ASR: (Barbecue lifehack, afraid tough meat, you can use this thing, which can make the meat fresh and tender.)

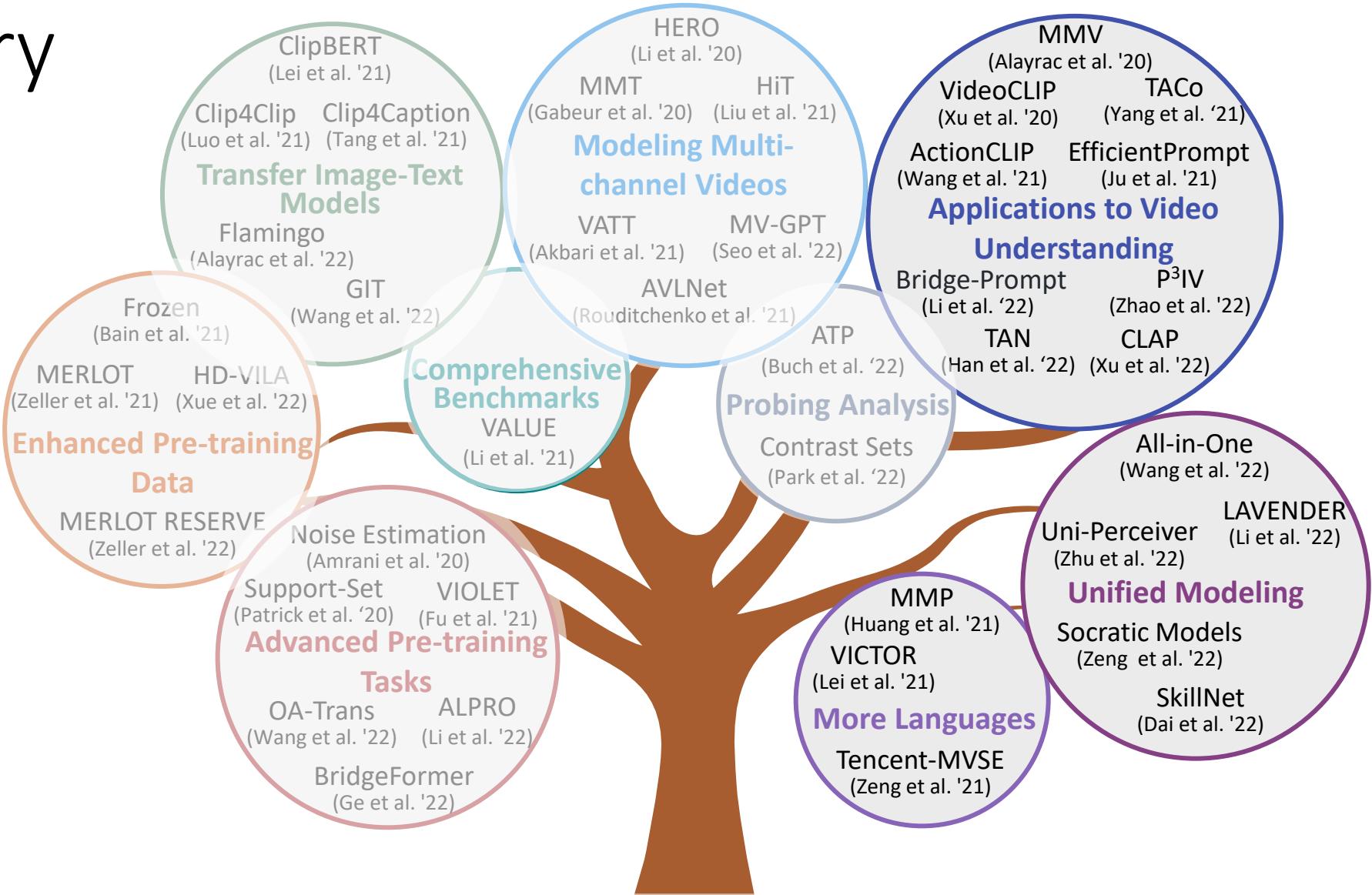
Category: 切一片洋葱圈。将肉馅放进去。能够保持肉馅的多汁

Tag: (Cut an onion ring, put the meat stuffing in, and it can keep the meat fresh and tender)

Category: 美食-菜谱 (food-menu)

Tag: (cooking lifehack; barbecue; food master; fresh and tender)

Summary



Pioneering work in Video-Text Pre-training

VideoBERT
(Sun et al. '19)

UniVL
(Luo et al. '20)

HTM
(Miech et al. '19)

MIL-NCE
(Miech et al. '20)