

# Agentic Context Engineering (ACE): Self-Improving LLMs via Evolving Contexts, Not Fine-Tuning

October 10, 2025

## Scaling AI with Haystack Enterprise: A Developer's Guide

Dr. Julian Risch, Engineering Team Lead, deepset  
Bilge Yücel, Developer Relations Engineer, deepset

When: October 15, 2025 | 10am ET, 3pm BST, 4pm CEST

[REGISTER NOW](#)



### Featured Speakers



Dr. Julian Risch  
Engineering Team Lead,  
Haystack



Bilge Yücel  
Developer Relations Engineer

**TL;DR:** A team of researchers from Stanford University, SambaNova Systems and UC Berkeley introduce [ACE framework](#) that improves LLM performance by **editing and growing the input context** instead of updating model weights. Context is treated as a living “playbook” maintained by three roles—**Generator, Reflector, Curator**—with small **delta items** merged incrementally to avoid brevity bias and context collapse. Reported gains: **+10.6%** on AppWorld agent tasks, **+8.6%** on finance reasoning, and **~86.9% average latency reduction** vs strong context-adaptation baselines. On the AppWorld leaderboard snapshot (Sept 20, 2025), **ReAct+ACE (59.4%) ≈ IBM CUGA (60.3%, GPT-4.1)** while using **DeepSeek-V3.1**.

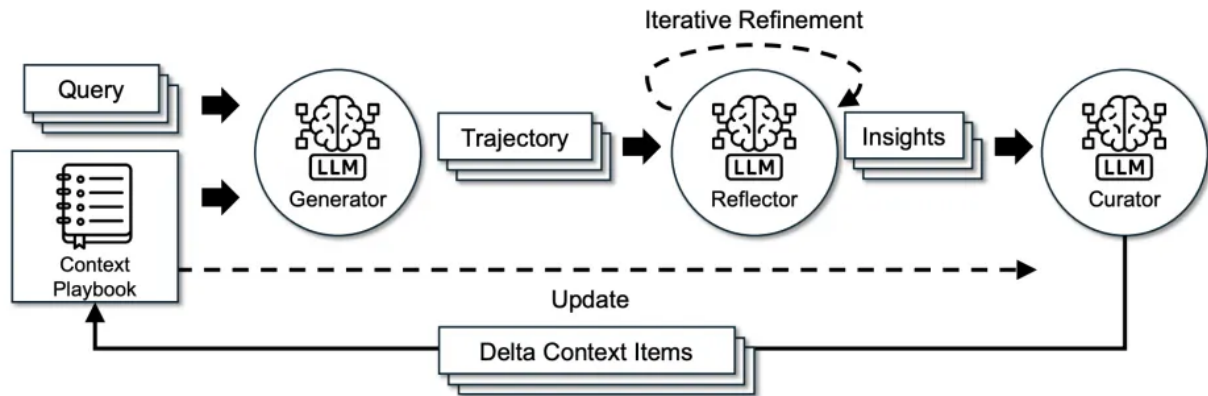
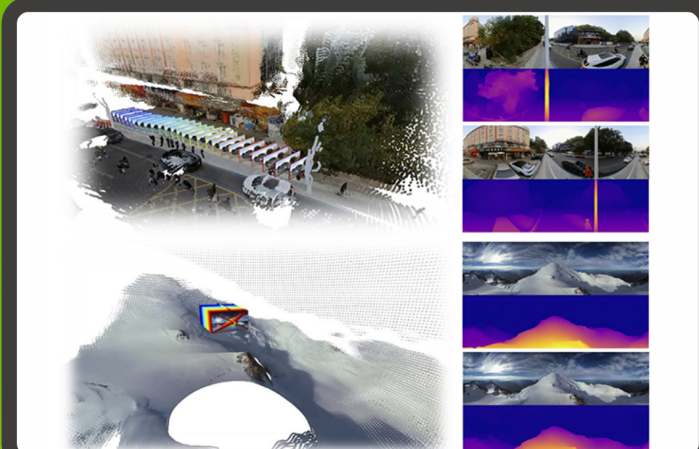


Figure 4: **The ACE Framework.** Inspired by Dynamic Cheatsheet, ACE adopts an agentic architecture with three specialized components: a Generator, a Reflector, and a Curator.


<https://arxiv.org/pdf/2510.04618>


## What ACE changes?

ACE positions “context engineering” as a first-class alternative to parameter updates. Instead of compressing instructions into short prompts, ACE **accumulates and organizes domain-specific tactics** over time, arguing that higher **context density** improves agentic tasks where tools, multi-turn state, and failure modes matter.



**Meet ViPE (Video Pose Engine)**  
**An Open-Sourced 3D Video**  
**Annotation Tool for Spatial AI**

 Fully Open Source



## Method: Generator → Reflector → Curator

- **Generator** executes tasks and produces trajectories (reasoning/tool calls), exposing helpful vs harmful moves.
- **Reflector** distills concrete lessons from those traces.

- **Curator** converts lessons into typed **delta items** (with helpful/harmful counters) and merges them deterministically, with de-duplication and pruning to keep the playbook targeted.

Two design choices—**incremental delta updates** and **grow-and-refine**—preserve useful history and prevent “context collapse” from monolithic rewrites. To isolate context effects, the research team fixes **the same base LLM (non-thinking DeepSeek-V3.1)** across all three roles.

## Benchmarks

**AppWorld (agents):** Built on the official ReAct baseline, **ReAct+ACE** outperforms strong baselines (ICL, GEPA, Dynamic Cheatsheet), with **+10.6% average** over selected baselines and **~+7.6%** over Dynamic Cheatsheet in online adaptation. On the **Sept 20, 2025 leaderboard**, **ReAct+ACE 59.4% vs IBM CUGA 60.3% (GPT-4.1)**; ACE surpasses **CUGA** on the harder **test-challenge** split, while using a smaller open-source base model.

**Finance (XBRL):** On **FiNER** token tagging and **XBRL Formula** numerical reasoning, ACE reports **+8.6% average** over baselines with ground-truth labels for offline adaptation; it also works with execution-only feedback, though quality of signals matters.

Method	GT Labels	Test-Normal		Test-Challenge		Average
		TGC↑	SGC↑	TGC↑	SGC↑	
DeepSeek-V3.1 as Base LLM						
ReAct		63.7	42.9	41.5	21.6	42.4
Offline Adaptation						
ReAct + ICL	✓	64.3 <sup>+0.6</sup>	46.4 <sup>+3.5</sup>	46.0 <sup>+4.5</sup>	27.3 <sup>+5.7</sup>	46.0 <sup>+3.6</sup>
ReAct + GEPA	✓	64.9 <sup>+1.2</sup>	44.6 <sup>+1.7</sup>	46.0 <sup>+4.5</sup>	30.2 <sup>+8.6</sup>	46.4 <sup>+4.0</sup>
ReAct + ACE	✓	76.2 <sup>+12.5</sup>	64.3 <sup>+21.4</sup>	57.3 <sup>+15.8</sup>	39.6 <sup>+18.0</sup>	59.4 <sup>+17.0</sup>
ReAct + ACE	✗	75.0 <sup>+11.3</sup>	64.3 <sup>+21.4</sup>	54.4 <sup>+12.9</sup>	35.2 <sup>+13.6</sup>	57.2 <sup>+14.8</sup>
Online Adaptation						
ReAct + DC (CU)	✗	65.5 <sup>+1.8</sup>	58.9 <sup>+16.0</sup>	52.3 <sup>+10.8</sup>	30.8 <sup>+9.2</sup>	51.9 <sup>+9.5</sup>
ReAct + ACE	✗	69.6 <sup>+5.9</sup>	53.6 <sup>+10.7</sup>	66.0 <sup>+24.5</sup>	48.9 <sup>+27.3</sup>	59.5 <sup>+17.1</sup>

Table 1: **Results on the AppWorld Agent Benchmark.** "GT labels" indicates whether ground-truth labels are available to the Reflector during adaptation. We evaluate the ACE framework against multiple baselines on top of the official ReAct implementation, both for offline and online context adaptation. ReAct + ACE outperforms selected baselines by an average of 10.6%, and could achieve good performance even without access to GT labels.

Method	GT Labels	FINER (Acc↑)	Formula (Acc↑)	Average
DeepSeek-V3.1 as Base LLM				
Base LLM		70.7	67.5	69.1
Offline Adaptation				
ICL	✓	72.3 <sup>+1.6</sup>	67.0 <sup>−0.5</sup>	69.6 <sup>+0.5</sup>
MIPROv2	✓	72.4 <sup>+1.7</sup>	69.5 <sup>+2.0</sup>	70.9 <sup>+1.8</sup>
GEPA	✓	73.5 <sup>+2.8</sup>	71.5 <sup>+4.0</sup>	72.5 <sup>+3.4</sup>
ACE	✓	78.3 <sup>+7.6</sup>	85.5 <sup>+18.0</sup>	81.9 <sup>+12.8</sup>
ACE	✗	71.1 <sup>+0.4</sup>	83.0 <sup>+15.5</sup>	77.1 <sup>+8.0</sup>
Online Adaptation				
DC (CU)	✓	74.2 <sup>+3.5</sup>	69.5 <sup>+2.0</sup>	71.8 <sup>+2.7</sup>
DC (CU)	✗	68.3 <sup>−2.4</sup>	62.5 <sup>−5.0</sup>	65.4 <sup>−3.7</sup>
ACE	✓	76.7 <sup>+6.0</sup>	76.5 <sup>+9.0</sup>	76.6 <sup>+7.5</sup>
ACE	✗	67.3 <sup>−3.4</sup>	78.5 <sup>+11.0</sup>	72.9 <sup>+3.8</sup>

Table 2: **Results on Financial Analysis Benchmark.** "GT labels" indicates whether ground-truth labels are available to the Reflector during adaptation. With GT labels, ACE outperforms selected baselines by an average of 8.6%, highlighting the advantage of structured and evolving contexts for domain-specific reasoning. However, we also observe that in the absence of reliable feedback signals (*e.g.*, ground-truth labels or execution outcomes), both ACE and other adaptive methods such as Dynamic Cheatsheet may degrade, suggesting that context adaptation depends critically on feedback quality.

## Cost and latency

ACE’s **non-LLM merges** plus localized updates reduce adaptation overhead substantially:

- **Offline (AppWorld): −82.3% latency** and **−75.1% rollouts** vs GEPA.
- **Online (FiNER): −91.5% latency** and **−83.6% token cost** vs Dynamic Cheatsheet.

Method	GT Labels	Test-Normal		Test-Challenge		Average
		TGC↑	SGC↑	TGC↑	SGC↑	
DeepSeek-V3.1 as Base LLM						
ReAct		63.7	42.9	41.5	21.6	42.4
Offline Adaptation						
ReAct + ACE w/o Reflector or multi-epoch	✓	70.8 <sup>+7.1</sup>	55.4 <sup>+12.5</sup>	55.9 <sup>+14.4</sup>	38.1 <sup>+17.5</sup>	55.1 <sup>+12.7</sup>
ReAct + ACE w/o multi-epoch	✓	72.0 <sup>+8.3</sup>	60.7 <sup>+17.8</sup>	54.9 <sup>+13.4</sup>	39.6 <sup>+18.0</sup>	56.8 <sup>+14.4</sup>
ReAct + ACE	✓	76.2 <sup>+12.5</sup>	64.3 <sup>+21.4</sup>	57.3 <sup>+15.8</sup>	39.6 <sup>+18.0</sup>	59.4 <sup>+17.0</sup>
Online Adaptation						
ReAct + ACE	✗	67.9 <sup>+4.2</sup>	51.8 <sup>+8.9</sup>	61.4 <sup>+19.9</sup>	43.2 <sup>+21.6</sup>	56.1 <sup>+13.7</sup>
ReAct + ACE + offline warmup	✗	69.6 <sup>+5.9</sup>	53.6 <sup>+10.7</sup>	66.0 <sup>+24.5</sup>	48.9 <sup>+27.3</sup>	59.5 <sup>+17.1</sup>

Table 3: **Ablation Studies on AppWorld.** We study how particular design choices of ACE (iterative refinement, multi-epoch adaptation, and offline warmup) could help high-quality context adaptation.

Method	Latency (s)↓	# Rollouts↓	Method	Latency (s)↓	Token Cost (\$)↓
ReAct + GEPA	53898	1434	DC (CU)	65104	17.7
ReAct + ACE	9517 <sup>(-82.3%)</sup>	357 <sup>(-75.1%)</sup>	ACE	5503 <sup>(-91.5%)</sup>	2.9 <sup>(-83.6%)</sup>

(a) **Offline** (AppWorld). (b) **Online** (FiNER).

Table 4: **Cost and Speed Analysis.** We measure the context adaptation latency, number of rollouts, and dollar costs of ACE against GEPA (offline) and DC (online).

<https://arxiv.org/pdf/2510.04618>

## Key Takeaways

- **ACE = context-first adaptation:** Improves LLMs by incrementally editing an evolving “playbook” (delta items) curated by Generator→Reflector→Curator, using the *same* base LLM (non-thinking DeepSeek-V3.1) to isolate context effects and avoid collapse from monolithic rewrites.
- **Measured gains:** ReAct+ACE reports **+10.6%** over strong baselines on AppWorld and achieves **59.4%** vs **IBM CUGA 60.3% (GPT-4.1)** on the Sept 20, 2025 leaderboard snapshot; finance benchmarks (FiNER + XBRL Formula) show **+8.6%** average over baselines.
- **Lower overhead than reflective-rewrite baselines:** ACE reduces adaptation latency by **~82–92%** and rollouts/token cost by **~75–84%**, contrasting with Dynamic Cheatsheet’s persistent memory and GEPA’s Pareto prompt evolution approaches.

## Conclusion

ACE positions context engineering as a first-class alternative to weight updates: maintain a persistent, curated playbook that accumulates task-specific tactics, yielding measurable gains on AppWorld and finance reasoning while cutting adaptation latency and token rollouts versus reflective-rewrite baselines. The approach is practical—deterministic merges, delta items, and long-context-aware serving—and its limits are clear: outcomes track feedback quality and task complexity. If adopted, agent stacks may “self-tune” primarily through evolving context rather than new checkpoints.

Check out the [PAPER here](#). Feel free to check out our [GitHub Page for Tutorials, Codes and Notebooks](#). Also, feel free to follow us on [Twitter](#) and don't forget to join our [100k+ ML SubReddit](#) and Subscribe to [our Newsletter](#). Wait! are you on telegram? [now you can join us on telegram as well.](#)



Asif Razzaq is the CEO of Marktechpost Media Inc.. As a visionary entrepreneur and engineer, Asif is committed to harnessing the potential of Artificial Intelligence for social good. His most recent endeavor is the launch of an Artificial Intelligence Media Platform, Marktechpost, which stands out for its in-depth coverage of machine learning and deep learning news that is both technically sound and easily understandable by a wide audience. The platform boasts of over 2 million monthly views, illustrating its popularity among audiences.



[Follow MARKTECHPOST: Add us as a preferred source on Google.](#)