

# From Observation to Abstractions: Efficient In-Context Learning from Human Feedback and Visual Demonstrations for VLM Agents

Gabriel Sarch<sup>1</sup> Lawrence Jang<sup>1</sup> Michael J. Tarr<sup>1</sup>  
William W. Cohen<sup>1,2</sup> Kenneth Marino<sup>2</sup> Katerina Fragkiadaki<sup>1</sup>  
<sup>1</sup>Carnegie Mellon University <sup>2</sup>Google DeepMind

## Abstract

We propose an efficient method, *In-Context Abstraction Learning (ICAL)*, to improve in-context VLM agents from sub-optimal demonstrations and human feedback. Specifically, given a noisy demonstration for a task in a new domain, LLMs/VLMs are used to fix inefficient actions and annotate four types of cognitive abstractions. These abstractions are then refined by executing the trajectory in the environment, guided by natural language feedback from humans. We demonstrate that this method rapidly learns useful experience abstractions. Our ICAL agent improves on the state-of-the-art when tested in dialogue-based instruction following in household environments in TEACH, action anticipation in Ego4D, and in multimodal autonomous web agents in VisualWebArena. In TEACH, we improve on the state-of-the-art by 12.6% in goal-condition success, outperforming LLM agents that use the raw visual demonstrations as in context examples without abstraction learning. In VisualWebArena, we improve on the state-of-the-art by an absolute 8.4% and relative 58.74% in task success, outperforming VLM agents that use hand-written examples. In Ego4D, we improve 6.4 noun and 1.7 action edit distance over few-shot GPT4V. Lastly, we find that weight fine-tuning and in-context abstraction learning complement each other, with their combination yielding the best performance.

## 1. Introduction

Humans acquire skills through language and observation, a model for automated systems. These systems must learn from verbal instructions and demonstrations to develop rapid learning technologies. This involves integrating linguistic feedback and demonstrative learning to refine knowledge across different contexts.

Research has used large language models (LLMs) and visual language models (VLMs) to derive insights from experiences, improving performance by adding these insights to prompts [7, 9, 10, 13]. However, there remains limitations in task transfer and underutilization of visual data.

We introduce a new method, In-Context Abstraction Learning (ICAL), for teaching VLMs using suboptimal demonstrations and feedback. ICAL helps VLMs create

and refine multimodal abstractions, aiding them in understanding task dynamics and critical knowledge [1, 2, 5, 14].

## 2. In-Context Abstraction Learning (ICAL)

ICAL starts by obtaining a noisy trajectory. It has two phases: (1) the abstraction phase  $F_{abstract}$ , where a VLM corrects the trajectory and adds language comments in isolation (Section 2.1), and (2) the human-in-the-loop phase  $F_{hitl}$ , where the trajectory is executed with human feedback to refine it (Section 2.2). Each corrected trajectory is stored as a contextual reference for learning and inference.

### 2.1. VLM-driven Abstraction Generation

The abstraction function  $F_{abstract}$  processes trajectory  $\xi_{noisy}$  into an optimized sequence  $\xi_{opt}$  with language abstractions  $L$  based on the instruction  $I$  and previous successful examples  $\{e^1, \dots, e^k\}$ .

$$F_{abstract} : (\xi_{noisy}, I, \{e^1, \dots, e^k\}) \rightarrow (\xi_{opt}, L) \quad (1)$$

The VLM is prompted to annotate subgoals [2], causal relationships [14], state changes [1], and relevant state [5], highlighting important demonstration aspects.

### 2.2. Human-in-the-loop Abstraction Verification

Human-in-the-loop learning involves executing the optimized trajectory  $\xi_{opt}$  in the environment. A human monitors and provides feedback  $H(a_t, o_t)$  on failures. The VLM is then prompted to revise the trajectory and annotations:

$$\Xi_{update}(\xi_{opt}, H(a_t, o_t), L, I, \{e^1, \dots, e^k\}) \rightarrow \xi'_{opt}, L' \quad (2)$$

The environment is reset after feedback, and the process repeats until the task is successful or a limit is reached.

### 2.3. Agent Deployment After Example Learning

Once examples are learned, the agent uses them to perform new tasks. The VLM generates actions based on the new instruction  $I$ , visual and textual state, and retrieving the top  $K$  ICAL examples from the learned set  $E$  to guide action generation, with similarity scores based on input instruction, textual, and visual state features. Implementation uses `gpt-4-1106-vision-preview` for the text generation, unless otherwise noted.

### 3. Experiments

#### 3.1. Environments

**TEACH** [11] This dataset includes over 3,000 dialogues in AI2-THOR [8], focusing on agents deducing actions from dialogue for tasks such as MAKE COFFEE. The data is split into training, seen, and unseen validations. Agents receive image observations and perform actions like `pickup(X)` based on given dialogue-based instructions. Using HELPER modules [12], agents navigate and manipulate environments, and are evaluated on fulfilling all task conditions. 250 human demonstrations in TEACH training were used for ICAL. We use ground truth semantic segmentation and depth for TEACH evaluations.

**VisualWebArena** [6] This dataset consists of 910 episodes of web tasks on sites like Shopping and Reddit. Agents receive visual and textual instructions and interact with web-pages using image screenshots, html text, and a fixed action API. Success is measured by task completion based on user instructions. 30 human demonstrations and 62 GPT4V-collected demonstrations were abstracted using ICAL.

**Ego4D** [4] Ego4D is a daily life activity video dataset of hundreds of scenarios. We focus on the long-term action anticipation task to predict the future user actions given an RGB egocentric video. 100 demonstrations from validation set were abstracted using ICAL, without human-in-the-loop (VLM-driven Abstraction Generation only). We take 200 unseen validation videos for evaluation. All GPT4V evaluations use DEVA tracking [3] + Set-of-Marks [17] for image inputs. Supervised baseline uses SlowFast with MViT [4].

#### 3.2. TEACH Evaluation

ICAL outperforms baseline approaches significantly, achieving a 17.9% absolute improvement in task success rate over unprocessed demonstrations. ICAL outperforms the handwritten examples used by the existing state-of-the-art in TEACH by 12.6% in goal condition success and 0.6% in task success (Table 1).

**Continual Improvement** ICAL shows progressive improvement in task success as more examples are learned, highlighting the benefits of continual learning and example accumulation (Figure 1).

**Improving with Fine-Tuning** Fine-tuning the LLM on ICAL examples further improves performance, especially when combined with retrieval-augmented generation, indicating the utility of integrating learned examples in training (Table 4).

**Ablation Studies** Ablation studies confirm that each component of ICAL—from the abstraction phase to the human-in-the-loop phase—is crucial for achieving the observed improvements in performance (Table 1).

Table 1. Evaluation of the TEACH unseen validation set using GPT-3.5-1106. Visual demos utilize an inverse dynamics model, while Kinesthetic demos are labeled with ground truth actions. GC = goal-condition success

	Success	GC
HELPER handwritten [12]	34.5	36.7
Zero-Shot CoT [7]	11.8	24.6
Raw Visual Demos	17.2	26.6
Raw Kinesthetic Demos	26.5	29.5
ICAL (ours)	<b>35.1</b>	<b>49.3</b>
w/o abstraction phase	29.4	44.9
w/o human-in-the-loop	29.9	41.0
w/ re-ranking [15]	<b>35.3</b>	<b>51.7</b>
w/ GPT4	41.7	63.6

Table 2. Evaluation on Ego4D Long Term Action Anticipation unseen validation subset. ICAL does not use human-in-the-loop due to the passive nature of this task.

	ED@(Z=20)			Success	GC
	Verb	Noun	Action		
Supervised [4]	0.7251	0.7393	0.9235	zero-shot retrieval	11.8
639x more data				finetuned	35.1
Few-shot CoT [16]	0.7877	0.7575	0.9414	finetuned + retrieval	23.2
ICAL (ours)	<b>0.7802</b>	<b>0.6934</b>	<b>0.9242</b>		40.3
					54.2

Table 3. Evaluation results on VisualWebArena. Ablations are done on a reduced 257 episodes.

	Seen	Unseen	Avg.
GPT4V+SoM [6]	16.3	14.1	14.3
ICAL (ours)	<b>38.8</b>	<b>20.9</b>	<b>22.7</b>

Ablations			
GPT4V+SoM [6]	11.5	12.9	12.7
ICAL (ours)	28.0	<b>21.6</b>	22.2
w/o image	28.0	17.3	19.0
w/ full trajectory	<b>57.7</b>	<b>21.6</b>	<b>25.5</b>

Table 4. Results on finetuning the LLM on the ICAL demonstrations. The model used is GPT3.5-turbo-1106.

	ED@(Z=20)			Success	GC
	Verb	Noun	Action		
Supervised [4]	0.7251	0.7393	0.9235	zero-shot retrieval	11.8
639x more data				finetuned	35.1
Few-shot CoT [16]	0.7877	0.7575	0.9414	finetuned + retrieval	23.2
ICAL (ours)	<b>0.7802</b>	<b>0.6934</b>	<b>0.9242</b>		40.3
					54.2

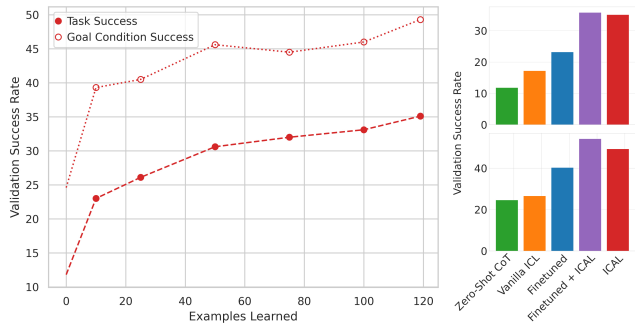


Figure 1. ICAL enables continual learning with more examples. Filled markers ● indicate full task success, while open markers ○ denote partial task (goal-condition) success. A bar plot comparing ICAL at 120 examples to baselines shows notable advancement over the Vanilla ICL, which represents the Raw Visual Demos baseline.

#### 3.3. Ego4D Evaluation

See Table 2. ICAL demonstrates superior few-shot performance on Ego4D action anticipation compared to handwritten few-shot GPT4V that uses chain of thought [16] by 6.4 noun and 1.7 action edit distance. ICAL also remains competitive with the fully supervised baseline [4] despite using 639x less training data. We find GPT4V video processing to have the most trouble with verb prediction.

#### 3.4. Visual Web Navigation Evaluation

ICAL demonstrates state-of-the-art performance on VisualWebArena, outperforming previous best methods by an absolute 8.4% (relative 58.74%) in success rate (Table 3).

## References

- [1] Lisa Feldman Barrett and Moshe Bar. See it with feeling: affective predictions during object perception. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 364(1521):1325–1334, 2009. [1](#)
- [2] Matthew M Botvinick, Yael Niv, and Andrew G Barto. Hierarchically organized behavior and its neural foundations: A reinforcement learning perspective. *cognition*, 113(3):262–280, 2009. [1](#)
- [3] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *ICCV*, 2023. [2](#)
- [4] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. [2](#)
- [5] Mark K Ho, David Abel, Carlos G Correa, Michael L Littman, Jonathan D Cohen, and Thomas L Griffiths. People construct simplified mental representations to plan. *Nature*, 606(7912): 129–136, 2022. [1](#)
- [6] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024. [2](#)
- [7] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022. [1](#), [2](#)
- [8] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli Vanderbilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017. [2](#)
- [9] Bodhisattwa Prasad Majumder, Bhavana Dalvi Mishra, Peter Jansen, Oyvind Tafjord, Niket Tandon, Li Zhang, Burch Callison-Burch, and Peter Clark. Clin: A continually learning language agent for rapid task adaptation and generalization. *arXiv*, 2023. [1](#)
- [10] Jesse Mu, Victor Zhong, Roberta Raileanu, Minqi Jiang, Noah Goodman, Tim Rocktäschel, and Edward Grefenstette. Improving intrinsic exploration with language abstractions (2022). *URL <https://arxiv.org/abs/2202.08938>*. [1](#)
- [11] Aishwarya Padmakumar, Jesse Thomason, Ayush Srivastava, Patrick Lange, Anjali Narayan-Chen, Spandana Gella, Robinson Piramuthu, Gokhan Tur, and Dilek Hakkani-Tur. Teach: Task-driven embodied agents that chat, 2021. [2](#)
- [12] Gabriel Sarch, Yue Wu, Michael Tarr, and Katerina Fragkiadaki. Open-ended instructable embodied agents with memory-augmented large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023. [2](#)
- [13] Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023. [1](#)
- [14] Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. How to grow a mind: Statistics, structure, and abstraction. *science*, 331(6022):1279–1285, 2011. [1](#)
- [15] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022. [2](#)
- [16] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022. [2](#)
- [17] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v, 2023. [2](#)