

# Continual Learning with Language Agents

**Bodhisattwa Majumder**

Allen Institute for AI

 @mbodhisattwa





When *not* watching  
seaplanes out of my  
office window

I dabble with

**Interactive Systems**  
**Language Agents**  
**Dialog Models**  
**Multi-Agent Systems**  
**Scientific Discovery**  
**Social Science**

# Outline

## Background

CLIN: Continual Learning from Interactions

Proposed Architecture

What does CLIN learn over time?

Results on ScienceWorld & ALFWorld

SSO: Skill Set Optimization

Skills

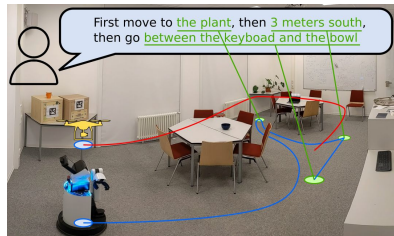
Skill Set Optimization

Results on ScienceWorld & NetHack



# Sequential Decision-making (SDM)

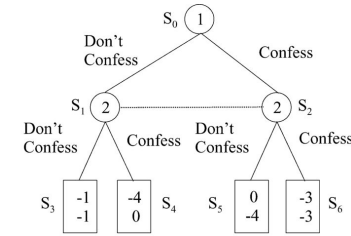
Real world decision-making tasks are **sequential** in nature



navigation

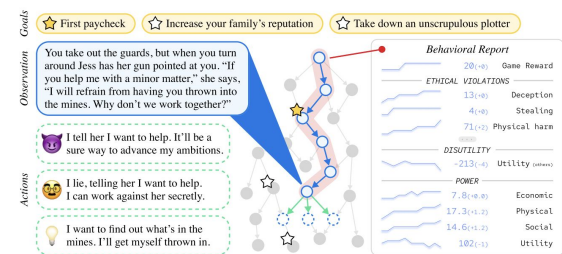
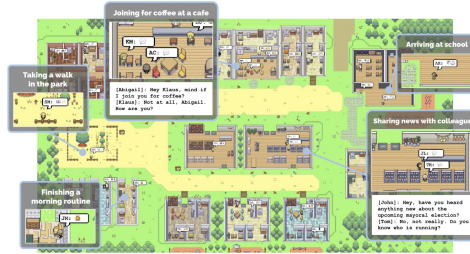


shopping



communication

**Text-worlds** are great simulations of real-world phenomena for testing systems ability to **long-term planning, sequential reasoning, and decision making**



# TextWorlds: E.g., ScienceWorld

## Task Description:

Your task is to boil water.



### > look around

This room is called the **hallway**.  
In it, you see: a picture, a substance called air, A door to the art studio (that is open) A **door to the kitchen** (that is open)...

### > go to kitchen

You move to the kitchen.

### > look around

This room is called the kitchen.  
In it, you see:  
a cupboard, a sink, a fridge, ...

. . . .

### > pick up thermometer

You move the thermometer to the inventory.

### > open cupboard

The cupboard is now open.

### > pick up metal pot

You move the metal pot to the inventory.

### > look around

This room is called the kitchen. ...

### > move metal pot to sink

You move the metal pot to the sink.

### > activate sink

The sink is now activated.

. . . .



### > focus on substance in metal pot

You focus on the water.

### > pick up metal pot

You move the metal pot to the inventory.

### > move metal pot to stove

You move the metal pot to the stove.

### > activate stove

The stove is now activated.

### > examine substance in metal pot

a substance called water

### > use thermometer in inventory

on substance in metal pot  
the thermometer measures a temperature of 13 degrees celsius

### > use thermometer in inventory on substance in metal pot

the thermometer measures a temperature of 102 degrees celsius  
(Task Completed)

our task is to boil

**ate:** This room is called  
way. In it, you see: ...  
**an episode**

- 
- k trials**

The thermometer  
shows a temperature of  
degrees celsius

# Existing Approaches for SDM

Not adaptive

Model classes	Learning	Interpretability	Generalization
RL (DRRN, CALM, KG-A2C)	Policies from environment feedback	Low	Low
Supervised (TDT)	Behavior cloning from gold trials	Low	Low
Generative (GPT-4)	Pre-training + Instruction tuning	Low	Moderate
Hybrid (SwiftSage)	Mix of Supervised + Generative	Low	Moderate

Adaptive

Meta RL (AdA)	Online RL on previous trials	Low	High
Reflexion	Mistakes from previous trials	High	Moderate
<b>What we want</b>	<b>More than mistakes</b>	High	High

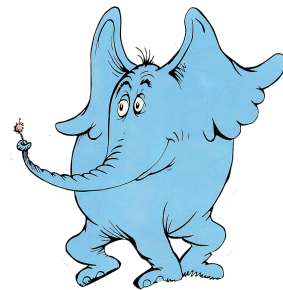
# Research Questions

Can SDM environments and tasks be  
**continually learnt**  
from interacting and observing world changes?

Can we build an agent that can  
**quickly adapt and generalize**  
to a new task or environment at the test time?



# CLIN: Continually Learning From INteractions



**Bodhi**, Bhavana, Peter Jansen, Oyvind,  
Niket, Harry, Chris, Pete



# Outline

Background

## **CLIN: Continual Learning from Interactions**

### **Proposed Architecture**

What does CLIN learn over time?

Results on ScienceWorld & ALFWorld

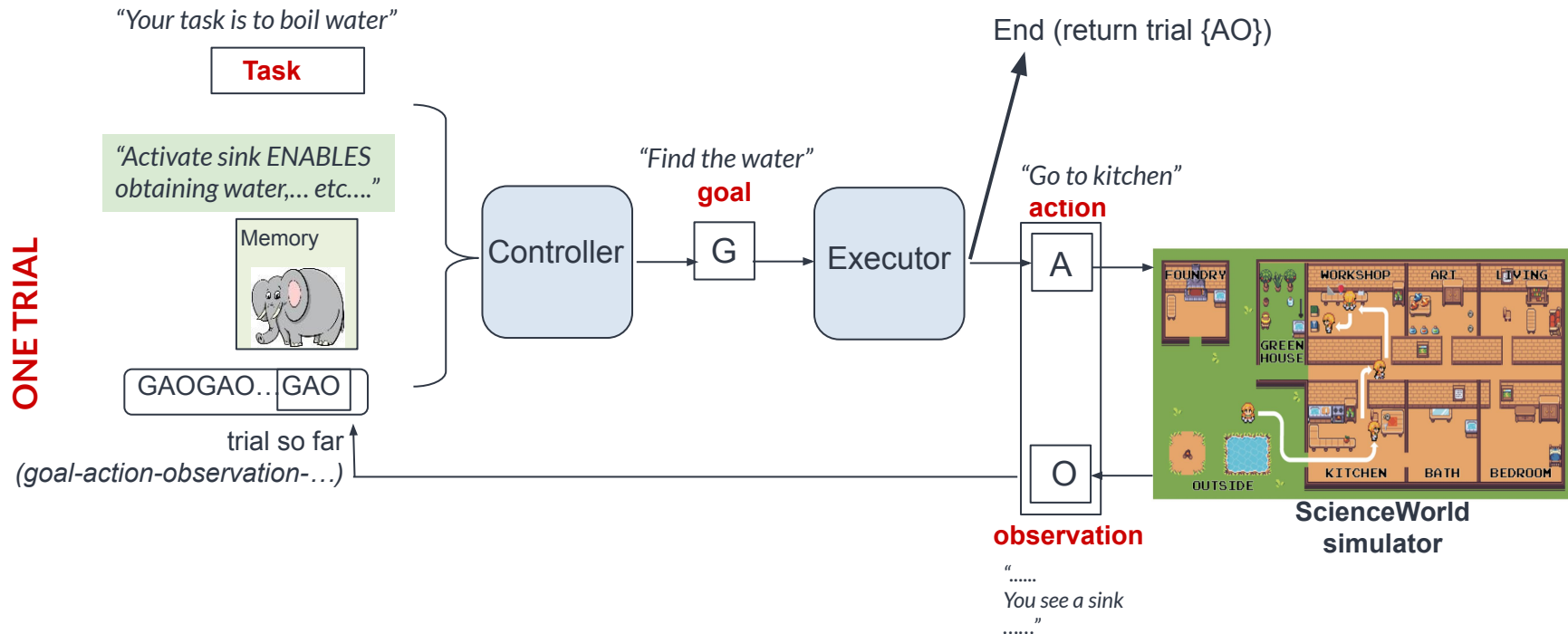
## **SSO: Skill Set Optimization**

Skills

Skill Set Optimization

Results on ScienceWorld & NetHack

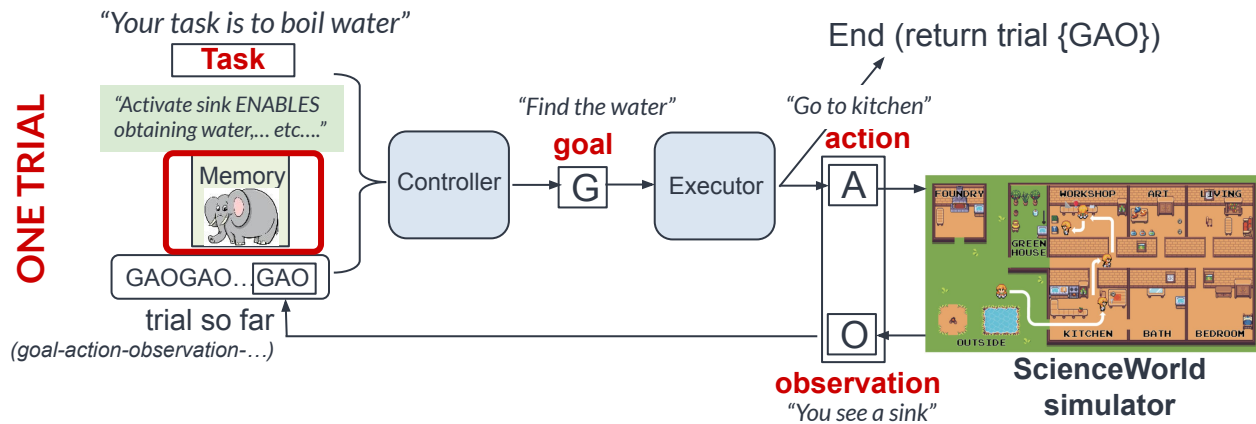
# CLIN: Continually Learning from INteractions



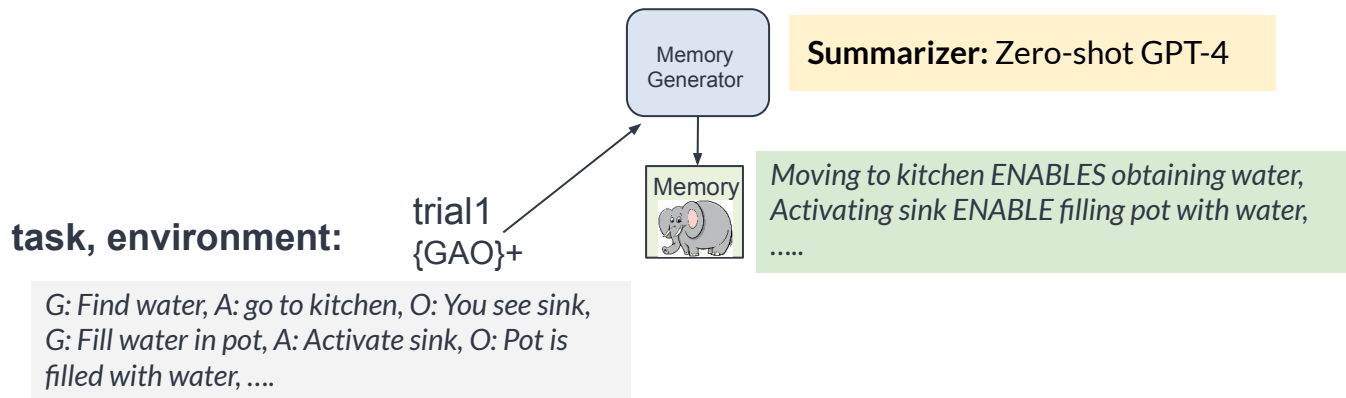
**\*\* Controller + Executor: Zero-shot GPT-4**  
(unlike Reflexion/ReAct, we do not use any task-specific few-shot examples)

# CLIN: Continually Learning from INteractions

ONE TRIAL

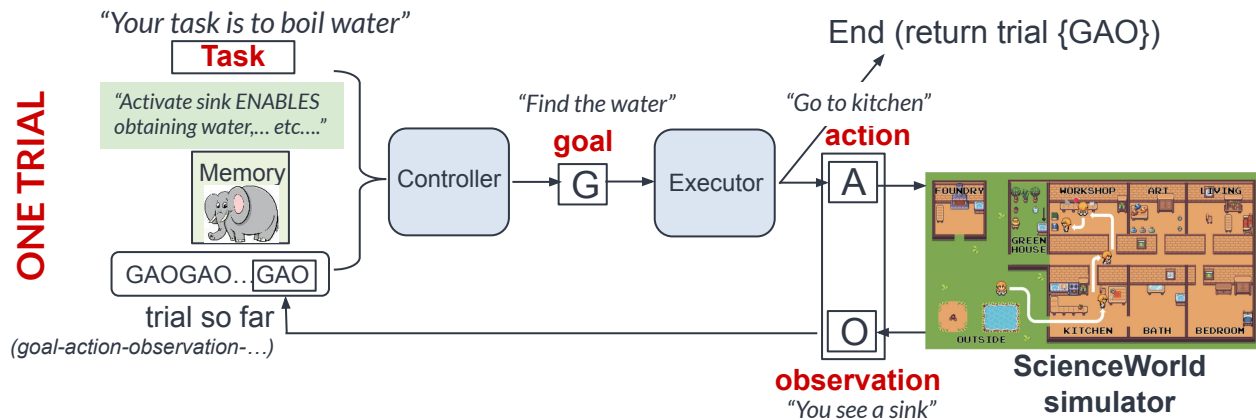


LEARNING

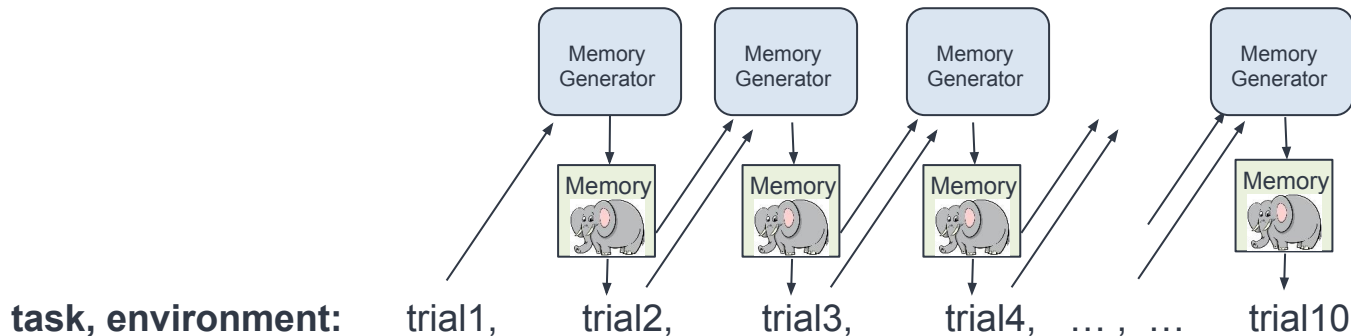


# CLIN: Continually Learning from INteractions

ONE TRIAL



LEARNING



Each trial refines the memory



# Outline

Background

## **CLIN: Continual Learning from Interactions**

Proposed Architecture

**What does CLIN learn over time?**

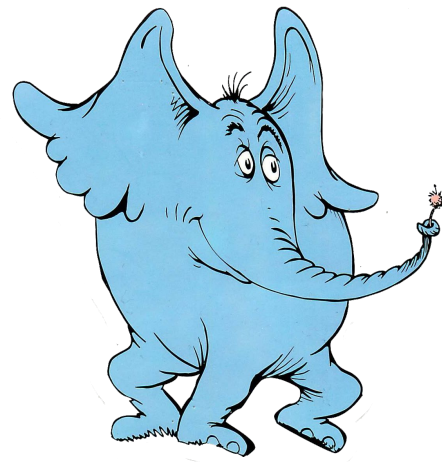
Results on ScienceWorld & ALFWorld

SSO: Skill Set Optimization

Skills

Skill Set Optimization

Results on ScienceWorld & NetHack



# Memory

Learning **state transitions** is essential for SDM

1. actions enabling **desired** state transitions
2. actions producing **undesired** or no changes
3. state transitions **contributing to the task**

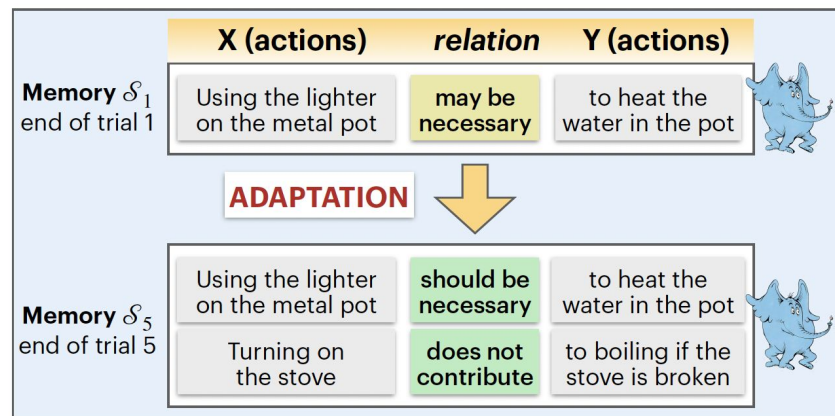
A collection of natural language statements capturing **causal abstractions of action-effects**  
*favorable to exploit at test-time like hindsight experience replay*

**Good effects:**  $X \rightarrow$  is necessary to  $\rightarrow Y$

**Bad effects:**  $X \rightarrow$  does not contribute  $\rightarrow Y$

**Uncertainty**

**Low:** may be; **High:** should be



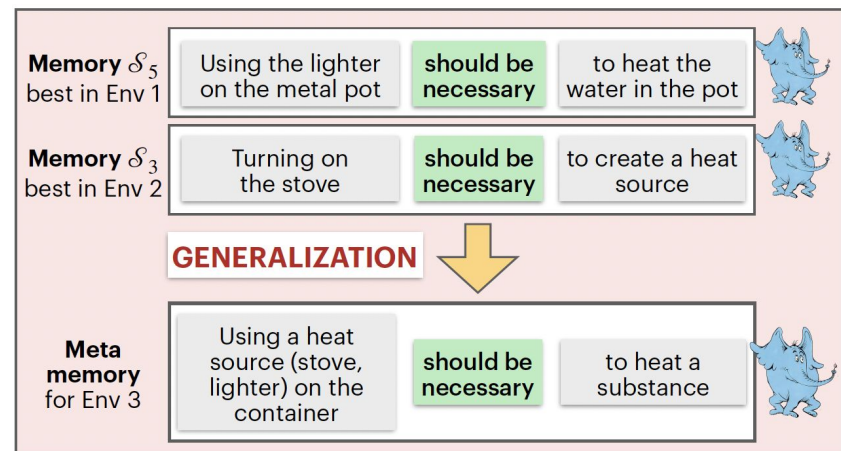
# Meta-Memory

Task- and environment-specific memory cannot help generalize such as knowing how to boil water may not help knowing how to boil cadmium unless *generalized abstractions*.

Select the best memories from past attempts across diverse environments/tasks  
*auto-curriculum selection*

Meta-memory with *generalized* instruction:

“Generate insights to solve the same task in a new environment configuration”



# CLIN: Summary

## CREATION (Env1, Trial1)

Task: **Grow an orange**

Goal: Find seeds

Action: Go to the bedroom

Observation: ...(no seeds)...

Action: Go to the garden

Observation: ...(no seeds)...

Action: Go to the kitchen

Observation: You see seeds

Action: Pick up seeds

Goal: Plant the seeds ...

**Trial failed!**

### MEMORY:

Going to the kitchen **may be necessary** to find seeds



## ADAPTATION (Env1, Trial2)

Task: **Grow an orange**

Goal: Find seeds

Memory:

Going to the kitchen **may be necessary** to find seeds

Action: Go to the kitchen

Observation: you see seeds

Action: pick up seeds...

Goal: plant the seeds

Action: move seeds to the pot ...

**Trial completed successfully!**

### MEMORY:

Going to the kitchen **may be necessary** to find seeds  
Moving seeds to the pot **may be necessary** for planting seeds



Env1, Trial3  
(etc)



## GENERALIZATION (for new Env3)

Env1



Trial1



Trial2



Trial3

Env2



Trial1



Trial2



Trial3

Env3, Trial1

### MEMORY:

Moving to different rooms **may be necessary** to find an object



# Outline

Background

## **CLIN: Continual Learning from Interactions**

Proposed Architecture

What does CLIN learn over time?

**Results on ScienceWorld & ALFWorld**

SSO: Skill Set Optimization

Skills

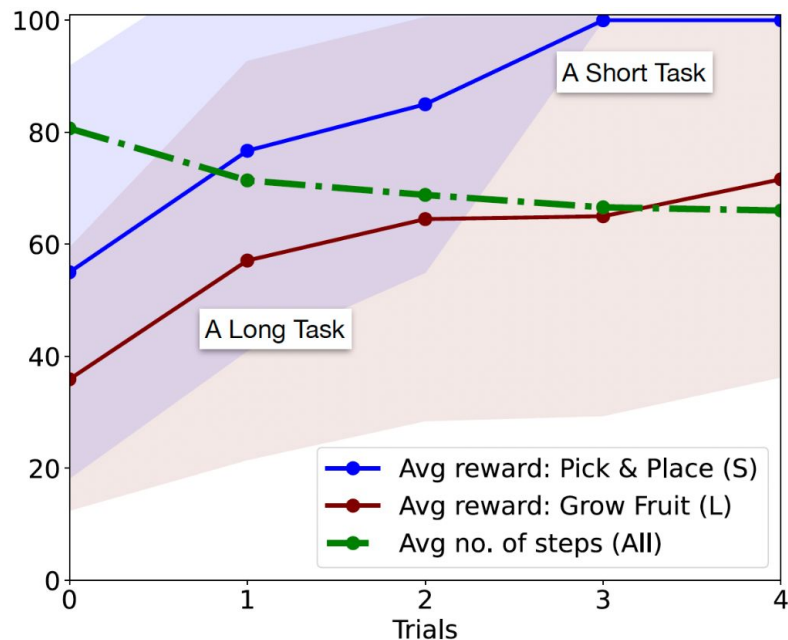
Skill Set Optimization

Results on ScienceWorld & NetHack

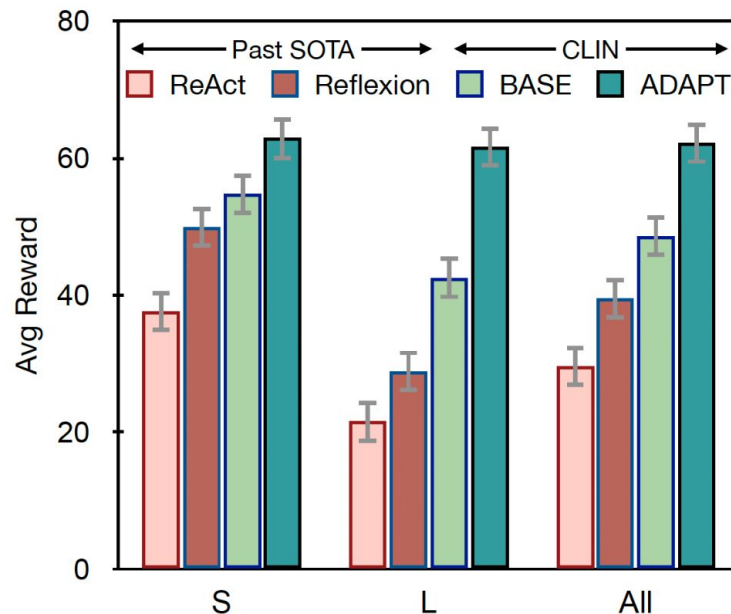


# CLIN Exhibits Rapid Task Adaptation

Quick adaptation, improved efficiency



CLIN beats reflective SOTA



# CLIN Generalizes to Novel Environments

Train:

Boil water

Boil chocolate

Test:

Boil Cadmium

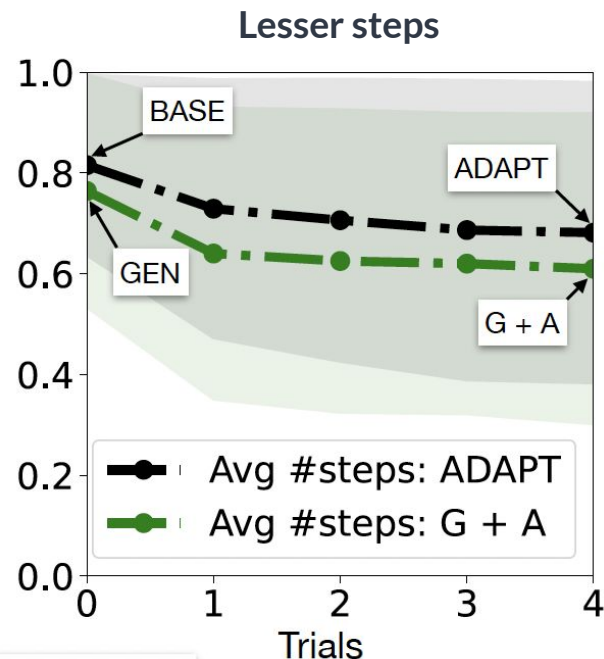
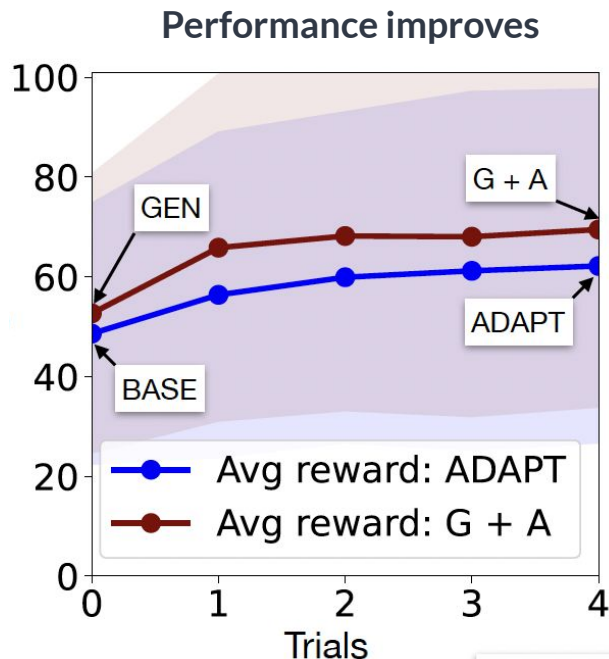
CLIN even beats  
imitation learning  
baselines (that uses  
gold trajectories) in  
most lengthy,  
complex tasks

		RL Methods			Generative Language Agents			CLIN (ours)		
Task	Type	DRRN	KGA2C	CALM	SayCan	ReAct	Reflexion	BASE	GEN-ENV	G+A
Temp <sub>1</sub>	S	6.6	6.0	1.0	<b>26.4</b>	7.2	5.9	25.2	15.7	13.8
Temp <sub>2</sub>	S	5.5	11.0	1.0	8.0	6.1	28.6	53.2	49.7	<b>58.2</b>
Pick&Place <sub>1</sub>	S	15.0	18.0	10.0	22.9	26.7	64.9	92.5	59.2	<b>100.0</b>
Pick&Place <sub>2</sub>	S	21.7	16.0	10.0	20.9	53.3	16.4	55.0	<b>100.0</b>	<b>100.0</b>
Chemistry <sub>1</sub>	S	15.8	17.0	3.0	47.8	51.0	<b>70.4</b>	44.5	42.2	51.7
Chemistry <sub>2</sub>	S	26.7	19.0	6.0	39.3	58.9	70.7	56.7	85.6	<b>93.3</b>
Lifespan <sub>1</sub>	S	50.0	43.0	6.0	80.0	60.0	<b>100.0</b>	85.0	65.0	<b>100.0</b>
Lifespan <sub>2</sub>	S	50.0	32.0	10.0	67.5	67.5	84.4	70.0	75.0	<b>90.0</b>
Biology <sub>1</sub>	S	8.0	10.0	0.0	16.0	8.0	8.0	10.0	32.0	<b>32.0</b>
Boil	L	3.5	0.0	0.0	<b>33.1</b>	3.5	4.2	7.0	4.4	16.3
Freeze	L	0.0	4.0	0.0	3.9	7.8	7.8	<b>10.0</b>	8.9	<b>10.0</b>
GrowPlant	L	8.0	6.0	2.0	9.9	9.1	7.3	10.2	10.9	<b>11.2</b>
GrowFruit	L	14.3	11.0	4.0	13.9	18.6	13.0	35.9	70.8	<b>94.5</b>
Biology <sub>2</sub>	L	21.0	5.0	4.0	20.9	27.7	2.6	70.0	42.8	<b>85.6</b>
Force	L	10.0	4.0	0.0	21.9	40.5	50.6	53.5	70.0	<b>100.0</b>
Friction	L	10.0	4.0	3.0	32.3	44.0	<b>100.0</b>	56.5	70.0	94.0
Genetics <sub>1</sub>	L	16.8	11.0	2.0	67.5	25.7	50.9	77.4	84.5	<b>100.0</b>
Genetics <sub>2</sub>	L	17.0	11.0	2.0	59.5	16.8	23.7	62.3	61.4	<b>100.0</b>
<b>S</b>		22.1	19.1	5.2	36.5	37.6	49.9	54.7	58.3	<b>71.0</b>
<b>L</b>		11.2	6.2	1.9	29.2	21.5	28.9	42.5	47.1	<b>68.0</b>
<b>All</b>		16.7	12.7	3.6	32.9	29.6	39.4	48.6	52.7	<b>69.5</b>

# Efficient Generalization

Performance drops 6.2 point and in 10% episodes if we do not use **causal format** for memory insights

**Controller** adds 18 points to a base (ReAct) performance improving 44% episodes



# CLIN Generalizes to Novel Tasks

Train (in Env 1):

Boil water

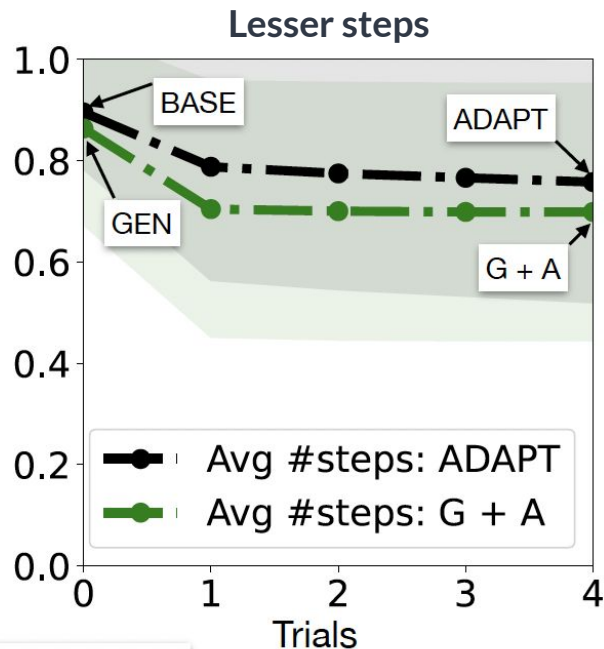
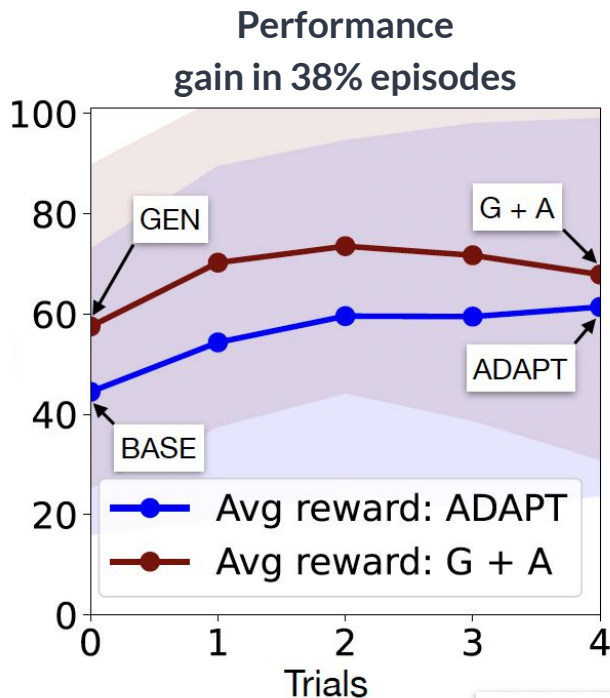
Boil apple juice

Test (in Env 1):

Freeze Water

The improvement  
attributes to *critical  
learning about the  
environment*

(apple juice is in the fridge)



# Memory Precision

Natural selection of good memory items over time shows CLIN can auto-correct when the starting memory is not applicable due to loss of specificity or lack of information.

CLIN converges to a more precise representation of the world

	GEN-ENV (Trial 0)	GEN-ADAPT (Best Trial)
No. of insights	100	105
Correct insights	72.0%	<b>91.4%</b>
Final score	39.1	<b>55.9</b>

	GEN-TASK (Trial 0)	GEN-ADAPT (Best Trial)
No. of insights	98	107
Correct insights	73.9%	<b>91.1%</b>
Final score	43.7	<b>58.1</b>



# Is Causal Abstraction Helpful?

Memory with no structure is generic (“Be clear with your actions”), often contains ungrounded information (“use a food processor”), and does not naturally abstract causal relations towards a world model (“this is unnecessary and wastes time”)

## Ablations for CLIN

Ablation Setup	$\Delta$ avg score ( $\downarrow$ )	%ep. drop. ( $\uparrow$ )
Abl-Causal-Memory	-6.2	10
Abl-Controller-BASE	-18.1	44.8

# ALFWorld

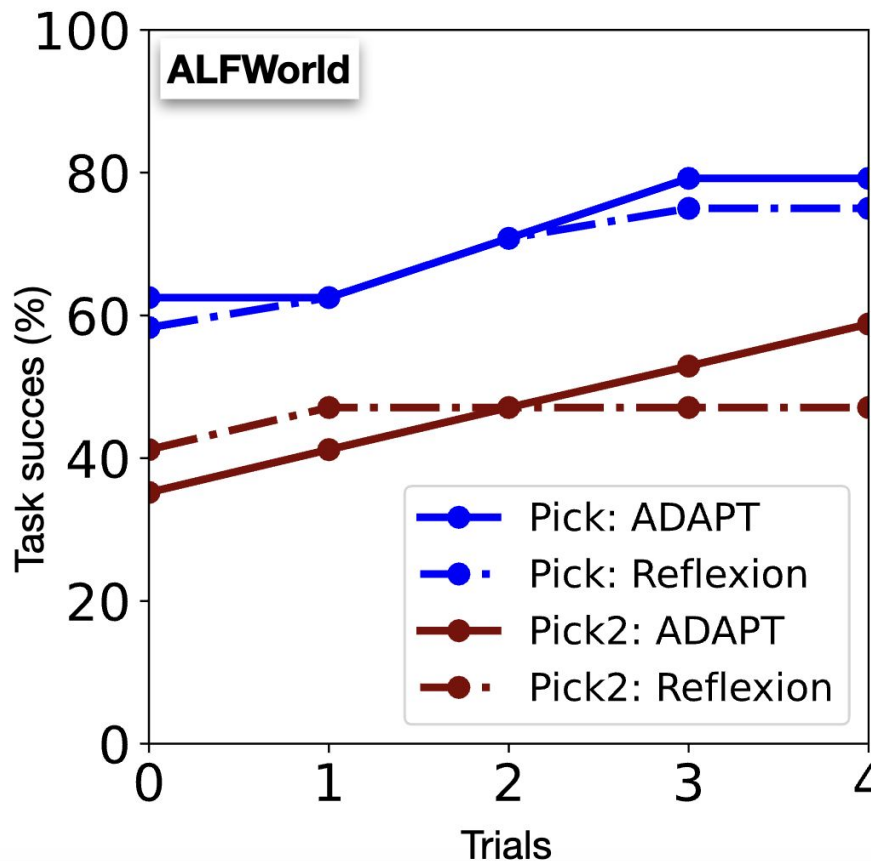


You are in the middle of a room. Looking quickly around you, you see a drawer 2, a shelf 5, a drawer 1, a shelf 4, a sidetable 1, a drawer 5, a shelf 6, a shelf 1, a shelf 9, a cabinet 2, a sofa 1, a cabinet 1, a shelf 3, a cabinet 3, a drawer 3, a shelf 11, a shelf 2, a shelf 10, a dresser 1, a shelf 12, a garbagecan 1, a armchair 1, a cabinet 4, a shelf 7, a shelf 8, a safe 1, and a drawer 4.

Your task is to: *put some vase in safe.*

> go to shelf 6

You arrive at loc 4. On the shelf 6, you see a vase 2.



# Failures

CLIN is able to **compose insights**

No stove, use furnace (Env 1) + Go to Kitchen for apple juice (Env 2)

But when it **fails**, it is due to:

1. **Lack of exploration**

If it has never visited an art studio,  
it will never “explore” to reach art studio for collecting paints

2. **Poor memory retrieval**

It knows to use stove for heating OR use furnace when stove is broken  
BUT to boil cadmium it needs to use furnace even if the stove is working

# Failures



	Transferable		Scalable	
	Multi-task	Modular	Lossless	Sublinear
Fewshot	X	X	✓	X
Reflexion	X	X	X	X
ExpeL	✓	X	✓	✓
Voyager	X	✓	✓	X
CLIN	✓	✓	X	✓

1. Shinn, Noah, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. "Reflexion: Language agents with verbal reinforcement learning." arXiv preprint arXiv:2303.11366 (2023).
2. Zhao, Andrew, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. "ExpeL: LLM Agents Are Experiential Learners." arXiv preprint arXiv:2308.10144 (2023).
3. Wang, Guanzhi, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. "Voyager: An open-ended embodied agent with large language models." arXiv preprint arXiv:2305.16291 (2023).
4. Majumder, Bodhisattwa Prasad, Bhavana Dalvi Mishra, Peter Jansen, Oyvind Tafjord, Niket Tandon, Li Zhang, Chris Callison-Burch, and Peter Clark. "CLIN: A Continually Learning Language Agent for Rapid Task Adaptation and Generalization." arXiv preprint arXiv:2310.10134 (2023).

# SSO: Skill Set Optimization

Kolby, **Bodhi**, Bhavana,  
Sameer, Pete, Roy



UCI



# Outline

Background

CLIN: Continual Learning from Interactions

Proposed Architecture

What does CLIN learn over time?

Results on ScienceWorld & ALFWorld

**SSO: Skill Set Optimization**

**Skills**

Skill Set Optimization

Results on ScienceWorld & NetHack

# Skills

- World model information should:
  - Be general, composable, editable, and retrievable
  - Contribute to LLM agent's knowledge of the world model (state & action transitions)



# Skill Definition

## Target:

- goal state feature

## Prerequisites:

- initial state features
- used for retrieval

## Instructions:

- generic actions to execute

### Example generated Skill

**Target:** agent is in the 'target location'

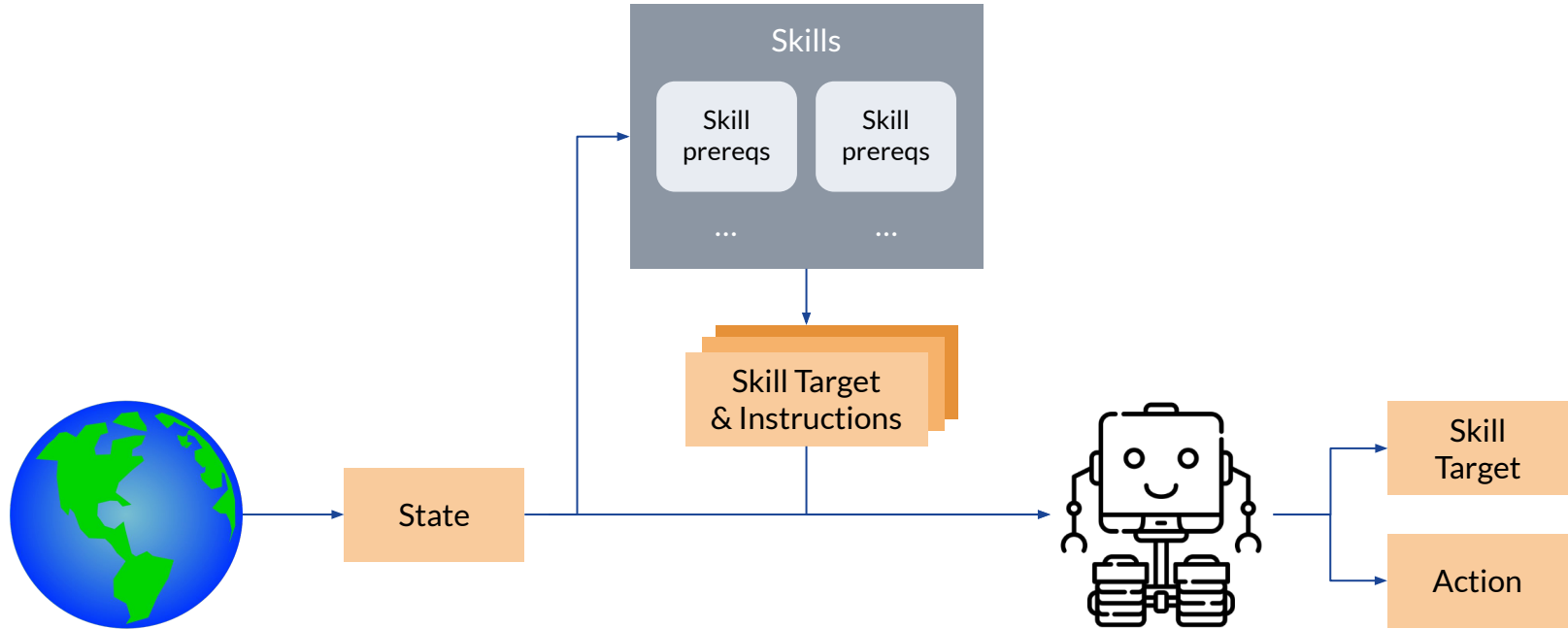
**Prereqs:**

1. agent is in a location that has a door leading to a hallway
2. there exists a known target location to which agent needs to move
3. agent is able to move (not restricted or blocked)

**Instructions:**

1. go to hallway
2. go to 'target location'

# Using Skills



# Outline

Background

CLIN: Continual Learning from Interactions

Proposed Architecture

What does CLIN learn over time?

Results on ScienceWorld & ALFWorld

**SSO: Skill Set Optimization**

Skills

**Skill Set Optimization**

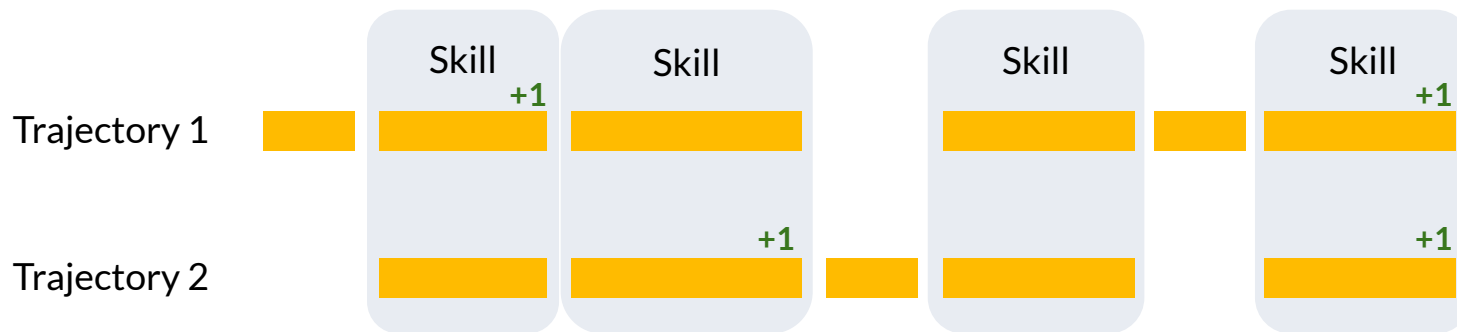
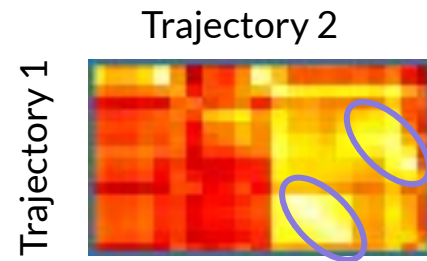
Results on ScienceWorld & NetHack

# Skill Set Optimization

We want to extract reusable sequences that lead to rewards

## 1. Find common sub trajectories

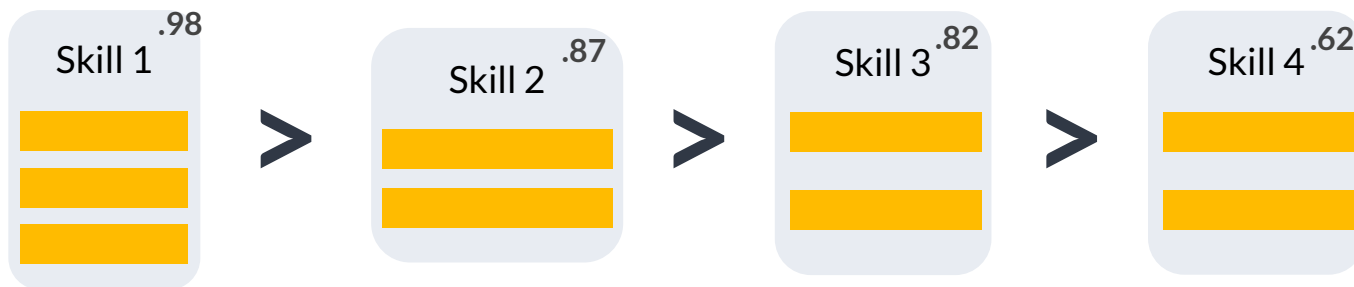
- Trim trajectories to end in positive rewards
- Align sub trajectories using state, action embedding from LLM



# Skill Set Optimization

We want to extract reusable sequences that lead to rewards

1. Find common sub trajectories
2. Score and sort skills
  - a. Similarity
  - b. Reward
  - c. Length

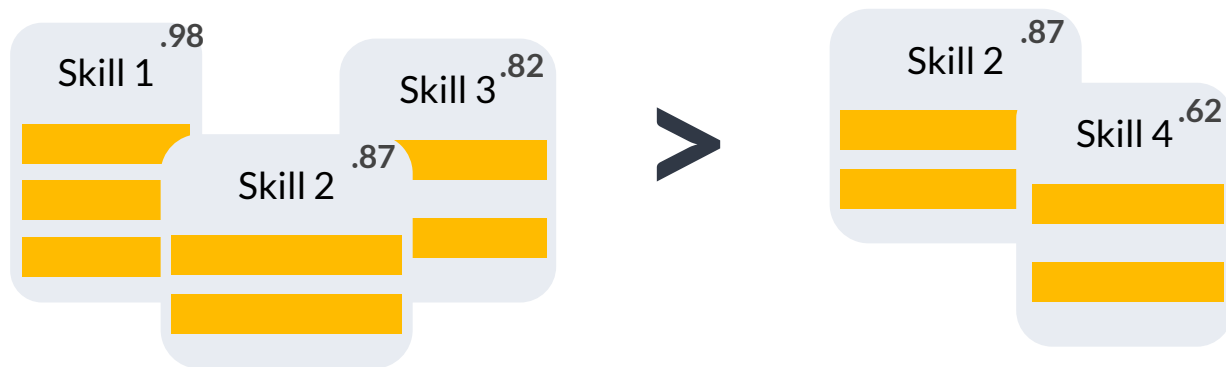




# Skill Set Optimization

We want to extract reusable sequences that lead to rewards

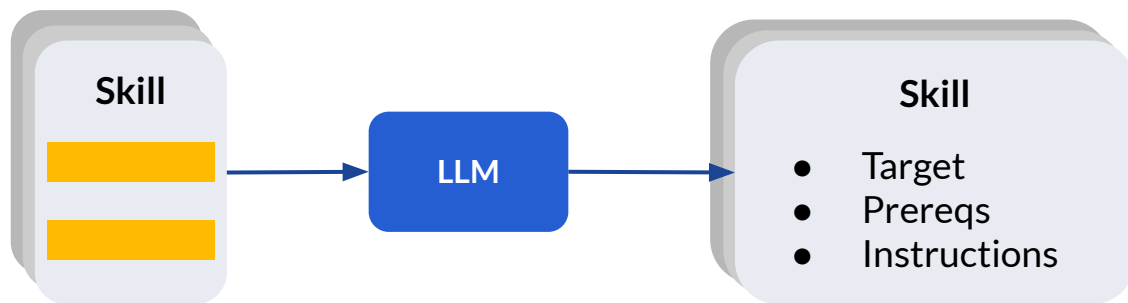
1. Find common sub trajectories
2. Score and sort skills
3. Construct skill set using beam search
  - a. Do not allow skill sub trajectories to overlap
  - b. Select best beam based on sum of scores: similarity, reward, length



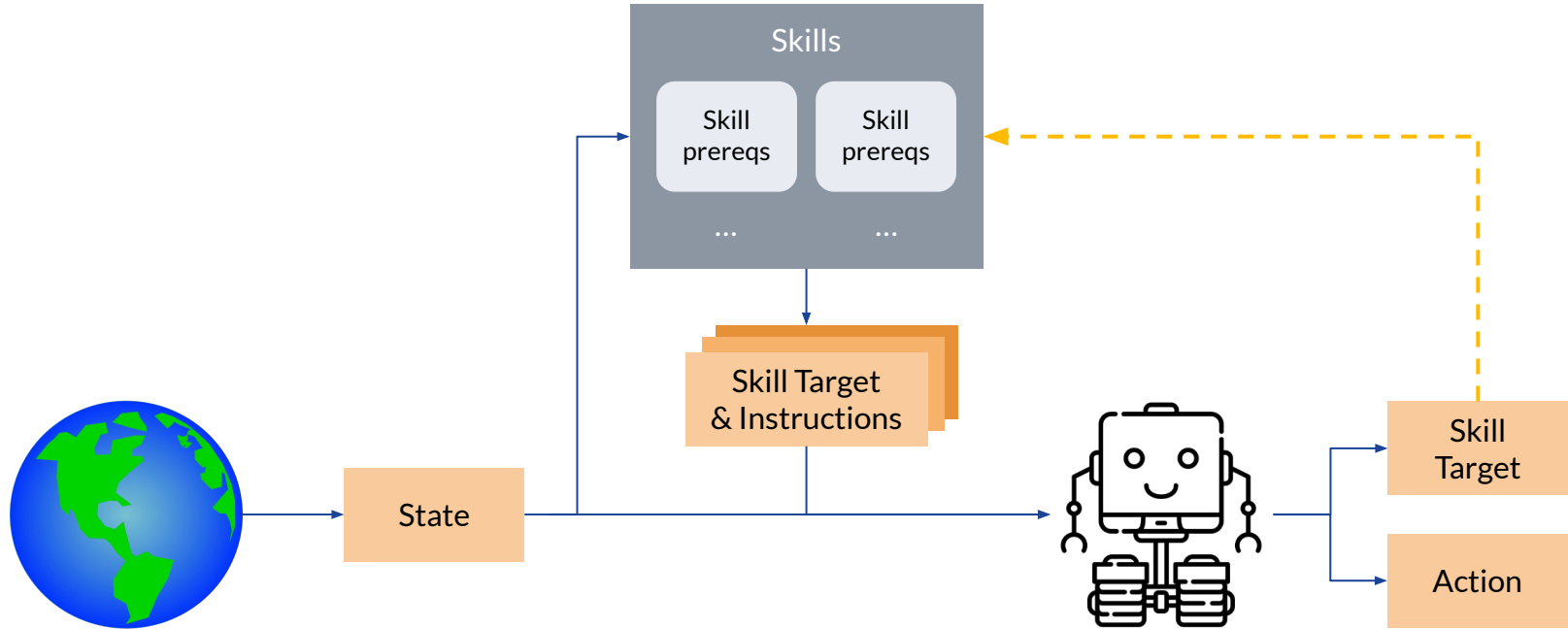
# Skill Set Optimization

We want to extract reusable sequences that lead to rewards

1. Find common sub trajectories
2. Score and sort skills
3. Construct skill set using beam search
4. Generate skill target, prerequisites, and instructions

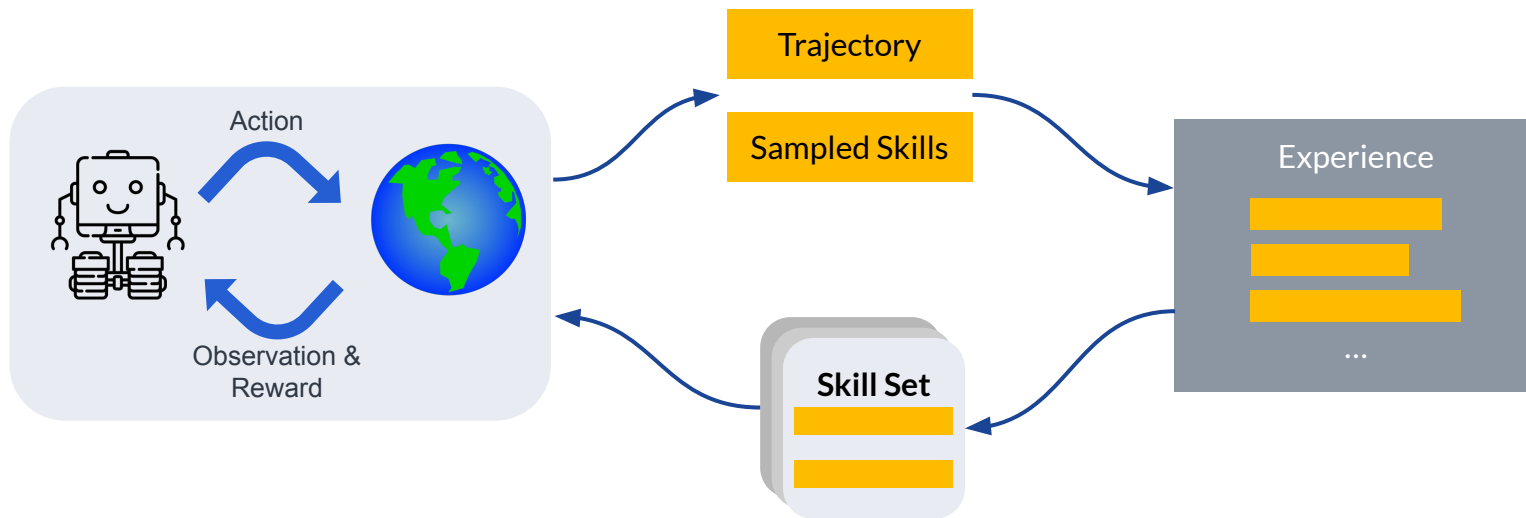


# Skill Set Optimization



# Skill Set Optimization

- Prioritize sampled skills that lead to positive reward
- Black list sampled skills that lead to negative reward
- After every trajectory, extract skills from last N trajectories



# Outline

Background

CLIN: Continual Learning from Interactions

Proposed Architecture

What does CLIN learn over time?

Results on ScienceWorld & ALFWorld

**SSO: Skill Set Optimization**

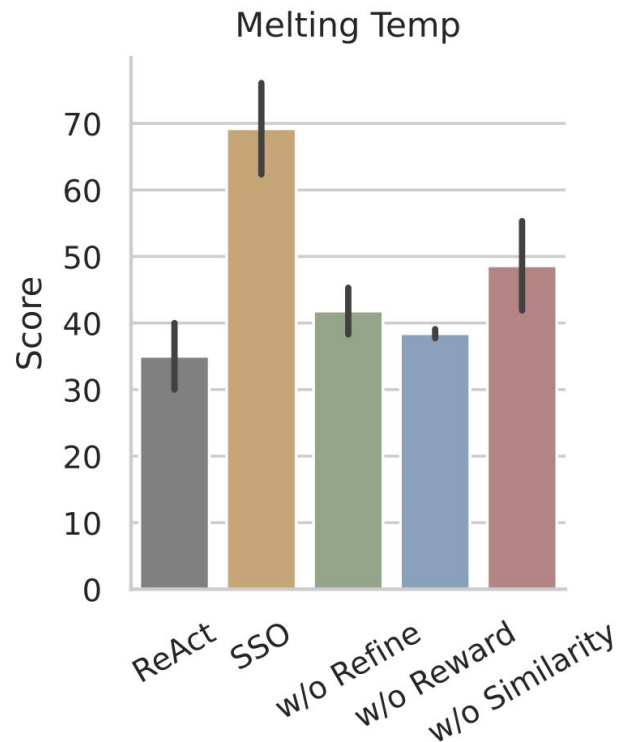
Skills

Skill Set Optimization

**Results on ScienceWorld & NetHack**

# SSO improves over CLIN

ScienceWorld Task	ReAct	Adaptation			Transfer	
		Reflexion	CLIN	SSO	CLIN	SSO
Temperature	7.2	5.9	14.3	<b>100</b>	15.7	<b>71.6</b>
Melting Temp	6.1	28.6	51.8	<b>97.3</b>	49.7	<b>69.2</b>
Find Plant	26.7	64.9	<b>100</b>	<b>100</b>	59.2	<b>100</b>
Find Living	53.3	16.4	<b>100</b>	96.7	<b>100</b>	90
Chemistry	51	70.4	44.4	<b>82.6</b>	42.2	<b>48</b>
Color Mixing	58.9	70.7	56.7	<b>81.1</b>	<b>85.6</b>	71.1
Lifespan, Longest	61	<b>100</b>	<b>100</b>	<b>100</b>	65	<b>90</b>
Lifespan, Shortest	67.5	84.4	90	<b>100</b>	75	<b>80</b>
Life Stages, Plant	8	<b>8</b>	<b>8</b>	6.2	<b>32</b>	3.4
Life Stages, Animal	27.7	2.6	81	<b>100</b>	42.8	<b>77</b>
Boil	3.5	4.2	15.2	<b>81.7</b>	4.4	<b>48.7</b>
Freeze	7.8	7.8	10	<b>74.3</b>	8.9	<b>38.9</b>
Grow Plant	9.1	7.3	11	<b>86.6</b>	10.9	<b>61.2</b>
Grow Fruit	18.6	13	71.6	<b>78</b>	<b>70.8</b>	28.3
Gravity	40.5	50.6	<b>100</b>	<b>100</b>	70	<b>74</b>
Friction	44	<b>100</b>	72.5	94	<b>70</b>	67.5
Genetics, Known	25.7	50.9	<b>100</b>	78.5	<b>84.5</b>	42.5
Genetics, Unknown	16.8	23.7	<b>92.6</b>	48.7	<b>61.4</b>	20.3
Average	29.6	39.4	62.2	<b>83.7</b>	52.7	<b>60.1</b>

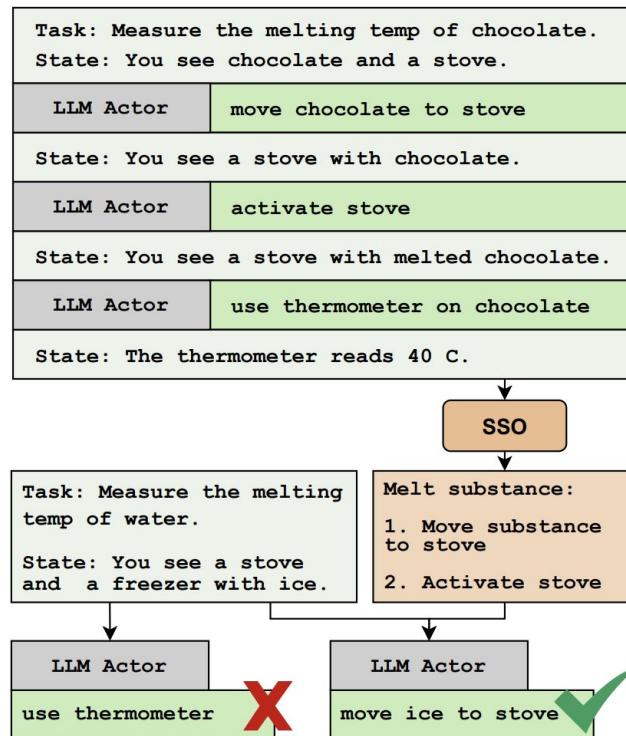


# An example skill

## ScienceWorld Melting Temp Task

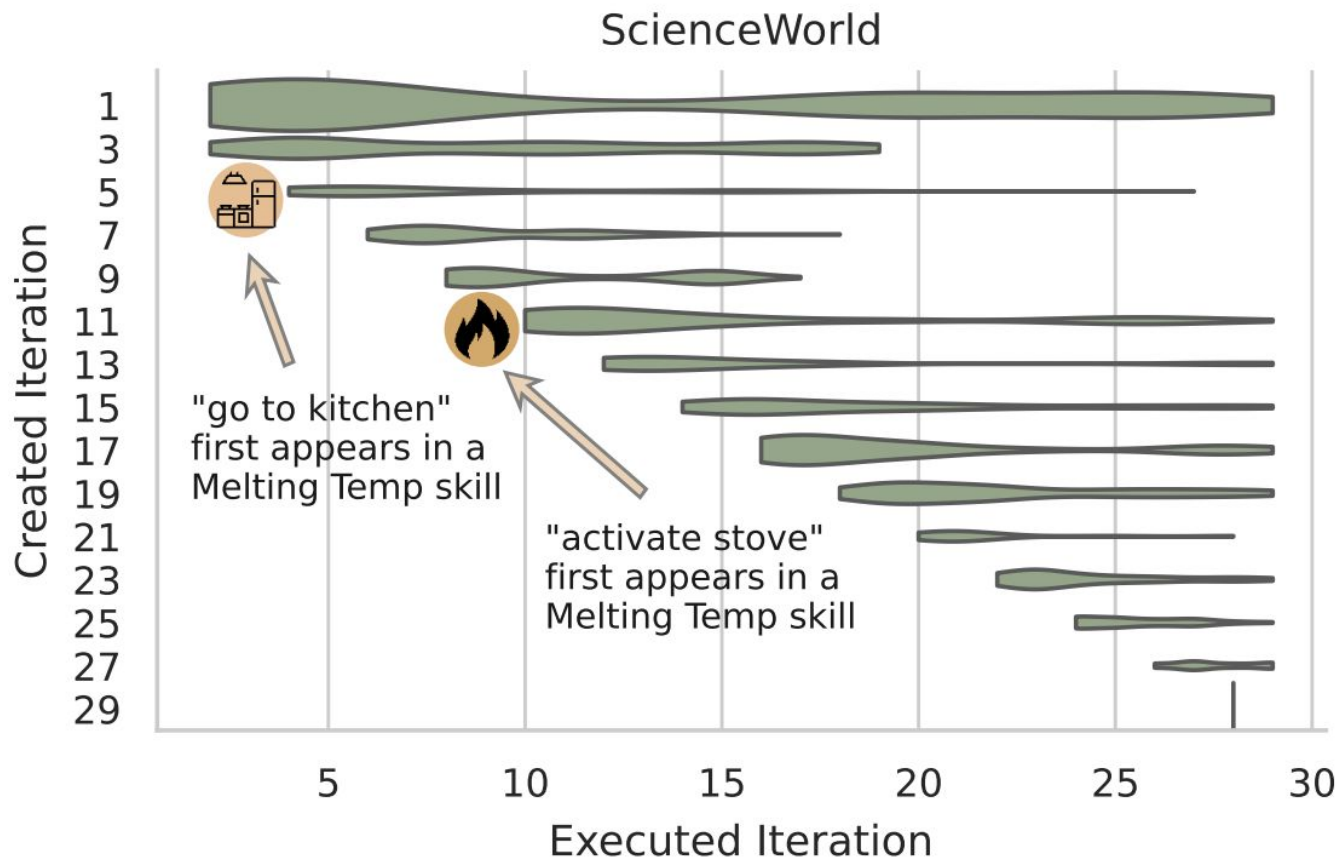
Subgoal: The stove is turned on. on the stove is:  
a substance called liquid [substance].

1. Focus on the thermometer
2. Focus on the substance you want to heat
3. Move the focused substance to the stove
4. Activate the stove

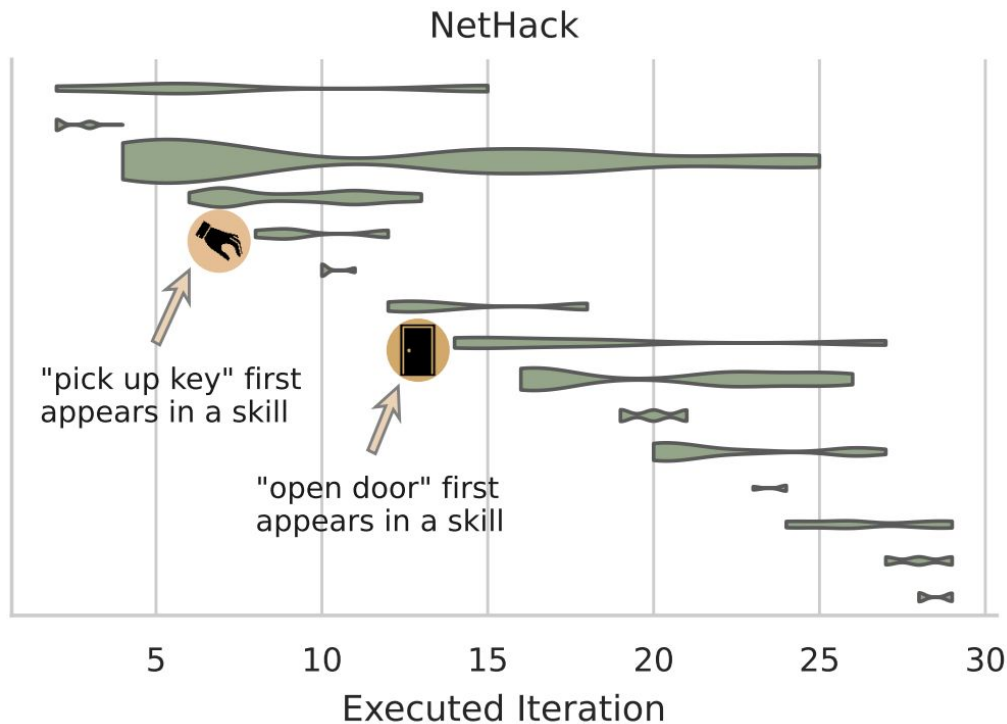
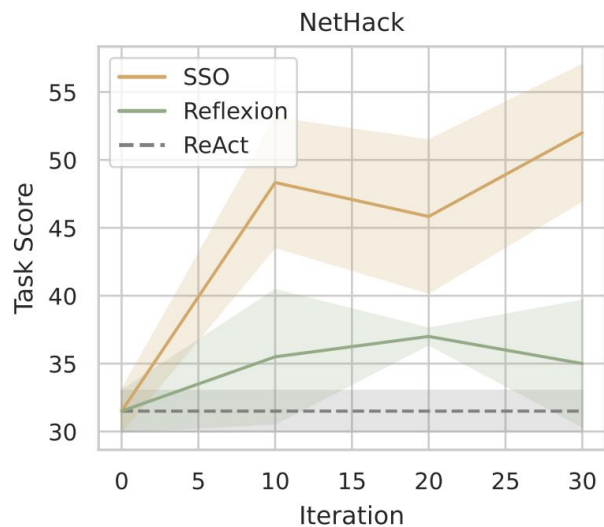
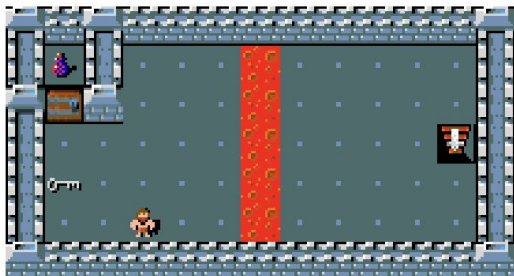




# Skill Lifecycle



# NetHack



# Conclusion

CLIN: <https://allenai.github.io/clin/>

CLIN: A CONTINUALLY LEARNING LANGUAGE AGENT  
FOR RAPID TASK ADAPTATION AND GENERALIZATION

Bodhisattwa Prasad Majumder<sup>1</sup>, Bhavana Dalvi Mishra<sup>1</sup>,  
Peter Jansen<sup>1,2</sup>, Oyvind Tafjord<sup>1</sup>, Niket Tandon<sup>1</sup>, Li Zhang<sup>3</sup>,  
Chris-Callison Burch<sup>3</sup>, Peter Clark<sup>1</sup>

<sup>1</sup>Allen Institute of AI

<sup>2</sup>University of Arizona

<sup>3</sup>University of Pennsylvania

SSO: <https://allenai.github.io/sso/>

---

**Skill Set Optimization:**  
**Reinforcing Language Model Behavior via Transferable Skills**

---

Kolby Nottingham<sup>1</sup> Bodhisattwa Prasad Majumder<sup>\*2</sup> Bhavana Dalvi<sup>\*2</sup>  
Sameer Singh<sup>1</sup> Peter Clark<sup>2</sup> Roy Fox<sup>1</sup>

- Dynamic memory in the form of causal abstractions or skills helps in generalization
- But current execution is greedy best, can we improve?
- Memory helps exploiting world knowledge, but how to incentivize exploration?

## Thank you!