



Condition: When ReAct fails to return an answer within a predefined number of steps.

• Action: The system falls back to using CoT-SC to attempt problem-solving.

• A) ReAct → CoT-SC (Switch from ReAct to CoT-SC):

- Step Thresholds: Set at 7 steps for HotpotQA and 5 steps for FEVER, as further steps did not improve ReAct's performance. Rationale: If ReAct struggles to find an answer using external search (indicating limitations in external knowledge utilization), the system leverages CoT-SC's strong internal reasoning capabilities to try and solve the problem. B) CoT-SC \rightarrow ReAct (Switch from CoT-SC to ReAct): Condition: When the majority answer among 'n' CoT-SC samples occurs less than 'n/2' times (i.e., the model's internal knowledge does not confidently support the
- Action: The system falls back to using ReAct to attempt problem-solving. • Rationale: If CoT-SC, despite multiple internal reasoning attempts, cannot converge on a confident answer (indicating uncertainty from internal knowledge alone), the system switches to ReAct to verify facts and augment information using external search, aiming for a more reliable solution.
- FineTuning => used **Bootstrapping Approach**

• Used **3,000 trajectories** (sequences of thoughts, actions, observations) that included correct answers. These trajectories were generated by ReAct (and other baselines) using larger, more capable models. • This automatically generated data was then used to **finetune smaller language models** (specifically, PaLM-8/62B)

ReAct Vs CoT

- Objective of Finetuning: To enable these smaller models to decode (generate) full ReAct-like trajectories (all thoughts, actions, and observations) in response to input questions or Benefit: Allows the advantages of ReAct to be transferred to more efficient, smaller models without requiring costly manual data annotation.

CoT's Hallucination issue

• ReAct outperforms CoT on FEVER (60.9 vs. 56.3).

• ReAct slightly lags CoT on HotpotQA (27.4 vs. 29.4).

• positive rate (14%) compared to ReAct (6%) in successful outputs. It constitutes CoT's major failure mode(56%).

	Type	Definition	ReAct	CoT	and making recovery difficult.This highlights a potential trade-off
Success	True positive False positive	Correct reasoning trace and facts Hallucinated reasoning trace or facts	94% 6%	86% 14%	between factuality (ReAct) and flexibility (CoT which motivated the proposed hybrid strategies. Hallucinations is a serious problem
Failure	Reasoning error	Wrong reasoning trace (including failing to recover from repetitive steps)	47%	16%	
	Search result error	Search return empty or does not contain useful information	23%	-	
	Hallucination	Hallucinated reasoning trace or facts	0%	56%	
	Label ambiguity	Right prediction but did not match the label precisely	29%	28%	
					for CoT

Decision Making Tasks

ALFWorld serves as a challenging benchmark for evaluating an AI's capacity for complex reasoning, planning, and acting in a dynamic, textonly environment by leveraging commonsense knowledge.

ALFWorld is a synthetic, text-based virtual environment game designed to test an AI agent's ability to perform complex, multi-step tasks within a simulated household. **Key Characteristics:**

ALFWorld

• Text-Based Interface: All interactions, environment descriptions, and actions are communicated purely through text commands. There are no visual elements. • **Simulated Household:** Features common rooms (kitchen, living room, bedroom) and objects found in a typical home. • Goal-Oriented: Agents are given high-level objectives (e.g., "examine paper under desklamp," "make coffee with the coffee machine"). • Complex Interactions: Achieving goals requires a sequence of precise text commands, such as navigating (go to table), manipulating objects (take paper), and using tools (use coffee machine). • Commonsense Reasoning: Success often relies on applying everyday knowledge (e.g., "books are usually on bookshelves," "refrigerators store food") to infer object locations or appropriate actions. • Long-Horizon Planning & Sparse Rewards: Tasks typically involve many steps, requiring long-term planning. Rewards are "sparse," meaning the agent only receives positive feedback upon successful completion of the entire goal, not for intermediate steps.

ALFWorld

How ALFWorld utlized:

Test ReAct's Reasoning and Acting Capabilities:

Selected ALFWorld to perform well in complex, Langague-based decision making environments

• Used as a complex, multi-step environment to evaluate ReAct's ability to plan and execute long action sequences, going beyond simple Q&A to actual decision-making. • Ideal for demonstrating how ReAct's iterative Thought-Action-Observation cycle functions effectively. Validate Commonsense Knowledge Utilization: • Leveraged ALFWorld's reliance on everyday commonsense (e.g., where items are typically found) to show how

ReAct, through its Thought component, effectively utilizes the LLM's pre-trained knowledge for efficient action planning. **Demonstrate the Importance of "Thoughts":** • ReAct was prompted with trajectories containing "sparse thoughts" (goal decomposition, subgoal tracking,

• Crucially, it was compared against an "Act" baseline (same trajectories but without thoughts). • This comparison aimed to empirically prove that explicitly incorporating reasoning ("thoughts") significantly enhances performance in complex environments like ALFWorld, quantifying the critical contribution of internal deliberation. **Assess Long-Term Planning and Exploration:**

• ALFWorld's multi-step, long-horizon tasks with sparse rewards were used to evaluate ReAct's capacity for strategic planning, subgoal tracking, and systematic exploration, facilitated by its Thought process.

WebShop

WebShop is a simulated online shopping environment designed to test how well AI can shop like **Key Points:** It uses real product data (1.18M items) and real human shopping instructions. • AI agents must **understand user requests** (e.g., "find a nightstand with drawers, • They then **interact with a simulated website** (search, click buttons, select options, • Unlike simpler environments, WebShop features **noisy**, **real-world text** (product descriptions, reviews) and requires **complex reasoning** to navigate. • Its goal is to evaluate AI's ability to understand complex language, make multi-step decisions, and perform practical tasks in a realistic web environment.

1 Benchmarking ReAct's Generalization: • It serves as a more challenging environment than ALFWorld (which is cleaner and more structured) to see if ReAct's Thought-Action mechanism can generalize to the complexities of real web interfaces and unstructured text. **Testing Robustness to Real-World Noise and Variety:** • WebShop's vast amount of real product data, diverse text types (titles, descriptions, user reviews), and complex user instructions

create a "noisy" environment. The paper uses this to test ReAct's robustness in extracting relevant information and making decisions amidst **Evaluating Complex Instruction Following and Information Extraction:** utilizes WebShop's detailed human instructions (e.g., "I am looking for a nightstand with drawers.

Findings : ALFWorld Findings: WebShop

outperformed both Act (no thoughts) and BUTLER (imitation learning), achieving a much higher success rate (ReAct 71% vs. Act 45% vs. BUTLER 37%). Crucial Role of 'Thoughts': Even the worst ReAct performance surpassed the best Act performance, clearly indicating that the presence or absence of 'thoughts' has an absolute, decisive impact on performance. Consistent Performance Gain: ReAct showed an average 62% relative performance gain over Act, proving that 'thoughts' are essential for goal decomposition and state tracking in complex environments. Conclusion: ALFWorld unequivocally demonstrated that ReAct's 'thought'-based reasoning leads to superior performance in complex simulated environments and is critical for task success.

ReAct's Dominant Performance: ReAct significantly

• **Performance Leap with 'Thoughts':** While Act already performed comparably to existing IL/IL+RL methods, ReAct, with improvement in success rate.

its "sparse reasoning," achieved a significant absolute 10% • Evolved Role of 'Thoughts': In WebShop, ReAct's 'thoughts' were crucial for identifying relevant information amidst noise and bridging the gap between user instructions and actual product/option details (e.g., reasoning how to search and filter for specific attributes). • Gap with Human Performance: Despite its improvements, ReAct still falls short of expert human **performance.** Humans perform more extensive product exploration and query reformulation, which remains challenging for current prompting-• Conclusion: WebShop showed that **ReAct significantly** enhances information identification and action decision-making in complex, noisy real-world web environments through 'thoughts,' but it still has limitations in flexible exploration compared to human capabilities.

• **Proof of ReAct's Practical Value:** WebShop confirmed

that ReAct is effective even in noisy, real-world web environments.

• Non-informative search accounts for

23% of ReAct's error cases, derailing its reasoning

— and making recovery difficult.