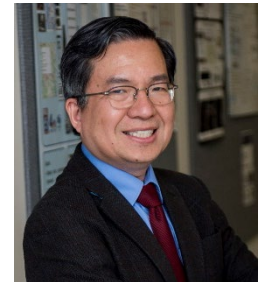
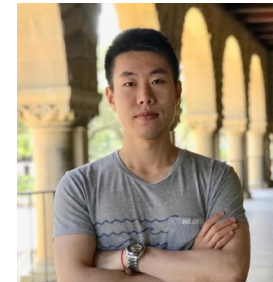
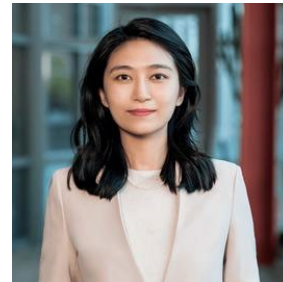
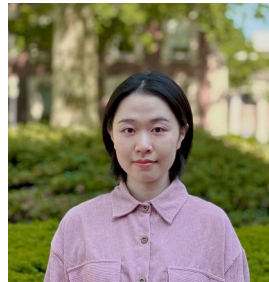




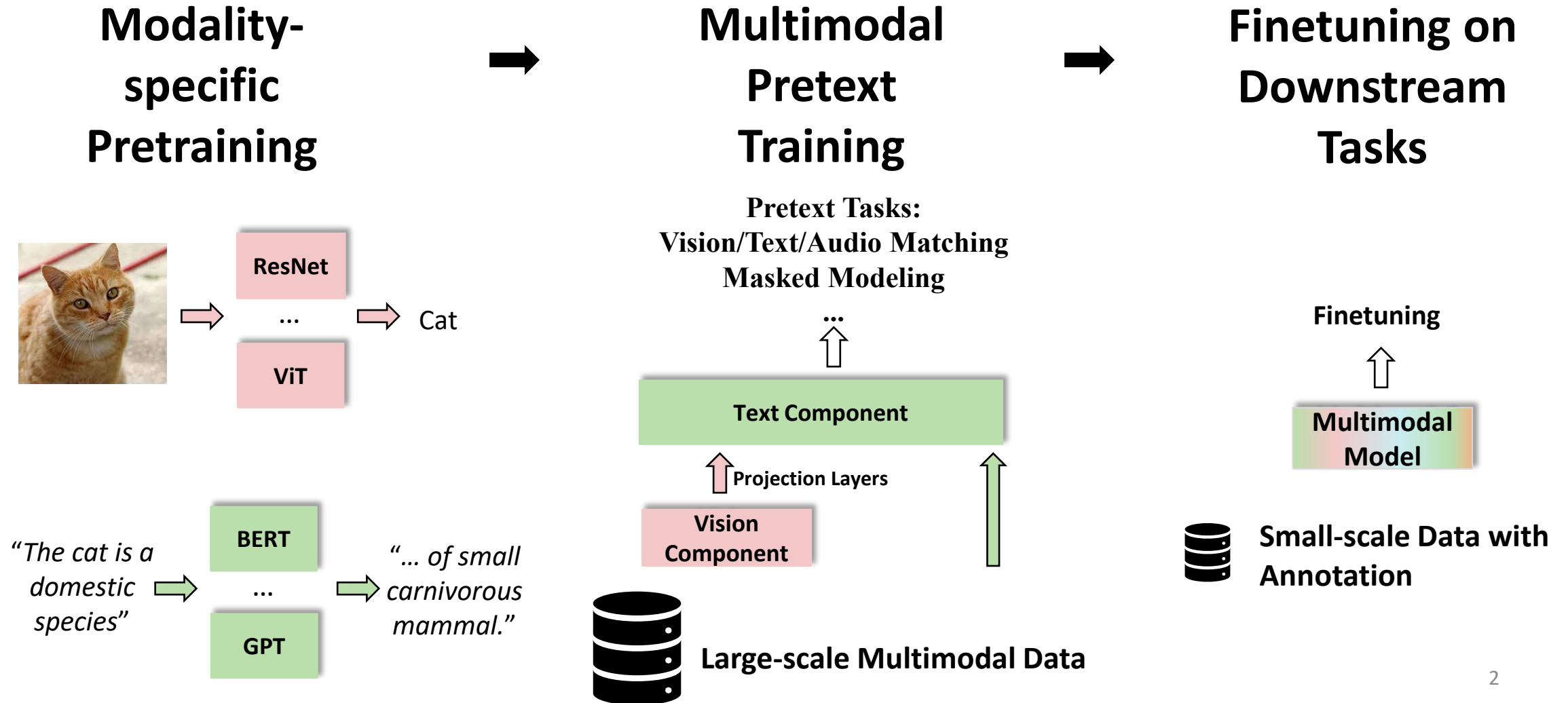
Towards Fast Adaptation of Pretrained Contrastive Models for Multi-Channel Video-Language Retrieval

Xudong Lin¹, Simran Tiwari¹, Shiyuan Huang¹, Manling Li²,
Mike Zheng Shou³, Heng Ji², Shih-Fu Chang¹

¹Columbia University ²UIUC ³NUS



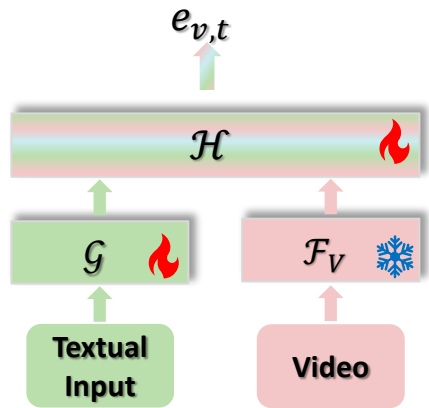
Tired of costly multimodal pretext training?



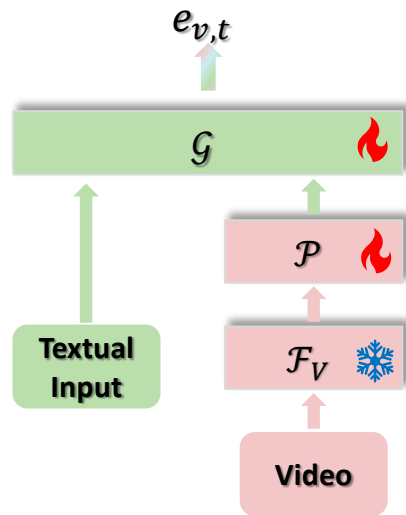
Towards Fast Adaptation of Pretrained Contrastive Models for Multi-Channel Video-Language Retrieval



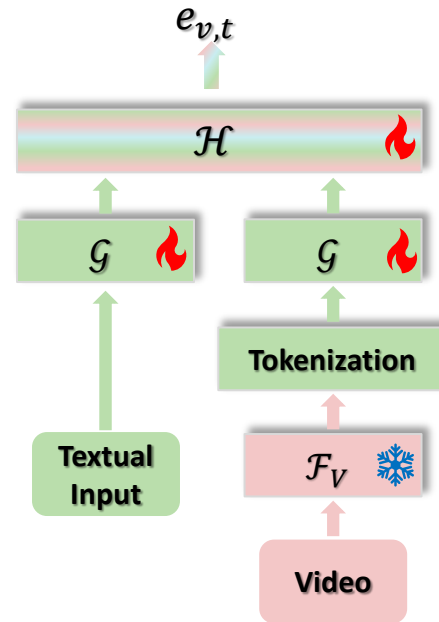
Video-language learning without costly large-scale pretext training!



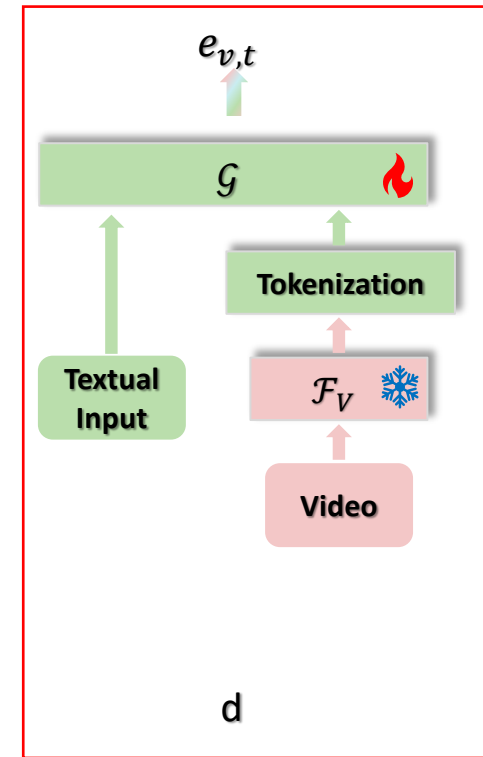
a



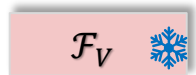
b



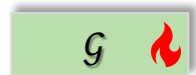
c



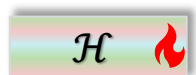
d



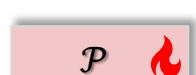
Pretrained Video Encoder



Pretrained Text Model

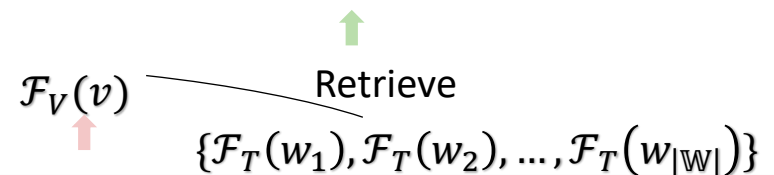


Multimodal Transformer



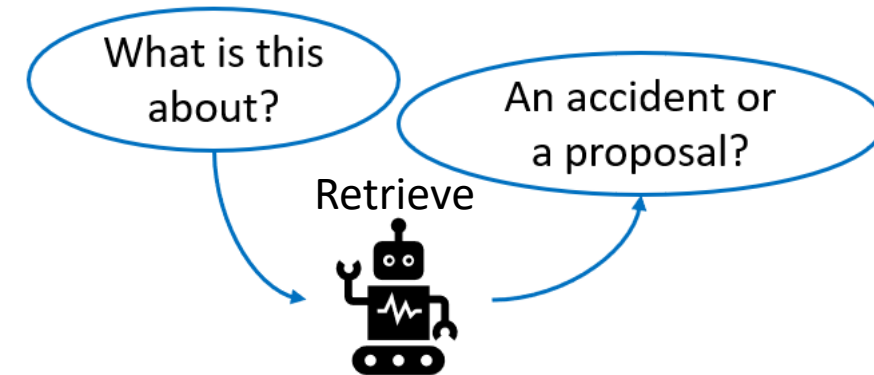
Projector

Tokenization



Multi-Channel Video-Language Retrieval

- What is Multi-Channel Video-Language Retrieval?
 - Given a multi-channel sample (e.g., video + text), retrieve a text response, or vice versa.
 - For example:
 - Open-ended video question answering
 - Text-video retrieval
- Key requirements:
 - Fusing multi-channel information
 - Discriminative representations



Video

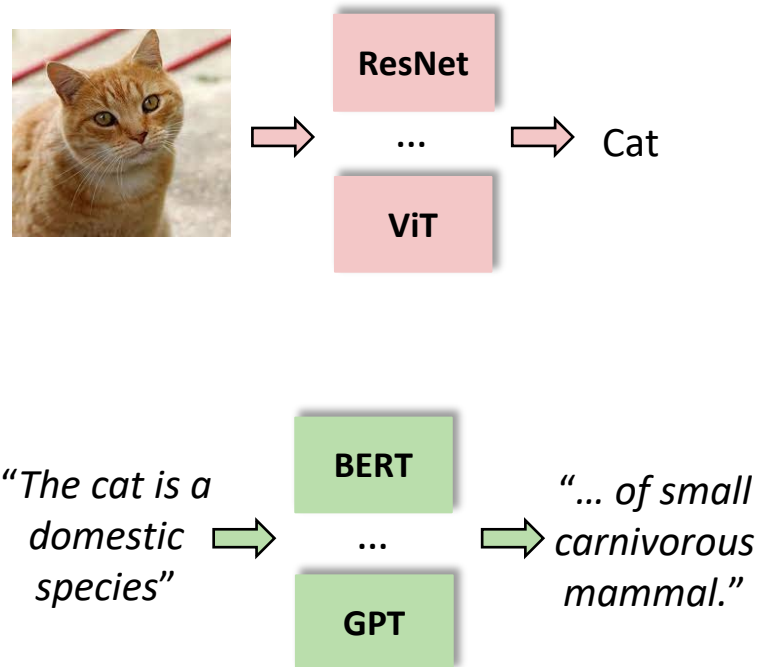
Speech *"...I have an important question to ask my girlfriend..."*

"...Robin, the ring..."

"...oh my..."

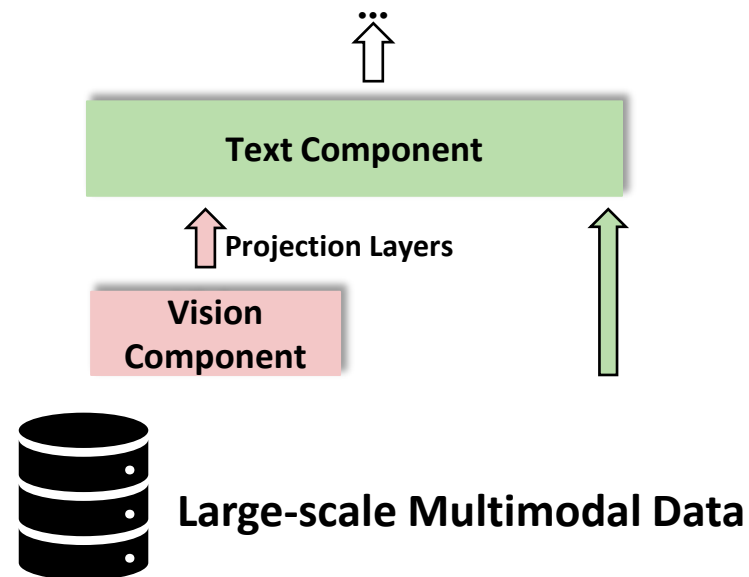
Existing Paradigm

Modality-specific Pretraining

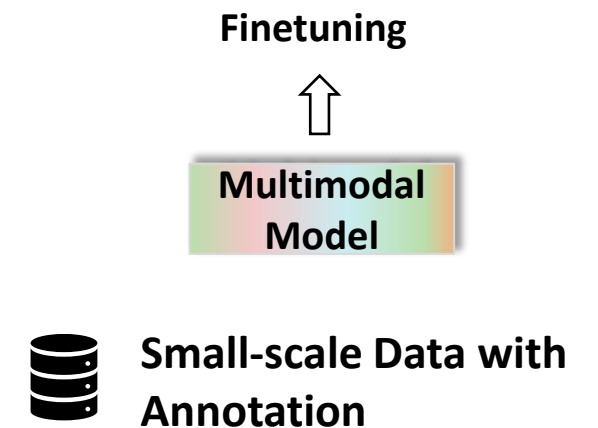


Multimodal Pretext Training

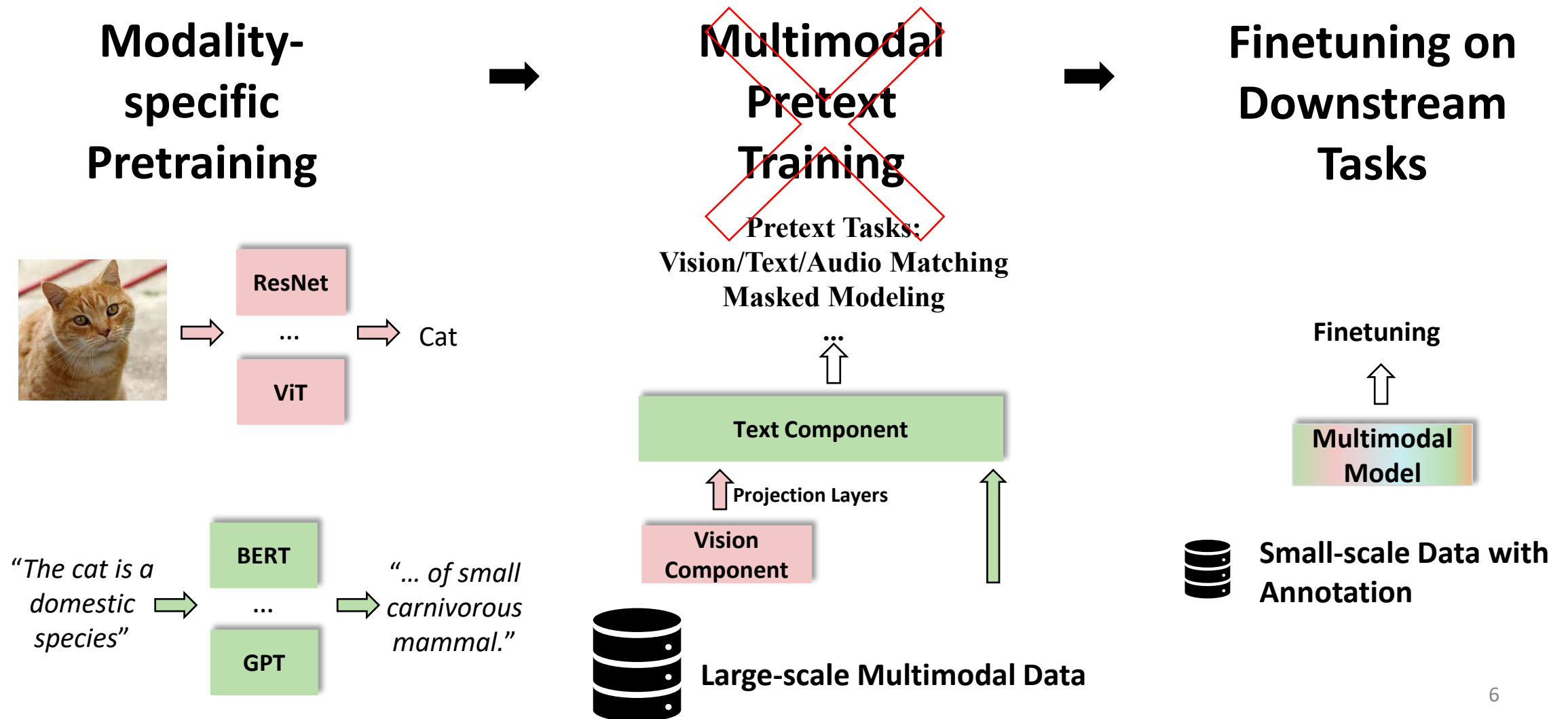
Pretext Tasks:
Vision/Text/Audio Matching
Masked Modeling



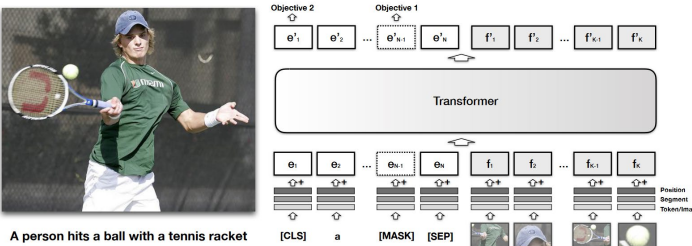
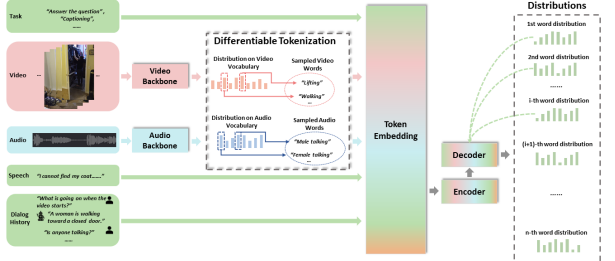
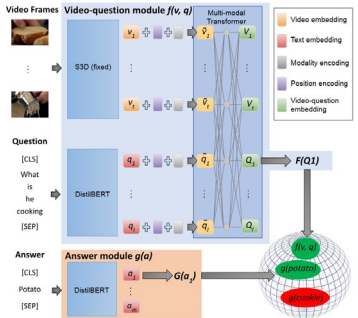
Finetuning on Downstream Tasks



Research Question: Can we remove the costly multimodal pretext training stage?



Design Space of Multimodal Models

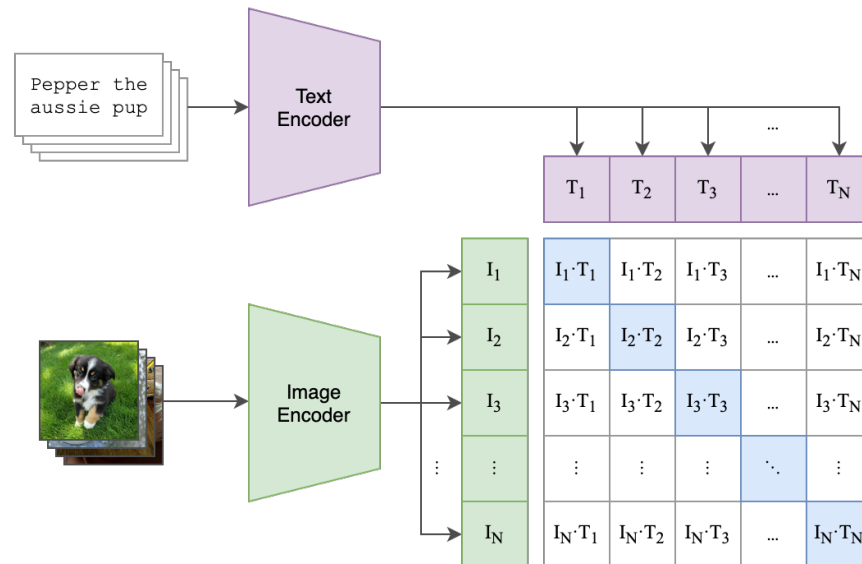
<p>Model Design Space</p>	<p>What is the representation of non-textual modalities, Continuous features or text tokens?</p>	
<p>What is the fusion strategy, Text Transformer Or Multimodal Transformer?</p>	<p>Continuous + Text Trans.</p>  <p>(Li, Liunian Harold, et al., 2019)</p>	<p>Text tokens + Text Trans.</p>  <p>(Lin, Xudong, et al., 2021)</p>
	<p>Continuous + MM Trans.</p>  <p>(Yang, Antoine, et al., 2021)</p>	<p>Text tokens + MM Trans.</p>

Li, Liunian Harold, et al. "Visualbert: A simple and performant baseline for vision and language." *arXiv preprint arXiv:1908.03557* (2019).
 Lin, Xudong, et al. "Vx2Text: End-to-End Learning of Video-Based Text Generation From Multimodal Inputs." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021.
 Yang, Antoine, et al. "Just ask: Learning to answer questions from millions of narrated videos." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021.

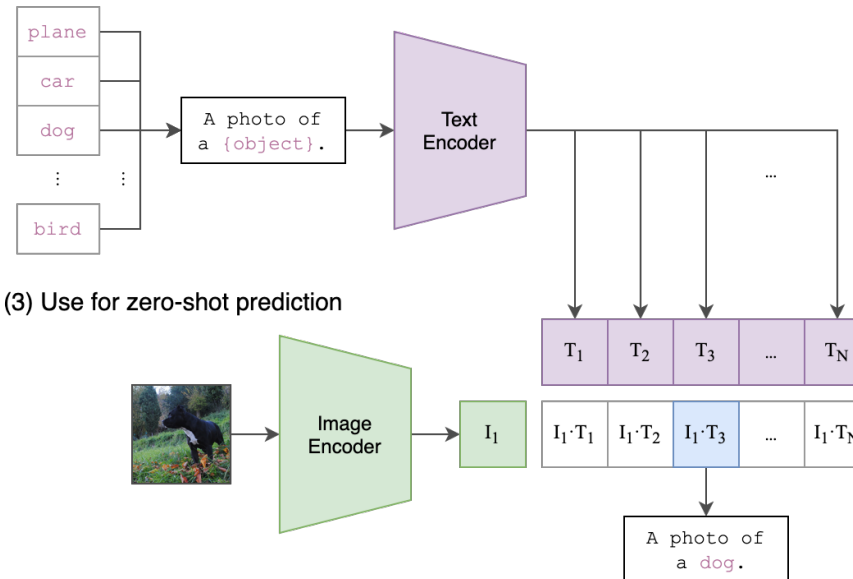
Background: Multimodal Contrastive Models

- Assumption: we have access to multimodal contrastive models like CLIP (Radford, Alec, et al., 2021) to perform zero-shot visual concept extraction.

(1) Contrastive pre-training

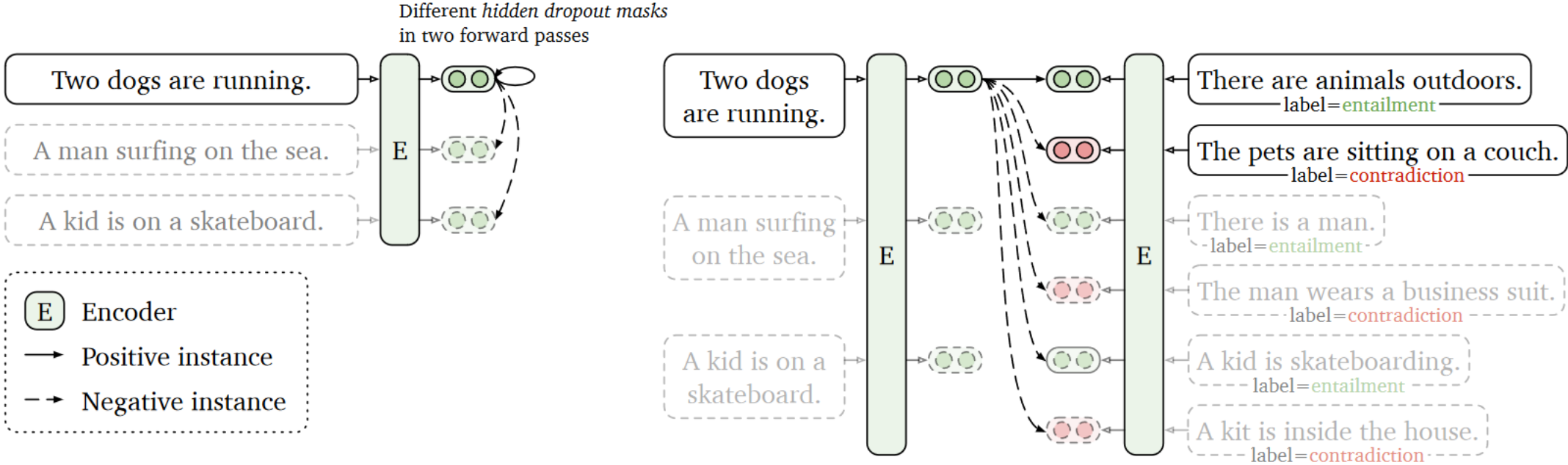


(2) Create dataset classifier from label text



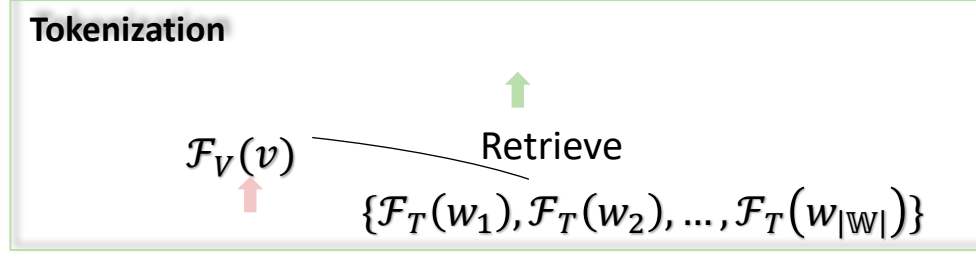
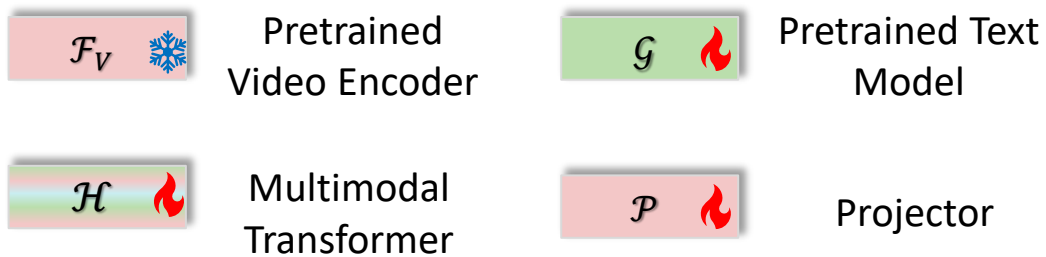
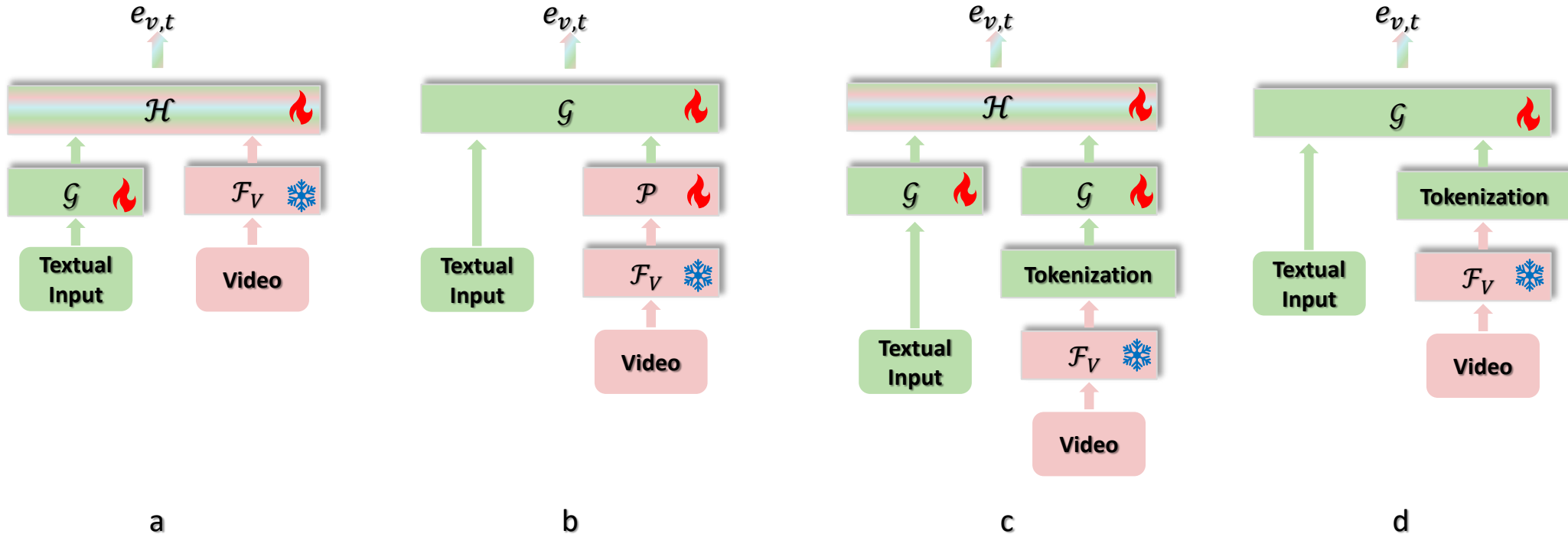
Background: Contrastive Text Model

- Assumption: we have access to pretrained contrastive text models like SBERT/SimCSE to produce discriminative representations.



Reimers, Nils, and Iryna Gurevych. "Sentence-bert: Sentence embeddings using siamese bert-networks." *arXiv preprint arXiv:1908.10084* (2019).
 Gao, Tianyu, Xingcheng Yao, and Danqi Chen. "Simcse: Simple contrastive learning of sentence embeddings." *arXiv preprint arXiv:2104.08821* (2021).

Towards Fast Adaptation of Pretrained Contrastive Models



Experimental Results

- Setting

- Task type:

- Open-ended, given question and video, retrieve the most relevant answer from the collection of all answers in the dataset.
 - Multiple choice
 - Text-video retrieval

Model	Open-ended Acc(%)		Multi-Choice Acc(%)	Retrieval AveR (%)	
	iVQA	ActivityNet-QA	How2QA	YouCook II	VATEX
Conti. + Multi.	22.4	36.9	79.2	41.9	69.4
Conti. + Text	23.2	37.3	80.4	46.2	72.7
Text + Multi.	23.4	37.1	79.4	40.4	67.5
Text + Text	31.6	38.7	82.9	49.7	74.8

Experimental Results

- Setting

- Task type:

- Open-ended, given question and video, retrieve the most relevant answer from the collection of all answers in the dataset.
 - Multiple choice
 - Text-video retrieval

Model	Open-ended Acc(%)		Multi-Choice Acc(%)	Retrieval AveR (%)	
	iVQA	ActivityNet-QA	How2QA	YouCook II	VATEX
Conti. + Multi.	22.4	36.9	79.2	41.9	69.4
Conti. + Text	23.2	37.3	80.4	46.2	72.7
Text + Multi.	23.4	37.1	79.4	40.4	67.5
Text + Text	31.6	38.7	82.9	49.7	74.8

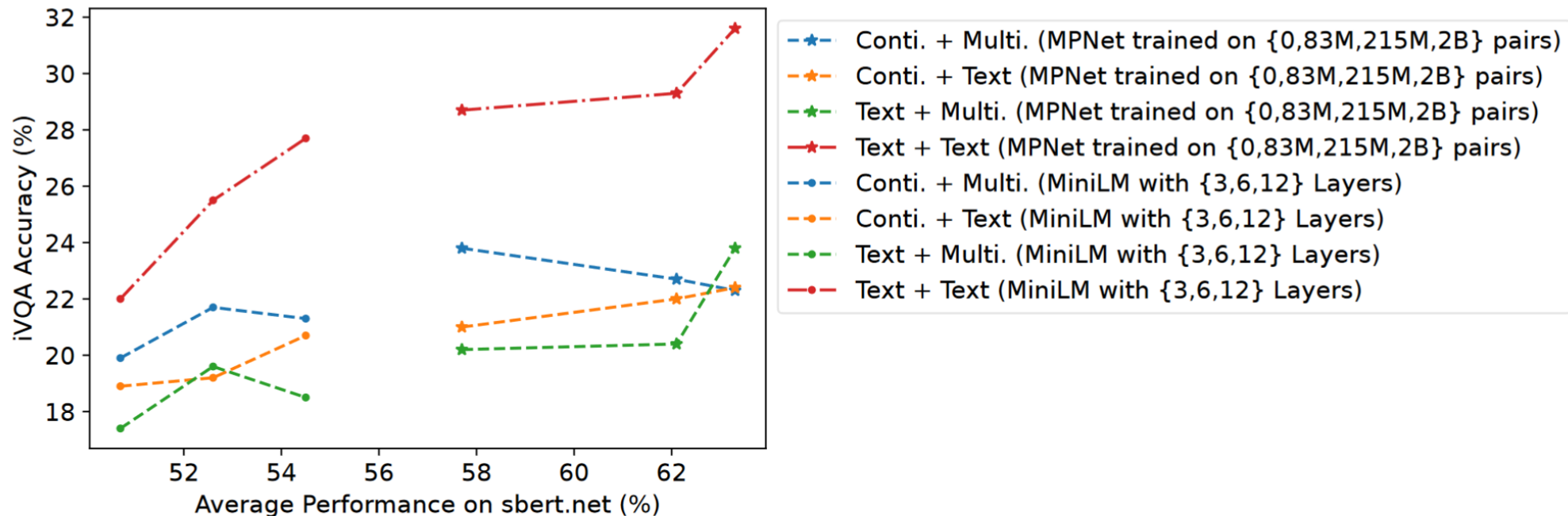
The Text + Text variant consistently achieves superior performance!

Comparison to SOTA

- Comparable or even better performance without extra multimodal pretext samples for training!
- Easily upgraded w.r.t. multimodal contrastive models!

Model	$\mathcal{F}_V, \mathcal{F}_T$	Extra MM Samples	Δ GPU hours	iVQA	ActivityNet	How2QA
MERLOT [38]	-	180M	-	-	41.4	-
SiaSamRea [36]	-	5.6M + 80K	-	-	39.8	84.1
VQA-T [33]	S3D [15]	69M + 3M	350 + 30	35.2	39.0	85.3
Conti. + Multi.	S3D [15]	69M	400	35.4	38.9	84.4
Conti. + Multi. (+ ASR)	S3D [15]	69M	400	36.0	38.9	84.8
Text + Text (Ours)	S3D [15]	0	0	31.6	38.7	82.9
Text + Text (+ ASR, Ours)	S3D [15]	0	0	36.8	38.8	84.6
FrozenBiLM [34]	CLIP [19]	10M	160	39.7	43.2	81.5
FrozenBiLM [34] (+ ASR)	CLIP [19]	10M	160	39.6	43.2	86.7
Text + Text (Ours)	CLIP [19]	0	0	36.9	41.4	92.4
Text + Text (+ ASR, Ours)	CLIP [19]	0	0	40.2	41.4	93.2

Ablation: how much do we gain by using better pretrained textual transformers?



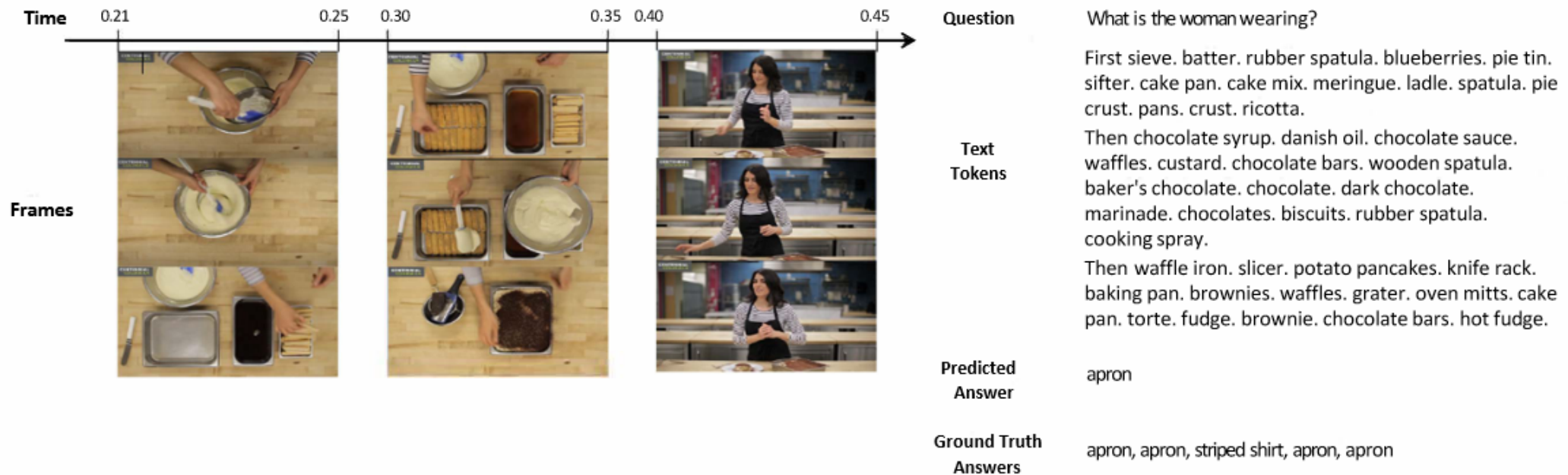
Ablation: what is the best vocabulary?

- Answer words produce the best results.


Vocabulary Source	iVQA Accuracy (%)
60K frequent English words	27.3
Visual Genome Objects and Attributes (VG-OA)	29.3
Words from all answers in iVQA training set (AW)	31.6
AW & VG-OA	29.5

Qualitative Results

- For more than 64% of iVQA samples, retrieved words include the ground-truth answer;
- When not directly retrieving the answer word, the model still learns to reason the correct answer.



Take-home messages

- Using text tokens of visual concepts to represent the visual information and using the pretrained text model as the modality fuser, a comparable or even better multimodal model can be obtained without costly multimodal pretraining!
- We hope this sheds light to upgradable multimodal intelligence.
- *We release our code in the repo.* 
- *Welcome to our poster session on Wednesday afternoon.*
- *Thanks!*