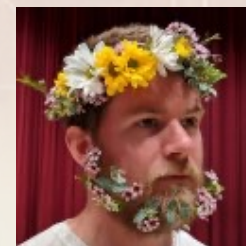
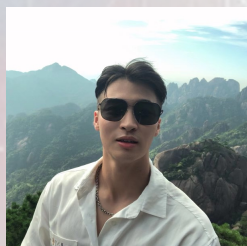
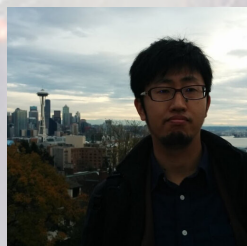


# Adapting Large Vision-Language Models to Visually-Aware Conversational Recommendation

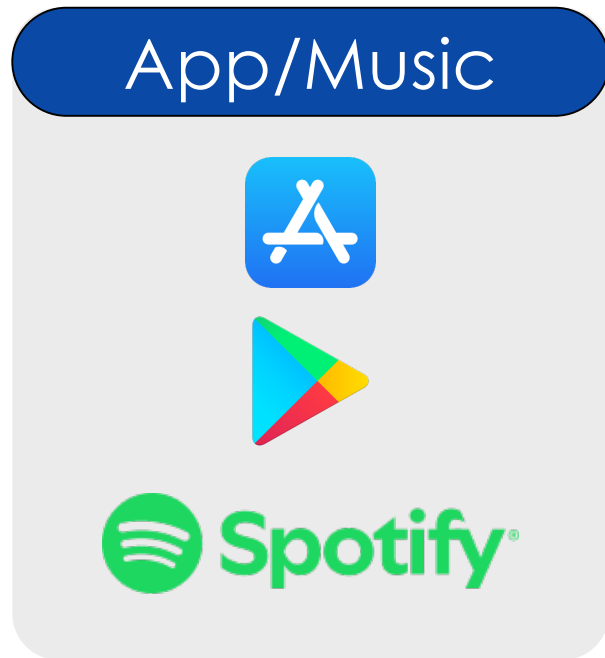
**Hyunsik Jeon**<sup>1</sup>, Satoshi Koide<sup>2</sup>, Yu Wang<sup>1</sup>, Zhankui He<sup>3</sup>, Julian McAuley<sup>1</sup>

*<sup>1</sup> UC San Diego, <sup>2</sup> Toyota Research, <sup>3</sup> Google DeepMind*



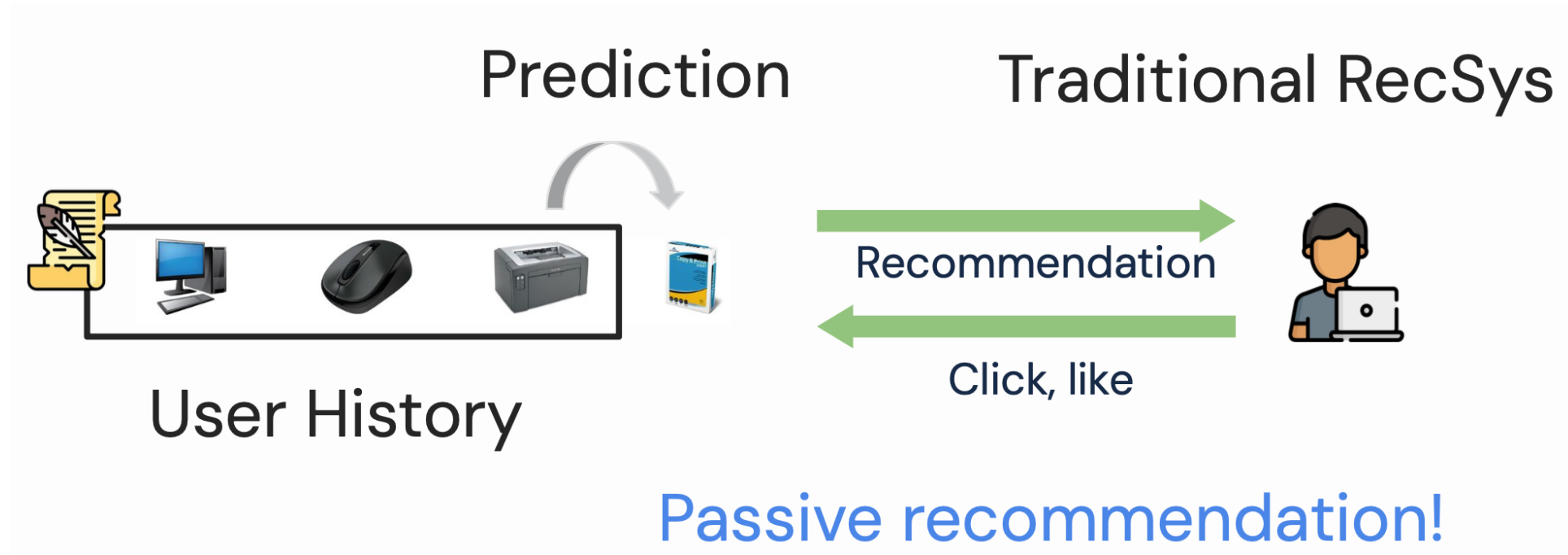
# Recommender Systems are Everywhere

- They deeply influence our daily choices and experiences



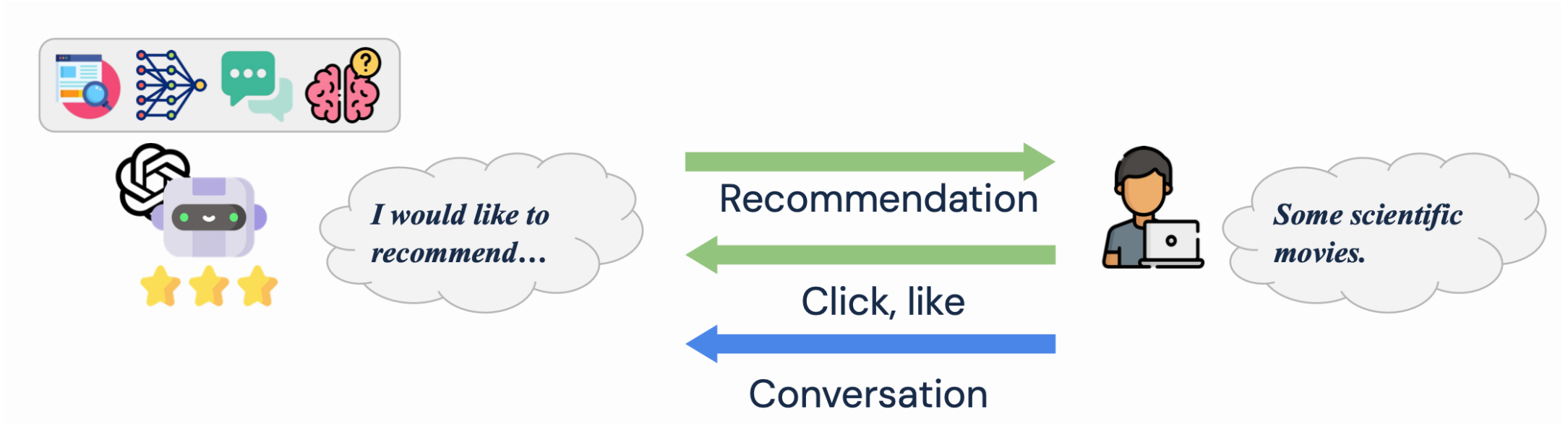
# Traditional Recommendation

- **Traditional RecSys** leverages passive signals



# Conversational Recommendation


- **Conversational RecSys** engages users in an interactive loop




# Conversation Brings Revolution

- Interactive recommendation using natural language dialogues

Can you find me a **mobile phone** on Amazon?  
Sure, what **operating system** do you prefer?  
I want an **Android** one.  
OK, and any preference on **screen size**?  
Better larger than **5 inches**.  
Do you have requirements on **storage capacity**?  
I want it to be at least **64 Gigabytes**.  
And any preference on **phone color**?  
**Not particularly**.  
Sure, then what about the following choices?



I don't like them very much...  
OK, do you have any preference on the **brand**?  
Better be **Samsung or Huawei**.  
Any requirement on **price**?  
Should be **within 700 dollars**.  
OK, then what about these ones?



Great, I want the first one, can you order it for me?  
Sure, I have placed the order for you, enjoy!

Hi

Hi! What kind of movies do you like?

I am looking for a movie recommendation. When I was younger, I really enjoyed the *A Nightmare on Elm Street (1984)*

I also enjoyed watching *The Last House on the Left(1972)*.

Oh, you like scary movies? I recently watched *Happy Death Day(2017)*. It was good for a new "scary movie".

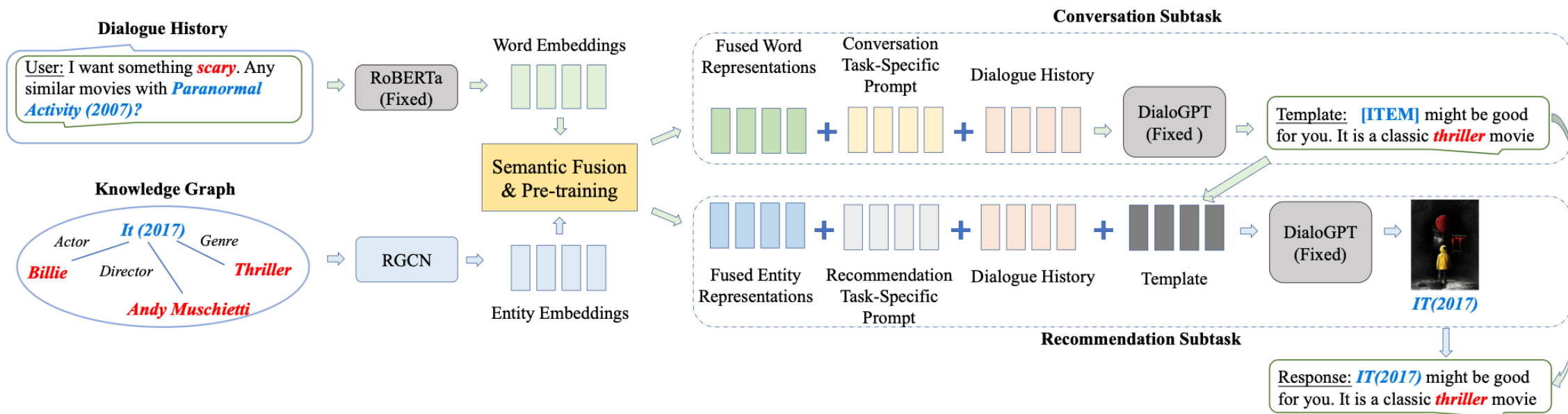
I do enjoy some of newer horror movies that I have seen as well.

I heard that *A Quiet Place (2018)* is good. It is still in theaters though.



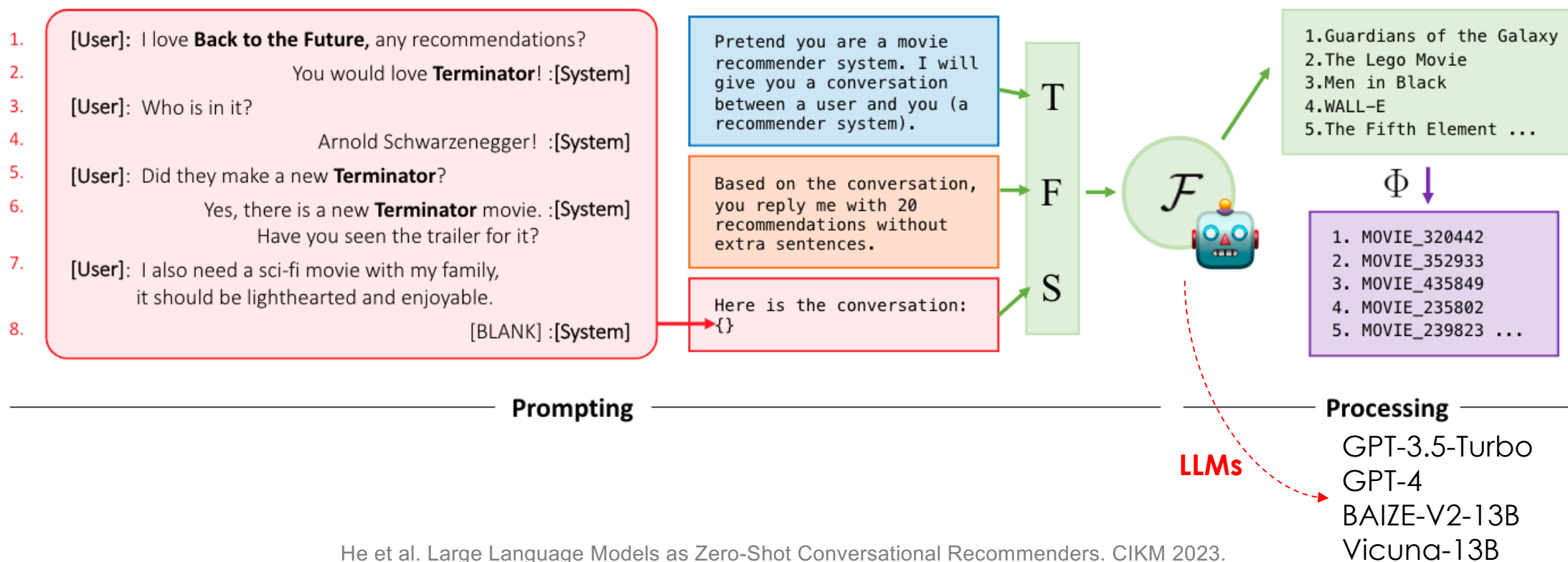
# Early Works Before the Era of LLM

- **Knowledge graph-based approaches (e.g., UniCRS)**
  - Consists of conversation/recommendation modules while utilizing external knowledge



# LLM as Conversational Recommender

- **LLM as a zero-shot recommender**: it is a single unified body of conversation, recommendation, and knowledge



# Research Motivation

- **Users can request recommendations based on visual features**

**[Seeker]:** I need help to find a similar jacket.

I got a red/burgundy jacket from Shopko a year or 2 before they closed.

The fabric was more like a hoodie but had button chest pockets on the shoulders.

The inner fabric was a super soft and warm fleece. It was the perfect jacket for fall or spring.

I'm looking for a hoodie-like military-style jacket with chest pockets.

- However, existing CRS cannot provide relevant items unless they have a visual ability



# Research Motivation

- **Visual information (image) matters**

- Images capture crucial details (design, color), especially in visually-driven domains
- Leveraging images can greatly enhance recommendation performance
- However, existing CRS rely on textual features (title)

**[Title]:** Levi's Men's Soft Shell Two Pocket Sherpa Lined Hooded Trucker Jacket

**[Images]:**



An example of  
Amazon product

# Research Motivation

**How can we design a visually-aware CRS solution?**

# Visually-Aware Conversational Recommendation

- **Given**

- A multi-turn dialogue with user preferences
- A set of candidate items, each with textual + visual data (images)

- **Goal**

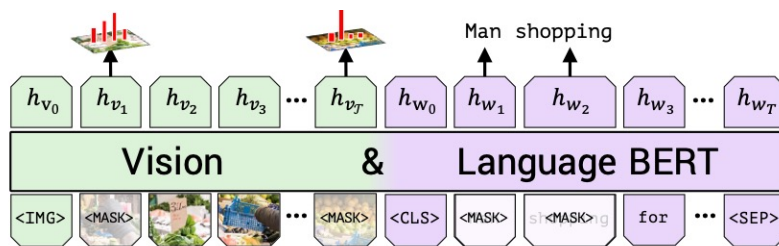
- Recommend the relevant item(s) that match user's needs

**Utilize vision-language models!**

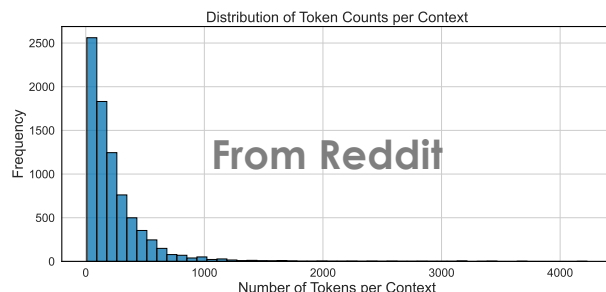
# Earlier Vision-Language Models

- Two types of early approaches

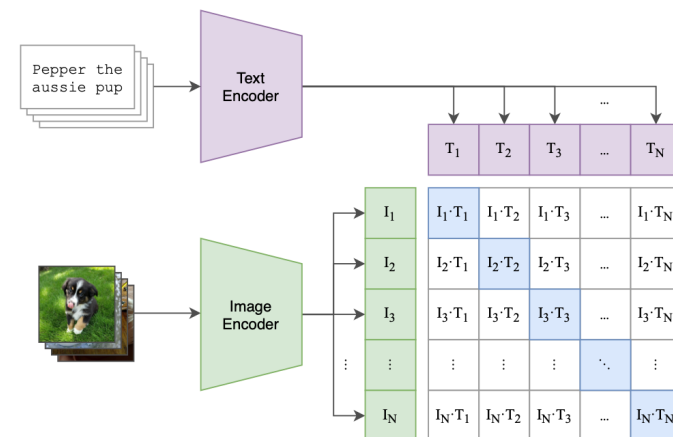
BERT-based models  
(ViBERT, VisualBERT, VL-BERT, UNITER, etc.)



Mostly trained on image captioning tasks: # tokens < 100  
Difficult to apply to long/complex conversations



Dual-encoder contrastive models  
(CLIP, ALIGN, CoCa, Florence, etc.)

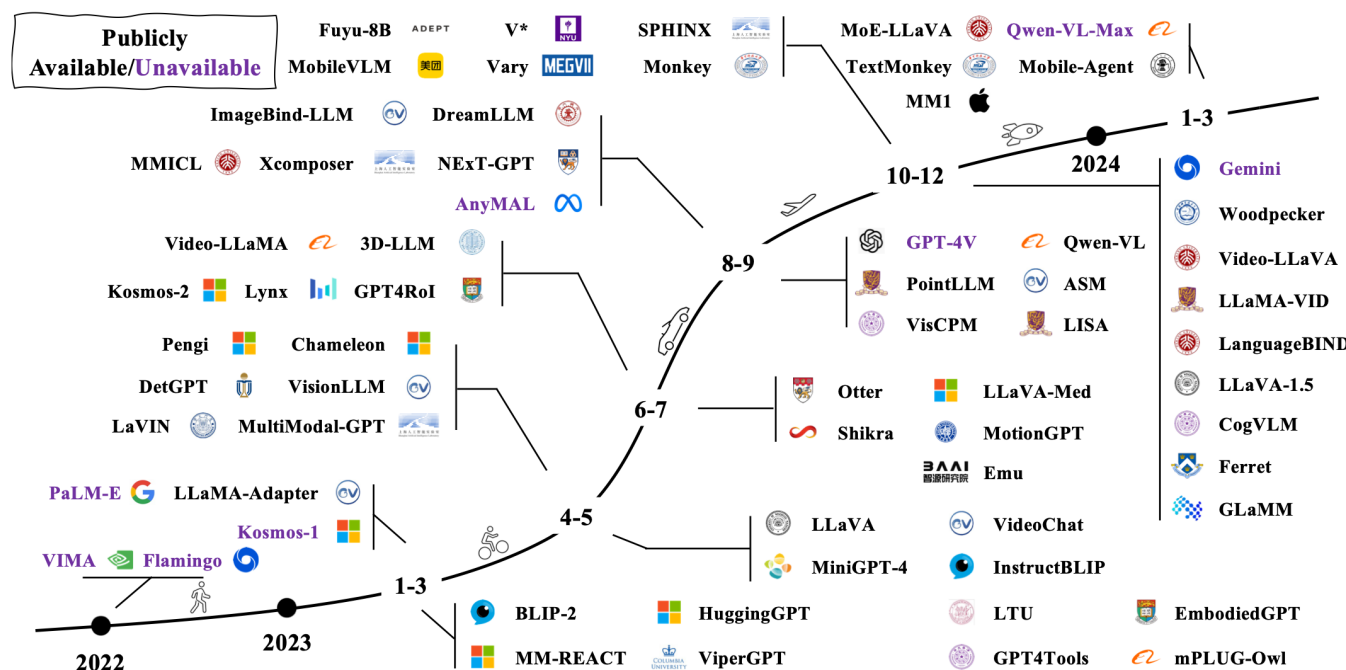
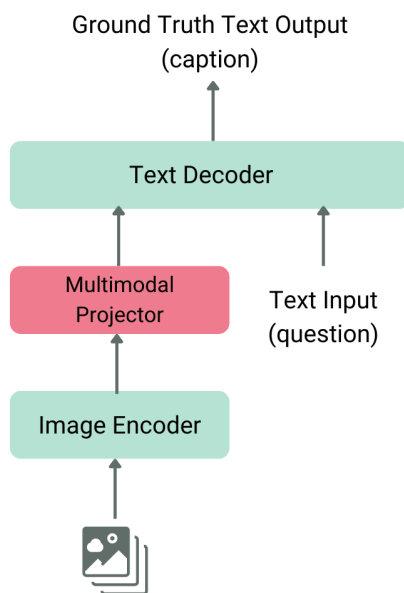


Lu et al. ViBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. NeurIPS 2019.

Radford et al. Learning Transferable Visual Models From Natural Language Supervision. ICML 2021.

# Large Vision-Language Models (VLMs)

- Leveraging LLMs for more complex downstream tasks
  - Vision encoder + LLM backbone

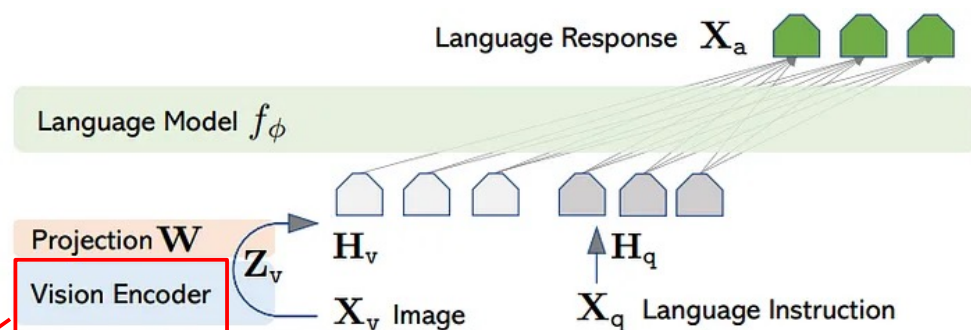


Noyan et al. Vision Language Models Explained. Hugging Face Blog 2024.

Yin et al. A Survey on Multimodal Large Language Models. arXiv 2024.

# Large Language and Vision Assistant (LLaVA)

- One of the most widely used VLMs in various of downstream tasks
  - Open sourced, strong performance, and handleable models available (7B)



LLaVA's Architecture



LLaVA-Next further splits each image into multiple sub-images

**Fine-tune LLaVA for visually-aware CRS!**

Liu et al. Visual Instruction Tuning. NeurIPS 2023.

Liu et al. Improved Baselines with Visual Instruction Tuning. CVPR 2024.

Dosovitskiy et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. ICLR 2021.



# Challenges

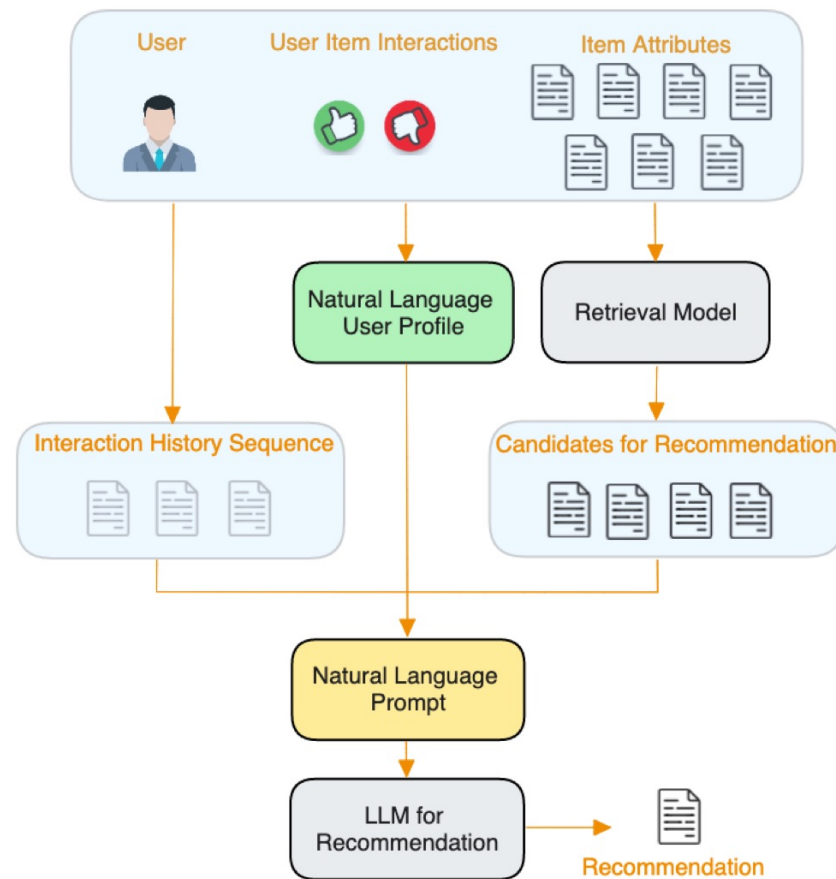
- Absence of adequate **datasets for visually-aware CRS**
  - Ideal: natural conversation  $\cap$  image features

List of CRS datasets

Datasets	#Conv.	#Turns	#Items	Domain	Source
FacebookRec [9]	1M	6M	-	Movies	Synthetic
ReDial [31]	10K	182K	6.2K	Movies	Crowd-sourced
GoRecDial [26]	9K	170K	-	Movies	Crowd-sourced
OpenDialKG [47]	15K	91K	-	Movies, music, etc.	Crowd-sourced
TG-ReDial [70]	10K	129K	-	Movies	Synthetic
DuRecDial 2.0 [44]	16.5K	255K	-	Movies, music, etc.	Crowd-sourced
CCPE-M [49]	502	11K	-	Movies	Crowd-sourced
INSPIRED [15]	1K	35K	1.9K	Movies	Crowd-sourced
Reddit-Movie <sub>base</sub> [18]	85K	133K	24.3K	Movies	Natural
Reddit-Movie <sub>large</sub> [18]	634K	1.6M	51.2K	Movies	Natural
U-NEED [43]	7K	53K	-	E-commerce	Natural
E-ConvRec [24]	25K	775K	-	E-commerce	Natural
HOOPS [13]	-	11.6M	-	E-commerce	Synthetic
MGConvRec [63]	7K	73K	-	Restaurant	Crowd-sourced
MMConv [34]	5K	39K	-	Travel	Crowd-sourced
MobileConvRec [45]	12.2K	156K	1.7K	Music, sports, etc.	Synthetic

# Challenges

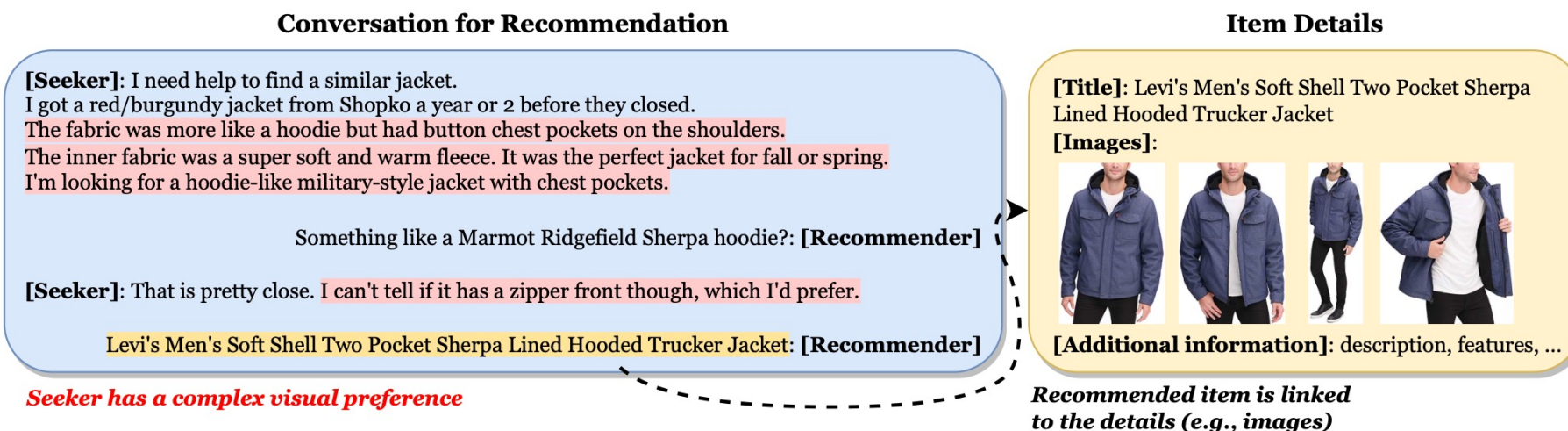
- **Candidate-based recommendation**
  - Effectively narrow down the search space
  - Suitable for domains that LLM is not familiar with
- **Token explosion for multiple images**
  - VLMs convert each image into hundreds/thousands of tokens
  - Multiple candidates cause token explosion
  - LLaVA-v1.6
    - Context length: 4K tokens, image: 2.8K tokens



Candidate-based recommendation

# Reddit-Amazon Dataset

- An example



- Statistics

Dataset	# Conv.	# Turns	# Items	# Images
Beauty	7,672	22,966	5,433	28,082
Fashion	8,039	21,831	6,716	31,162
Home	3,701	6,675	3,077	18,505

# Proposed Method: LaViC

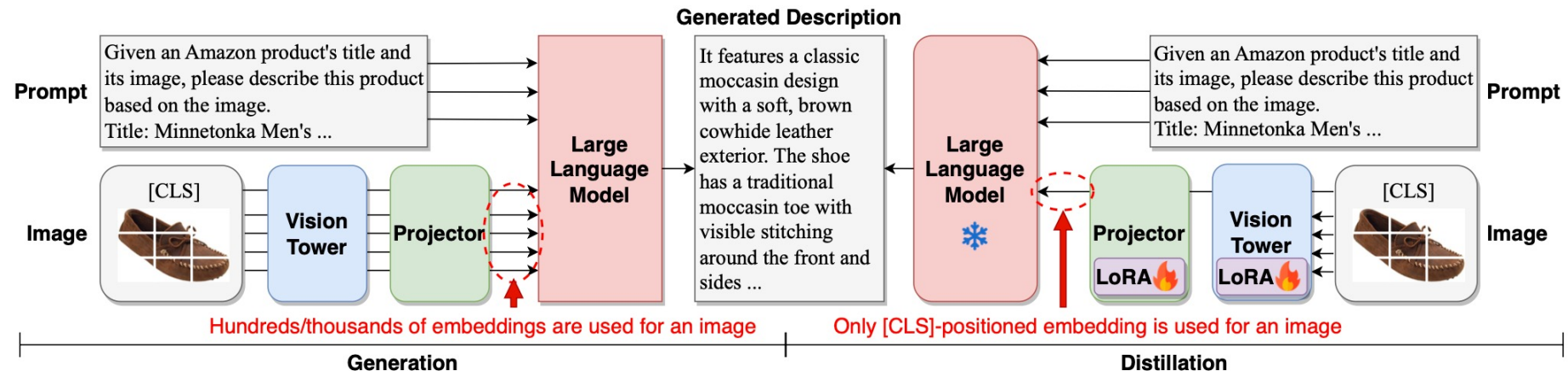
- Designed **LaViC** (Large Vision-Language Conversational Recommendation Framework)
- Overview of **LaViC**
  - 1) Visual Knowledge Self-Distillation
    - Compress thousands image tokens into minimal embeddings
  - 2) Recommendation Fine-Tuning
    - Incorporate visual embeddings + dialogue context into the LLMs

# Visual Knowledge Self-Distillation



- How to avoid token explosion of using multiple images?
  - LLaVA-V1.6 splits each image into 577 x 5 tokens
  - Maximum context length: 4K
- Self-distillation of the visual knowledge
  - Compress each image into 5 [CLS] tokens

Minimize negative log likelihood

$$\min_{\Omega_{\text{vision}}} \sum_i -\log P_{\Omega_{\text{LM}} + \Omega_{\text{vision}}} (D_i \mid \mathcal{T}_{\text{desc}}, \{\text{cls}_{i,r}\}_{r=1}^5)$$



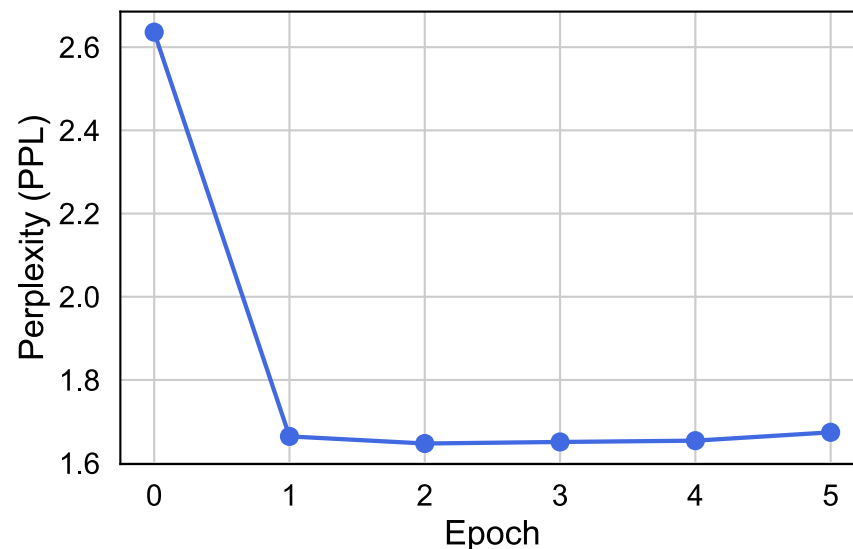
# Example of Description Generation

Image	Domain	Title	Description generated by LLaVA-v1.6
	Fashion	LYANER Women's Pajamas Set 4pcs Satin Silk Cami Top Button Down Loungewear Pjs Set	<p>The product is a women's pajama set made of satin silk, consisting of four pieces. The set includes:</p> <ol style="list-style-type: none"> <li>1. A button-down long-sleeve top with a collar and black piping details.</li> <li>2. A cami top with thin spaghetti straps.</li> <li>3. A pair of long pants with an elastic waistband.</li> <li>4. A pair of shorts with an elastic waistband.</li> </ol> <p>The fabric has a smooth, shiny finish and features a pink base color with a tropical print, including leaves and flowers in shades of blue, white, and black. The overall design is elegant and comfortable, suitable for lounging and sleeping.</p>
	Fashion	Super Mario Bros Classic Retro NES Design Bi-fold Wallet	<p>The product is a bi-fold wallet designed with a retro theme inspired by the classic Super Mario Bros game for the NES (Nintendo Entertainment System). The wallet features a design that mimics the appearance of an NES game cartridge, with the iconic Super Mario Bros label prominently displayed on the front. The label includes pixelated graphics reminiscent of the original game, along with the Nintendo logo and other classic design elements. The wallet appears to be made of a material that resembles the texture and color of an NES cartridge, giving it a nostalgic and vintage look.</p>



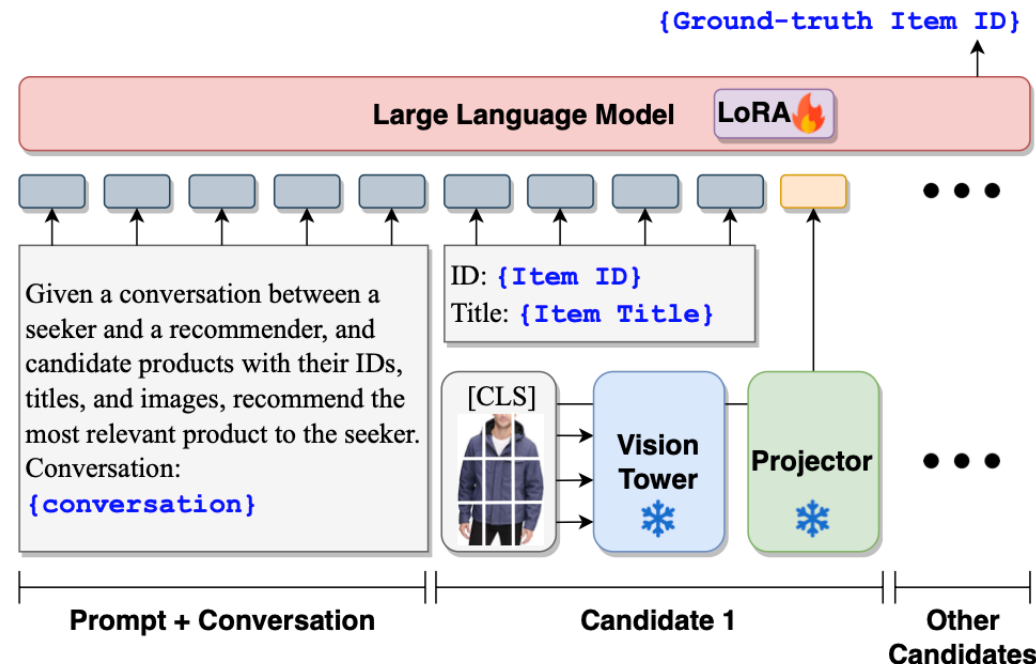
# Perplexity on Visual Knowledge Self-Distillation

- The validation PPL reaches a plateau after 1-2 epochs



# Recommendation Fine-Tuning

- Candidate-based pipeline
  - Retrieved top-10 items using SBERT
    - Title-dialogue matching
  - For each candidate, we provide {ID, Title, [CLS], ...} as input
  - LLM is trained to answer the ground-truth item ID
- Training only LLM parameters
  - Freeze vision tower and projector
  - LoRA is adapted for LLM



$$\min_{\Omega_{\text{LM}}} \sum_{(\mathcal{T}_{\text{conv}}, \mathcal{I}_{\text{cand}})} -\log P_{\Omega_{\text{LM}} + \Omega_{\text{vision}}} \left( ID_{i^*} \mid \mathcal{T}_{\text{conv}}, \{\mathbf{X}_{i_j}\}_{j=1}^{10} \right)$$

Minimize negative log likelihood

# Experimental Setup

- **Datasets**

- Reddit-Amazon (Beauty, Fashion, Home)
  - Each item has title + image(s)

- **Baselines**

- Retrieval-based: BM25, SBERT, RoBERTa, SimCSE, BLaIR
- Generative: Vicuna, LLaVA, GPT-3.5, GPT-4o

- **Implementation**

- Single A100 40GB GPU, LoRA, batch size=4 (distill), 1 (prompt)

- **Evaluation Metrics**

- HitRatio@1 (HR@1): correct item selection rate
- ValidRatio (VR): valid item ID output (no hallucination)

# Experimental Results

- Comparison w/ open-source methods

Method	<i>Beauty</i>		<i>Fashion</i>		<i>Home</i>	
	HR@1	VR	HR@1	VR	HR@1	VR
<i>Retrieval Baselines (item title)</i>						
BM25	0.0169	-	0.0140	-	0.0479	-
SBERT	0.0551	-	0.0681	-	<u>0.2166</u>	-
RoBERTa <sub>large</sub>	0.0640	-	0.0631	-	0.1814	-
SimCSE <sub>large</sub>	0.0326	-	0.0301	-	0.0957	-
BLaIR <sub>base</sub>	0.0371	-	0.0441	-	0.1335	-
<i>Generative Baselines (item title) + SBERT</i>						
Vicuna-v1.5	0.0533	0.9870	0.0481	0.9903	0.1184	1.0000
LLaVA-v1.5	0.0476	0.9896	0.0441	0.9855	0.0932	1.0000
LLaVA-v1.6	<u>0.0770</u>	0.9870	<u>0.0827</u>	0.9867	0.2030	0.9919
<i>Generative Baselines (item title and image) + SBERT</i>						
LLaVA-v1.5	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
LLaVA-v1.6	0.0584	0.9741	0.0459	0.9843	0.1089	0.9919
<i>Proposed Method (item title and image) + SBERT</i>						
<b>LaViC (ours)</b>	<b>0.1187</b>	0.9702	<b>0.1232</b>	0.9298	<b>0.3197</b>	0.9892
<b>Improvement</b>	<b>+54.2%</b>	-	<b>+49.0%</b>	-	<b>+47.6%</b>	-

7B parameters {

→ Outperforms the others

# Experimental Results

- Comparison w/ proprietary methods

Method	<i>Beauty</i>		<i>Fashion</i>		<i>Home</i>	
	HR@1	VR	HR@1	VR	HR@1	VR
<i>Generative Baselines (item title) + SBERT</i>						
GPT-3.5-turbo	0.0968	0.9935	0.0977	0.9903	0.2343	1.0000
GPT-4o-mini	0.1213	1.0000	0.1160	0.9927	0.3258	0.9973
GPT-4o	0.1271	0.9987	0.1278	0.9976	0.3350	1.0000
<i>Generative Baselines (item title and image) + SBERT</i>						
GPT-4o-mini	0.1081	0.9974	0.1098	0.9927	0.2861	0.9946
GPT-4o	0.1160	0.9974	0.1231	0.9939	0.3308	0.9973
<i>Proposed Method (item title and image) + SBERT</i>						
<b>LaViC (ours)</b>	0.1187	0.9702	0.1232	0.9298	0.3197	0.9892

150-200B parameters {

7B parameters ←

→ Comparable performance

# Experimental Results

- **Ablation study**

- “Entire tokens” use all tokens for each image
- “w/o images” use only item titles
- “w/o self-distillation” use [CLS] tokens but without self-distillation

Method	<i>Beauty</i>		<i>Fashion</i>		<i>Home</i>	
	HR@1	VR	HR@1	VR	HR@1	VR
Entire tokens ( $5 \times 577$ )	0.0256	0.9456	<i>o.o.m.</i>		<i>o.o.m.</i>	
w/o images	<u>0.0972</u>	0.9767	0.1022	0.9358	<u>0.2944</u>	0.9946
w/o self-distillation	0.0842	0.9793	<u>0.1084</u>	0.9649	0.2861	0.9973
<b>LaViC (ours)</b>	<b>0.1187</b>	0.9702	<b>0.1232</b>	0.9298	<b>0.3197</b>	0.9892



# Experimental Results

## • Case study

- (a) LaViC captures subtle visual attributes (color, design) not evident in the item title
- (b) LaViC captures additional details such as extra straps or shape using compressed image tokens

**[Seeker]:** I just landed an internship. It's business casual, and I already have OCBDs and chinos, but what shoes should I buy? I only have sneakers, and I'm not really looking to spend a couple hundred dollars for a two-month internship.

LLaVA-v1.6  
(w/ title)



*G.H. Bass & Co. Men's Burlington Oxford*



LLaVA-v1.6  
(w/ title & image)



*Dockers Men's Gordon Leather Oxford Dress Shoe*



LaViC (ours)



*G.H. Bass & Co. Men's Buckingham Oxford*



(a) Same brand but different style

**[Seeker]:** Skates backpack that doesn't pull the center of gravity too far back? My fiancée wants: 1) two separate shoulder straps, 2) capacity to fit both gear and skates, 3) minimal backward pull on center of gravity, and 4) a cross-chest strap.

LLaVA-v1.6  
(w/ title)



*Gutezeiten Kick Scooter Shoulder Strap*



LLaVA-v1.6  
(w/ title & image)



*Mares Cruise Backpack Pro Bag - Black White*



LaViC (ours)



*Atom Skates Back Pack - Sport Backpack*



(b) Specific requirements

# Conclusion

- Address the visually-aware conversational recommendation
- Propose a novel framework: LaViC
  - Large vision-language model-based framework
  - Visual knowledge self-distillation
  - Recommendation fine-tuning
- Open-source the benchmark datasets



**Thank You!**  
[hyjeon@ucsd.edu](mailto:hyjeon@ucsd.edu)