

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

Yue Yang, Joongwon Kim, Artemis Panagopoulou, Mark Yatskar, Chris Callison-Burch

## Abstract

Schemata are structured representations of complex tasks that can aid artificial intelligence by allowing models to break down complex tasks into intermediate steps. We propose a novel system that induces schemata from web videos and generalizes them to capture unseen tasks with the goal of improving video retrieval performance.

## Methods

Our system proceeds in 3 major phases: 1) Schema Induction; 2) Schema Editing; 3) Schema-Guided Video Retrieval. See figure in the middle for details.

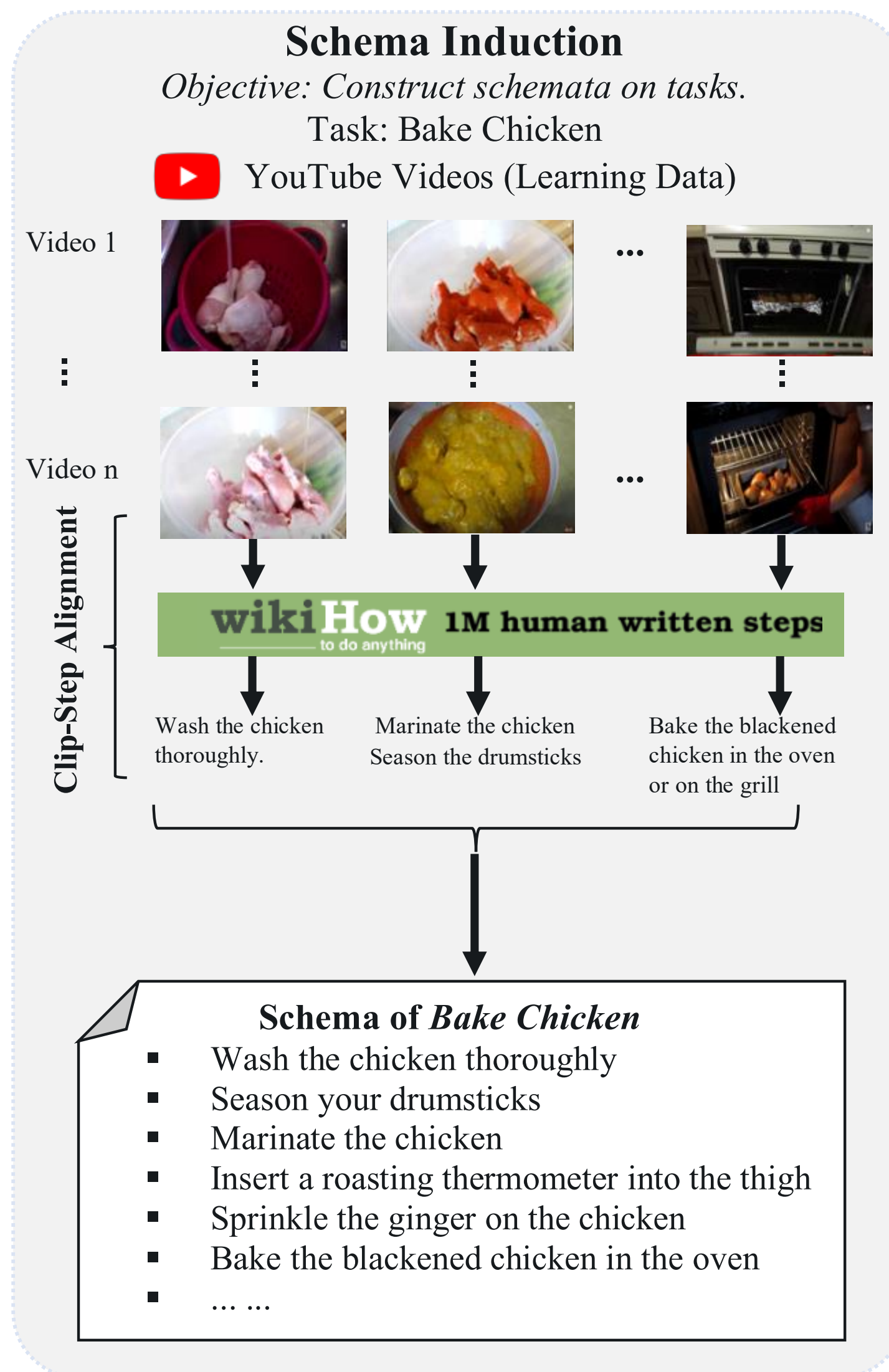
## Experiments

### Datasets:

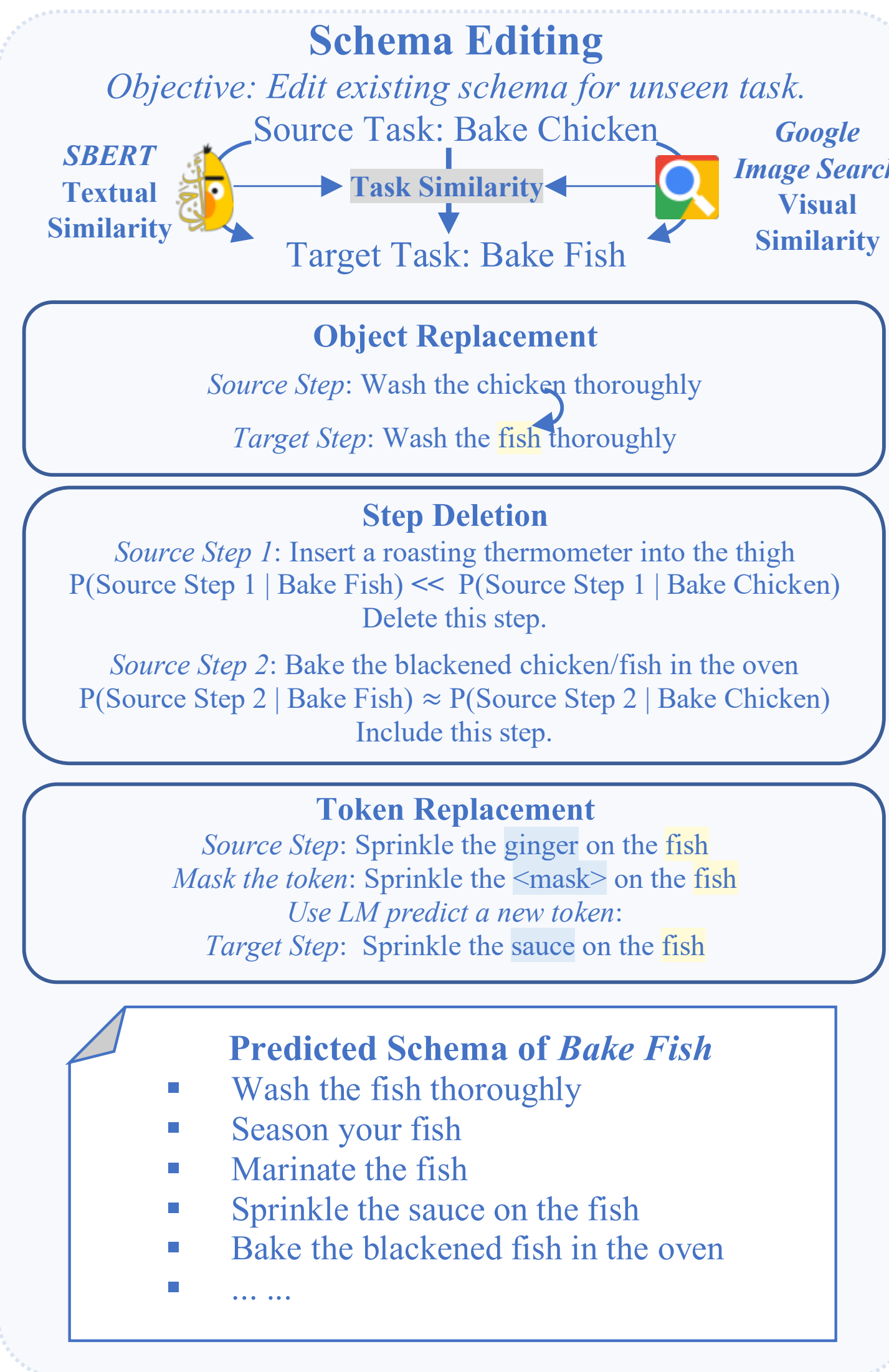
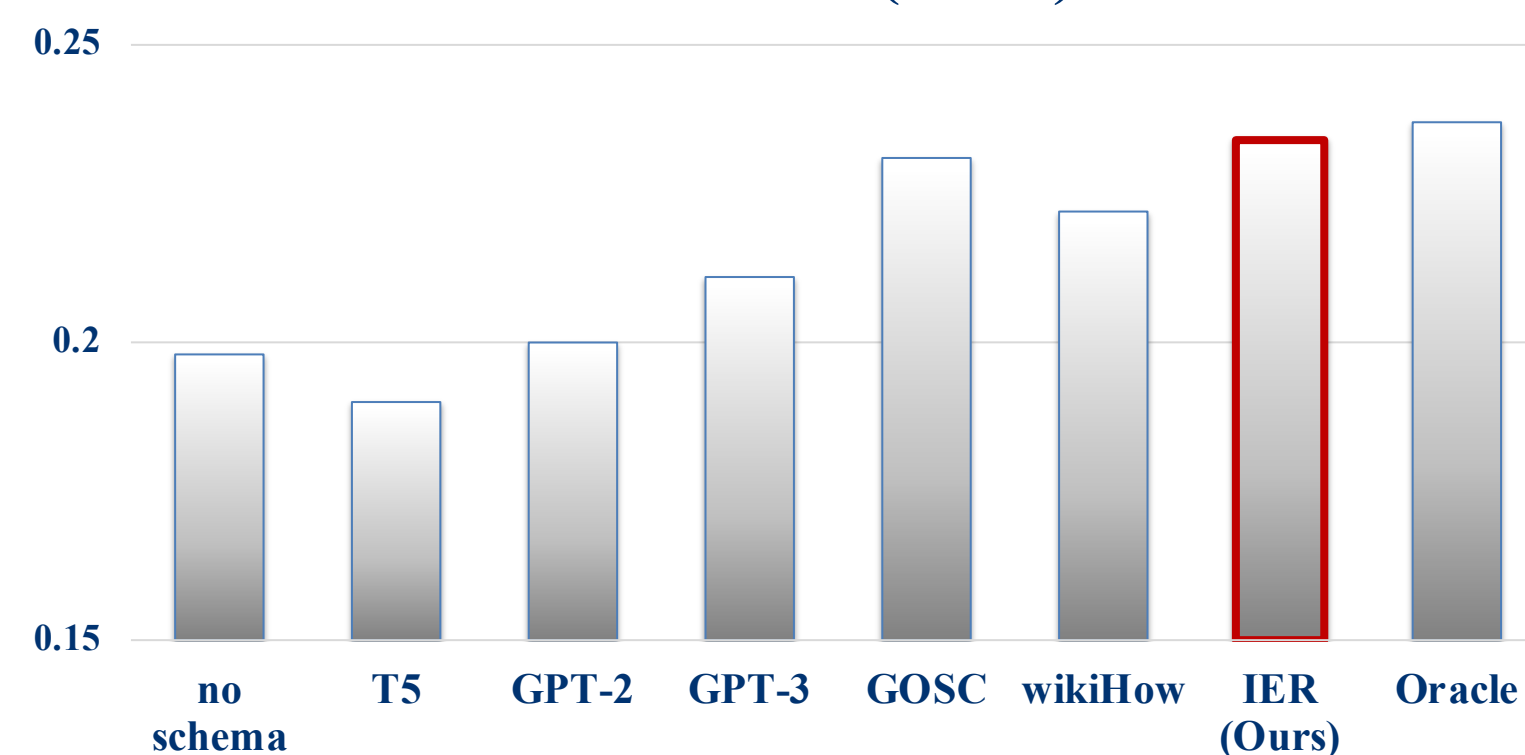
- Howto100M: 1.2M videos, 23K tasks. We induce 21,299 schemata using it.
- Howto-GEN: 3,365 from Howto100M. Train/Val/Test: 500/500/1088.
- COIN: 11,827 videos for 180 tasks.
- Youcook2: 2,000 videos for 89 tasks.

### Baselines (Other induction methods):

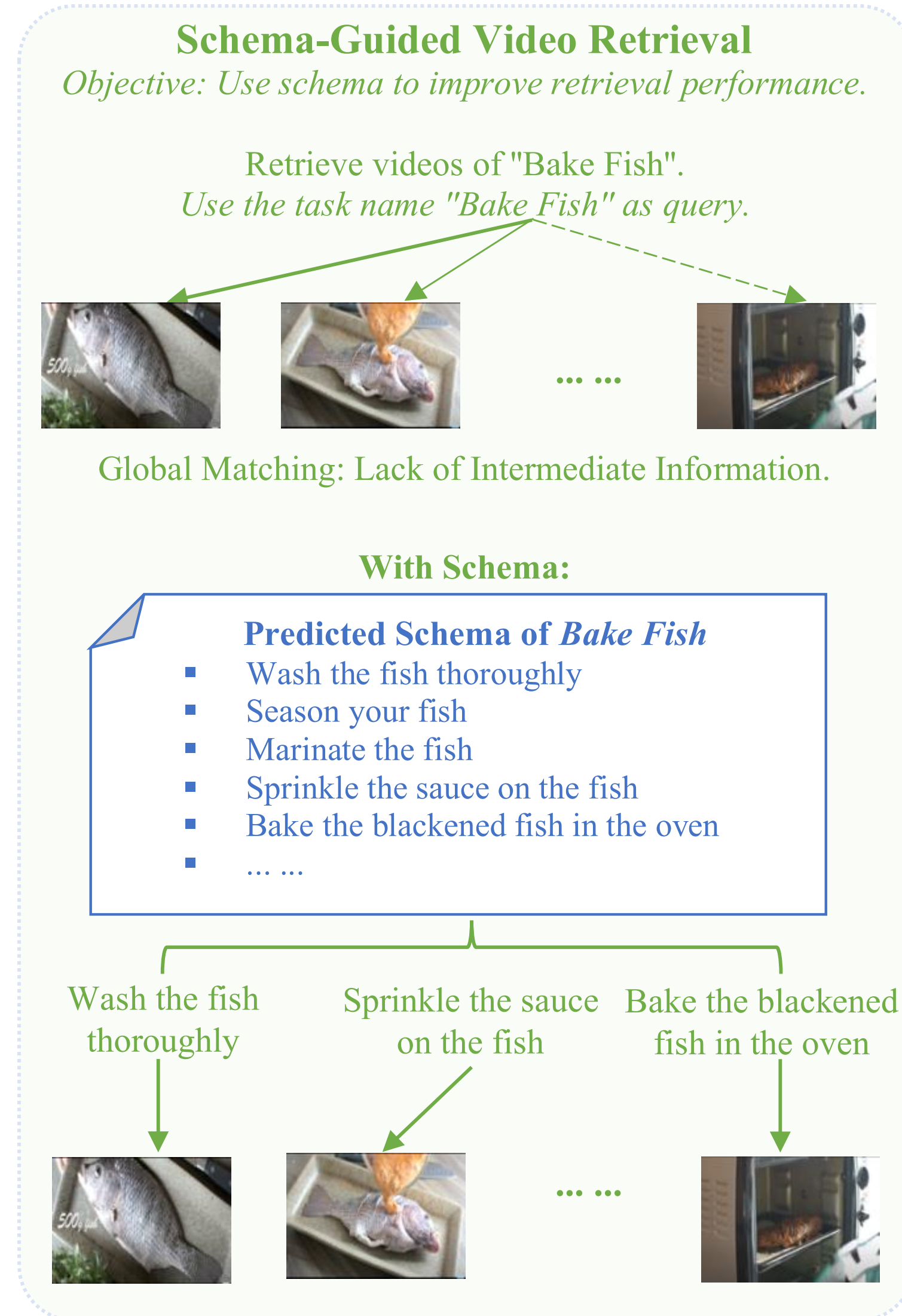
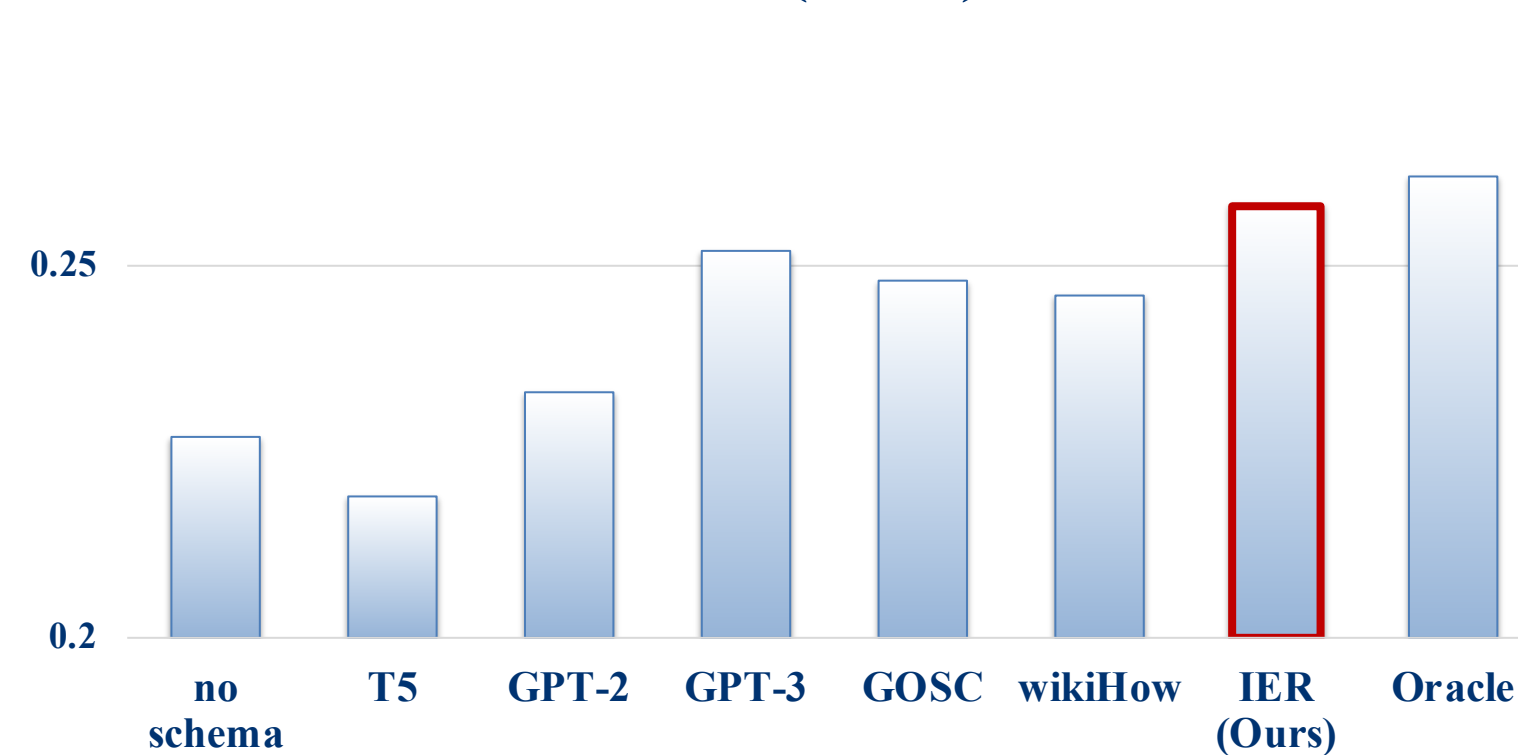
- Language Models: T5, GPT-2, GPT-3.
- Goal Oriented Script Construction.
- Human written: wikiHow and Oracle.



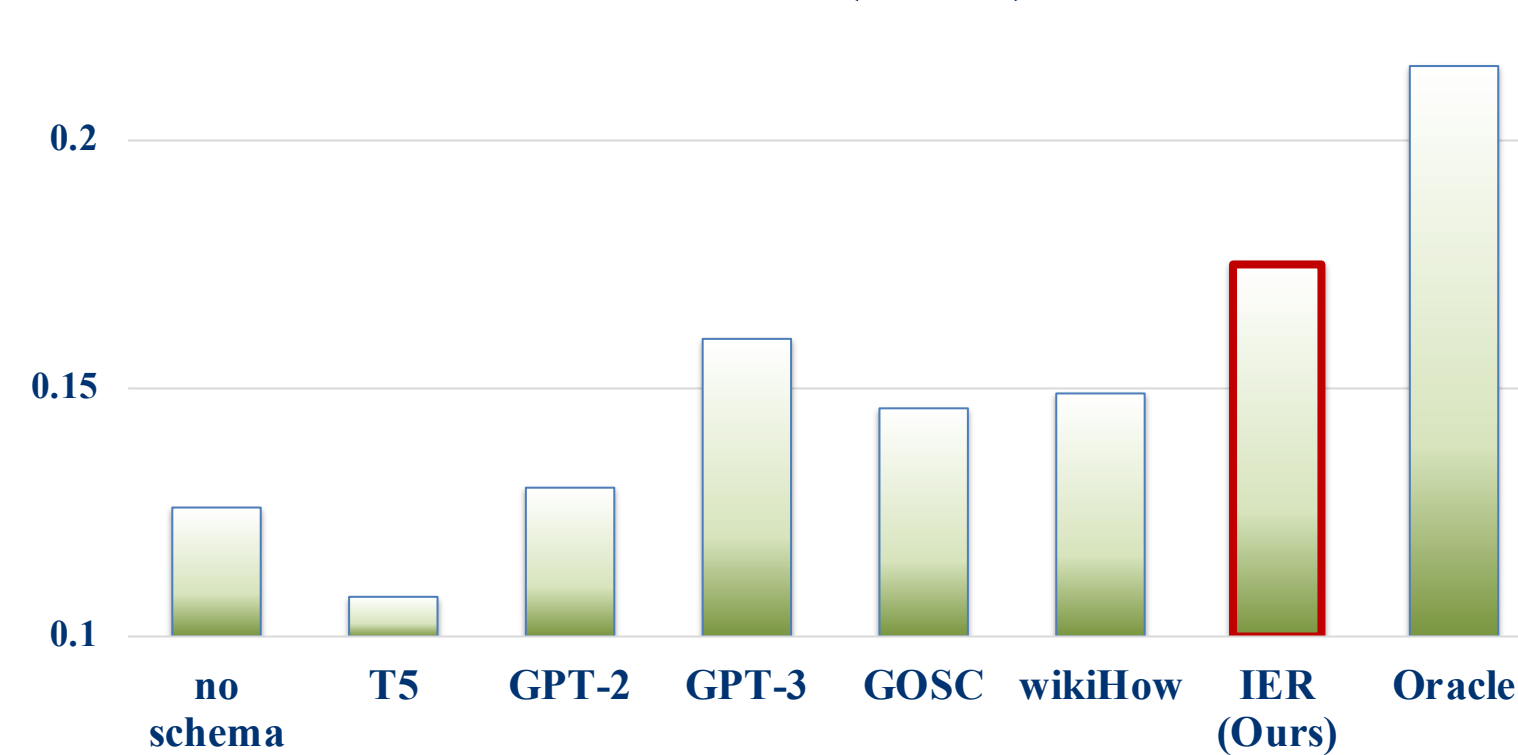
Howto-GEN (MRR)



COIN (MRR)

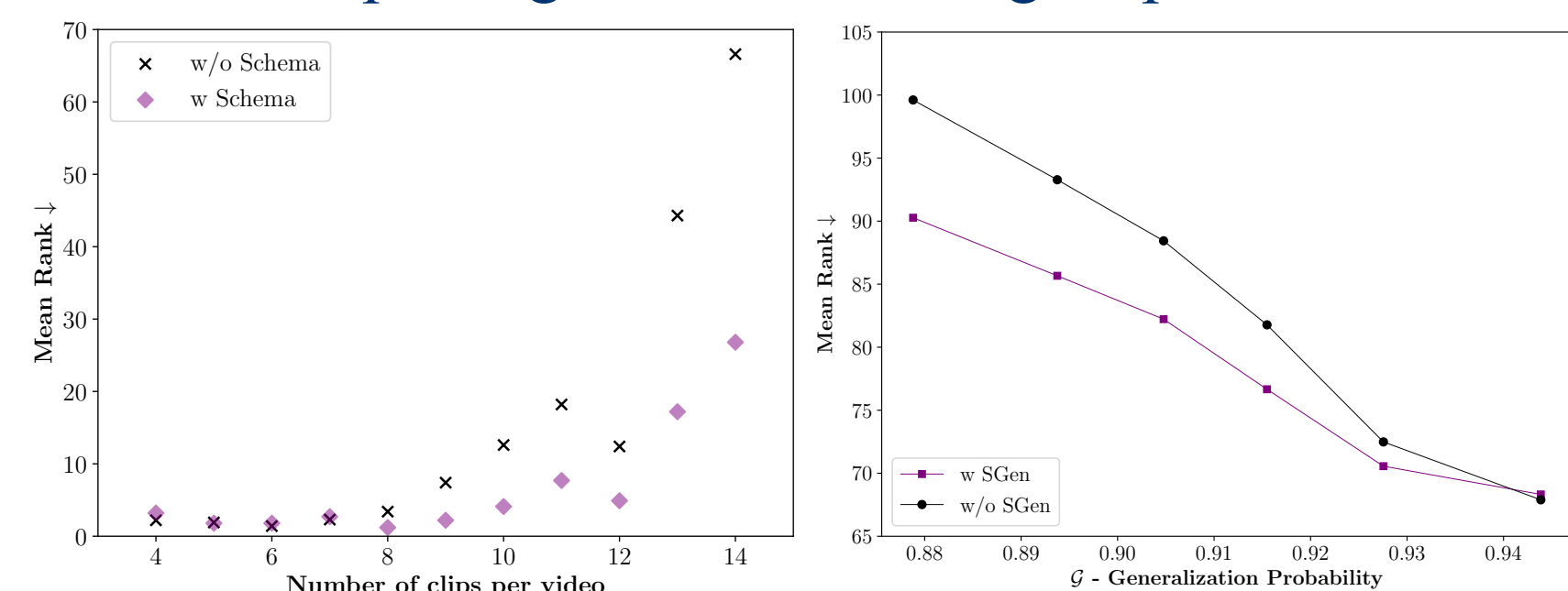


Youcook2 (MRR)



## Results

Schema helps long videos: Editing helps unseen tasks:



All three editing modules are beneficial:

	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	53.7	36.3	49.3	11.0	.229
	— deletion	53.6	36.9	49.8	11.0	.230
	— replacement	51.5	34.9	47.3	12.0	.220
	— all	45.5	31.0	43.1	15.0	.199

Schemata can be used by other vision-text models:

	Model	P@1↑	R@5↑	R@10↑	Med r↓	MRR↑
MIL-NCE	full	48.3	37.1	52.8	9.5	.227
	+schema	57.2	42.2	57.8	7.0	.256
CLIP [38]	full	58.9	44.9	58.8	6.0	.264
	+schema	65.0	47.4	60.8	5.5	.282

## 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, such as the temporal order, and discover other applications of schemata.