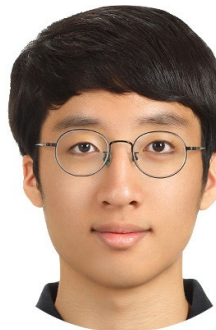


# Induce, Edit, Retrieve: Language-grounded Multimodal Schemata for instructional Video Retrieval

Yue Yang



Joongwon Kim



Artemis Panagopoulou



Mark Yatskar

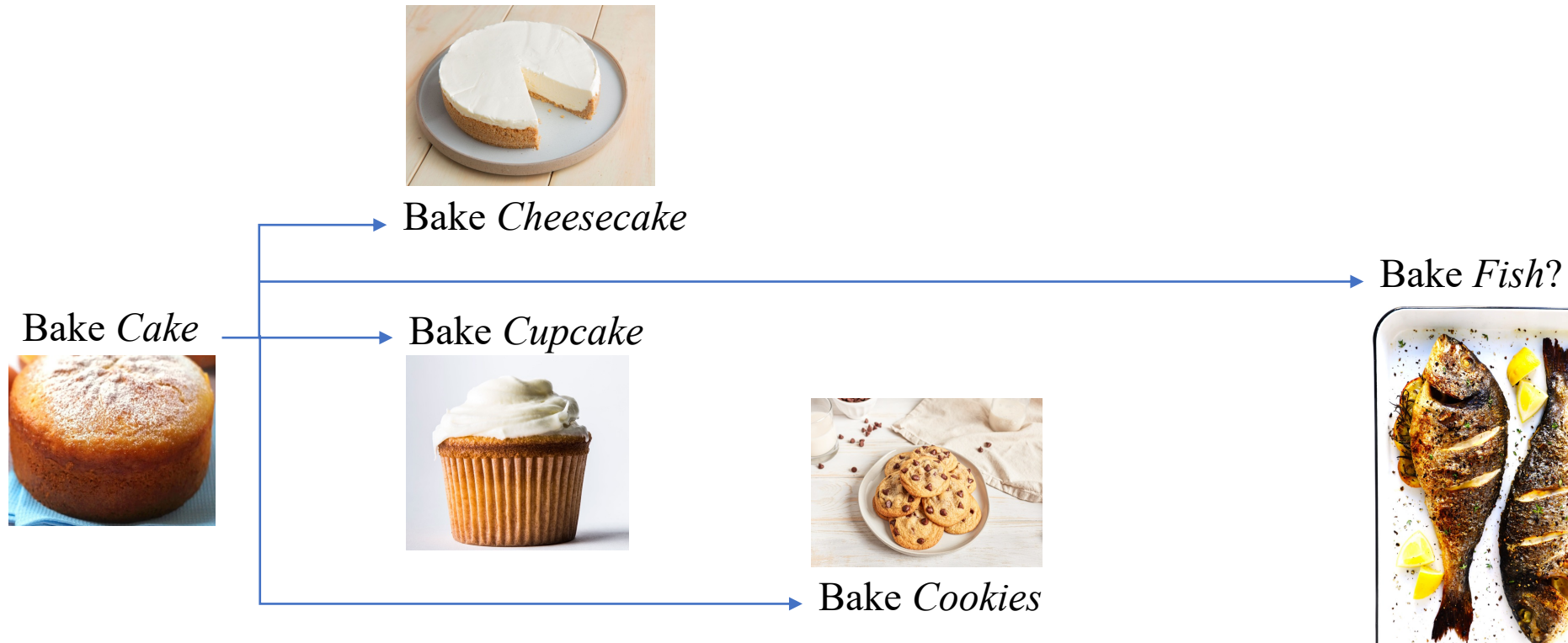


Chris Callison-Burch



# Motivation

- Schema: a set of rules people use to perform everyday tasks.
- Schema can be generalized.
- When facing new tasks, people use prior knowledge. (Chen et al., 2004)



# IER Overview

- Our **Induce, Edit, Retrieve (IER)** system:
- **Induce:**
  - Input: a task name, a set of related videos.
  - Output: a set of sentences as the schema.
- **Edit:**
  - Given an unseen task.
  - Use language models to modify schema.
- **Retrieve:**
  - Improve video retrieval using schema.



Figure 1. An example from our IER system, which first induces a schema for *Bake Chicken* using a set of videos. Then it edits the steps in the schema to adapt to the unseen task *Bake Fish* (the tokens that have been edited are highlighted). Finally, IER relies on the edited schema to help retrieve videos for *Bake Fish*.

# Schema Induction

## Schema Induction

*Objective: Construct schemata on tasks.*

Task: Bake Chicken

 YouTube Videos (Learning Data)

Video 1



...



⋮

⋮

⋮

⋮

Video n



...



Clip-Step Alignment

**wikiHow** to do anything 1M human written steps

Wash the chicken thoroughly.

Marinate the chicken  
Season the drumsticks

Bake the blackened chicken in the oven  
or on the grill

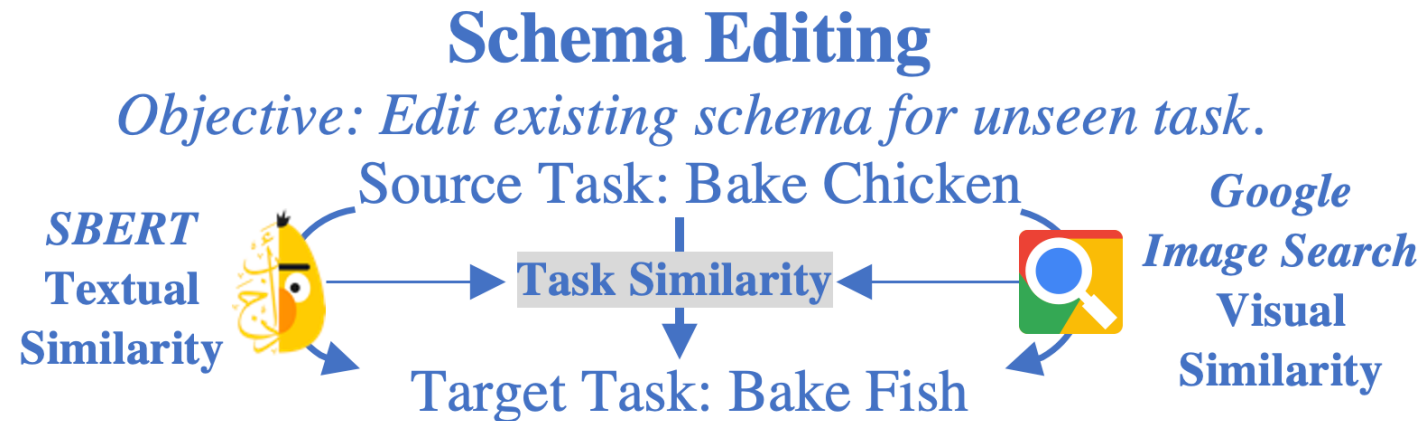
### Schema of *Bake Chicken*

- Wash the chicken thoroughly
- Season your drumsticks
- Marinate the chicken
- Insert a roasting thermometer into the thigh
- Sprinkle the ginger on the chicken
- Bake the blackened chicken in the oven
- ... ..

- Align wikiHow steps with videos.
- MIL-NCE as the video-text model.
- We induced 21,299 schemata using 1.2 M videos from Howto100m.

# Schema Editing

- Given an **unseen** task without videos, edit **existing** schema.
- Find the most similar task in the schema library:



- Textual Similarity = cosine similarity of SBERT embeddings
- Visual Similarity (Google Image Search)
- Task Similarity =  $\max(\text{Textual Similarity}, \text{Visual Similarity})$

# Schema Editing

- Editing Module 1: **Object Replacement**
  - Every task has a main object, e.g., “chicken” of “Bake Chicken”.
  - Use POS tagger to find the 1st occurred noun as main object.
  - Replace the objects in all steps.

Object Replacement	
Cook Ham $\xrightarrow{0.86}$ Cook Lamb	
Put the ham in the oven.	
↓	
Put the lamb in the oven.	
Clean a Guitar $\xrightarrow{0.84}$ Build a Violin	
Use a polish for particularly dirty guitars.	
↓	
Use a polish for particularly dirty violins.	
Trap a Rat $\xrightarrow{0.84}$ Trap a Rabbit	
Bait and set snap rat traps.	
↓	
Bait and set snap rabbit traps.	

# Schema Editing

- Editing Module 2: **Step Deletion**
  - Delete the steps no longer suitable for the new target task.
  - “Insert a roasting thermometer into the thigh” of “Bake Chicken” **×** “Bake Fish”
  - Sentence BERT pretrained on QA pairs.
    - Compute the score of (task, step).
    - if (source task, step)  $\gg$  (target task, step), delete, otherwise include

Step Deletion	
Transplant a Young Tree	$\xrightarrow{0.89}$ Remove a Tree
Fill your pot with a balanced fertilizer.	
↓delete	
<del>Fill your pot with a balanced fertilizer.</del>	
Fix a Toilet	$\xrightarrow{0.85}$ Remove a Toilet
Test out the new flapper.	
↓delete	
<del>Test out the new flapper.</del>	
Brush a Cat	$\xrightarrow{0.87}$ Brush a Long Haired Dog
Comb and groom your pet.	
↓include	
Comb and groom your pet.	

# Schema Editing

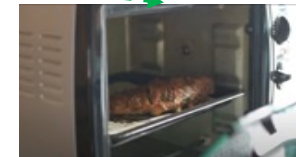
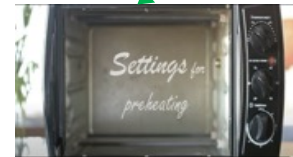
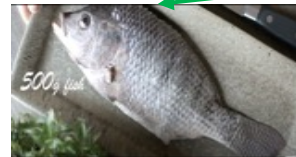
- Editing Module 3: **Token Replacement**
  - Use *masked language model* to replace the token with the lowest probability.
  - “Season the **drumstick**” in “Bake Chicken”
  - Mask the token “Season the <mask>”.
  - Use a prompt: How to [TASK]? [STEP]
    - How to Bake Fish? “Season the <mask>”.
  - Predict a new token from vocabulary.
    - <mask> --> fish, “Season the fish”

Token Replacement	
Prepare Fish	$\xrightarrow{0.82}$ Prepare Crabs
Cut the <b>fins</b> from the <b>fish</b> using <b>kitchen shears</b> .	
↓	
Cut the <b>shells</b> from the <b>crabs</b> using <b>steel scissors</b> .	
Make Healthy Donuts	$\xrightarrow{0.88}$ Bake Healthy Cookies
Slice your <b>donuts</b> into <b>disks</b> .	
↓	
Slice your <b>cookies</b> into <b>squares</b> .	
Wash Your Bike	$\xrightarrow{0.84}$ Wash a Motorcycle
Clean the <b>bike chain</b> with a <b>degreaser</b> .	
↓	
Clean the <b>motorcycle</b> <b>thoroughly</b> with a <b>towel</b> .	

# Schema-Guided Video Retrieval

- Query: Task Name (short) Retrieve long multi-minute videos.
- Global Matching (use task name only)
- Step Aggregation (use schemata to expand task name)

*Use the task name "Bake Fish" as Query*



*With Schema*

*Wash the fish*



*Sprinkle the sauce  
on the fish*



*Preheat the oven.*



*Bake the blackened  
fish in the oven*



# Experiment-Datasets

---

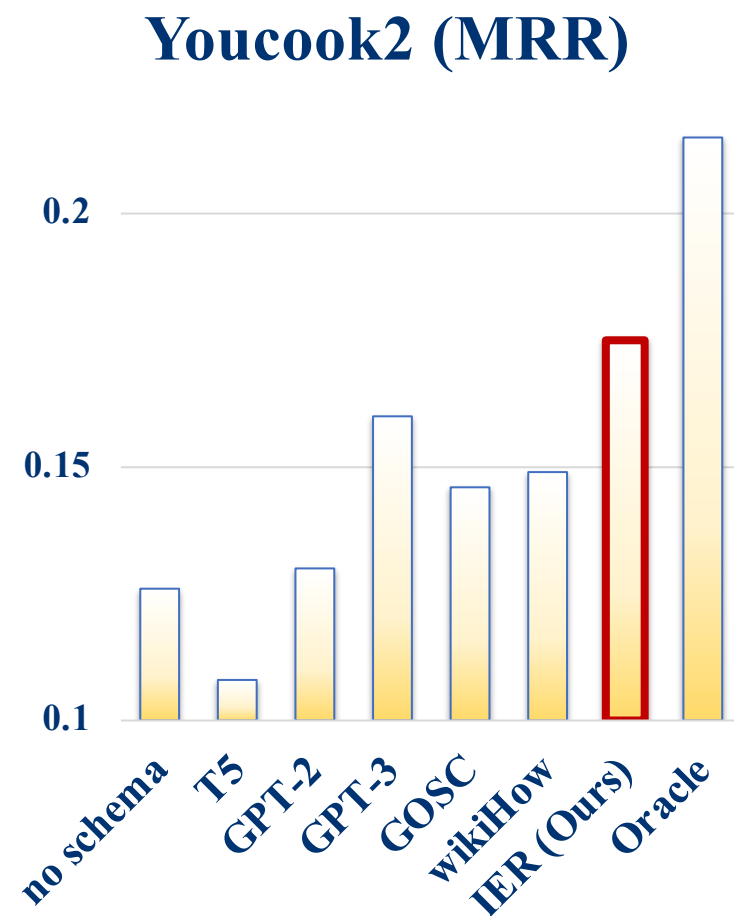
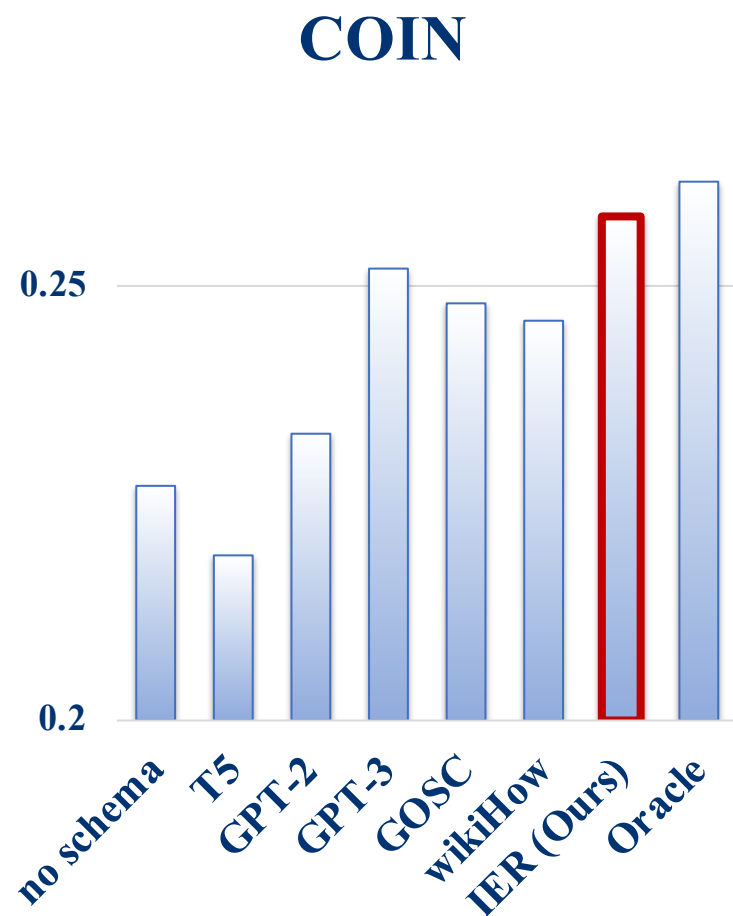
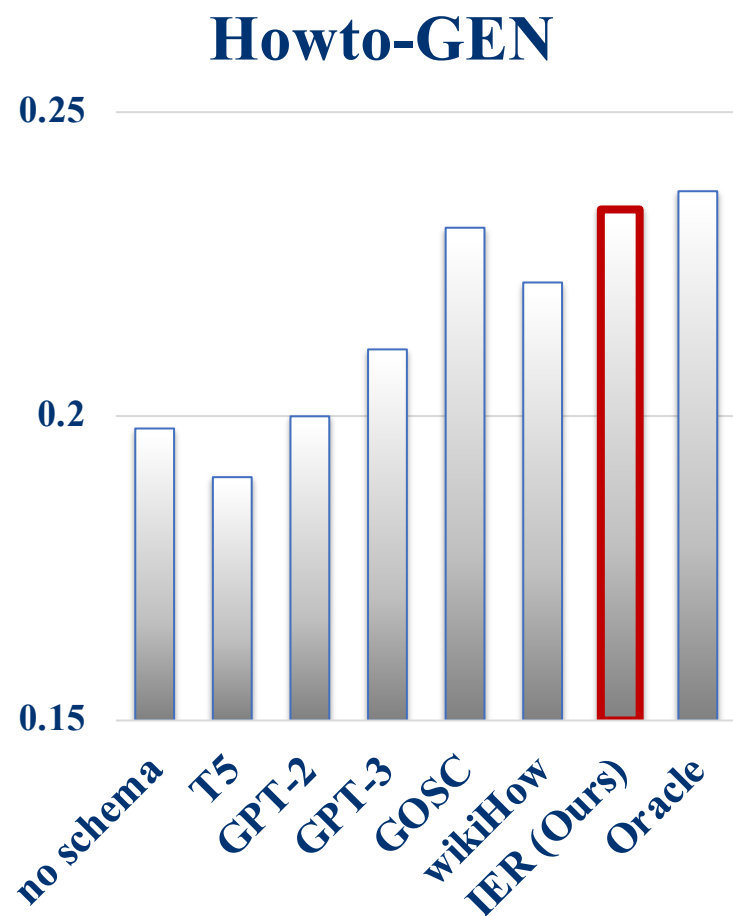
- **Howto-GEN** (a new split of Howto100M)
  - 3365 tasks, 2,184 unique main objects.
  - Random select 500 tasks for training, 500 for validation, 2365 for test.
- **COIN** (**CO**mprehensive **IN**structional video analysis)
  - 180 tasks, 11,827 videos.
  - Unseen tasks, e.g., “Blow Sugar”, “Make Youtiao”, etc.
- **Youcook2**
  - 89 recipes – tasks, 2,000 long videos.

# Experiment - Baselines

---

- **Generation Models:**
  - T5, GPT-2-large.
  - GPT-3: Zero-shot generation (How to [Task Name]? Give me several steps.)
- **Goal-Oriented Script Construction**
  - Given the input task name, retrieve the set of desired steps from wikiHow.
- Oracle Schemata (human written, upper bound)
  - Howto-GEN (from wikiHow)
  - COIN/Youcook2 (Human annotation)

# Results:



# Results:

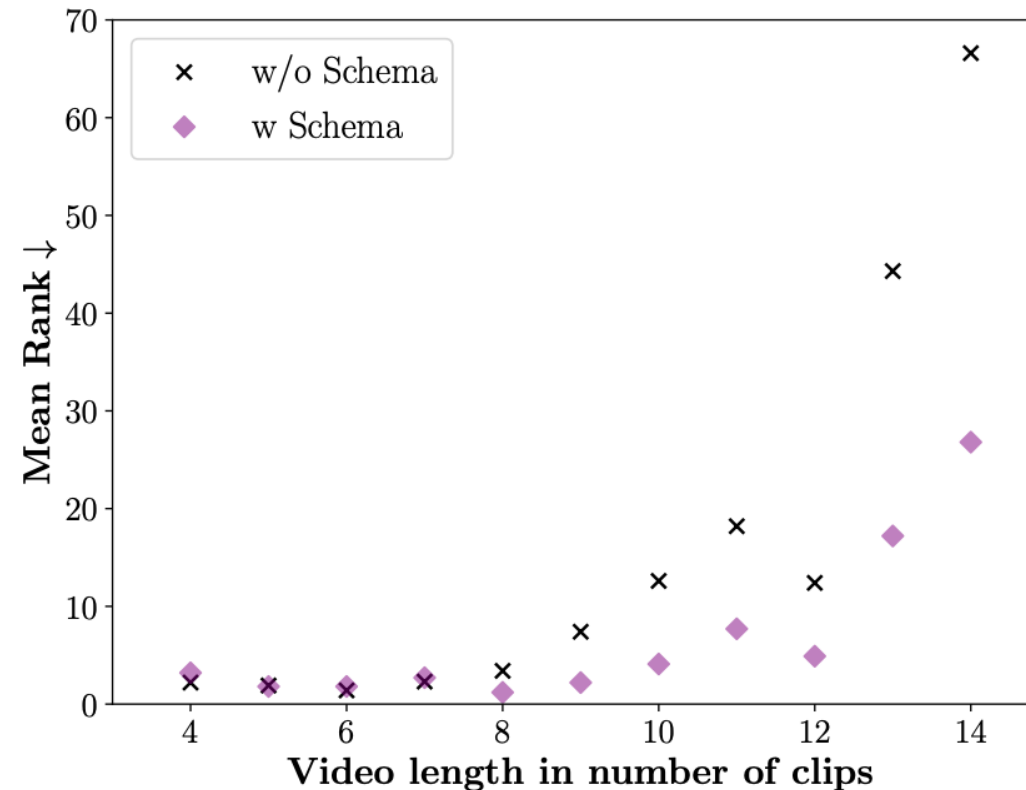


Figure 3. Retrieval performance by video length (in the number of clips). We group the test videos of Youcook2 by the number of clips per video and compute the mean rank for each group.

**IER helps more for longer videos**

# Results:

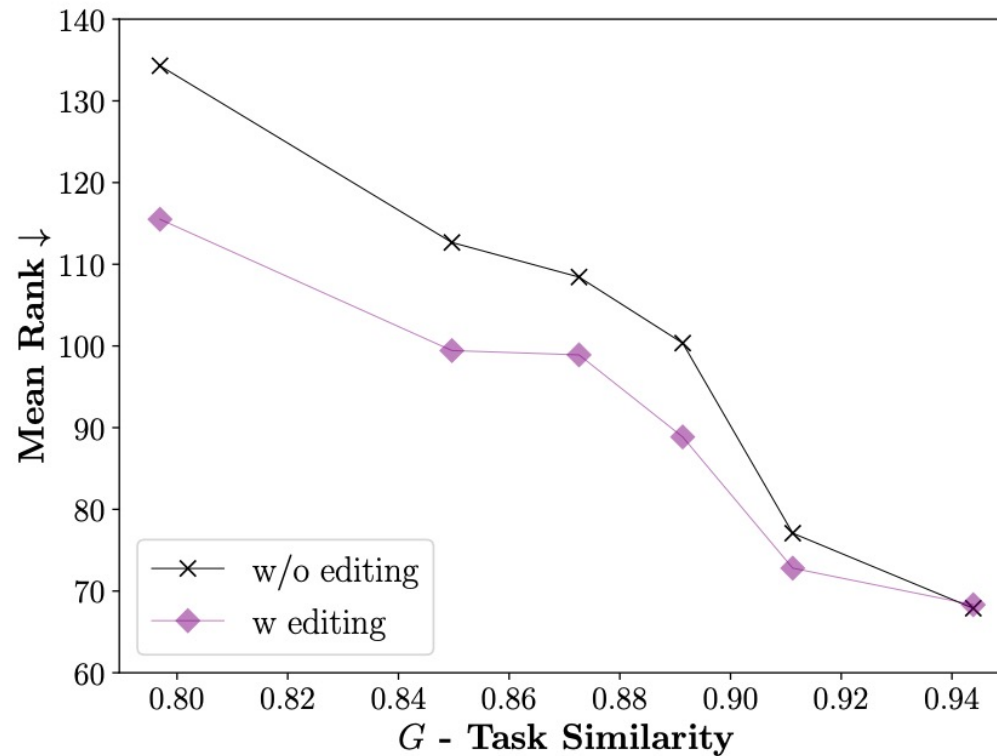


Figure 4. Retrieval performance by task similarity. We sort the test tasks of Howto-GEN based on their task similarity ( $G$ ) and compute their mean rank for every batch of 400 tasks.

**Editing helps more when task similarity is low.**

# Results:

	Method	P@1↑	R@5↑	R@10↑	Med r↓	MRR↑
Howto-GEN	full	<b>54.4</b>	<b>37.3</b>	<b>50.1</b>	<b>10.0</b>	<b>.231</b>
	– mask	<u>53.7</u>	<u>36.3</u>	<u>49.3</u>	<u>11.0</u>	<u>.229</u>
	– deletion	<u>53.6</u>	<u>36.9</u>	<u>49.8</u>	<u>11.0</u>	<u>.230</u>
	– replacement	<u>51.5</u>	<u>34.9</u>	<u>47.3</u>	<u>12.0</u>	<u>.220</u>
	– all	<u>45.5</u>	<u>31.0</u>	<u>43.1</u>	<u>15.0</u>	<u>.199</u>

**All three editing modules are beneficial.**

# Results:

<b>Model</b>	<b>P@1<math>\uparrow</math></b>	<b>R@5<math>\uparrow</math></b>	<b>R@10<math>\uparrow</math></b>	<b>Med r<math>\downarrow</math></b>	<b>MRR<math>\uparrow</math></b>
MIL-NCE	48.3	37.1	52.8	9.5	.227
+schema	57.2	42.2	57.8	7.0	.256
CLIP [38]	58.9	44.9	58.8	6.0	.264
+schema	65.0	47.4	60.8	5.5	.282

Table 5. Retrieval performance on COIN using MIL-NCE and CLIP as the matching functions. +schema represents using schema induced by IER (MIL-NCE as matching function) for retrieval.

**Schemata can be reused by different video-text model.**

# Conclusion

---

- We propose a schema induction and generalization system that improves instructional video retrieval performance.
- We demonstrate that the induced schemata benefit video retrieval on unseen tasks, and our IER system outperforms other methods.
- In the future, we plan to investigate the structure of our schemata.

Thank you!