

Multi-Agentive Recommender Systems: Foundations, Design Patterns, and E-Commerce Applications – An Industrial Tutorial



Yashar Deldjoo¹



Reza Yousefi Maragheh²



Chi Wang³



Jason Cho²



Derek Zhiyuan Cheng³

ACM Recommender Systems Conference, (RecSys'25)

Prague, Czech Republic, 22-26 Sep 2025

¹ Polytechnic University of Bari, Italy

² Walmart Global Tech, USA

³ Google DeepMind, USA

Outline

Introduction (Yashar- 30 min)

- **Classical recommendation** → **generative models**
 - Context to traditional (classical) models
 - Definition and types of generative models
 - Applications Area
 - **Generative** → **Agentic models**
 - Evolution
 - Definitions
 - Different Capabilities of RS
 - Level of autonomy
 - Language agents vs RL agents
 - Use Cases
-



RecSys 2025
Prague

Walmart Global Tech

Google DeepMind



Politecnico
di Bari

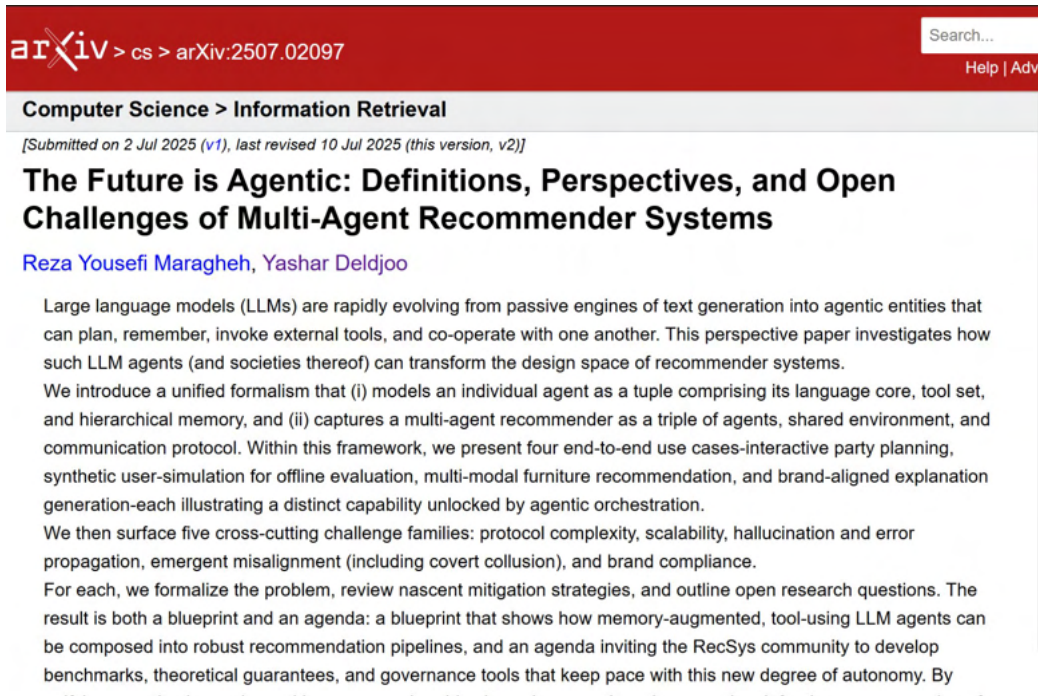
PART I- Introduction, Foundation, Alphabets

Agentic Recommender Systems (Agentic-RecSys)

GPTs and Assistants are precursors to agents. They will gradually be able to plan and perform more complex actions on your behalf.
These are our first step toward AI Agents



Check out our ACM TORS 2025 Perspective paper!



The screenshot shows the arXiv interface for a paper. At the top left, the arXiv logo is followed by the text '> cs > arXiv:2507.02097'. On the right, there is a search bar with the text 'Search...' and a 'Help | Adv' link. Below this, the category 'Computer Science > Information Retrieval' is displayed. The submission information reads: '[Submitted on 2 Jul 2025 (v1), last revised 10 Jul 2025 (this version, v2)]'. The title of the paper is 'The Future is Agentic: Definitions, Perspectives, and Open Challenges of Multi-Agent Recommender Systems'. The authors are listed as 'Reza Yousefi Maragheh, Yashar Deldjoo'. The abstract text follows, starting with 'Large language models (LLMs) are rapidly evolving from passive engines of text generation into agentic entities that can plan, remember, invoke external tools, and co-operate with one another. This perspective paper investigates how such LLM agents (and societies thereof) can transform the design space of recommender systems. We introduce a unified formalism that (i) models an individual agent as a tuple comprising its language core, tool set, and hierarchical memory, and (ii) captures a multi-agent recommender as a triple of agents, shared environment, and communication protocol. Within this framework, we present four end-to-end use cases-interactive party planning, synthetic user-simulation for offline evaluation, multi-modal furniture recommendation, and brand-aligned explanation generation-each illustrating a distinct capability unlocked by agentic orchestration. We then surface five cross-cutting challenge families: protocol complexity, scalability, hallucination and error propagation, emergent misalignment (including covert collusion), and brand compliance. For each, we formalize the problem, review nascent mitigation strategies, and outline open research questions. The result is both a blueprint and an agenda: a blueprint that shows how memory-augmented, tool-using LLM agents can be composed into robust recommendation pipelines, and an agenda inviting the RecSys community to develop benchmarks, theoretical guarantees, and governance tools that keep pace with this new degree of autonomy. By

Link: <https://arxiv.org/abs/2507.02097>

Check out our ACM TORS 2025 Perspective paper!

arXiv > cs > arXiv:2507.02097

Search...
Help | Adv

Computer Science > Information Retrieval

[Submitted on 2 Jul 2025 (v1), last revised 10 Jul 2025 (this version, v2)]

The Future is Agentic: Definitions, Perspectives, and Open Challenges of Multi-Agent Recommender Systems

Reza Yousefi Maragheh, Yashar Deldjoo

Large language models (LLMs) are rapidly evolving from passive engines of text generation into agentic entities that can plan, remember, invoke external tools, and co-operate with one another. This perspective paper investigates how such LLM agents (and societies thereof) can transform the design space of recommender systems.

We introduce a unified formalism that (i) models an individual agent as a tuple comprising its language core, tool set, and hierarchical memory, and (ii) captures a multi-agent recommender as a triple of agents, shared environment, and communication protocol. Within this framework, we present four end-to-end use cases-interactive party planning, synthetic user-simulation for offline evaluation, multi-modal furniture recommendation, and brand-aligned explanation generation-each illustrating a distinct capability unlocked by agentic orchestration.

We then surface five cross-cutting challenge families: protocol complexity, scalability, hallucination and error propagation, emergent misalignment (including covert collusion), and brand complexity.

For each, we formalize the problem, review nascent mitigation strategies, and outline open research questions. The result is both a blueprint and an agenda: a blueprint that shows how memory-augmented, tool-using LLM agents can be composed into robust recommendation pipelines, and an agenda inviting the RecSys community to develop benchmarks, theoretical guarantees, and governance tools that keep pace with this new degree of autonomy. By

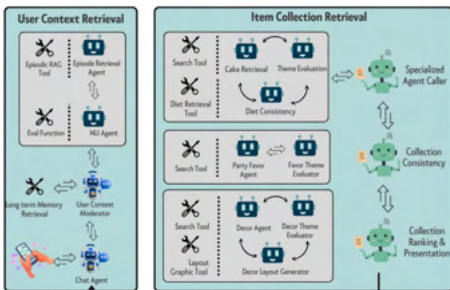
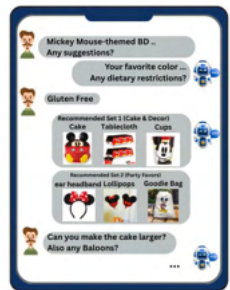
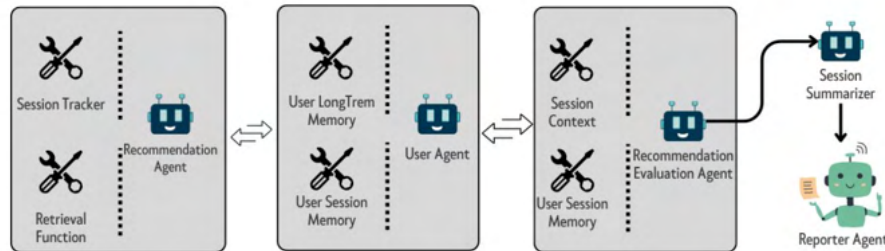
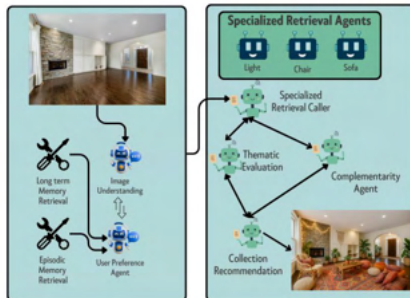


Table 1. Core Capabilities of Agentic Recommender Systems

Icon	Capability	Description
	Planning & Task Decomposition	Breaks down complex goals into sub-tasks, enabling multi-step reasoning and planning (e.g., candidate selection, bundle retrieval, re-ranking) [12, 44, 79].
	Tool Use & Action Execution	Invokes external tools/APIs or performs real-world actions (e.g., database queries, web searches, retrieving live inventory) [76].
	Memory & State Management	Maintains state across steps using memory modules (semantic, vector, knowledge graph), enabling context retention and personalized, multi-turn recommendations [4, 31, 34, 49, 53, 72, 74].
	Autonomy & Goal-Driven Behavior	Operates autonomously in a closed-loop, observing the environment, self-refining, and continuing until the objective is reached [42, 64].



Link: <https://arxiv.org/abs/2507.02097>

Also Take a look at our Gen-RecSys Work

Foundations and Trends® in
Information Retrieval
16:1-2

Recommendation with Generative Models

Y. Deldjoo, Z. He, J. McAuley, A.
Korikov, S. Sanner, A. Ramisa, R.
Vidal, M. Sathiamoorthy, A.
Kasrizadeh, S. Milano, and F. Ricci

now

the essence of knowledge



A Review of Modern Recommender Systems Using Generative Models (Gen-RecSys)



Yashar Deldjoo¹



Zhankui He²



Julian McAuley²



Anton Korikov³



Scott Sanner³



Arnau Ramisa⁴



René Vidal⁴



Mahesh
Sathiamoorthy⁵



Atoosa
Kasizadeh⁶



Silvia Milano⁷

¹ Polytechnic University of Bari, Italy

² University of California San Diego, USA

³ University of Toronto, Canada

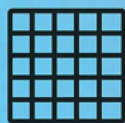
⁴ Amazon, USA*

⁵ BespokeLabs.AI, USA

⁶ CMU, USA

⁷ University of Exeter, UK & LMU
Munich, Germany

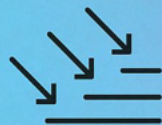
* This work does not relate to the author's position at Amazon



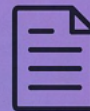
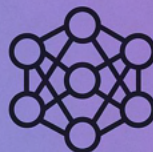
Matrices



Embeddings



Pattern Matching
(Classical)

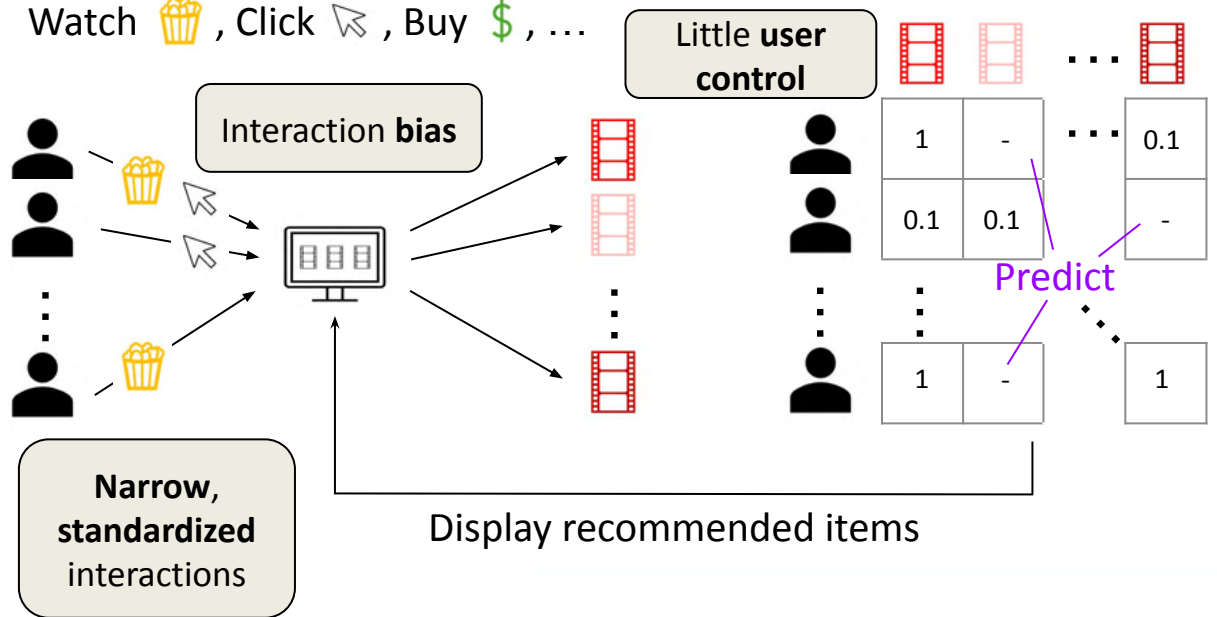


Reasoning
(Generative)

Traditional Recommendation (Recap)

Observe **fixed** user-item interactions:

Watch 🍿 , Click 🖱️ , Buy 💰 , ...

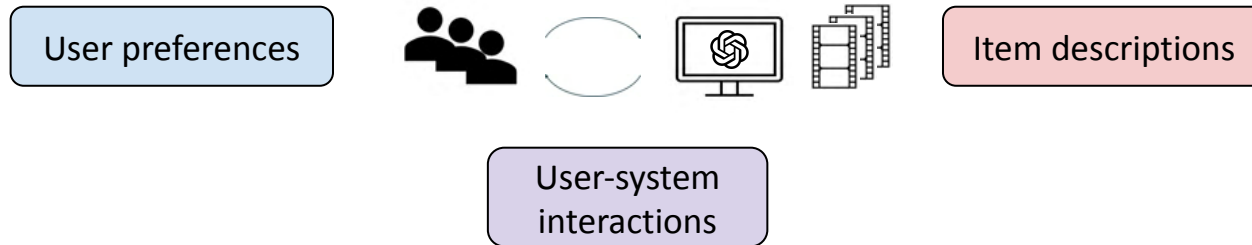


Task-specific optimization

Need lots of data

LLMs Unlock NL in Recommender Systems

LLMs **unlock NL** as a medium to represent:



Opportunities of LLM-Driven Recommendation

Rich NL data

+

LLM general reasoning abilities

=

Opportunities:

Nuanced **personalization** in **diverse contexts**

Interactive, real-time recommendation

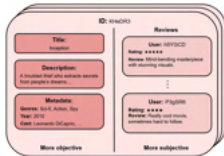
Faster system design and deployment

Long term preferences:
Love sci-fi movies with ethical dilemmas

Current preferences:
In a 1970s phase

Detailed Interactions:

- Watched inception (2010)
- Clicked The Bourne Identity (2002)
- Search for "1960s sci-fi thriller"
- ...

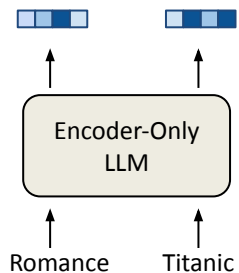


Pretrained LLMs:

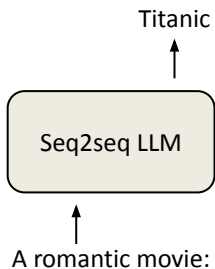
- **Internalized knowledge** about many items and human preferences
- Can reason about **diverse tasks** with **little or no new training data**

NL-based Models

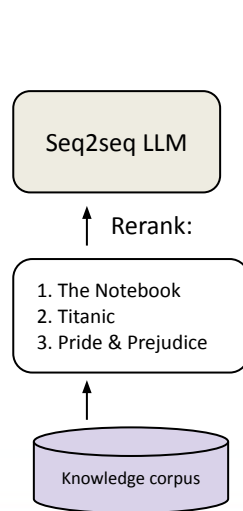
Encoder-Only LLMs



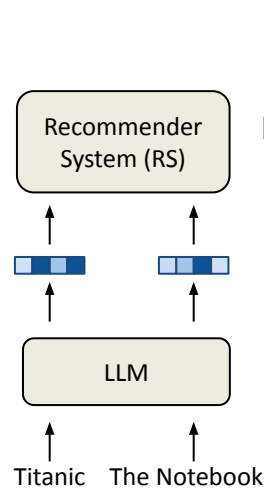
Sequence to Sequence (Seq2seq) LLMs



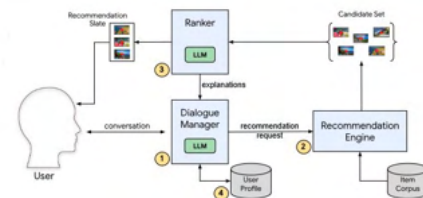
Retrieval-Augmented Generation (RAG)



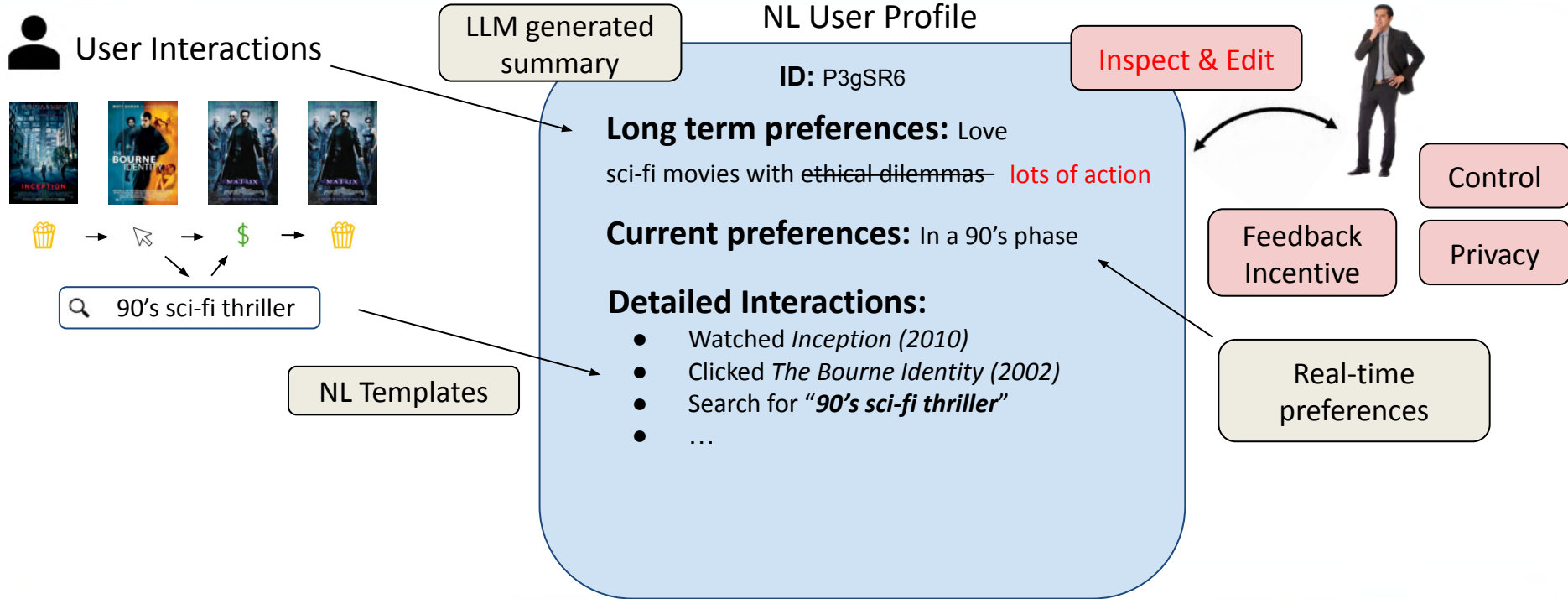
LLM Representation Generation



Conversational Recommendation



Editable NL User Profiles



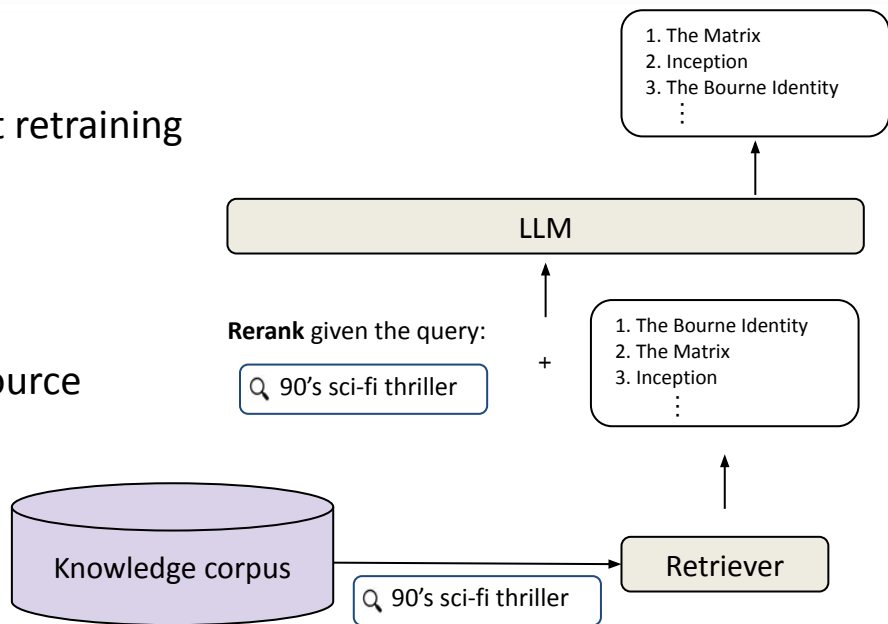
Retrieval Augmented Generation (RAG)

Limitations of LLM internal knowledge:

- Cannot add or revise knowledge without retraining
- More knowledge → more parameters
- No source attribution, hallucination

RAG:

- Retrieve information from an external source
- Prompt LLM with retrieved information
- Advantages:
 - Smaller LLM
 - Opportunity for source attribution

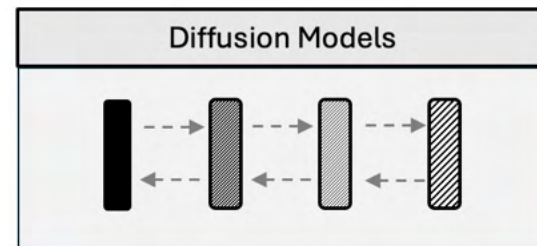
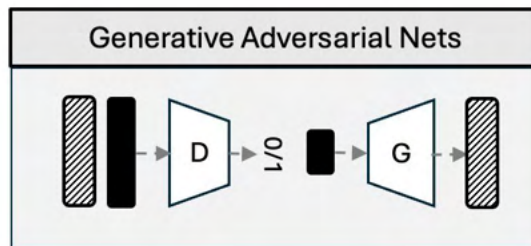
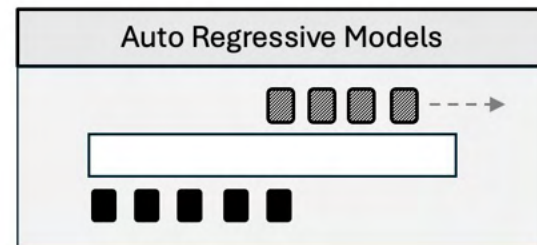
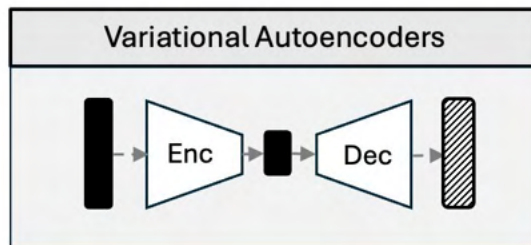


Example: RAG for recommendation reranking

Deep Generative Models (DGMs)

- DGMs

- VAEs
- AR Models
- GANs
- Diffusion Models
- Others



Other Models

What are Generative Recommender Systems? (Gen-RecSys)

Recommender systems enhanced by Generative AI



Structured List

generate lists, sequences or bundles tailored to user



Conversational & textual

Generate natural language responses, recommendations and explanations



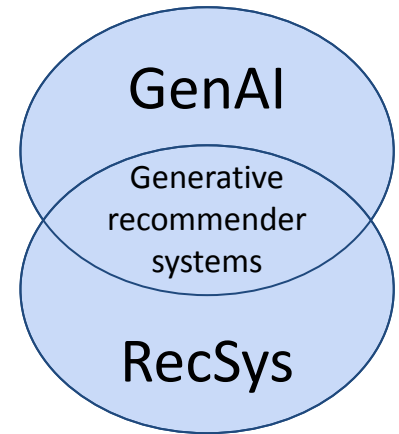
Visual & multimodal

Generate images or other media to enrich the recommendation experience

Recommender Systems with Generative Models (Gen-RecSys)

What are generative recommender systems?

- Systems that generate **structured outputs**, like bundles, lists, sets, sequences, etc.
- Systems that generate **text**, including conversational models, abstractive explanations, etc.
- Systems that generate **images**, including virtual try-on, fashion designs, image generation
- Existing generative models that involve **personalization**, e.g. personalized LMs or diffusion models



Essentially: any model that combines ideas from **Generative AI** and **recommender systems**

Deep Generative Models (DGMs)

- DGMs

Unconditional: $p_{\text{DGM}}(\mathbf{x}) \approx p(\mathbf{x})$

Conditional: $p_{\text{DGM}}(\mathbf{x}|C) \approx p(\mathbf{x}|C)$

- Attention!

1. **Different** \mathbf{X} in *RecSys w/ UI Data* high-dimensional
2. **Different** reasons to model $p(\mathbf{x}|\cdot)$ in *RecSys w/ UI Data*

Deep Generative Models (DGMs)

- **Different** reasons ?

Direct

- Generate an **interaction vector** as recommendations
- Generate **item lists** or **pages** as recommendations
-

Indirect

- Augment training data to train better **scoring function** $f_{\theta}(u, i)$
- Better sequential encoding for **next-item prediction** $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t}, \cdot)$
-

Some motivating context (from r/MovieSuggestions)

USER: Normally I watch horror and thrillers , but I've been really into comedies lately , specifically from the mid to late 90's (and some 00's) . I watched White Chicks last night and can't remember the last time I laughed so hard that my face hurt! Looking for recommendations on movies that will give me a good laugh , preferably 90's but anything will do !

SYSTEM: 90's era chuckles: *Friday, Groundhog Day, Kingpin, The Cable Guy, Rush Hour, and Office Space.*

Q: How should a recommender system generate these recommendations?

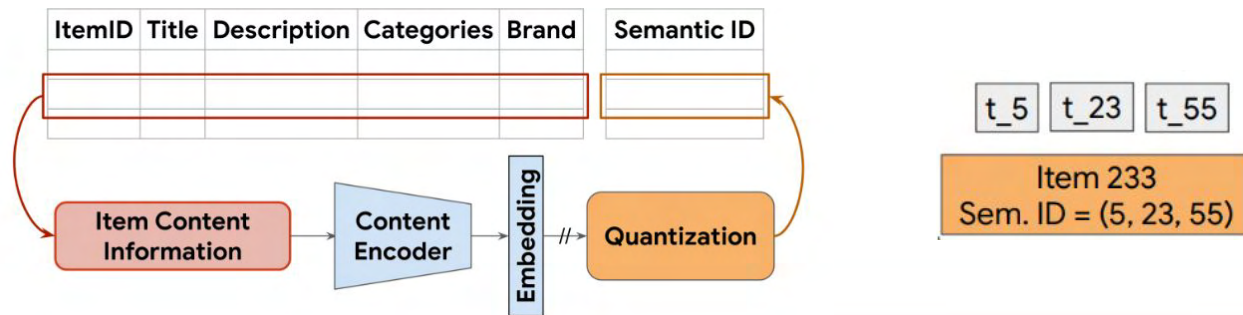
Some motivating context

Q: How should a recommender system generate these recommendations?

- **“Pure” language:** just generate some tokens!
 - Hallucination? Constrained beam search? Controllability? Fairness/bias?
- Add items as **new tokens** in the language model; maybe connect it with a pre-trained recommender
 - Scalability? Our Amazon review dataset has 48M item IDs! Cold-start?
- **Something else?** *Are there options for tokenization between language and items?*

Lately: Generative recommendation

- **Semantic IDs:** a few discrete tokens from a shared vocabulary, representing semantics for each item (e.g. quantized sentence embeddings)
- **Generative Recommendation:** models that autoregressively generate semantic IDs as recommendations



Lately: Generative recommendation

- “Recommendation” becomes a **seq**-to-**seq** generation problem:

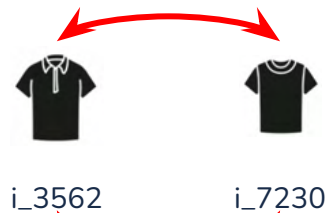
Input: user interacted items $\{c_{11}, c_{12}, c_{13}, c_{14}, c_{21}, c_{22}, \dots\}$

Output: next item $\{c_{t1}, c_{t2}, c_{t3}, c_{t4}\}$

Item IDs vs Semantic IDs

Item IDs:

- Contain **no prior knowledge** about items;
- **#Parameters** grows proportionally to #items



No correlation across IDs

Semantic IDs:

- **Prior knowledge**: semantically similar items have similar tokens;
- **#Parameters**: proportional to #tokens from a compact vocabulary;



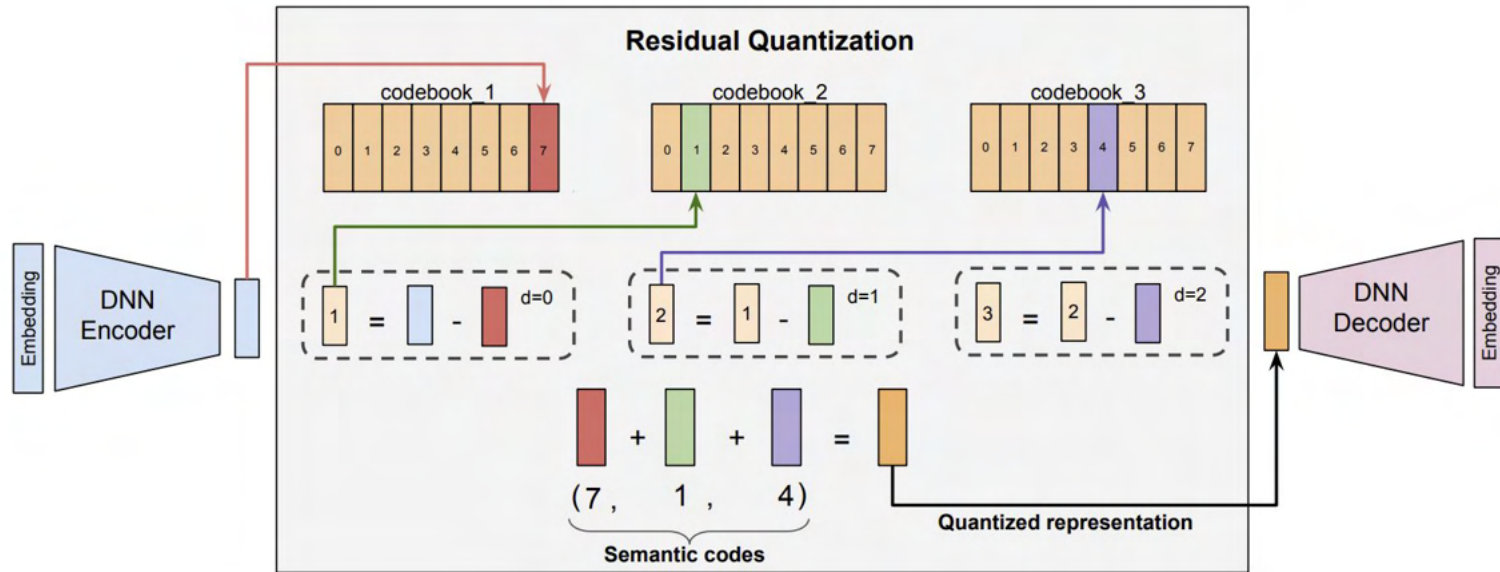
[t3, t65, t131, t243]



[t3, t65, t173, t243]

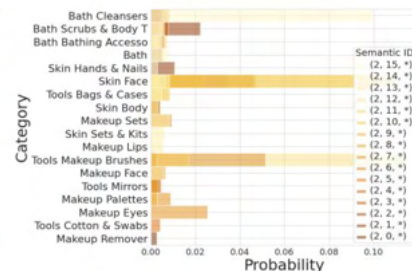
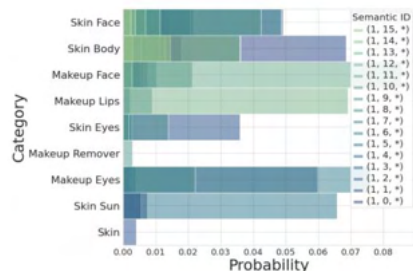
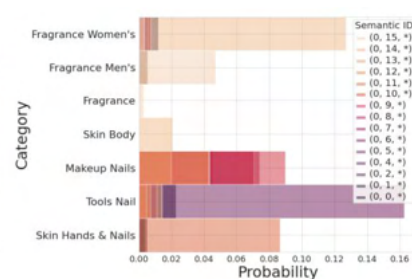
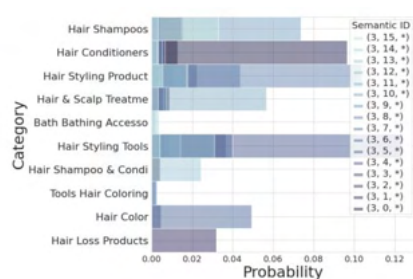
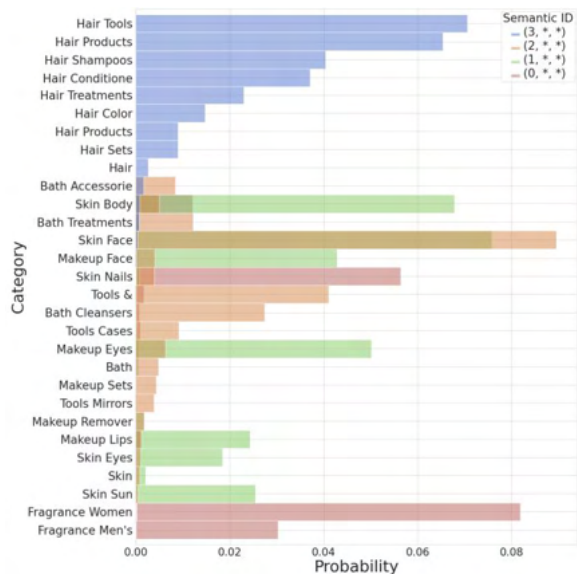
Lately: Generative Recommendation

- E.g. TIGER (from youtube):

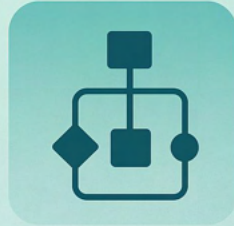


Lately: Generative Recommendation

- E.g. TIGER (from youtube):



Generative



Agentic



Pursues a
goal-oriented task

Comparison: Personalization

Classical

Pattern matching

Surface-level personalization

Generative (pLLMs)

Semantic Reasoning

Contains deeper understanding, infers deeper context

Agentic

Goal driven

aligns recommendations with user goal

Evolution of Recommendation Models

Comparison: Proactivity

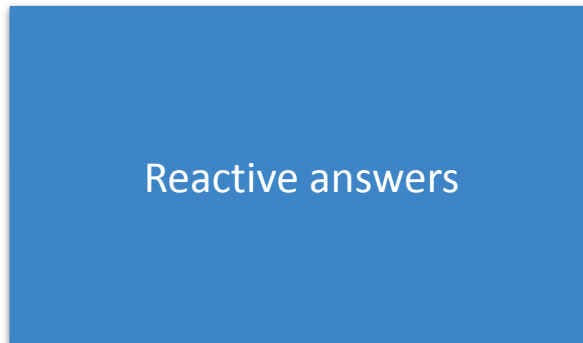


Classical



suggests only when asked

Generative (pLLMs)



answers specific prompts

Agentic



initiates suggestions and leads the conversation

—————>
Evolution of Recommendation Models

Comparison: Context Awareness

Classical

Past behavior

Generative (LLMs)

Prompted context

can also uses context explicitly given
in the query/prompt

Agentic

Situational, Multimodal
Context

→
Evolution of Recommendation Models

Comparison: Task Scope

Classical

Single step
recommendation

Generative (pLLMs)

Single-turn answer

Agentic

Multi-Step Planning

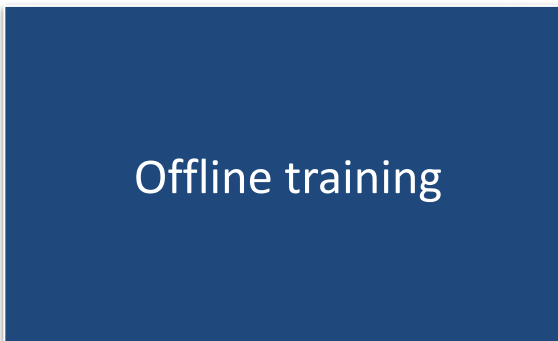
plans and executes complex goals

→
Evolution of Recommendation Models

Comparison: Adaptivity



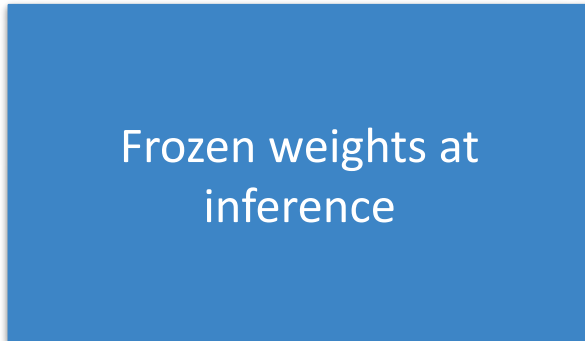
Classical



Offline training

Offline training and batch updates
(rarely refreshed)

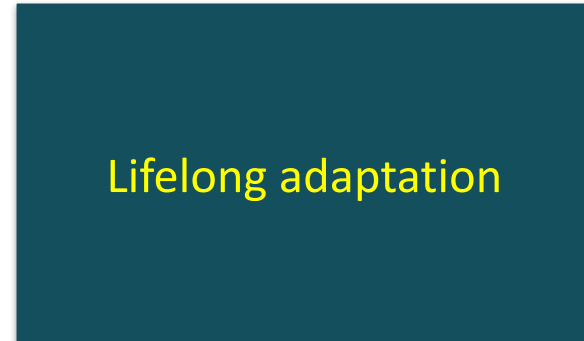
Generative (pLLMs)



Frozen weights at
inference

Pretrained model weights: static at
inference

Agentic



Lifelong adaptation

plans and executes complex goals

→
Evolution of Recommendation Models

Comparison: Memory

Classical

No explicit working
memory

No explicit short-term memory

Generative (pLLMs)

Context window

Session context limited by the LLM's
context window

Agentic

Structured long-term
memory

Structured long-term memory for
persistent user history

Evolution of Recommendation Models

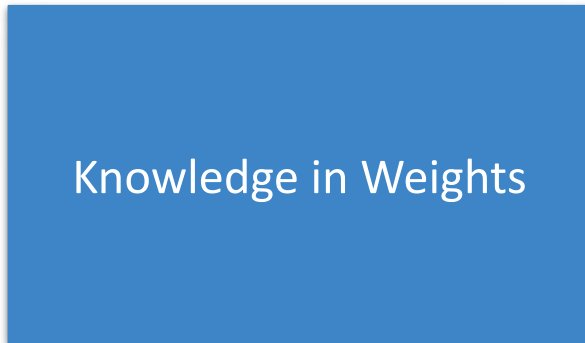
Comparison: Tools and Knowledge

Classical



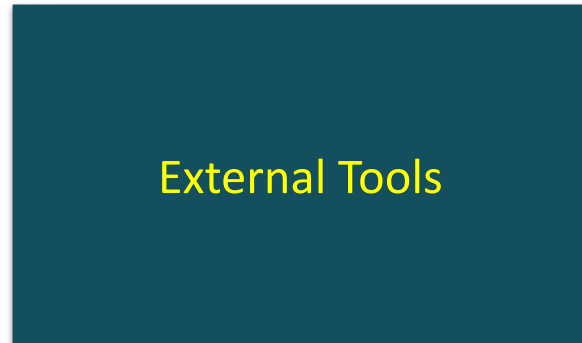
Fixed catalog or database only

Generative (pLLMs)



Relies on pre-training data for knowledge

Agentic



Employs external tools (search, APIs) during reasoning

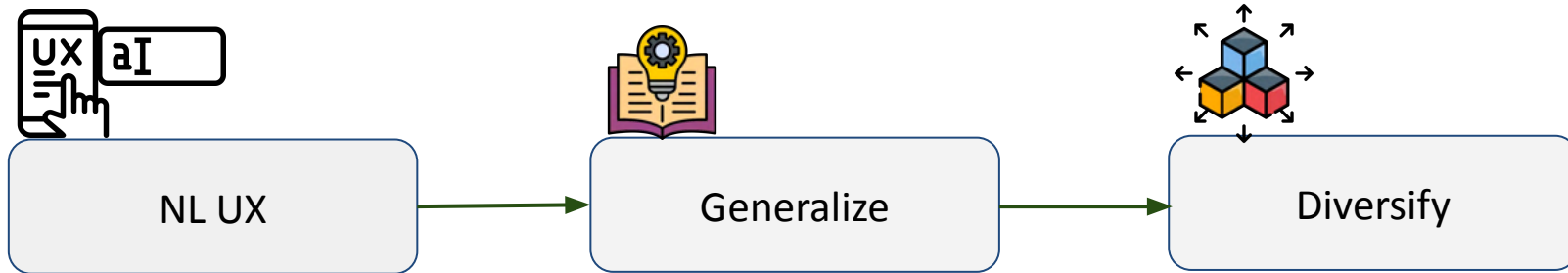
→
Evolution of Recommendation Models

Gen-RecSys Strength



Strength at a glance:

- Natural language interface accelerates preference elicitation.
- Pretrained knowledge improves zero-/few-shot generalization.
- Richer UX via explanations and diversified suggestions

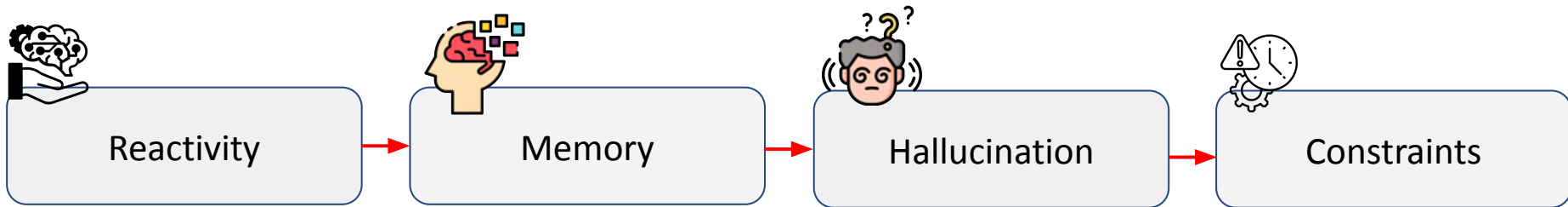


Gen-RecSys Limitations



Risks hotspots:

- Mostly reactive → it must be prompted with a **specific query** before generating a response
- Limited long-term memory.
- Constraint satisfaction can fail without explicit checkers.
- Hallucination risk without grounding + verification.



What is a Deep Research AI Agent

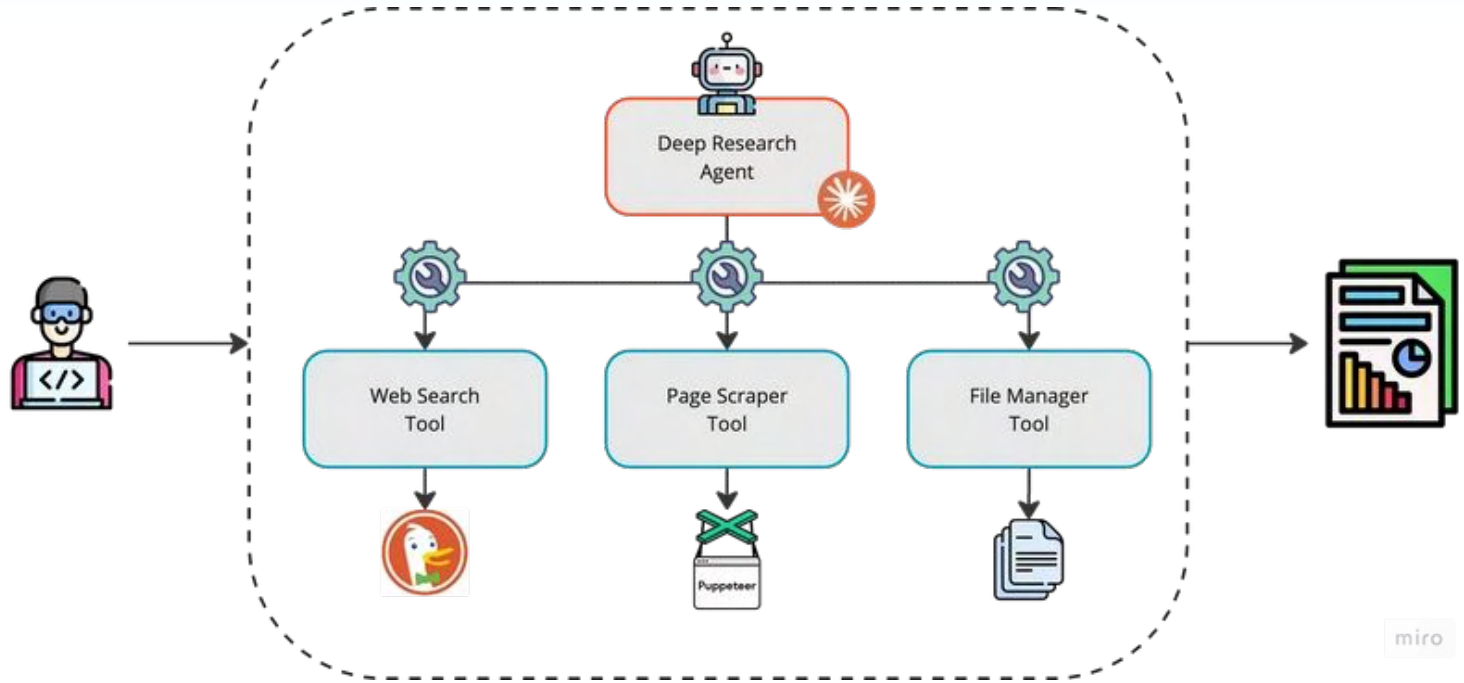
Definition: Autonomous system leveraging large language models (LLMs) to perform complex research tasks with minimal human intervention.

Purpose: Automate tasks such as literature reviews, competitive analysis, and technical documentation.

Key Features:

- Multi-step reasoning
 - Tool integration
 - Memory management
 - Autonomous decision-making
-

Simple Architecture



Why Agentic Recommender Systems?

Recommenders are evolving from **static predictors** to **adaptive agents** capable of reasoning, planning and interaction.

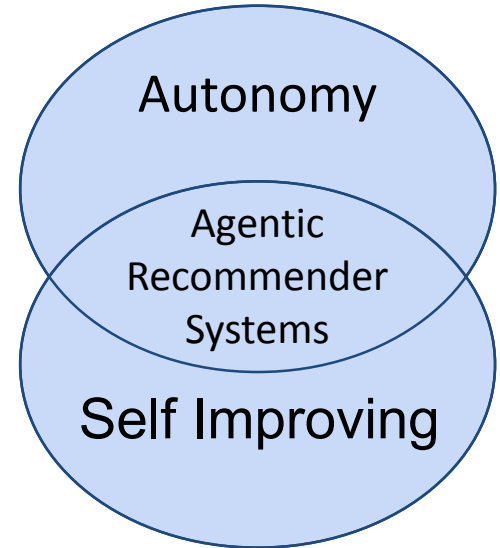
- Large language models (LLMs) enable natural language understanding and flexible tool usage.
- Agentic systems incorporate **memory**, **reflection** and **feedback** to personalize and align recommendations.



Agentic Recommender Systems (Agentic-RecSys)

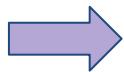
- **Autonomy.** Agents decide actions and adapt based on state and feedback
- **Self-Improvement.** Agents learn from experience, refine models and,

Agentic RS Combine autonomy and self-improvement to proactively gather context, plan, act and reflect.



Agentic Recommender Systems (Agentic-RecSys)

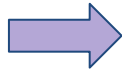
single-turn ranking



Autonomous goal-driven assistance

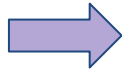
Agentic Recommender Systems (Agentic-RecSys)

single-turn ranking



Autonomous goal-driven assistance

“show items”



“solve tasks”

(bundle creation, constraint satisfaction, negotiation,
explanation).

Four Categories of Recommender Systems

Level 0

Data driven RS

Offline training &
static predictions



Level 1

Conversational

Interactive RS leveraging
user prompts



Level 2

Exploratory

Retrieval augmented
& task specific



Level 3

Fully Agentic

Autonomous, multimodal,
self evolving

Four Categories of Recommender Systems

Level 0

Data driven RS

Offline training & static predictions



Level 1

Conversational

Interactive RS leveraging user prompts



Level 2

Exploratory

Retrieval augmented & task specific



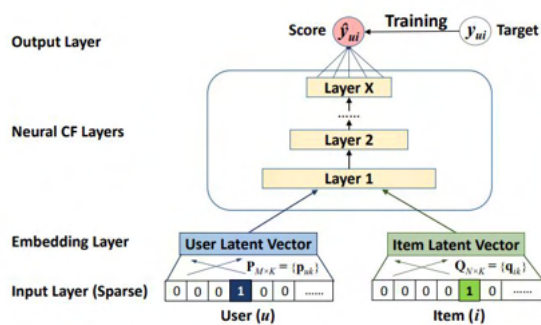
Level 3

Fully Agentic

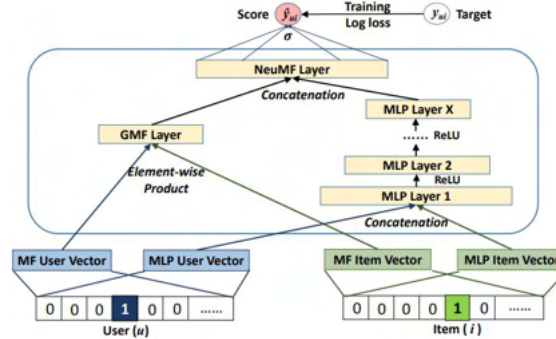
Autonomous, multimodal, self evolving

Level 0 - Data Driven RS

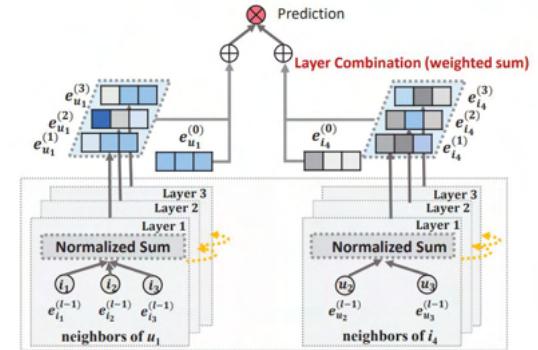
. Latent-space Modeling in Mainstream CF Models



Neural collaborative filtering framework (NCF)



Neural matrix factorization model (NeuMF)



Light Graph Convolution (LGC)

[1] He, X., Liao, L., Zhang, H., Nie, L., Hu, X., & Chua, T. S. (2017, April). Neural collaborative filtering. In Proceedings of the 26th international conference on world wide web (pp. 173-182).

[2] He, X., Deng, K., Wang, X., Li, Y., Zhang, Y., & Wang, M. (2020, July). Lightgcn: Simplifying and powering graph convolution network for recommendation. In Proceedings of

Four Categories of Recommender Systems

Characteristics

- Offline training and static predictions
- Learn from historical interactions; no dialog or reasoning

Representative models

- Collaborative Filtering (k NN, matrix factorisation: SVD, SVD++, Funk SVD);
- Neural models: Neural CF, Wide & Deep, BERT4Rec, SASRec, GNN
- Hybrid: content based + collaborative (e.g., Fab system)

Level 0

Data driven RS

Offline training & static predictions

Four Categories of Recommender Systems

Level 0
Data driven RS

Offline training &
static predictions



Level 1
Conversational

Interactive RS leveraging
user prompts



Level 2
Exploratory

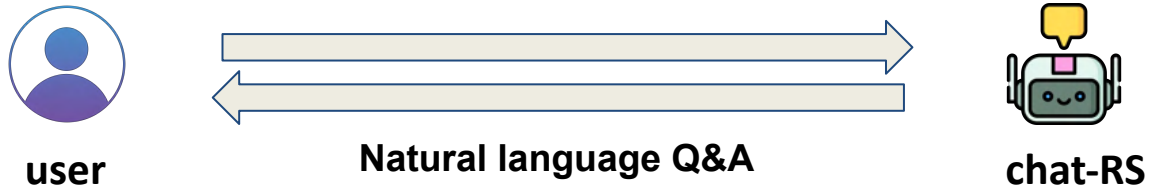
Retrieval augmented
& task specific



Level 3
Fully Agentic

Autonomous, multimodal,
self evolving

Level 1 - Conversational



- Interactive dialogue between user and system
- LLM interprets user intent and asks clarifying questions

Example: "Recommend a romantic movie like Titanic" → suggestions and follow up questions

Level 1 - Conversational

Characteristics

- Interactive dialogue with users via natural language
- Collect preferences on the fly and adjust recommendations

Representative models

- Early CRSs: ReDial (2018), KBRD, CR-MLS, etc.
- ChatRec and LLaMA Rec: LLM-based conversational recommenders
- MACRec & RecMind: multi agent and self inspiring architectures enabling conversation

Level 1

Conversational

Interactive RS leveraging user prompts

Four Categories of Recommender Systems

Level 0
Data driven RS

Offline training &
static predictions



Level 1
Conversational

Interactive RS leveraging
user prompts



Level 2
Exploratory

Retrieval augmented
& task specific

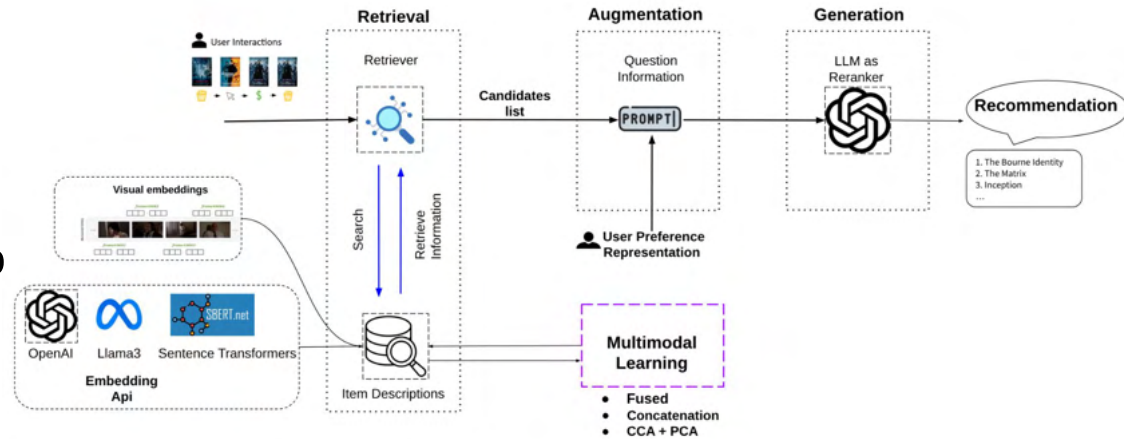


Level 3
Fully Agentic

Autonomous, multimodal,
self evolving

Level 2 - Exploratory (RAG)

- Augments LLM knowledge with external data
- Retrieves relevant documents/items on demand
- Enables source attribution and up to date answers



Example: "90s sci fi thriller" → retrieve candidates and let LLM re rank/explain

Four Categories of Recommender Systems

Characteristics

- Retrieval augmented reasoning: mix LLMs with search and database tools;
- Task specific modules for rating prediction, sequential and direct recommendations

Representative models

- RecMind: Self Inspiring planning with database and search tools
- MACRec: searcher retrieves external information for explanation;
- Self Ask + Search and RAG style agents: decompose tasks and query external sources

Level 2

Exploratory

Retrieval augmented &
task specific

Four Categories of Recommender Systems

Level 0

Data driven RS

Offline training &
static predictions



Level 1

Conversational

Interactive RS leveraging
user prompts



Level 2

Exploratory

Retrieval augmented
& task specific

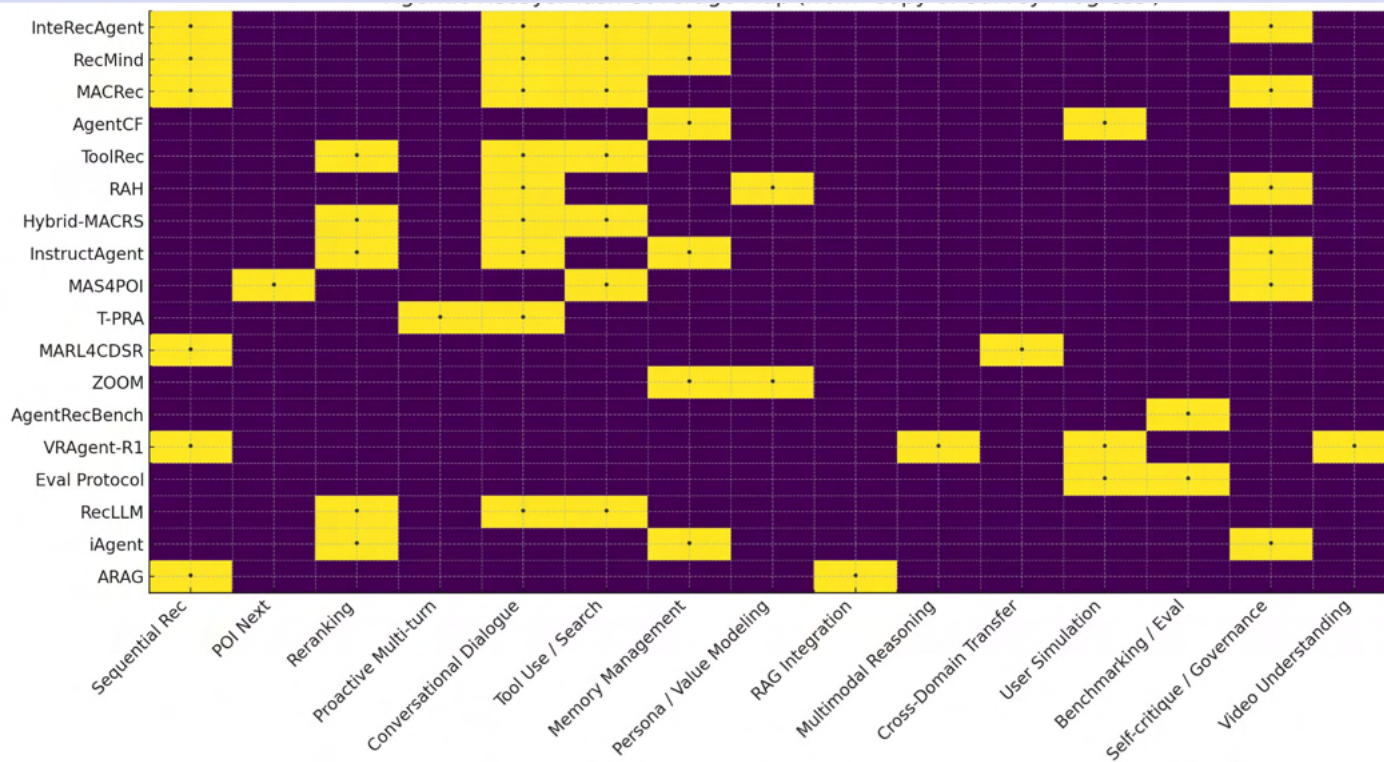


Level 3

Fully Agentic

Autonomous, multimodal,
self evolving

Notable Examples



Example 1: InteRecAgent

- A: recommender
- B: single agent
- C: plan-then-execute, reflection

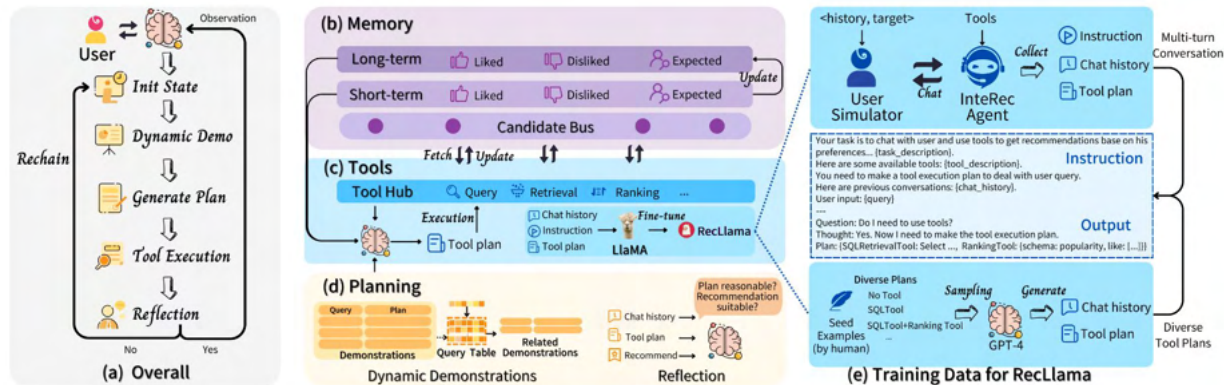


Figure 1: InteRecAgent Framework. (a) The overall pipeline of InteRecAgent; (b) The memory module, consisting of a candidate memory bus, a long-term and a short-term user profile; (c) Tool module, consisting of various tools, the plan-first execution strategy and the fine-tuning of RecLlama; (d) Planning module, involving the dynamic demonstrations and the reflection strategy; (e) Sources of fine-tuning data for RecLlama.

Example 2: MACRec

A: recommender
B: multi
C: interleaved plan and act

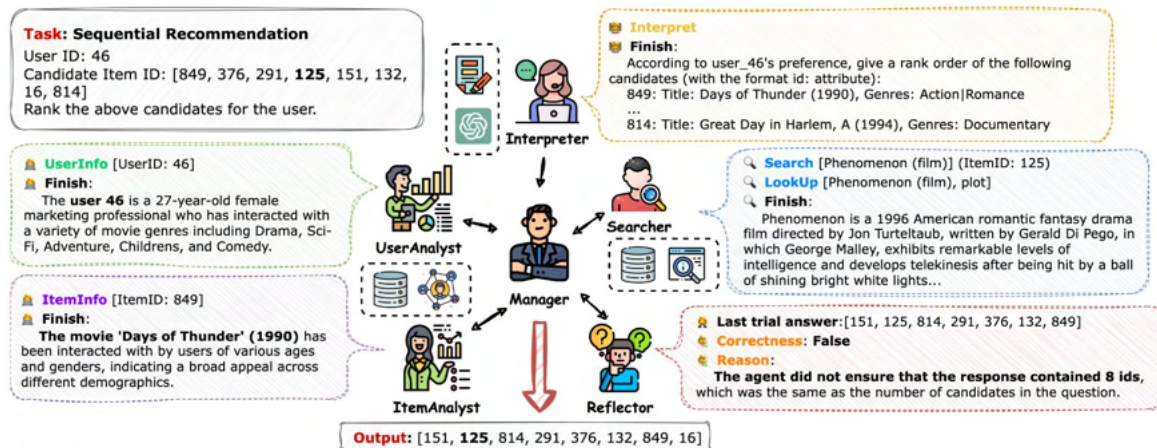


Figure 1: The Framework of MACRec. We take a sequential recommendation task as an example to show how these agents work collaboratively.

Example:3 AgentCF

A: recommender
B: multi agent
C: no plan

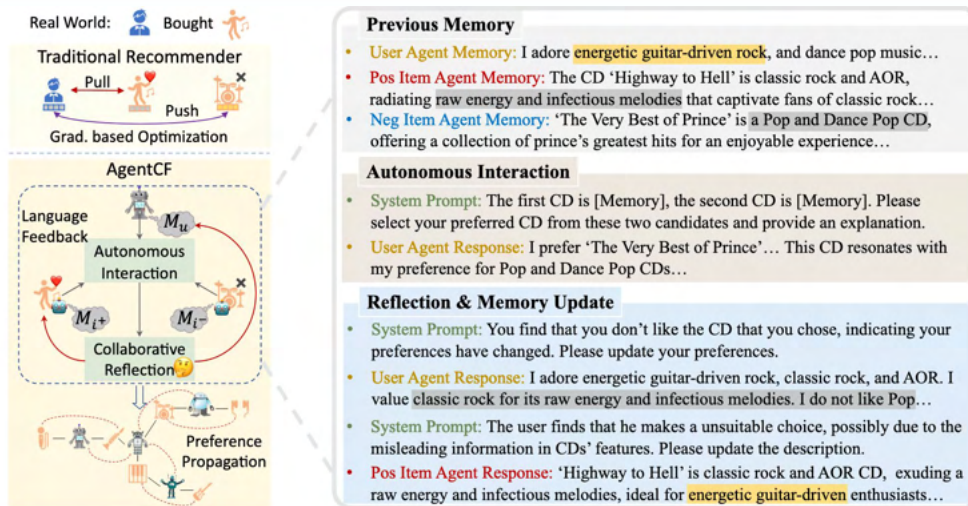


Figure 1: The overall framework of AgentCF and a case about the optimization process of agents: (1) The user and item agents are first prompted to autonomously interact. (2) These agents adjust the misconceptions in their memory, by reflecting on the disparities between their decisions and real-world interactions. In this process, the simulated preferences of user and item agents aggregate (as indicated by the highlighted content) and can propagate to other agents in subsequent interactions.

Multi-Agentive Recommender Systems: Foundations, Design Patterns, and E-Commerce Applications

An Industrial Tutorial

Reza Yousefi Maragheh¹ Yashar Deldjoo² Chi Wang³ Jason Cho¹ Derek Cheng³

¹Walmart Global Tech

²Polytechnic University of Bari

³Google DeepMind

RecSys 2025

Thought Process Evolution

Explain briefly Not Long!

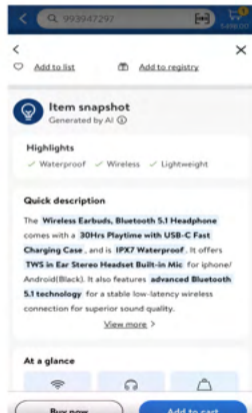
TLDR

Multi-stage modeling framework for **highlight extraction**, **Theme-Aware Keyword Extraction model** (LLM-TAKE) utilizing Large Language Models (LLMs), focused on scalability and integrated evaluation

Objective of LLM-TAKE framework

- Summarizing products through highlights that are context and theme aware, and are diverse
- **Helps customers** in their decision journey through **swiftly understanding product** characteristics
- Overcome LLM-specific **challenges**:
 - **Scalability**
 - **Evaluating the results**
 - **Hallucinations** in LLMs

Highlights



Thought Process Evolution

What's the deal with this item?

Hypothesis

Planned Solution

Why this is harder?

WE KNOW

Customers are looking for guidance and comparison between products they are interested in while browsing

IF WE

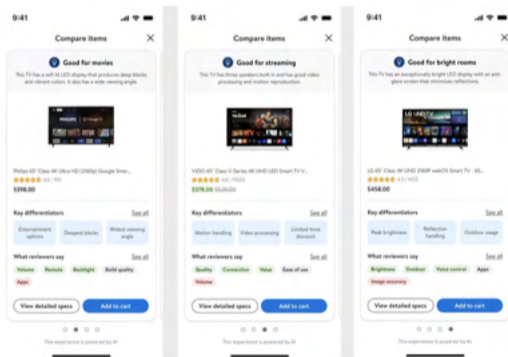
Show customers the most relevant use case of an item while they are comparing

THEN

We will make it easier for customers to make a decision

RESULTING IN

In an increase in the quality of ATC and higher conversion to purchases



But this one is harder. Why?
This falls upon largescale AI reasoning

Question:

If I give you item features, can You give me the usecases?
You should reason your way to generation and validation of usecases

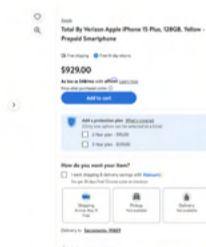
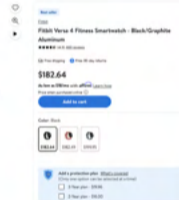
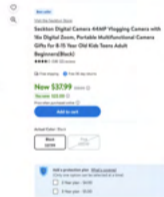
Thought Process Evolution

Myriad of Problems

The Relevance Problem

The Phrasing Problem

The Non-Informativeness Problem



Good for Professional Photography

Good for Fashion Accessory!

Good for Communication!

The question: How to deal with these edge cases in scale?

Thought Process Evolution

The Problem of Hallucination

- **Trustworthiness** of generated text by Language Models is an important element for large-scale implementations.
- Large Language Models have shown to **hallucinate** (Ji et al., 2023; Rawte et al., 2023)
- They generate responses that are **not consistent with instructions** (Duan et al., 2024).
- Thus, **evaluating the outcome of these models is crucial**.

Do LLMs Know about Hallucination? An Empirical Investigation of LLM's Hidden States

Hanyu Duan Yi Yang Kar Yan Tam

Department of Information Systems, Business Statistics and Operations Management
The Hong Kong University of Science and Technology
hduanac@connect.ust.hk {inyiyang, kytan}@ust.hk

Abstract

Large Language Models (LLMs) can make up answers that are not real, and this is known as hallucination. This research aims to see if, how, and to what extent LLMs are aware of hallucination. More specifically, we check whether and how an LLM reacts differently in its hidden states when it answers a question right versus when it hallucinates. To do this,

Mayer et al., 2020. Recently, this issue has triggered a notable discussion among thousands of AI researchers around the world, resulting in over 30,000 signatures on an open letter¹, calling for a six-month pause on "open AI experiment" (Pause Giant, 2023).

A growing body of NLP literature has examined hallucinations in LLMs, such as looking into the source of hallucinations from a training

Classic Approaches

- **n-gram based metrics**
- **model-based** such as
 - Rouge Lin (2004),
 - BLEU Papineni et al. (2002),
 - BERTScore Zhang et al. (2019)
 - BARTScore Yuan et al. (2021)
- But have proven to show **weak correlation with human evaluations** Wang et al. (2023).

BERTSCORE: EVALUATING TEXT GENERATION WITH BERT

Tianyi Zhang¹, Varsha Kishore¹, Felix Wu¹, Kilian Q. Weinberger^{1,2}, and Yoav Artzi¹

¹Department of Computer Science and ²Cornell Tech, Cornell University
{tvk352, zw243, kilian}@cornell.edu {yoav}@cs.cornell.edu

³ASAPP Inc.
tzhang@asapp.com

ABSTRACT

We propose BERTSCORE, an automatic evaluation metric for text generation. Analogously to common metrics, BERTSCORE computes a similarity score for each token in the candidate sentence with each token in the reference sentence. However, instead of exact matches, we compute token similarity using contextual embeddings. We evaluate using the outputs of 363 machine translation and image captioning systems. BERTSCORE correlates better with human judgments and

Neo-Classic Approaches

- Rapid evolution of LLMs, some researchers and practitioners have considered using LLMs for the evaluation purpose
 - Wang et al. (2023);
 - Zheng et al. (2024);
 - Maragheh et al. (2023);
 - Zhou et al. (2024).
- In this framework, due to high capability of **LLMs in mimicking human language capabilities**, they are used as a judge (Zheng et al., 2024)

Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena

Lianmin Zheng¹ Wei-Lin Chiang² Ying Sheng³ Siyuan Zhang⁴

Zhonghao Wu⁵ Yonghao Zhang⁶ Zi Lin⁷ Zhehan Li⁸ Dacheng Li¹

Eric P. Xing⁹ Hao Zhang¹⁰ Joseph E. Gonzalez¹¹ Ian Stoica¹

¹UC Berkeley ²UC San Diego ³Carnegie Mellon University ⁴Stanford ⁵MIT/CMU

Abstract

Evaluating large language model (LLM) based chat assistants is challenging due to their broad capabilities and the inadequacy of existing benchmarks in measuring human preferences. To address this, we explore using strong LLMs as judges to

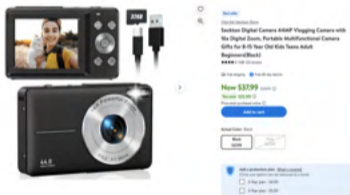
Thought Process Evolution

Images Credit: **Walmart.com**

Recommendation Explanation

For reasoning tasks,

- LLM as judge don't quite work really.
- Meaning it does not work.
- Yes, it does not work on complex tasks.



Is the usecase **“professional Photography”** relevant to this item given its features? **What is a professional camera? Really?**

Chat Completion

Metrics	Understandability		Groundedness	
	ρ	κ	ρ	κ
Vanilla	0.336	0.304	0.455	0.455
CoT	0.270	0.215	0.485	0.485
Self-Refine	0.167	0.144	0.412	0.412

Recommendation Explanation

Metrics	Relevance		Phraseness			
	ρ	κ	ρ	κ	ρ_{avg}	κ_{avg}
Vanilla	0.102	0.058	0.018	0.011	0.060	0.035
CoT	0.084	0.062	0.040	0.029	0.062	0.045
Self-Refine	0.075	0.046	0.026	0.012	0.050	0.029

Kappa, rho < 0.1 is not good at all!
Single Prompts Don't work!

Startups!

Startups

Image Credit: [dreamstime.com](https://www.dreamstime.com)



mindtrip.

for Creators for Business Get inspired

Get app Log in Get started

Travel differently

Mindtrip brings the world to you and empowers you to experience it **your way**.

Start chatting Play video

MindTrip

Image Credit: mindtrip.ai

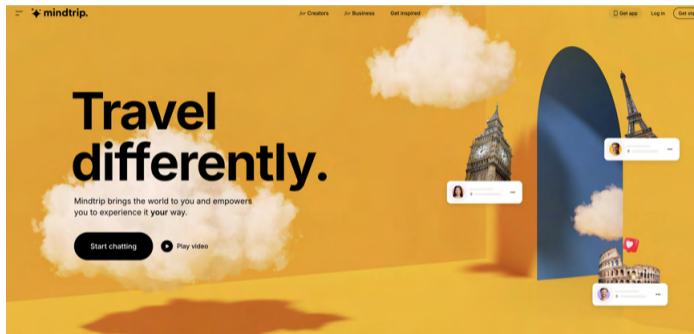



Figure: AI travel agent that recommends/optimizes itineraries, hotels, activities-Series A: Oct 2024

Tracxn Customers Offerings Company Pricing

Discover > Companies > Mindtrip > Funding & Investors

Navigate to

- Overview
- Funding & Investors
 - Funding Rounds
 - Investors
- Competitors Landscape 🔒
- Comparables 🔒
- Who's likely to invest? 🔒
- All Related Reports 🔒



Mindtrip - Funding & Investors

Last updated: August 28, 2025

[Claim Profile](#) [Suggest Edits](#) [↶](#) [⋮](#)

Mindtrip's Funding Rounds

Mindtrip has raised a total of \$19M over 2 funding rounds. One was Series A round of \$12M from Forerunner Ventures and Costanoa Venture and the other was Seed round of \$7M from Costanoa.

Here is the list of all funding rounds of Mindtrip:

Date of funding	Funding Amount	Round Name	Post money valuation	Revenue multiple	Investors
Sep 17, 2024	\$12M	Series A	\$120M	10x	Forerunner Ventur
Sep 08, 2023	\$7M	Seed	\$7M	1x	Costanoa

Note: Investors shown in bold are lead investors in that round

DayDream

Image Credit: daydream.ing

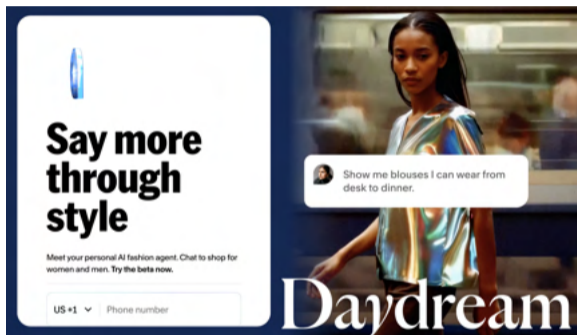


Figure: Chat-based fashion shopping agent for discovery and recommendations across thousands of brands. Seed: June 2024

The screenshot shows a news article on the PR Newswire website. The article title is "Daydream Secures \$50 Million In Seed Funding to Launch New AI-Powered Search and Discovery Shopping Platform". The article is dated June 20, 2024, at 08:00 ET. The text of the article states that e-commerce tech and retail veteran Julie Bornstein, along with co-founders Matt Fisher, Dan Cary, Lisa Green, and Richard Kim, have raised a \$50M seed round. The round is co-led by Forerunner Ventures and Index Ventures, with participation from GV (Google Ventures) and True Ventures. The funding is for the launch of Daydream, an AI-powered platform designed to change the way people shop online. The article also mentions that Daydream will be launching in Beta this fall, offering a highly personalized shopping experience powered by a built-in, superhuman search engine. This engine aims to provide a better way to find and discover women's and men's fashion. Additionally, Daydream will feature the largest high-quality branded fashion catalog found anywhere.

PR Newswire News Products Contact Search

News in Focus Business & Money Science & Tech Lifestyle & Health Policy & Public Interest People & Culture

Daydream Secures \$50 Million In Seed Funding to Launch New AI-Powered Search and Discovery Shopping Platform

Daydream

NEWS PROVIDED BY
Daydream
Jun 20, 2024, 08:00 ET

SHARE THIS ARTICLE

NEW YORK, June 20, 2024 /PRNewswire/ -- Today e-commerce tech and retail veteran Julie Bornstein along with co-founders Matt Fisher, Dan Cary, Lisa Green and Richard Kim announce they raised a \$50M seed round co-led by Forerunner Ventures and Index Ventures with participation from GV (Google Ventures) and True Ventures to launch Daydream, a powerful new, AI-powered platform that will change the way people shop online.

Launching in Beta this fall, Daydream will introduce its first shopping category, a highly personalized shopping experience powered by a built-in, superhuman search engine that offers a better way to find and discover women's and men's fashion. Daydream will also have the largest high-quality branded fashion catalog found anywhere.

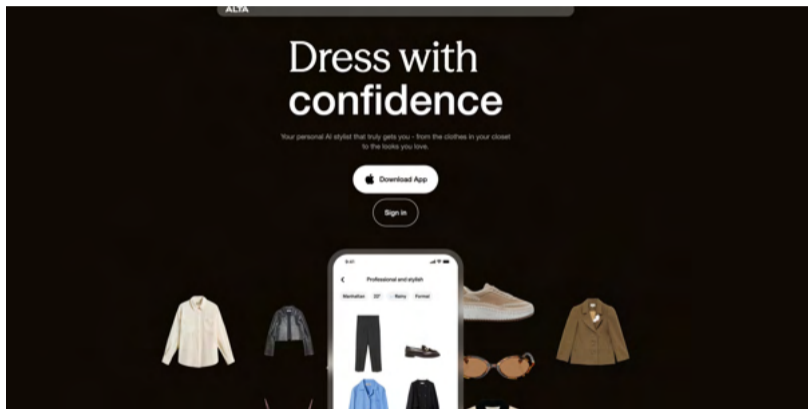


Figure: Agentic AI personal stylist that plans outfits and recommends products. Seed: June 2025

The screenshot shows a PR Newswire article page. At the top, there is a navigation bar with "PR Newswire" on the left, "News" (highlighted), "Products", and "Contact" on the right. A search icon is also present. Below this is a secondary navigation bar with categories: "News in Focus", "Business & Money", "Science & Tech", "Lifestyle & Health", "Policy & Public Interest", and "People & Culture". The main headline reads "Alta Raises \$11M Seed Round to Build the Future of Agentic Shopping". Below the headline, it says "NEWS PROVIDED BY Alta Daily" and "Jun 16, 2025, 09:30 ET". To the right of this is a "SHARE THIS ARTICLE" section with social media icons for Facebook, X, LinkedIn, Print, Email, and RSS. The article text begins with "NEW YORK, June 16, 2025 /PRNewswire/ -- In a \$185 billion U.S. apparel e-commerce industry saturated with choice and friction, Alta is crafting a brand new AI-native shopping experience by empowering shoppers with a personalized styling companion." The next paragraph states "Alta announced today it has raised \$11 million in seed funding to build the next generation of personal shopping and styling—powered by AI."

PR Newswire News Products Contact Search

News in Focus Business & Money Science & Tech Lifestyle & Health Policy & Public Interest People & Culture

Alta Raises \$11M Seed Round to Build the Future of Agentic Shopping

NEWS PROVIDED BY
Alta Daily
Jun 16, 2025, 09:30 ET

SHARE THIS ARTICLE

NEW YORK, June 16, 2025 /PRNewswire/ -- In a \$185 billion U.S. apparel e-commerce industry saturated with choice and friction, [Alta](#) is crafting a brand new AI-native shopping experience by empowering shoppers with a personalized styling companion.

Alta announced today it has raised **\$11 million in seed funding** to build the next generation of personal shopping and styling—powered by AI.



Figure: Aampe's Agentic AI learns what works for each customer. Then it instantly adapts your messaging and delivers at optimal times to drive better engagement, growth and unlock valuable insights. Series A-DEc 2024



Aampe deploys 100 million AI agents to power the next wave of personalization for consumer apps, as it raises \$18M

Aampe's agentic infrastructure enables marketing and product teams to deliver continuous personalization across channels and surfaces without having to build and maintain complex segments and campaigns across multiple tools.

10. Dezember 2024 10:00 ET | Quelle: [Aampe](#)

Folgen

San Francisco, Dec. 10, 2024 (GLOBE NEWSWIRE) – While companies building consumer apps and prosumer tools invest heavily in personalizing user experiences through product usage data, teams still manually craft the workflows that deliver those personalized moments. Today, [Aampe](#) reveals it has deployed over 100 million intelligent agents into consumer applications running across four continents. Businesses that have deployed Aampe agents include some of the leading food delivery and on-demand apps in South and Southeast Asia, top sports and fitness apps in Europe, as well as major fintech and entertainment apps

Constructor

Image Credit: constructor.com

The screenshot shows the Constructor website homepage. At the top left is the Constructor logo. To its right is a navigation menu with links for "Native Commerce Core™", "Solutions", "Customers", "Resource Hub", and "About". Further right are "Login" and a blue "Book a Demo" button. The main heading reads "Turn more ecommerce searches into sales with Constructor". Below this is a paragraph: "Constructor's AI-powered search and product discovery platform delivers unmatched KPI optimization and fast ROI for enterprise ecommerce brands — while giving customers a personalized, enjoyable shopping experience." The bottom section features a collection of seven circular images: a solid blue circle, a woman with glasses, a solid blue circle, a camera lens, a solid orange circle, a pair of black headphones, and a woman at a laptop with a plant.

Figure: AI Shopping Agent + search & recommendations for enterprise ecommerce. Series B: June 2024

Constructor

Credit: prnewswire.com-all funds: 85M

The screenshot shows a PR Newswire article. At the top, there is a navigation bar with 'PR Newswire' on the left and 'News', 'Products', and 'Contact' in the center. A search bar is on the right. Below the navigation bar, there is a secondary menu with categories: 'News in Focus', 'Business & Money', 'Science & Tech', 'Lifestyle & Health', 'Policy & Public Interest', and 'People & Culture'. The main headline of the article is 'Constructor Raises \$25M Series B Led by Sapphire Ventures, Tripling Valuation to \$550M'. To the right of the headline is the Constructor logo. Below the headline, there is a 'SHARE THIS ARTICLE' section with icons for Facebook, Twitter, LinkedIn, Email, and Print. The article text begins with 'SAN FRANCISCO, June 17, 2024 /PRNewswire/ - Constructor, the leading AI-powered product discovery and search platform for enterprise ecommerce companies, today announced that it has closed a \$25 million Series B round, bringing the company's valuation to \$550 million. This marks a nearly triple valuation since its 2021 Series A and brings total funds raised to date to more than \$85 million. Sapphire Ventures led the round with participation from existing investor SilverSmith Capital Partners. With this new capital, Constructor will continue to accelerate product development and innovation — applying the cutting-edge clickstream-based AI it pioneered to further improve ecommerce product discovery — and continue its rapid international expansion.'

Why so much \$\$\$?

What is an agent anyway?

Image Credit: Quite Like You! Andy Shauf, youtube.com

“I’ll call “Society of Mind” this scheme in which each mind is made of many smaller processes. These we’ll call agents. Each mental agent by itself can only do some simple thing that needs no mind or thought at all. Yet when we join these agents in societies-in certain very special ways-this leads to true intelligence”.

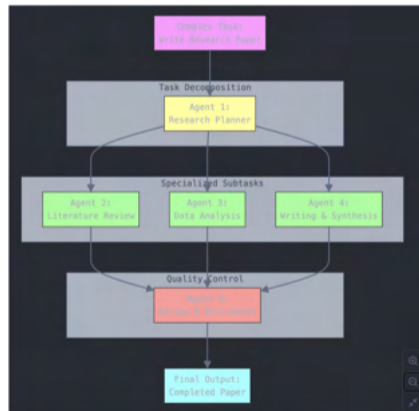
From: **The Society of Mind,**
Marvin Minsky



Why Multi-Agents?

Task Chunking

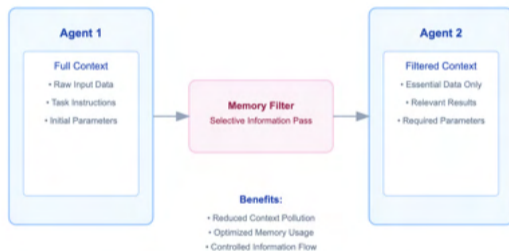
- **Reducing Reasoning Load:** Complex tasks often involve multiple interdependent steps that can overwhelm a single agent's decision-making capacity. Breaking these down into smaller, manageable chunks.
- **Failure Isolation:** Failures can be isolated to specific subtasks rather than requiring a complete restart of the entire process.
- **Specialized Agents:** Different agents can be optimized for different types of subtasks.



Why Multi-Agents?

Memory Moderation

- **Reducing Context Pollution:**
Controlled information flow between agents prevents context pollution and reduces the likelihood of conflicting or irrelevant information affecting decision-making processes.
- **Reducing Computational Overhead:**
Selective memory passing allows for efficient resource utilization by only sharing essential information.



Why Multi-Agents?

Tool Calling

- **Beyond Language Processing:** Agents can be designed to interface with external tools like databases, APIs, or specialized software, expanding their capabilities beyond pure language processing.
- **Less expensive compute:** Task-specific tools can handle computationally intensive operations (like complex mathematical calculations or data processing) more efficiently than language models alone.



Why Multi-Agents?

We practice what we present:

- Past three slides (and only past three of them) are researched, written, drawn, evaluated, corrected by an agentic pipeline!





Formalism

Formal Definition: LLM Agent

Definition

An **LLM Agent** is an intelligent system whose core decision-making and interaction capabilities are powered by one or more large language models.

$$A_{\text{LLM}} = (\mathcal{M}, \mathcal{I}, \mathcal{O}, \mathcal{F}, \Omega)$$

Formal Definition: LLM Agent

- \mathcal{M} : base model(s) for *understanding, reasoning, generation* (e.g., instruction-tuned transformer; optionally a mixture-of-experts).
- \mathcal{I} : *input space* observable by the agent (user text, images, traces, features, tool responses).
- \mathcal{O} : *output space* producible by the agent (natural language, structured JSON, function/tool calls, actions).
- $\mathcal{F} = \{f_k\}_{k=1}^K$: *external tools/APIs* the agent may invoke (retrieval, DB lookup, calculators, vision encoders, planners).
- Ω : *state & memory* enabling persistence across turns/sessions (short-term/working, episodic, semantic, procedural).

Multi-Agent System (MAS): Formal Definition

Definition

A **Multi-Agent System** is an ordered triple

$$\text{MAS} = (\mathcal{A}, \mathcal{E}, \Pi),$$

where:

- $\mathcal{A} = \{A_1, \dots, A_n\}$: finite set of agents; each A_i may be an LLM agent or another modular service.
- \mathcal{E} : shared *environment* exposing percepts/resources (APIs, UIs, simulators). Formally, a partially observable state space from which each agent receives observations and executes actions.
- $\Pi = (\mathbf{C}, \Gamma)$: *interaction protocol* constraining who may talk to whom and how.

Execution view: a run of MAS is a sequence of environment states and inter-agent messages that respect Π .

Interaction Protocol $\Pi = (\mathbf{C}, \Gamma)$

Communication matrix.

$\mathbf{C} \in \{0, 1\}^{n \times n}$, $\mathbf{C}_{ij} = 1 \iff$ messages of type $\gamma \in \Gamma$ are permitted from $A_i \rightarrow A_j$.

- **Preset (static) routing:** \mathbf{C} fixed at design time.
- **Autonomous (dynamic) routing:** channels toggled by guard functions $g_{ij} : \text{state} \rightarrow \{0, 1\}$, e.g. $\mathbf{C}_{ij}(t) = g_{ij}(\mathcal{E}_t, \Omega_t^i, \Omega_t^j)$.

Message schemata Γ .

- *Performatives* (semantics): INFORM, REQUEST, ACCEPT, REJECT, ...
- *Serialization rules* (syntax): JSON schema, field names/types for machine readability.
- *Timing constraints:* ordering, deadlines, rate limits, retry/timeout policies.

Each $\gamma \in \Gamma$ defines both **syntax** and **semantics** so a receiver can parse and act on messages deterministically.

Minimal Example (Rec \leftrightarrow Eval)

Agents: $\mathcal{A} = \{A_{\text{rec}}, A_{\text{eval}}\}$.

Environment:

$\mathcal{E} = (\text{UserProfileDB}, \text{ProductCatalogue}, \text{BusinessRulesKB})$.

Protocol:

$\Pi = (\mathbf{C}, \Gamma)$, $\mathbf{C} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, $\Gamma = \{\text{candidate_list}, \text{compliance_report}\}$.

Message flow.

- $A_{\text{rec}} \rightarrow A_{\text{eval}}$: `candidate_list` (ranked JSON array $\{(s_1, \text{score}_1), \dots, (s_L, \text{score}_L)\}$).
- $A_{\text{eval}} \rightarrow A_{\text{rec}}$: `compliance_report` (violations w.r.t. `brand/policy` \in `BusinessRulesKB`).

Off-diagonal ones in \mathbf{C} permit bidirectional messaging; zeros on the diagonal disallow self-messaging.

Autonomous moderation of links.

$$C_{\text{eval,rec}}(t) = \begin{cases} 1, & \text{NOT Halt}(t) \\ 0, & \text{Halt}(t) \end{cases} \quad \text{with} \quad \text{Halt}(t) \iff (\text{violations} = \emptyset) \vee (t - t_0 > \text{timeout}).$$

Protocol-level guards in Γ :

- *Stop condition*: end dialogue when all items pass compliance.
- *Retry/backoff*: bounded retries on transient failures.
- *Escalation*: on persistent violation, route to A_{policy} or human-in-the-loop.

Takeaway. Π can be preset yet *autonomously reconfigured* according to system state which balances accuracy with efficiency.

Let's Recap



Industrial Agentic RecSys-Usecases

Industrial Agentic RecSys-Usecases

Planning → Simulation → Multi-modal → Explanation

Scenarios we will walk through

- **Interactive conversational Recommendation**
 - sub-agents coordinate (planner, retrieval, constraints) → curated “Mickey-Mouse party” bundles.
- **User-simulation-Reco Eval:**
 - synthetic users + logging/summarization → stress-test policies pre-deployment.
- **Contextual & multi-modal Recommendation:**
 - vision+text (room photo → Boho set) → complementarity & aesthetic coherence.
- **Recommendation Explanation:**
 - brand-consistent narratives that expose latent ranking logic (“inspired by your Disney purchases”).

Interactive Recommendation

Interactive Recommendation

Agentic conversational planning for context-aware bundles

- Multi-turn conversation \Rightarrow refine constraints & intent
- Sub-agents for retrieval, validation, ranking, layout
- Memory-aware (STM/SEM/EPI/PROC) & tool-using



Interactive Recommendation

Task at a Glance

Definition (informal)

Interactive Recommendation engages in a **multi-turn** dialogue to identify, refine, and present suitable items, adaptively incorporating user feedback, contextual constraints, and memory/tool signals.

Why interactive?

- Real-world user intents are **open-ended** and **multi-step**
- Preferences emerge through **clarification** and **trade-offs**

Running example

“Mickey-Mouse themed birthday”: decorations, gluten-free cake, favors; one conversation, many sub-goals.

Interactive Recommendation

Formal Definition (I)

Let \mathcal{D} be conversation turns and \mathcal{S} the item universe (SKUs). d_i and a_i are user and agent's utterances. A session at step t has transcript

$$\mathcal{C}_{1:t} = (d_1, a_1, \dots, d_{t-1}, a_{t-1}, d_t), \quad d_i, a_i \in \mathcal{D}.$$

Define the interactive recommendation mapping

$$\Phi : (\mathcal{C}_{1:t}, \mathcal{E}) \rightarrow \langle s_{(1)}, \dots, s_{(L)} \rangle, \quad s_{(j)} \in \mathcal{S},$$

where, \mathcal{E} is the shared space among agents. Φ outputs a ranked list of length L .

Objective with implicit constraints

$$\langle s_{(1)}, \dots, s_{(L)} \rangle = \arg \max_{\langle s_1, \dots, s_L \rangle} \sum_{j=1}^L \text{Rel}(s_j \mid \mathcal{C}_{1:t})$$

subject to constraints (theme, dietary, budget) encoded in $\mathcal{C}_{1:t}$ and available tools/environment \mathcal{E} .

- $\text{Rel}(\cdot)$ may incorporate user history (SEM/EPI) and session cues (STM)
- Feasibility checked via tools (e.g., inventory, price, layout fit)

Interactive Recommendation

High-Level Goal

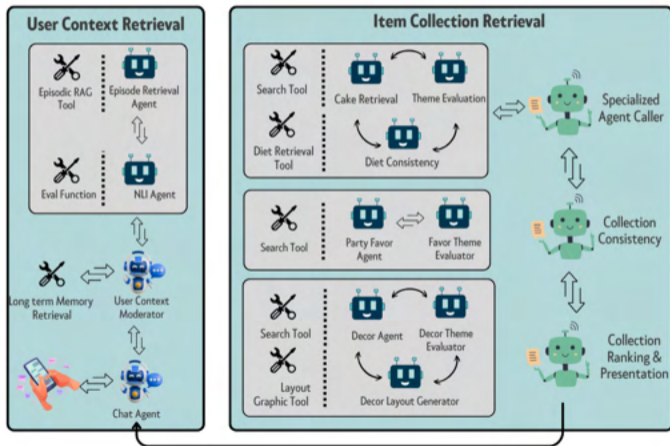
Objective

Dynamically adapt the Retrieval/Ranking function Φ across turns to minimise user effort while satisfying latent sub-goals (e.g., *gluten-free chocolate cake*, *palette-consistent decor*).

- Spawn specialised sub-agents on demand (category retrievers, validators)
- Maintain **coherence** across turns via short/long-term memory
- Present **holistic** bundles (ranked collection + layout preview)

Interactive Recommendation

The Visual



Interactive Recommendation

System Architecture: Agents & Environment

Agent set \mathcal{A}

A_{chat} (chat), A_{epi} (episodic retrieval \mathcal{Q}), A_{nli} (episode vetting), A_{SAC} (specialised-agent caller),

A_{cake} , A_{decor} , A_{favor} (category retrieval), $A_{\text{col_check}}$ (collection consistency), A_{rank} (ranking/presentation).

Environment & Protocol

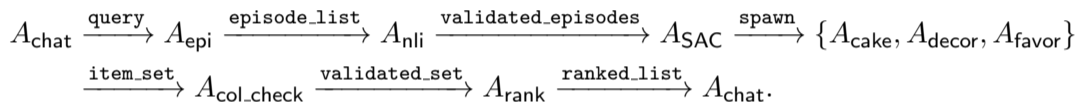
$\mathcal{E} = (\text{ProductCatalogue}, \text{UserProfileDB}, \text{VectorDB});$

$\Pi = (\mathbf{C}, \Gamma)$ with $\mathbf{C}_{ij}=1$ iff A_j is spawned by A_i or $A_i=A_{\text{chat}};$

$\Gamma = \{\text{query}, \text{episode_list}, \text{tool_call}, \text{item_set}, \text{ranked_list}\}.$

Interactive Recommendation

Message Flow & Orchestration



- A_{SAC} spawns micro-MAS blocks per category
- A_{nli} and $A_{\text{col_check}}$ act as *gates* (reduce hallucinations, enforce theme/diet)
- A_{rank} personalises final presentation; LayoutTool renders a board

Combination of workflow and autonomy: A_{SAC} determines which specialized agents to call.

Interactive Recommendation

Memory Requirements (I)

Stores

- **STM** (working): A_{chat} keeps last L turns of $C_{1:t}$
- **EPI**: A_{epi} manages episode store via U/Q
- **SEM**: stable traits (colors, dietary rules, budget)
 - This may have its update pipeline.
- **PROC**: reusable prompts/templates, DB schema for autonomous tool use and DB read.
 - What are the internal tables and when should the agent call them?

Interactive Recommendation

Memory Requirements (II): Operators & Example

Retrieve Retrieval Function Q for sub-goal τ of cake retrieval:

$$\begin{aligned}\hat{C} &= Q(\Omega_t, \tau) \\ &= \{\text{Flavor_pref: chocolate,} \\ &\quad \text{query: "Mickey-Mouse Cake"}\}.\end{aligned}$$

Feeds SearchCakeAPI with gluten-free, chocolate filter.

Update Update function for memory \mathcal{U} :
User: "gluten-free restriction."

- Episode Distillation:

$$\tilde{C}_t = \text{"\{allergy: gluten\}"}$$

- Episodic Memory Update:

$$\Omega_{t+1}^{\text{EPI}} = \Omega_t^{\text{EPI}} \cup \{\text{allergy: gluten}\}.$$

Interactive Recommendation

Tooling

- Tools \mathcal{F} : SpecializedSearchAPI, (graphic board of decor set).
- Sparse \mathbf{C} limits communication.
- Observability: per-agent logs for `query/tool_call/ranked_list`
- Timeouts & fallbacks:
 - Control on hit rate,
 - Failure back up

Interactive Recommendation

Discussion: Benefits in Production

- **Modularity**: category blocks are micro-MAS units, reusable across themes
- **Memory-aware personalisation**: injects user- and session-specific constraints , SEM/EPI yield relevant recommendation.
- **Error containment**: NLI & collection checks act as fact gates
- **Autonomy**: specialised-agent caller spawns workers on demand.
- **Richer User Experience**: interactive, using tools like LayoutTool offers better experience than a static recommendation.

Agentic Recommendation Evaluation

Task Definition

Simulation-based evaluation (informal)

Generate **synthetic user behavior** to stress-test recommendation policies *before* live deployment.

Core mappings

$R_\phi : \mathcal{X} \rightarrow \mathcal{S}^L$ (recommender maps state x to L items)

$\mathcal{M}_{\theta,\omega} : \mathcal{X} \times \mathcal{S}^L \rightarrow \mathcal{A}_{\text{user}}$ (user simulator maps (x, list) to an action)

$\mathcal{A}_{\text{user}} \in \{\text{Select, Not Select}\}$ or richer (e.g. $\{\text{Click, Pass, Purchase}\}$).

Agentic Recommendation Evaluation

Fall Seasonal Recommendation

The image displays a website interface for a home catalog. On the left, a large featured image shows a living room with a green armchair, a wooden sideboard, and a blue door. Text overlay reads "Introducing our fall home catalog" with a "Shop now" link. The main content area is a grid of product recommendations under the heading "Home Catalog Shop All". The grid includes items such as a green armchair, a patterned throw pillow, a brown blanket, a mirror, a candle, a beige pillow, a wooden box, a plant, a vase, and a framed picture. Each item has a price tag, a "Shop now" button, and a "View details" link. The website has a navigation bar at the top with various category icons and filters.

User Simulation

What is a new trend?

Results for "labubu" (76)
Uses item details. Price when purchased online

in 30+ people's carts Not in cart in 30+ people's carts

Item 1: Pop Mart Labubu The Monsters Coca Cola Series Vinyl Face Single Blind Box, from StockX. Price: ~~\$81.00~~ \$69.00. Free shipping, arrives in 3-5 days. Only 1 left.

Item 2: Pop Mart Labubu The Monsters Big into Energy Series Luck Vinyl Plush Pendant, from StockX. Price: \$71.00. Free shipping, arrives in 3-5 days. Only 1 left.

Item 3: Pop Mart Labubu The Monsters Big into Energy Series Love Vinyl Plush Pendant, from StockX. Price: \$71.00. Free shipping, arrives in 3-5 days. Only 1 left.

Item 4: Pop Mart Labubu The Monsters Coca Cola Series Figure Single Blind Box, from StockX. Price: \$97.00. Free shipping, arrives in 3-5 days. Only 1 left.

Item 5: Pop Mart Labubu The Monsters Big into Energy Series Vinyl Face Single Blind Box, from StockX.

Item 6: Pop Mart Labubu The Monsters Big into Energy Series Vinyl Plush Pendant, from StockX.

Item 7: Pop Mart Labubu The Monsters Big into Energy Series Love Vinyl Plush Pendant, from StockX.

Item 8: Pop Mart Labubu The Monsters Coca Cola Series Figure Single Blind Box, from StockX.

User Simulation

What is a new trend?

Stitch Toys in Lilo and Stitch (565)
Uses item details. Price when purchased online

Item	Price	Description	Rating	Shipping
Disney Stitch Plush, Kids Toys for Ages 2 and up	Now \$9.00 (was \$10.00)	Disney Stitch Plush, Kids Toys for Ages 2 and up	4.5 stars (273)	Free shipping, arrives in 3-5 days
Disney Stitch Many Moods Stitch Sounds and Phrases Interactive Plush, 14 in, Ages 3 and up	\$39.97	Disney Stitch Many Moods Stitch Sounds and Phrases Interactive Plush, 14 in, Ages 3 and up	4.5 stars (142)	Free pickup today Delivery today Free shipping, arrives tomorrow
Disney Stitch Plush Toy, 14 in, Ages 2 and up	Now \$19.97 (was \$25.00)	Disney Stitch Plush Toy, 14 in, Ages 2 and up	4.5 stars (23)	Free pickup today Delivery today Shipping, arrives today
Disney Stitch Angel Plush Stuffed Animal, 14 in, Ages 2 and up	\$19.97	Disney Stitch Angel Plush Stuffed Animal, 14 in, Ages 2 and up	4.5 stars (10)	Free pickup today Delivery today Shipping, arrives today

High-Level Goals (I)

- **De-risk A/B:** approximate user responses when real traffic is scarce or expensive
 - Sometimes, duration of the recommendation being live is short-lived (no time for AB),
 - Example: Seasonal Recommendation.
- **Fast iteration:** offline loops to evaluate design choices, ranking features, and guardrails
 - AB tests by design are slower (specially if we wait for stat sig results on niche general merchandise items)
- **Coverage:** explore edge cases (long-tail preferences, rare journeys) at scale
 - Hard to perform AB test on user cohorts,
 - Example: all new users with account opening in the last year.
- **Safety:** detect regressions, policy violations, and mode collapse before serving
 - Better control on the edge cases because we can simulated edge cases at scale.

High-Level Goals (II)

- **Measurement:** CTR/CVR, diversity, novelty, calibration, fairness, robustness
 - Should optimize on the above metrics,
 - Hard to track for customer cohorts.
- **Session granularity:** short (*quick browse*) vs. long (*deliberative*)
 - In each AB test we want to analyse short term effect and long term effect of recommendation,
 - Short term: purchase within session,
 - Longer term: product return patterns,
 - Longerer term: User loyalty.
- **Population realism:** sample (θ, ω) to reflect heterogeneous users
 - θ : can represent user profiles.
 - ω : we can add controlled noise to better simulate user variability withing cohorts.
- **Explainability:** log rationales and session summaries for root-cause analysis

Formal Definition (I)

In Recommendation Evaluation and User simulation task

- We define a user/platform state \mathcal{X} ,
- Recommend \mathcal{S}^L ,
- Observe user behavior $\mathcal{A}_{\text{user}}$.

$$R_\phi : \mathcal{X} \rightarrow \mathcal{S}^L, \quad x_t \mapsto s_t^{(L)}$$
$$\mathcal{M}_{\theta, \omega} : \mathcal{X} \times \mathcal{S}^L \rightarrow \mathcal{A}_{\text{user}}, \quad (x_t, s_t^{(L)}) \mapsto a_t^{\text{user}}$$

Iterated for $t = 1, \dots, T$ to generate trajectories $\{x_t, s_t^{(L)}, a_t^{\text{user}}\}_{t=1}^T$.

Formal Definition (II)

Task objective

Task objective is to estimate any relevant metrics for user interaction,

$$\Psi(R_\phi, \mathcal{M}_{\theta, \omega}) = \mathbb{E} \left[\sum_{t=1}^T g(x_t, s_t^{(L)}, a_t^{\text{user}}) \right],$$

where g can encode Select/Purchase reward, diversity/penalty, fairness, etc.

- Estimate Ψ under controlled distributions of (θ, ω)
- Simulate and estimate results via Monte Carlo simulation.

Formal Definition (III): Population & Estimation

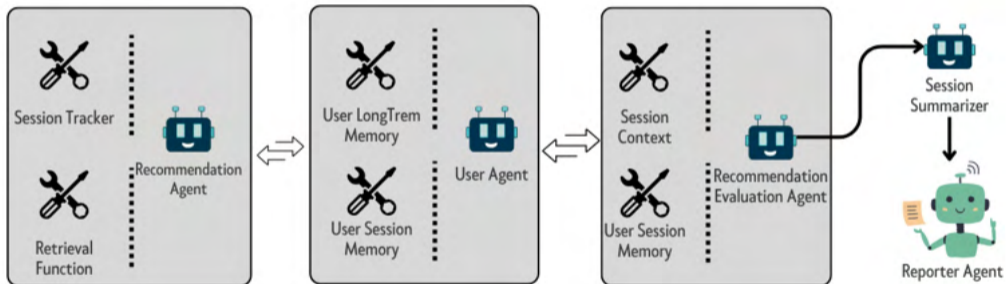
- **Population:** $(\theta, \omega) \sim P$ (PreferenceSampler)
- **Monte Carlo estimator:**

$$\hat{\Psi} = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T g(x_t^{(n)}, s_t^{(L,n)}, a_t^{\text{user}(n)})$$

- **Risk-aware variants:** Measure metrics like
 - CVaR@ k ,
 - worst-case across cohorts,
 - per-segment fairness gaps.

Architecture

The Visual



$$\text{MAS}_{\text{sim}} = (\mathcal{A}, \mathcal{E}, \Pi)$$

$$\mathcal{A} = \{A_{\text{rec}}, A_{\text{user}}, A_{\text{note}}, A_{\text{eval}}, A_{\text{summ}}, A_{\text{report}}\}$$

- A_{rec} : wraps R_ϕ ; calls SearchAPI, SessionTracker
- A_{user} : implements $\mathcal{M}_{\theta, \omega}$ with structured memory
- A_{note} : durable logging of $(x_t, s_t^{(L)}, a_t^{\text{user}})$
- A_{eval} : computes g_t ; emits eval_event
- A_{summ} : *compresses* sessions via \mathcal{R}_{sum}
- A_{report} : aggregates, tests, produces PDF/CSV dashboards

Environment

$\mathcal{E} = (\text{ProductDB}, \text{PreferenceSampler}, \text{LogStore})$

Protocol

$A_{\text{rec}} \rightleftarrows A_{\text{user}} \rightarrow A_{\text{eval}} \rightarrow A_{\text{summ}} \rightarrow A_{\text{report}},$

$\Gamma = \{\text{rec_list}, \text{user_action}, \text{log_entry}, \text{eval_event}, \text{session_summary}, \text{final_report}\}.$

C permits minimal paths to reduce chatter; logs are append-only.

End-to-End Flow

- 1 A_{rec} : produce $s_t^{(L)}$ for state x_t
- 2 A_{user} : act $a_t^{\text{user}} \sim \mathcal{M}_{\theta, \omega}(x_t, s_t^{(L)})$
- 3 A_{note} : log $(x_t, s_t^{(L)}, a_t^{\text{user}})$
- 4 A_{eval} : compute g_t and `eval_event`
- 5 A_{summ} : periodic \mathcal{R}_{sum} summaries
- 6 A_{report} : aggregate \Rightarrow KPIs, cohort analyses

Memory Requirements (I)

- A_{user} :
 - Ω^{SEM} (latent taste vector, price sensitivity, stable user features)
 - Ω^{STM} (current trajectory)
 - Ω^{EPI} (prior sessions / exposure effects)
 - Ω^{PROC} (navigation/click heuristics)
 - Encodes how the user proceeds in a session
- A_{eval} :
 - Raw text log + vector store for fast Q by item/action
 - Sliding-window stats (CTR, entropy) in STM

Tools & Data Access (II)

- `SearchAPI`, `ProductDB.query` for candidate context
- `PreferenceSampler` for (θ, ω) populations
- `LogStore.append/LogStore.scan` for \mathcal{U}/\mathcal{Q}
- Offline analytics: aggregation, stratified metrics, significance tests

User Simulation

Fall Seasonal Recommendation

The image shows a user simulation interface for a home catalog. On the left, there is a featured image of a living room with a green armchair and a wooden coffee table. The text "Introducing our fall home catalog" is overlaid on this image, with a "Shop now" link below it. To the right, a navigation bar includes categories like Living Room, Bed, Dining Room, etc. Below the navigation bar is a grid of product recommendations. A red arrow points from the green armchair in the featured image to a specific product in the grid: a brown blanket. The product card for the blanket is highlighted with a red border and includes a price tag of \$99. Other products in the grid include pillows, a wooden box, a plant, a vase, and a framed picture.

User Simulation

Fall Seasonal Recommendation

The screenshot shows a Walmart product page for a 'Better Homes & Gardens Cherille Knit Super Soft Oversized Throw Blanket, Terracotta Clay'. The main image shows the blanket draped over a white chair. To the left is a vertical gallery of smaller images showing the blanket in different colors and settings. The product title is 'Better Homes & Gardens Cherille Knit Super Soft Oversized Throw Blanket, Terracotta Clay' with a 4.5-star rating and 42 reviews. The price is \$19.96. Below the price are three color swatches: Terracotta Clay (selected), Dark Green, and Dark Grey. The 'About this item' section lists features: 100% Recycled Polyester, 50" x 72" Throw, Textured Throw, Adult, Unisex, Warm, and Oversized. The 'At a glance' section provides a summary of the product's attributes in a grid format.

At a glance		
Brand Better Homes & Gardens	Size 50" x 72"	Color Terracotta Clay
Weight 2.05-lb.	Pattern Ribbed	Decor style Farmhouse, French Country, Traditional,...

Additional details on the right side of the page include the price (\$19.96), a 'Free 90-day returns' badge, an 'Add to cart' button, and delivery options. The shipping section indicates 'Arrives Sep 10' with a 58-minute delivery time. The location is 'Bentonville, 72718' and the item is sold and shipped by Walmart.com.

User Simulation

Scroll Down-Click on second item third color

Similar items you might like

Based on what customers bought

Best seller



Options

Sponsored

\$22.79

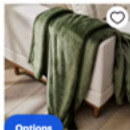
More options from \$13.99

Exclusivo Mezcla Queen Size Flannel Fleece Velvet Plush Bed Blanket as Bedspread, Coverlet, Bed Cover...

★★★★★ 623

Save with W+

Best seller



Options

● ● ● ● +10

\$18.24

Better Homes & Gardens Solid Velvet Plush Soft Fleece Throw Blanket, Oversized, Sea Turtle

★★★★★ 1000

Save with W+

Shipping, arrives today

Best seller



Options

Sponsored

Now \$16.99 \$20.99

Options from \$16.99 - \$32.99

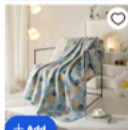
Nestl Cut Plush Fleece Blanket, Soft Lightweight Fuzzy Luxury Throw Size Bed Blankets for Bed, Throw, Gray

★★★★★ 2070

Save with W+

Shipping, arrives in 2 days

500+ bought since yesterday



+ Add

● ● ● ● ●

\$4.22

Options from \$4.22 - \$56.64

Mainstays Cozy Fleece Throw Blanket, Gray Paw 50" x 60", All Ages

★★★★★ 5538

Save with W+

Shipping, arrives today



User Simulation

Color Out of Stock!

Best offer

Better Homes & Gardens

Better Homes & Gardens Textured Velvet Plush Soft Fleece Throw Blanket, Oversized, Mauve Gem

★★★★☆ (4.7) | 300 reviews

Actual Color: Violet Gem - **Out of stock**

\$18.24	\$18.24	\$18.24	\$18.24	\$18.24	\$18.24
\$18.24	\$18.24	\$18.24	\$18.24	\$18.24	\$18.24
\$18.24	\$18.24				

About this item

- Velvet Plush Throw
- 100% Polyester
- 50" x 72" Throw
- Textured Fleece
- Unisex
- Oversized
- Adult
- Soft ...

View more ▾

\$18.24
Price when purchased online

Not Available

How you'll get this item:

 Shipping Not available	 Pickup Not available	 Delivery Not available
-------------------------------	-----------------------------	-------------------------------

Not available

[Add to list](#) [Add to registry](#)

Sponsored

\$18.00
Elegant Comfort Velvet Touch Ultra Plush Halloween Hobble...

★★★★☆ (4.7)

2-day shipping

[+ Add](#)

User Simulation

Add to correct the color you like

The screenshot shows a Walmart product page for a "Better Homes & Gardens Solid Velvet Plush Soft Fleece Throw Blanket, Oversized, Copper Pipe". The main image shows the blanket draped over a white chair. To the left is a vertical gallery of color swatches, with the "Copper Pipe" color selected. Below the main image is a grid of 18 color swatches, each with a price of \$18.24. The "Copper Pipe" swatch is highlighted with a red border. The product title and price are displayed prominently. The right side of the page shows the price "\$18.24", a "Free 90-day returns" badge, and an "Add to cart" button. Below this, there are options for shipping, pickup, and delivery. The bottom of the page features a "Get free delivery, shipping and more" banner.

Save on 5G phones. Get one at T-Mobile.

35+ bought alone yesterday | In 30+ people's carts

Bestseller

Better Homes & Gardens

Better Homes & Gardens Solid Velvet Plush Soft Fleece Throw Blanket, Oversized, Copper Pipe

★★★★★ (4.7) | 1000 reviews

Actual Color: Copper Pipe

\$18.24	\$18.24	\$18.24	\$18.24	\$18.24	\$18.24
\$18.24	\$18.24	\$18.24	\$18.24	\$18.24	\$18.24
\$18.24	\$18.24				

About this item

- Velvet Plush Throw
- 100% Polyester
- 50" x 72" Throw
- Adult
- Fleece
- Unisex
- Soft
- Insulating...

View more

\$18.24

Price when purchased online

Free 90-day returns

Add to cart

How you'll get this item:

I want delivery savings with Walmart+
You get \$0 down from Choose a plan at checkout.

 Shipping Arrives today Order within the 30 min	 Pickup As soon as 3pm today	 Delivery As soon as 1 hour
---	---	--

Bentonville, 72713 Change

Arrives by **Today**. Order within the 30 min

Sold and shipped by Walmart.com

Free 90-day returns Details

This item is gift eligible Learn more

Add to list Add to registry

Walmart+ Get free delivery, shipping and more*


*Restrictions apply. Start 30-day free trial

User Simulation

Added to cart-Rest of the Journey

Tyson Try Asian-inspired flavor Korean BBQ & hot orange. Sponsored

Added to cart!


 Better Homes & Gardens Solid Velvet Plush Soft Fleece Throw Blanket, Oversized, Copper Pipe **\$18.24 ea** \$18.24/ea

View cart (2)

Subtotal	\$34.23
Taxes	Calculated at checkout
Estimated total	\$34.23

Customers also bought these products

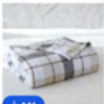
Best seller

 **+ Add**

Sponsored
\$34.97 \$42.32
Better Homes & Gardens 4-Piece 400 Thread Count Arctic White...
★★★★★ 3805
Save with W+

Shipping, arrives in 2 days


Best seller

 **+ Add**

\$17.96
Mainstays Super Soft Plush Blanket, Brown Plaid, Twin, Adult/Teen
★★★★★ 5738
Save with W+

Pickup today Delivery today Shipping, arrives today


Best seller

 **+ Add**

\$24.72 \$35.54
Better Homes & Gardens Luxury Velvet Plush Blanket, Dark Grey...
★★★★★ 1836
Save with W+

Shipping, arrives in 2-3 days

Best seller

 **+ Add**

\$29.96
Better Homes & Gardens Cozy Knit Blanket, Beige, Full/Queen
★★★★★ 520
Save with W+

Pickup today Delivery today Shipping, arrives today

Top picks to explore

Benefits (I)

- **Cost-effective experimentation:** run $N \gg 1$ users in parallel; compare ϕ variants safely
 - Can study cohort groups better,
 - more granularity.
- **Cold-start mitigation:** sample sparse/unseen θ to probe failure modes early
 - new trends,
 - new items,
 - new users
- **Repeatability:** deterministic seeds \Rightarrow reproducible KPI deltas
 - With different language models,
 - Same LLM to analyze variations in user journey trajectories

User Simulation

What is a new trend?

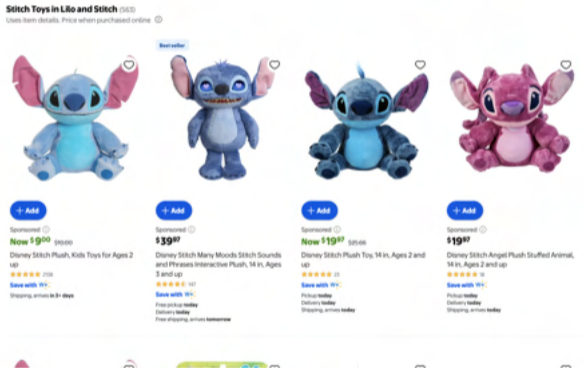
Results for "labubu" (76)
Uses item details. Price when purchased online

The screenshot displays a search results page for 'labubu' toys. At the top, it says 'Results for "labubu" (76)' and 'Uses item details. Price when purchased online'. Below this, there are four items in a row, each with a 'to 30+ people's carts' or 'Best seller' label, a heart icon, and an 'Add' button. The first item is a 'Pop Mart Labubu The Monsters Coca Cola Series Vinyl Face Single Blind Box, from StockX' priced at \$81.00 (was \$88.00). The second is a 'Pop Mart Labubu The Monsters Big into Energy Series Luck Vinyl Plush Pendant, from StockX' priced at \$71.00. The third is a 'Pop Mart Labubu The Monsters Big into Energy Series Love Vinyl Plush Pendant, from StockX' priced at \$97.00. The fourth is a 'Pop Mart Labubu The Monsters Coca Cola Series Figure Single Blind Box, from StockX'. Below these, there are four more items in a row, including another blind box, a vinyl figure in a box, and two different plush pendants.

User Simulation

What is a new trend?

Stitch Toys in Lilo and Stitch (5/5)
Uses item details. Price when purchased online



Product Name	Price	Rating	Shipping
Disney Stitch Plush, Kids Toys for Ages 2 and up	Now \$9.00 (was \$10.00)	4.5/5 (273)	Free shipping, arrives in 3-5 days
Disney Stitch Many Moods Stitch Sounds and Phrases Interactive Plush, 14 in, Ages 3 and up	\$39.97	4.5/5 (142)	Free pickup today, Delivery today, Free shipping, arrives tomorrow
Disney Stitch Plush Toy, 14 in, Ages 2 and up	Now \$19.97 (was \$25.00)	4.5/5 (23)	Free pickup today, Delivery today, Shipping, arrives today
Disney Stitch Angel Plush Stuffed Animal, 14 in, Ages 2 and up	\$19.97	4.5/5 (10)	Free pickup today, Delivery today, Shipping, arrives today

Benefits (II)

- **Bias diagnosis:** aggregate trends (price sensitivity, popularity bias, demographic skew)
 - Informed edge case handling,
 - Informed recommendation algorithm design.
- **Memory-driven realism:** SEM/EPI/PROC capture carry-over and bounded user attention
 - What is the effect on other recommendation algorithm?
 - what is the effect on long term user loyalty?
 - what is the effect of previous purchases in current user session behavior?
- **Modular extensibility:** add fairness/robustness evaluators without retraining R_ϕ
 - Custom metric evaluation.

Contextual & Multi-Modal Recommendation

fill the image with Boho style furniture



Task Definition

Contextual & multi-modal recommendation (Informal)

Recommend concept-oriented, while adapting to different modalities.

Contextual & multi-modal recommendation (formal)

Let $\mathbf{x} \in \mathcal{T}$ be textual query (e.g. “Bohemian, earthy colors”), $\mathbf{v} \in \mathcal{V}^K$ a set of visuals, and $\mathbf{u} \in \mathcal{N}$ a long-term user profile. We model

$$\mathcal{R}_\phi : (\mathbf{x}, \mathbf{v}, \mathbf{u}) \longrightarrow \hat{\mathbf{y}} = \langle s^1, s^2, \dots \rangle, \quad s^i \in \mathcal{S},$$

yielding a ranked set subject to aesthetic coherence & complementarity.

High-Level Goals (I)

- **Holistic curation:** move beyond single-item retrieval to a *cohesive* set
 - Focus can be on curation of non-similar bundles
 - Topic-oriented Recommendation,
 - Aesthetic Complementarity.
- **Context grounding:** leverage spatial cues (layout, lighting) and palette from images.
Ground Concepts:
 - What does “Boho” Style exactly mean?
 - What does “Industrial” house furniture exactly mean?
- **Personalisation:** incorporate long-term preferences & budget constraints
 - Be aware of user long term features,
 - User’s behavior in similar sessions,
 - Other similar users behavior
- **Effort reduction:** “one-shot designer” experience vs. piecemeal searching
 - Reduce the time and effort of collection building by user.

High-Level Goals (II)

- **Brand/style awareness:** enforce on-brand style guides and user's declared themes
 - The model should be aware of high level queries
 - The model should be aware of brand to style relations
- **Consistency at scale:** ensure cross-item compatibility (materials, colors, sizes)
 - This is hard,
 - constructing cohesive bundles: a lot of edge cases,
- **Interactive refinement:** accept quick corrections
 - Design feature of this task
- **Ready-to-buy boards:** present ranked set with room mock-ups for fast decisions
 - Make life easier for the user.

How to implement?

Sample pseudo-algo

Goal: turn image + text + history into a coherent, feasible, on-budget bundle.

Key primitives in this module:

- 1 **Palette vector** \mathbf{p} : compact color/material signature of the image.
- 2 **Layout affordance** ℓ : spatial constraints (free space, obstacles, anchors).
- 3 **Compatibility kernel** $\kappa(\cdot, \cdot; \mathbf{p})$: pairwise harmony in a bundle.
- 4 **Total score** S_{tot} : personalized relevance + coherence – cost/penalties.

What is a *palette vector* \mathbf{p} ?

Definition

A **palette vector** $\mathbf{p} \in \mathbb{R}^d$ is a compact, normalized descriptor of the dominant *colors & materials* present in the scene images $\mathbf{v} \in \mathcal{V}^K$:

$$\mathbf{p} = [p_1, \dots, p_d], \quad p_i \geq 0, \quad \sum_{i=1}^d p_i = 1.$$

Instantiations (common in practice):

- **Hybrid** (color + material): concatenate color bins with material tags (e.g., rattan, linen, leather) and train a classifier.

Palette Vector: Role in the Pipeline

Usage: \mathbf{p} conditions retrieval/ranking and the *compatibility function* $\kappa(\cdot, \cdot; \mathbf{p})$ so that selected items harmonize with the scene palette.

How to train: Maybe on co-purchase data! Or distill a large VLM to classifier.

Where it enters:

- *Category targeting:* prefer categories whose canonical palettes align with \mathbf{p} .
- *Per-item filtering:* discard items with palettes far from \mathbf{p} .
- *Bundle scoring:* $\kappa(s^i, s^j; \mathbf{p})$ rewards pairwise harmony under the observed scene palette.

What is a layout affordance ℓ ?

Definition

Layout affordance ℓ encodes *spatial constraints* derived from the image including

- geometry: extract the floor plan
- free space dimensions,
- anchor elements in image like window, door, etc

to determine the test feasibility of placing a set $\hat{\mathbf{y}}$ of items.

Feasibility predicate:

$\text{Layout}(\hat{\mathbf{y}}; \ell) = 1 \iff \exists$ non-overlapping placement satisfying extracted geometry, free space di

CAL-RAG: Retrieval-Augmented Multi-Agent Generation for Content-Aware Layout Design

Najmeh Forouzandehmehr, Reza Yousefi Maragheh, Sriram Kollipara, Kai Zhao, Topojoy Biswas, Evren Korpeoglu, Kannan Achan

Walmart Global Tech, Sunnyvale, California, USA

{najmeh.forouzandehmehr, reza.yousefimaragheh, sriram.kollipara, kai.zhao, topojoy.biswas, evren.korpeoglu, kannan.achan}@walmart.com

ABSTRACT

Automated content-aware layout generation—the task of arranging visual elements such as text, logos, and underlays on a background canvas—remains a fundamental yet underexplored problem in intelligent design systems. While recent advances in deep generative models and large language models (LLMs) have shown promise in structured content generation, most existing approaches lack grounding in contextual design exemplars and fall short in handling semantic alignment and visual coherence. In this work, we introduce CAL-RAG, a Retrieval-Augmented, Agentic framework for content-aware layout generation that integrates multimodal retrieval, large language models, and collaborative agentic reasoning. Our system retrieves relevant layout examples from a structured knowledge base and invokes an LLM-based layout recommender to propose structured element placements. A vision-language grader agent evaluates the layout based on visual metrics, and a feedback agent provides targeted refinements, enabling iterative improvement. We implement our framework using LangGraph and evaluate on the PKU PosterLayout dataset, a benchmark rich in semantic and structural variability. CAL-RAG achieves state-of-the-art performance across multiple layout metrics—including underlay effectiveness, element alignment, and overlap—substantially outperforming strong baselines such as LayoutPrompter. Our results demonstrate that combining retrieval augmentation with agentic multi-step reasoning provides a scalable, interpretable, and high-fidelity solution for automated layout generation.

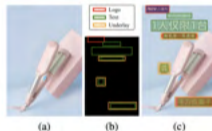


Figure 1: Example for content-aware layout generation: (a) Input canvas with background and content elements; (b) Generated layout based on visual and textual content awareness; (c) Final rendered presentation using the layout from (b).

offer limited generalization, particularly in diverse or semantically rich design contexts. The advent of deep generative models—such as GANs, VAEs, and transformers—has enabled more expressive layout synthesis by learning from annotated layout corpora [1, 2, 7, 9, 10]. However, these models are typically data-hungry, brittle to out-of-distribution inputs, and often struggle to incorporate visual-semantic alignment at inference time.

More recently, Large Language Models (LLMs) and Large Vision-Language Models (LVLMs) have demonstrated remarkable zero-

506.21934v1 [cs.IR] 27 Jun 2025

CAL-RAG

check out this paper

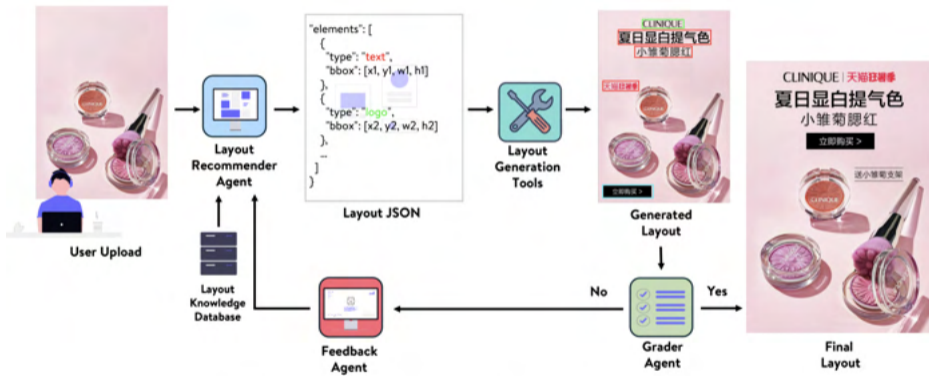


Figure 2: System Architecture Diagram of CAL-RAG.

Compatibility Kernel $\kappa(\cdot, \cdot; \mathbf{p})$ (Examples)

$$\kappa(s^i, s^j; \mathbf{p}) = \underbrace{\text{sim}(\text{pal}(s^i), \text{pal}(s^j) \mid \mathbf{p})}_{\text{color harmony}} + \rho \underbrace{\text{sim}_{\text{material}}(s^i, s^j)}_{\text{e.g., GloVe/CLIP on material tags}} + \tau \underbrace{\text{sim}_{\text{style}}(s^i, s^j)}_{\text{style embedding cosine}} .$$

Bundle coherence:

$$S_{\text{comp}}(\hat{\mathbf{y}}) = \frac{2}{L(L-1)} \sum_{i < j} \kappa(s^i, s^j; \mathbf{p}).$$

Notes:

- Color harmony condition by \mathbf{p} to favor harmony around scene-dominant hues.
- Learned sim heads can be trained to match designer “goes-well-with” labels.

Explaining the *Total Score* S_{tot}

Decomposition

$$S_{\text{tot}}(\hat{\mathbf{y}}) = \underbrace{\sum_{i=1}^L \alpha \text{Rel}(s^i | \mathbf{x}, \mathbf{u})}_{\text{personal relevance}} + \underbrace{\beta S_{\text{comp}}(\hat{\mathbf{y}})}_{\text{coherence/complementarity}} - \underbrace{\gamma \text{Cost}(\hat{\mathbf{y}})}_{\text{budget}} - \underbrace{\lambda \mathcal{P}_{\text{brand}}(\hat{\mathbf{y}})}_{\text{style/brand penalty}} .$$

Hard vs. soft constraints:

Hard: $\text{Compat}(\hat{\mathbf{y}}; \mathbf{p}, \mathcal{G}) = 1$; Soft: add penalties in S_{tot} .

Choosing weights: $\alpha, \beta, \gamma, \lambda$ via grid/Bayes search or *learned* from designer judgments (pairwise ranking).

Optimizing S_{tot} under Constraints

Problem:

$$\hat{\mathbf{y}}^* = \arg \max_{\hat{\mathbf{y}}} S_{\text{tot}}(\hat{\mathbf{y}}) \quad \text{s.t. Budget} \leq B, \text{ Layout}(\hat{\mathbf{y}}; \ell) = 1.$$

Solvers (latency/quality trade-offs):

- *Greedy w/ lookahead*: fast, good for serving.
- *ILP/MIP relaxations*: exact/near-exact, best offline.

Caching: precompute per-item relevance and per-pair κ for top categories to accelerate online assembly.

Another Formal Definition

$$\mathbf{x} \in \mathcal{T}, \quad \mathbf{v} \in \mathcal{V}^K, \quad \mathbf{u} \in \mathcal{N}, \quad \hat{\mathbf{y}} = \langle s^1, \dots, s^L \rangle, \quad s^i \in \mathcal{S}.$$

$$\text{Objective: } \hat{\mathbf{y}}^* = \arg \max_{\hat{\mathbf{y}}} \sum_{i=1}^L \text{Rel}(s^i \mid \mathbf{x}, \mathbf{v}, \mathbf{u})$$

Subject to design constraints:

$$\text{Compatability}(\hat{\mathbf{y}}; \mathbf{p}, \mathcal{G}) = 1, \quad \text{Budget}(\hat{\mathbf{y}}) \leq B, \quad \text{Layout}(\hat{\mathbf{y}}; \ell) = \text{feasible},$$

where \mathbf{p} is a palette vector from vision, \mathcal{G} style rules, and ℓ layout affordances.

Another Formal Definition

History retrieval: $\hat{\mathcal{C}}_{\text{SEM}} = \mathcal{Q}(\Omega^{\text{SEM}}, \tau(\mathbf{x}))$

Category selection: $C = \text{Cat}(\mathbf{x}, \mathbf{p}, \hat{\mathcal{C}}_{\text{SEM}})$

Per-category retrieval: $\forall c \in C, \mathcal{F}_{\text{search}}(c, \mathbf{p}, \mathbf{u}) \rightarrow \mathcal{S}_c$

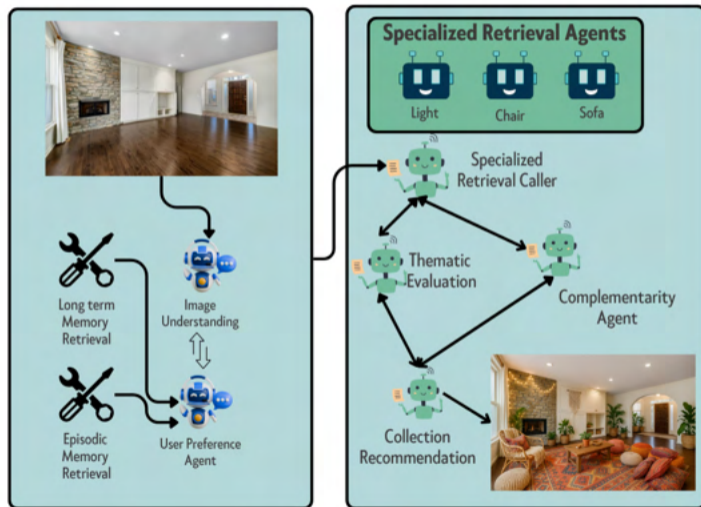
Bundle assembly: $\hat{\mathbf{y}} = \bigcup_{c \in C} \text{TopK}(\mathcal{S}_c)$

Semantic & complementarity gates: $\text{SemChk}(\hat{\mathbf{y}}) = 1, \text{CompChk}(\hat{\mathbf{y}}) = 1.$

End-to-End Algorithm (Pseudo-Code)

- 1 $(\mathbf{x}, \mathbf{v}) \leftarrow$ user input; $\hat{\mathcal{C}}_{\text{SEM}} \leftarrow \mathcal{Q}(\Omega^{\text{SEM}}, \tau(\mathbf{x}))$.
- 2 $\mathbf{p} \leftarrow f_{\text{img}}(\mathbf{v})$; $\ell \leftarrow$ affordance extractor(\mathbf{v}).
- 3 $C \leftarrow \text{Cat}(\mathbf{x}, \mathbf{p}, \hat{\mathcal{C}}_{\text{SEM}})$.
- 4 For $c \in C$: $\mathcal{S}_c \leftarrow \mathcal{F}_{\text{search}}(c, \mathbf{p}, \mathbf{u})$; keep TopK per category.
- 5 Assemble $\hat{\mathbf{y}}$ by maximizing S_{tot} under Budget, Layout.
- 6 Run SemChk, CompChk; if fail, revise C /constraints and re-solve.
- 7 Render mock-up; return bundle and visualization to A_{chat} .

Architecture: The Visual



$$\text{MAS}_{\text{mm}} = (\mathcal{A}, \mathcal{E}, \Pi)$$

$$\mathcal{A} = \{A_{\text{chat}}, A_{\text{image}}, A_{\text{history}}, A_{\text{cat}}, A_{\text{caller}}, A_{\text{semChk}}, A_{\text{compChk}}, A_{\text{collect}}\}$$

- A_{image} : palette & layout from images ($f_{\text{img}}, \mathcal{F}_{\text{layout}}$)
- A_{history} : \mathcal{Q} over Ω^{SEM} for stable tastes/budget
- A_{caller} : combines signals; calls specialized micro-MAS,
- $A_{\text{semChk}}/A_{\text{compChk}}$: rule compliance & harmony gates
- A_{collect} : rank & render; return set to A_{chat}

$$\begin{aligned} A_{\text{chat}} &\rightarrow A_{\text{image}}(\mathbf{v}) \ \& \ A_{\text{history}}(\mathbf{x}) \\ &\rightarrow A_{\text{caller}}(\mathbf{x}, \mathbf{p}, \hat{\mathcal{C}}_{\text{SEM}}) \\ &\rightarrow \{\text{micro-MAS}_c\}_{c \in \mathcal{C}} \rightarrow A_{\text{semChk}} \rightarrow A_{\text{compChk}} \rightarrow A_{\text{collect}} \rightarrow A_{\text{chat}} \end{aligned}$$

Autonomy can happen based in Agent calling to adjust the retrieval process per query.

Memory Requirements (I)

- **Short-term (STM):** live prompt, extracted palette \mathbf{p} , current candidate sets
- **Semantic (SEM):** durable preferences (colors, materials, budget caps), style affinities
- **Episodic (EPI):** past styling sessions (accepted/rejected bundles, returns)
- **Procedural (PROC):** routing templates, per-category query schemes, brand/style rules

Tools & Data Access (II)


- VisionEncoder f_{img} , LayoutRenderer \mathcal{F}_{render}
- VectorSearch over ProductDB for semantic retrieval
- UniversalSearch/PriceAPI for stock/price validation
- CompatEngine for pairwise/multi-item harmony checks

- **Designer-quality curation:** coherent, theme-consistent sets in one step
 - Now, near designer quality!
- **Lower cognitive load:** fewer iterations vs. item-by-item browsing
 - Easy user experience
- **Personal, grounded:** vision-aware + history-aware = higher acceptance
 - This was not possible before!
 - **Error containment:** semantic/complementarity gates reduce mismatch & drift
- **Modular extensibility:** add/replace micro-MAS for new styles or categories
 - Easier implementation.

Example

Minnie Mouse

✓ Added to cart!




6V Huffy Disney Minnie Mouse Battery-Powered Ride-On Car, Kids Ages 3+ - Pink
\$174.00 ea ~~\$174.00/ea~~

— 1 +

Explore Fun Minnie Mouse Gear for Kids

Generated by AI


Best seller



+ Add

\$31.10
Bell Disney Minnie Mouse Bike Helmet, Pink Flowers, Toddler 3 (48-52cm)
★★★★☆ 290
Save with W+
Shipping, arrives in 3+ days


Best seller



+ Add

\$5.97
Disney Minnie Mouse Girl's Brow Bar Sunglasses Pink
★★★★★ 432
Save with W+
Pickup tomorrow
Delivery today
Shipping, arrives tomorrow


Best seller



+ Add

\$16.95
Bell Disney Minnie Mouse Protective Pad and Glove Set
★★★★☆ 106
Save with W+
Shipping, arrives in 3+ days

Best seller



+ Add

\$5.97
Minnie Mouse Girl's Pink Heart Sunglasses
★★★★★ 103
Save with W+
Pickup tomorrow
Delivery today
Shipping, arrives tomorrow

Example

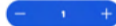
Summer Dress

✓ Added to cart!



MOSHU Ribbed Trim Tank Tops
for Women Flowy Round Neck
Women Shirts Loose Fit
Sleeveless Summer Tops

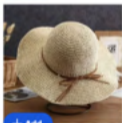
Now \$10.00 ea ~~\$25.99~~



Stay Stylish in the Sun with Summer Essentials

Generated by AI

Best seller



+ Add

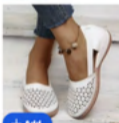
\$6.80

Set Of 4,Silicone Stainless
Steel Food Tongs-Black

★★★★☆ 53

Shipping, arrives in 3+ days

Best seller



+ Add

\$14.29

50% off Summer! TMOYZQ
Wedge Sandals for Women,
Summer Closed Round To...

★★★★☆ 31

Shipping, arrives in 3+ days

Best seller



+ Add

Now \$12.97 ~~\$17.00~~

Joopin Oversized Polarized
Sunglasses for Women
Vintage Lady UV Protectio...

★★★★★ 371

Save with W+

Shipping, arrives in 2 days

Best seller



+ Add

Now \$13.99 ~~\$19.99~~

BCOOSS Summer Sun Hat
for Women Wide Brim Sun
Protection Women Straw...

★★★★★ 317


Save with W+

Shipping, arrives in 2 days

Example

Laptop for Everyday Use

✓ **Added to cart!**




ASUS CX15
15.6 inch
Laptop
Now \$159.00 ea
~~\$219.99~~
~~\$159.00/ea~~

— 1 +


Discover Essential Accessories for Your Everyday Tech Needs
Generated by AI

Reduced price



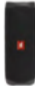
+ Add

Now \$13.40 ~~\$22.99~~
Logitech M185 Wireless
Computer Mouse
★★★★★ 1475
Save with W+
Shipping, arrives in 2 days




+ Add

\$133.12
Sit to Stand Mobile Laptop
Computer Stand with
Height Adjustable & Tiltab...
Shipping, arrives in 3+ days



+ Add

Now \$89.95 ~~\$119.99~~
JBL Flip 5 - Portable
Waterproof Speaker - Black
Matte
★★★★★ 4427
Save with W+
Shipping, arrives in 2 days

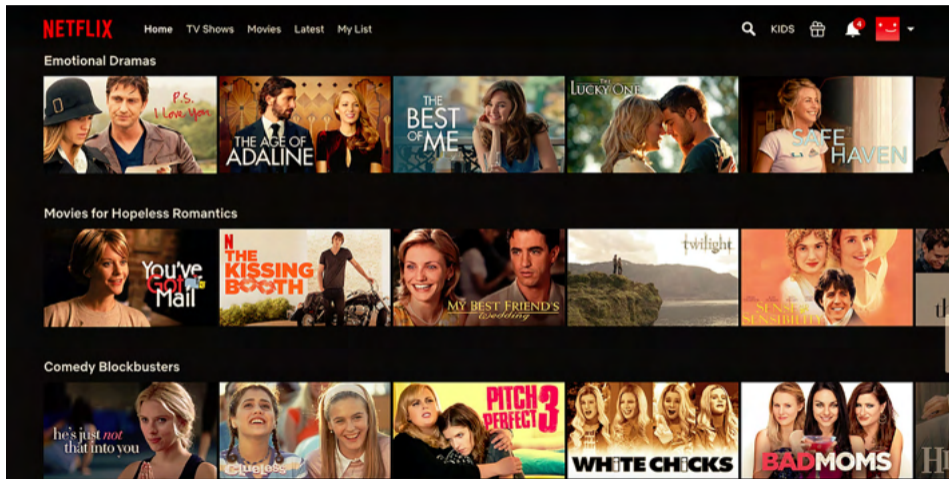


+ Add

Now \$13.49 ~~\$22.85~~
Logitech M185 Wireless
Mouse, 2.4GHz with USB
Mini Receiver, 12-Month...
★★★★★ 15
Save with W+
Shipping, arrives in 3+ days

Recommendation Explanation

Why are you recommending this?



Task Definition

Informal Definition

Recommendation Explanation is the process of generating intelligible and contextually relevant justifications for why particular items are recommended to a user.

Recommendation-Explanation maps items and user context to text:

$$\mathbf{r} \in \mathcal{S}^L, \quad \mathbf{u}_t = (\hat{\mathcal{C}}_t^{\text{SEM}}, \hat{\mathcal{C}}_t^{\text{EPI}}), \quad e_t = \Phi_\psi(\mathbf{r}, \mathbf{u}_t) \in \mathcal{E}_{\text{text}}.$$

Constraints (consistency predicate)

$$\text{Consistency}(e_t) = \underbrace{\text{Factual}(e_t : \mathbf{r}, \mathbf{u}_t)}_{\text{no hallucinations}} \wedge \underbrace{\mathcal{C}_P(e_t) = 1}_{\text{brand/style compliance}}$$

Notes. $\hat{\mathcal{C}}_t^{\text{SEM}}, \hat{\mathcal{C}}_t^{\text{EPI}}$ obtained via memory retrieval \mathcal{Q} ; \mathcal{C}_P is a policy checker (tone, vocabulary, legal).

High-Level Goal

- **Transparency:** expose the rationale behind recommended items in one concise narrative.
 - Explaining can help bridging the gap with what user wants and the recommendation,
 - makes the experience easier.
- **Trust & engagement:** increase user confidence and CTR with factual, user-aware explanations.
 - If the explanation explains what user is looking for $-j$ then the decision process is faster and conversion happens with more probability.
 - we have tested this and have evidence for it.
- **Brand consistency:** ensure on-tone language and regulatory alignment across languages & markets.
 - This is one of the most important issues, in large launches.
- **Low latency:** integrate into serving path with bounded overhead and revise-on-fail loop.
 - Cost benefit analysis, if doing so has a positive net impact!

System Architecture (Agents)

$$\text{MAS}_{\text{expl}} = (\mathcal{A}, \mathcal{E}, \Pi), \quad \mathcal{A} = \{A_{\text{journey}}, A_{\text{session}}, A_{\text{rec}}, A_{\text{expl}}, A_{\text{eval}}\}.$$

- A_{session} (User Session Summary): \mathcal{Q} over $\Omega^{\text{SEM}} \cup \Omega^{\text{EPI}}$ to form \mathbf{u}_t .
- A_{journey} (Journey Detector): tags live intent from clickstream (e.g., “holiday décor upgrade”).
- A_{rec} (Recommender): returns $\mathbf{r} \in \mathcal{S}^L$ given $(\mathbf{u}_t, \text{intent})$.
- A_{expl} (Explanation Writer): $e_t = \Phi_{\psi}(\mathbf{r}, \mathbf{u}_t)$.
- A_{eval} (Explanation Evaluator): verifies factuality & brand; issues revise if needed.

Environment \mathcal{E} : ProductDB, UserProfileDB, BrandPolicy, LogStore.

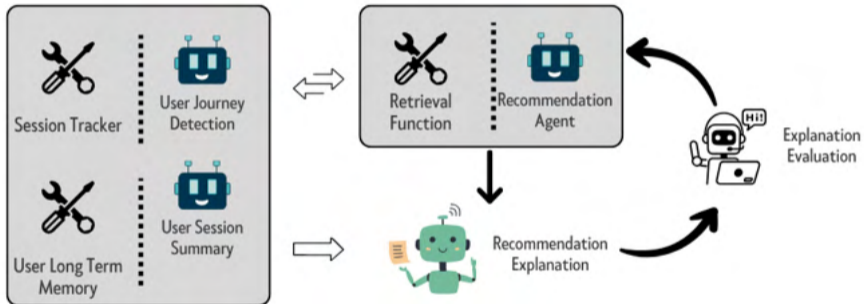
Allowed edges ($C_{ij} = 1$):

$$A_{\text{session}} \rightarrow A_{\text{rec}}, \quad A_{\text{journey}} \rightarrow A_{\text{rec}}, \quad A_{\text{rec}} \rightarrow A_{\text{expl}} \rightarrow A_{\text{eval}} \rightarrow A_{\text{rec}}.$$

Interaction loop:

- 1 A_{session} : $\mathbf{u}_t \leftarrow \mathcal{Q}(\Omega^{\text{SEM}} \cup \Omega^{\text{EPI}}, \tau)$
- 2 A_{journey} : infer intent from clicks.
- 3 A_{rec} : produce \mathbf{r} conditioned on $(\mathbf{u}_t, \text{intent})$.
- 4 A_{expl} : generate e_t .
- 5 A_{eval} : check Factual & $\mathcal{C}_{\mathcal{P}}$; if fail \Rightarrow revise \rightarrow (3–4).

Architecture



Formal Consistency Checks

Factuality (NLI/grounding):

$$\text{Factual}(e_t : \mathbf{r}, \mathbf{u}_t) = \mathbf{1}[\forall \text{fact} \in e_t : \text{fact} \in \text{Facts}(\mathbf{r}) \cup \text{Facts}(\mathbf{u}_t)]$$

Brand/style compliance:

$$\mathcal{C}_{\mathcal{P}}(e_t) = \begin{cases} 1, & \text{all rules in policy } \mathcal{P} \text{ satisfied} \\ 0, & \text{otherwise} \end{cases}$$

Joint predicate: $\text{Consistency}(e_t) = \text{Factual}(e_t) \wedge \mathcal{C}_{\mathcal{P}}(e_t)$.

Optional score:

$$J(e_t) = \lambda_1 \text{Readability}(e_t) + \lambda_2 \text{Specificity}(e_t) - \lambda_3 \text{Redundancy}(e_t),$$

maximized subject to $\text{Consistency}(e_t) = 1$.

Memory (by agent):

- A_{session} : Ω^{SEM} (stable prefs), Ω^{EPI} (recent sessions), Ω^{STM} (current turns).
- A_{expl} : Ω^{PROC} (brand tone exemplars, banned phrases, templates).
- A_{eval} : cached fact table for \mathbf{r} ; policy store \mathcal{P} .

Tools \mathcal{F} :

- `UserData.fetch`, `ProductDB.lookup` (explanation and item grounding).
- `NLI.verify` (trace facts).
- `PolicyCheck` (brand/legal), `StyleTransfer` (tone repair).


- **Trust & CTR uplift:** clear, user-aware rationale improves acceptance.
- **Brand safety at scale:** evaluator veto enforces \mathcal{P} with low integration cost.
- **Personal continuity:** explanations reuse the same $\Omega^{\text{SEM}}/\Omega^{\text{EPI}}$ driving ranking.
- **Modularity:** new style guides or fairness rules added via Ω^{PROC} and A_{eval} prompts.

Implementations

Intent Base Recommendation Model-post add to cart explore page

Image Credit: **Walmart.com**

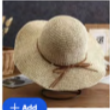



✓ Added to cart!



MOSHU Ribbed Trim Tank Tops for Women Flowy Round Neck Women Shirts Loose Fit Sleeveless Summer Tops
Now \$10.00 ea ~~\$19.99~~

— 1 +


Stay Stylish in the Sun with Summer Essentials
Generated by AI

Best seller	Best seller	Best seller	Best seller
			
\$6.80 Set Of 4,Silicone Stainless Steel Food Tongs-Black ★★★★☆ 53 Shipping, arrives in 3+ days	\$14.29 50% off Summer! TMOYZQ Wedge Sandals for Women, Summer Closed Round To... ★★★★☆ 31 Shipping, arrives in 3+ days	Now \$12.97 \$17.00 Joopin Oversized Polarized Sunglasses for Women Vintage Lady UV Protectio... ★★★★☆ 311 Save with W+ Shipping, arrives in 2 days	Now \$13.99 \$19.99 BCOOSS Summer Sun Hat for Women Wide Brim Sun Protection Women Straw... ★★★★☆ 317 Save with W+ Shipping, arrives in 2 days

Intent Base Recommendation Model-post add to cart explore page





Image Credit: [Walmart.com](https://www.walmart.com)

Added to cart



6V Huffy Disney Minnie Mouse Battery-Powered Ride-On Car, Kids Ages 3+ - Pink
\$174.00 ea ~~\$174.00 ea~~

Explore Fun Minnie Mouse Gear for Kids
Generated by AI





Best seller	Best seller		Best seller
			
+ Add	+ Add	+ Add	+ Add
\$31.10 Bell Disney Minnie Mouse Bike Helmet, Pink Flowers, Toddler 3 (48-52cm) ★★★★★ 290 Save with W+ Shipping, arrives in 3+ days	\$5.97 Disney Minnie Mouse Girl's Brow Bar Sunglasses Pink ★★★★★ 432 Save with W+ Pickup tomorrow Delivery today Shipping, arrives tomorrow	\$16.95 Bell Disney Minnie Mouse Protective Pad and Glove Set ★★★★★ 106 Save with W+ Shipping, arrives in 3+ days	\$5.97 Minnie Mouse Girl's Pink Heart Sunglasses ★★★★★ 903 Save with W+ Pickup tomorrow Delivery today Shipping, arrives tomorrow

Explain It to Me-item page

Cozy Home Scented Candles for Ambiance

Based on what customers bought





(A)

 + Add	Best seller  + Add	Best seller  + Add	Best seller  + Add
\$218 42.4 ct/oz Mainstays Beachside Linen Scented 3 Wick Candle, 11.5 oz ★★★★☆ 15 Save with WU Shipping, arrives in 3+ days	\$3.96 34.4 ct/oz Mainstays Garden Rain 3 Wick Candle, 11.5 oz ★★★★☆ 19 Save with WU Shipping, arrives in 3+ days	\$3.96 34.4 ct/oz Mainstays Fall Farmhouse 3 Wick Candle, 11.5 oz ★★★★☆ 20 Save with WU Shipping, arrives in 3+ days	\$3.96 34.4 ct/oz Mainstays Vanilla Scented 3 Wick Glass Jar Candle, 11.5 oz ★★★★☆ 408 Save with WU Shipping, arrives in 2 days

All-Season Performance and Durability Tires

Based on what customers bought

(B)

 + Add	Reduced price  + Add	 + Add	 + Add
\$178.88 Goodyear Wrangler Fortitude HT 255/65R17 10T All-Season Tire ★★★★☆ 203 Shipping, arrives in 3+ days	Now \$207.88 \$210.99 Goodyear Wrangler All-Terrain Adventure 255/65R17 10T All-Terrain Tire ★★★★☆ 220 Shipping, arrives in 3+ days	\$767.96 4 New Goodyear Wrangler Fortitude HT All-Season Tires - 255/65R17 10T Fits 2004-08... ★★★★☆ 1 Shipping, arrives in 3+ days	\$224.99 Goodyear Wrangler SteadFast HT All-Season 255/65R17 10T Light Truck Tire ★★★★☆ 1 Save with WU Shipping, arrives in 2 days

Explain It to Me-item page

Smartphone Trio: Power, Performance, Style

Based on what customers bought

In 200+ people's carts



+ Add

\$59.88

Cricket Wireless Moto G Stylus 2023, 128GB, 4GB RAM, 8MP FF Camera, Blue - Prepaid...

★★★★☆ 136

Save with W+

Pickup available
Delivery available
Shipping, arrives in 2 days



+ Add

Sponsored

\$199.00

Straight Talk Motorola Moto G Stylus 5G (2024), 128GB, Beige - Prepaid Smartphone (Locked...

★★★★☆ 22

Save with W+

Shipping, arrives in 3+ days



+ Add

\$59.88

Total by Verizon Motorola Moto G Stylus 4G (2023), 64GB, Blue - Prepaid Smartphone (Locked...

★★★★☆ 113

Save with W+

Pickup available
Delivery available
Shipping, arrives in 3+ days



+ Add

Now \$49.88 ~~\$199.99~~

AT&T Motivate Max 32GB, Celestial Blue - Prepaid Smartphone

★★★★☆ 101

Save with W+

Pickup available
Delivery available
Shipping, arrives in 2 days

(C)

Luxurious Bedding for Ultimate Comfort and Style

Based on what customers bought



+ Add

\$12.00

1500 Thread Count Hospitality Fitted Sheet, Queen Size, Gray

★★★★☆ 704

Save with W+

Shipping, arrives in 2 days



Options

Sponsored

Now \$18.99 ~~\$24.99~~

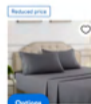
Options from \$18.99 - \$38.99

Shiucheng Cooling 4 Piece Luxury Bed Sheets Set, 1800 Series Microfiber Bed Sheets...

★★★★☆ 263

Save with W+

Shipping, arrives in 2 days



Options

• • • • • +10

Now \$17.99 ~~\$29.99~~

More options from \$16.99

Lux Decor Collection Twin Sheets Set, Deep Pocket 4 Pc Bed Sheets Set - Wrinkle, Fade...

★★★★☆ 167

Extra savings available

Shop it

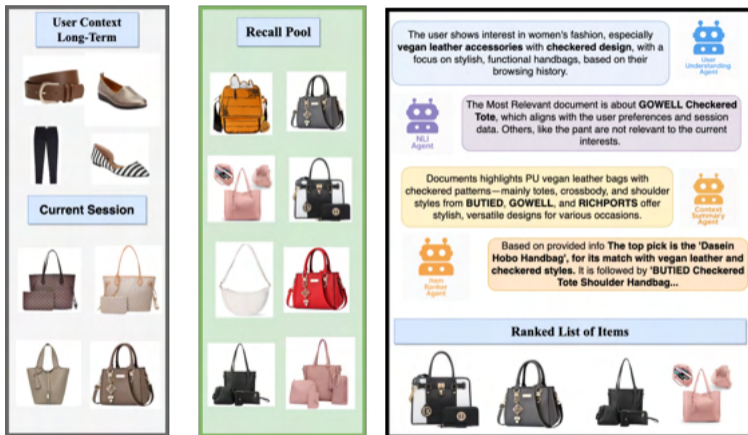
Shipping, arrives in 3+ days

(D)

Table: A/B testing evaluation results. The lift results are stat sig.

	CTR	GMV
Lift	>50pbs	>100pbs

Agentic RAG for Personalized Recommendation



MetaSynth for SEO

Credit: **Google.com**



Walmart

<https://www.walmart.com> > ... > Girls Tank Tops

Justice Girls Side Tie 2FER Tank, Sizes XS-XLP

This tank top is a size XS 5-6 and fit her perfect. It comes with a sports bra that is bright and colorful. The shirt is white with a matching colorful ...



Walmart

<https://www.walmart.com> > ... > Girls Tank Tops

Justice Girls Side Tie 2FER Tank, Sizes XS-XLP - Walmart.com

Get ready for summer with the Justice Girls Side Tie 2FER Tank! This colorful 2-piece set is perfect for sports practice, family outings, or lounging.



Table: A/B testing evaluation results. The lift results are stat sig.

	CTR	Traffic
Lift	+1026 bps	+7510 pbs

Search Compare Eval-MAREval

Credit: **Walmart.com**



Good for high-speed computing

Sleek • High-performance

GOOD



Apple MacBook Pro 15" Touchbar - Intel Core i7 2.9GHz - 16G...

\$699.00

★★★★☆ 4.5 | 11

Free shipping, arrives in 2 days

Table: online A/B test results

Metrics	Improvement	P-Value
ATC	118 bps	0.05
GMV	136 bps	< 0.01

Personalization Team

Acknowledgment



MCP in Multi-Agent RecSys: Why It Matters

- Agents must exchange **user context**, **item signals**, and **intermediate results** reliably.
- A robust **Model Context Protocol (MCP)** ensures message clarity, task negotiation, and state sync.
- Weak/adhoc protocols \Rightarrow *bottlenecks*, *misinterpretations*, and *integration debt*.
- RecSys adds pressure:
 - **low latency**,
 - **high update rate**,
 - **clear semantics-in explanasion usecases.**

Protocol Standardization Problem & Goals

Snapshot: FIPA.org

- Goal: **plug-and-play** agents across teams/vendors via a shared syntax & semantics. A standard communication protocol.
- Legacy lesson (e.g., FIPA ACL): semantics help, but **complexity can hinder adoption**.

FOUNDATION FOR INTELLIGENT PHYSICAL AGENTS

FIPA ACL Message Structure Specification

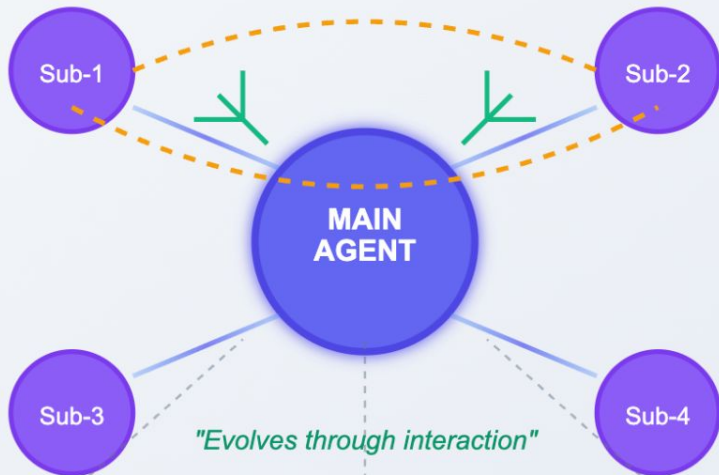
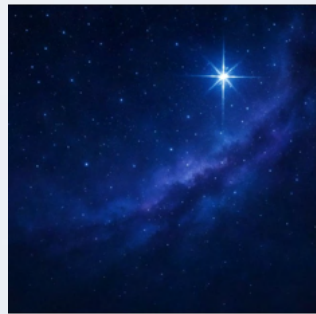
Document title	FIPA ACL Message Structure Specification		
Document number	SC00061G	Document source	FIPA TC Communication
Document status	Standard	Date of this status	2002/12/03
Supersedes	None		
Contact	fab@fipa.org		
Change history	See <i>Informative Annex A – ChangeLog</i>		

Frontiers of Agentic AI

RecSys'25 Tutorial
Agentic Recommendation Systems

✕ [@Chi_Wang_](#)





2023: AutoGen (AG2)
Foundation laid

2025: Industry Progress
What remains?

North Star: AI Agent That Grows With You

NATURAL INTERFACE
Sight • Sound • Context

STRONG CAPABILITY
Code • Action • Learning

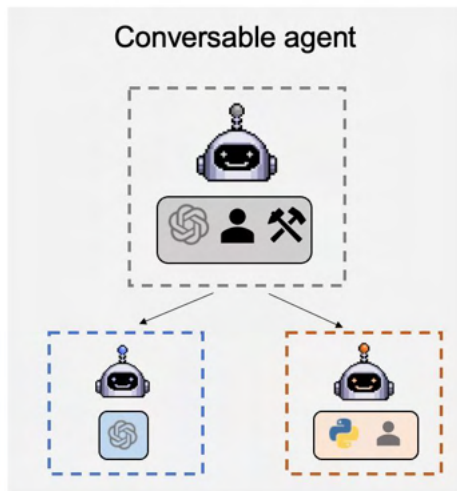
SCALABLE ARCHITECTURE
Multi-Agent • Parallel

AG2: Open-Source AgentOS

- Simplicity: Simplify developers' thought process
 - **Intuitive unified agentic abstraction**
- Capability: Enable advanced exploration & diverse needs
 - **Flexible multi-agent orchestration**
- Reusability: Minimize unnecessary complexity
 - **Composable agentic design patterns**

Agentic Abstraction

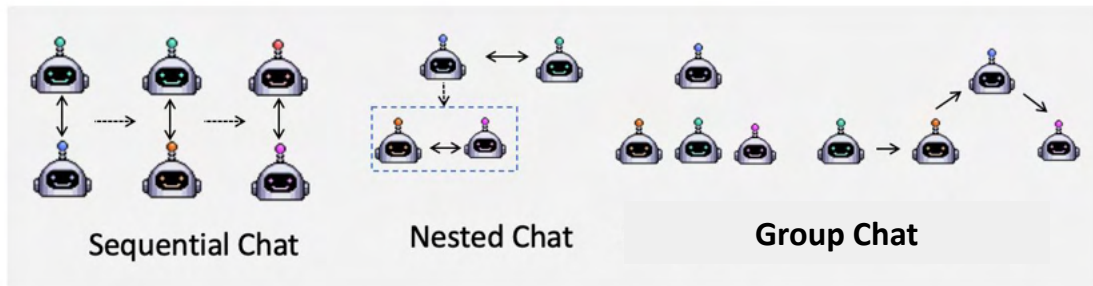
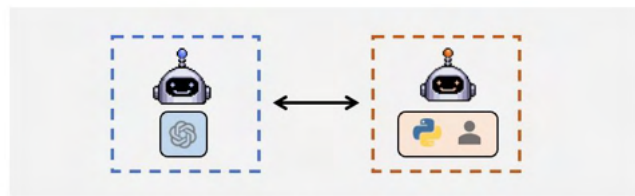
Step 1. Define AG2 agents:
Conversable & Customizable



Agent Customization

Multi-Agent Orchestration

Step 2. Get them to talk:
Conversation Programming

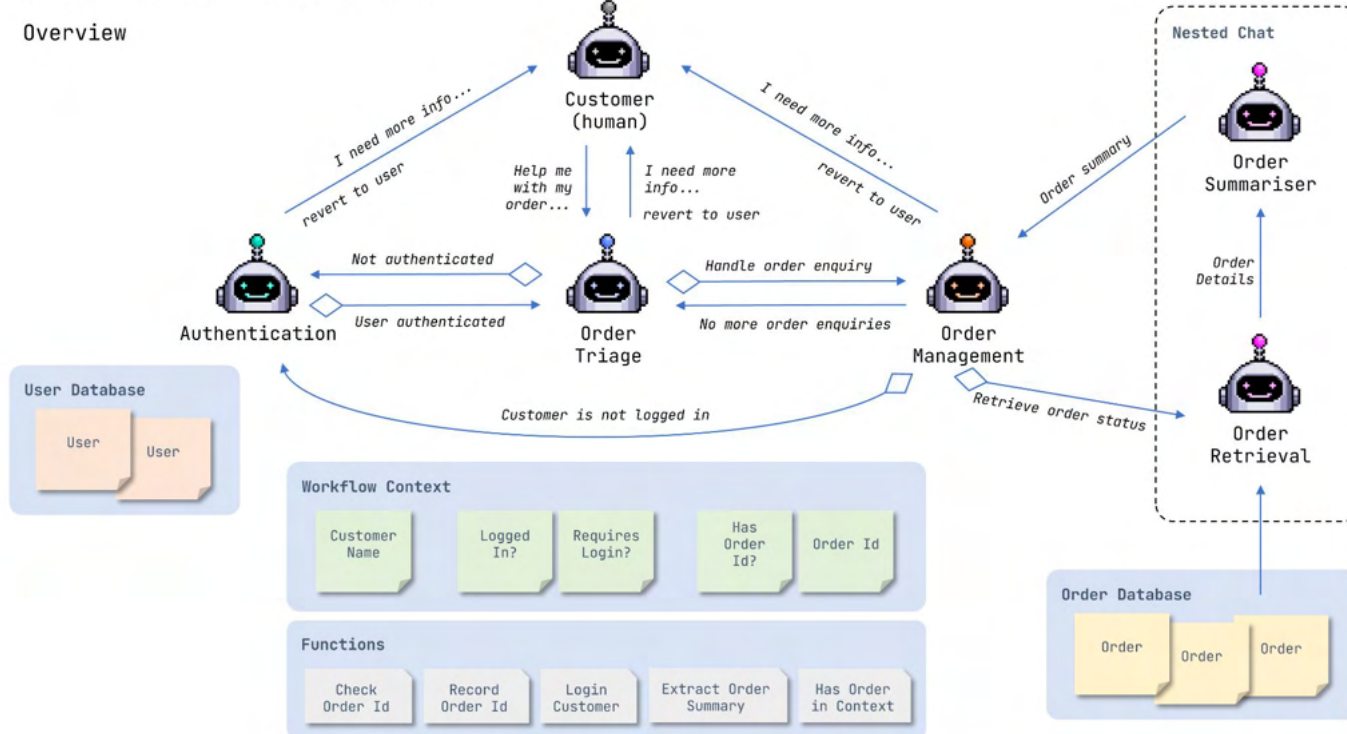


Flexible Conversation Patterns

Multi-Agent Orchestration Example

Customer Service

Overview





Top 100
Open source
achievements

5K
Forks

Forbes, The
Economist,
WIRED...

40K
Stars

08.2023: Research
paper

10.2023: Top
trending on GitHub

12.2023:
5 favorite AI
papers by The
Sequence

03.2023: Initial Prototype
- Flexible multi-agent
conversation framework
- Code/function execution

AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation Framework

Qingyun Wu¹, Gagan Bansal², Jerry Zhang², Yixia Wu¹, Shaohan Zhang¹, Erkang Zhu¹, Bozhi Li¹, Li Jiang², Xiaoyan Zhang², and Chi Wang²

¹Pennsylvania State University
²Microsoft
³University of Washington



Figure 1: AutoGen enables complex LLM-based workflows using multi-agent conversations. (Left) Multiple agents are customizable and can be based on LLMs, tools, humans and even a combination of them. (Top-right) Agents can converse to solve tasks. (Bottom-right) The framework supports many additional complex conversation patterns.

Abstract

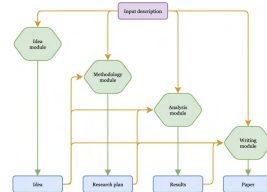
This technical report presents AutoGen, a new framework that enables development of LLM application using multiple agents that can converse with each other to solve tasks. AutoGen agents are customizable, conversable, and extensible other human participants. They can operate in various modes that employ combinations of LLMs, human inputs, and tools. AutoGen's design offers multiple advantages: (1) it generally requires the strong task-specific operations and reasoning abilities of these LLMs; (2) it leverages human understanding and intelligence, while providing reliable execution through conversation between agents; (3) it enables and scales the implementation of complex LLM workflow; (4) it allows human developers can easily integrate their coding, mathematics, open-domain knowledge, etc.

1 Introduction

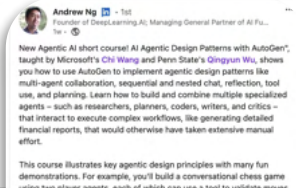
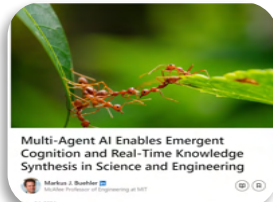
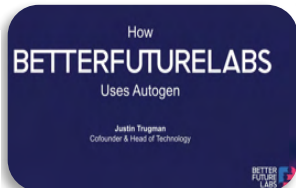
Large Language Models (LLMs), like GPT-4, GPT-4o, and others, are becoming a major part of AI applications.

https://arxiv.org/abs/2308.08155v1

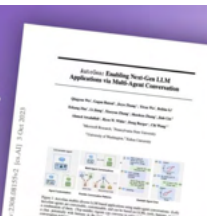
Denario: Modular Automation of Research



<https://github.com/AstroPilot-AI/Denario>



Best Paper
at ICLR 2024
LLM Agents
Workshop



Example Production Use Cases

Agent Platform (Google)

Autonomous Trading

Business Automation (Cegid)

Chip Design (Nvidia)

Customer Support (Parker Hannifin)

Cyber Security

Data Engineering (Nexla)

Farming

Investment Research (BFL)

Marketing (Walmart)

Patient Support

Recommendation (Walmart)

Software Engineering

Software Testing

Task Management

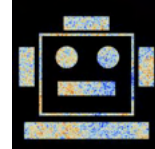
Web Automation (Emergence)

Automation of Research Workflows with AI Agents

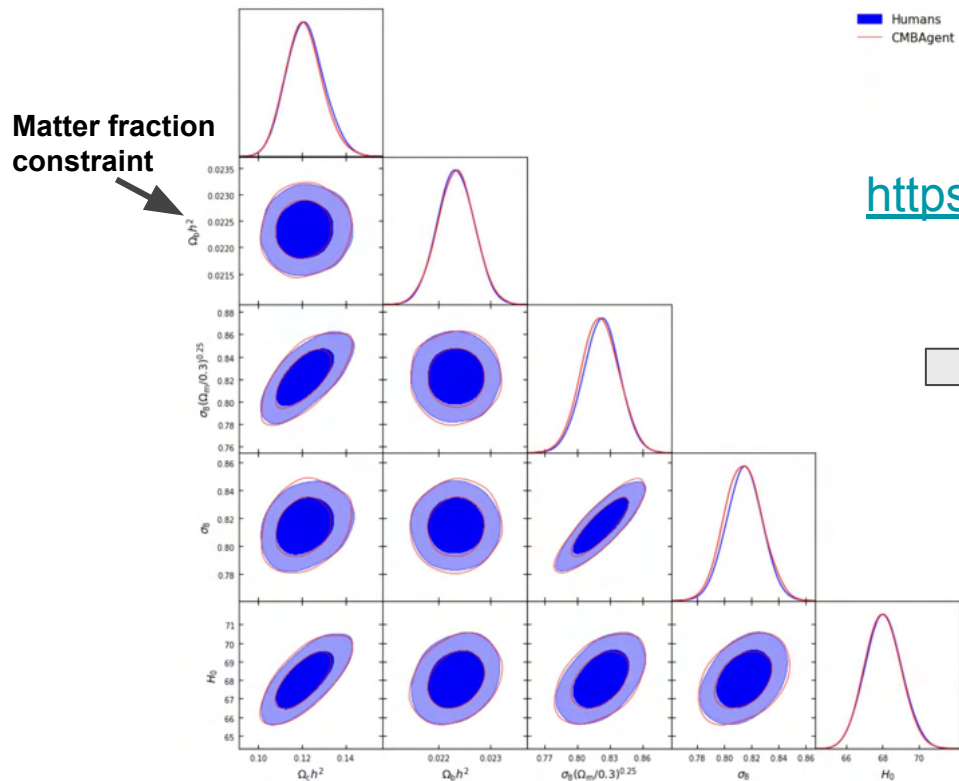
Boris Bolliet
University of Cambridge



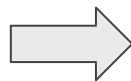
Cmbagent



Agents reproduces the lensing analysis, with minimal human input



<https://github.com/CMBAgents/cmbagent>



*Analysis fully reproduced in ~10'
Pipeline entirely written by cmbagent*

*Would have taken ~ a full day of work
to an experienced cosmologist.*

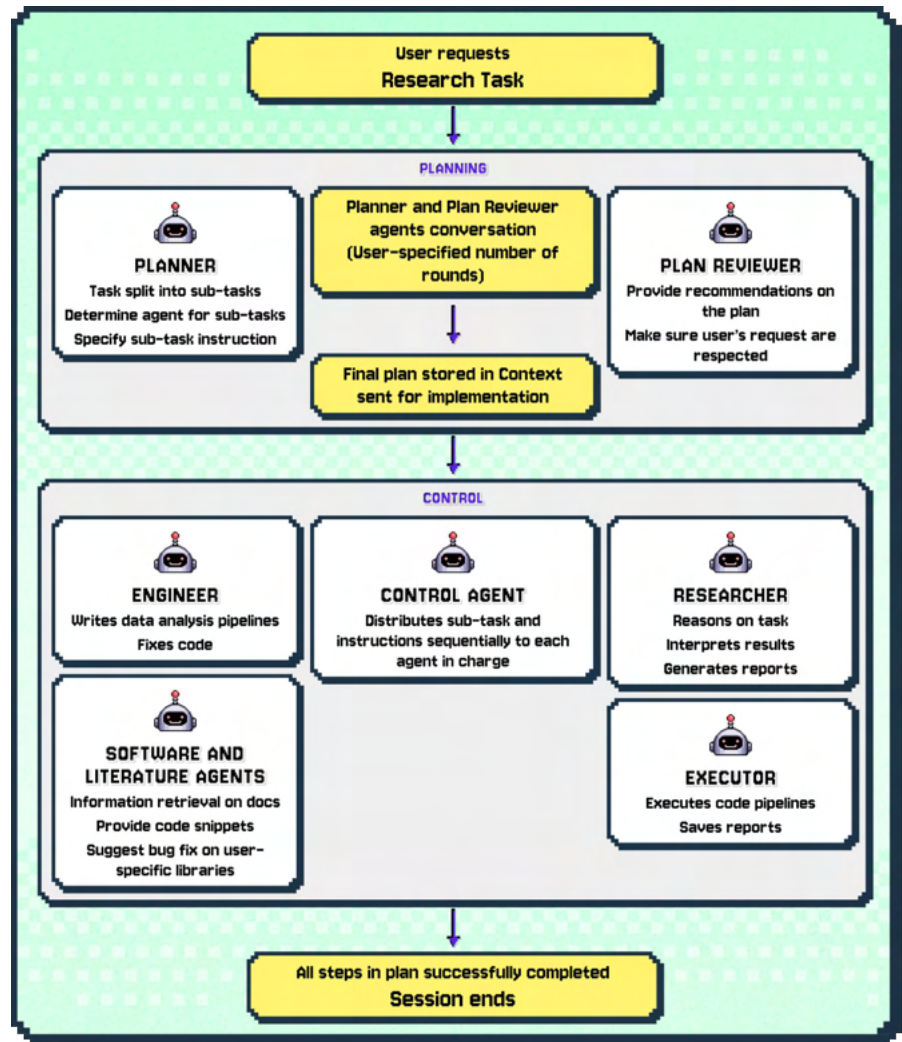


Hubble parameter constraint



Features That Helped:

- 🤖 Tool Calls
- 🤖 Multi-LLM Calls
- 🤖 Group Chat
- 🤖 Structured Output
- 🤖 Agentic RAG



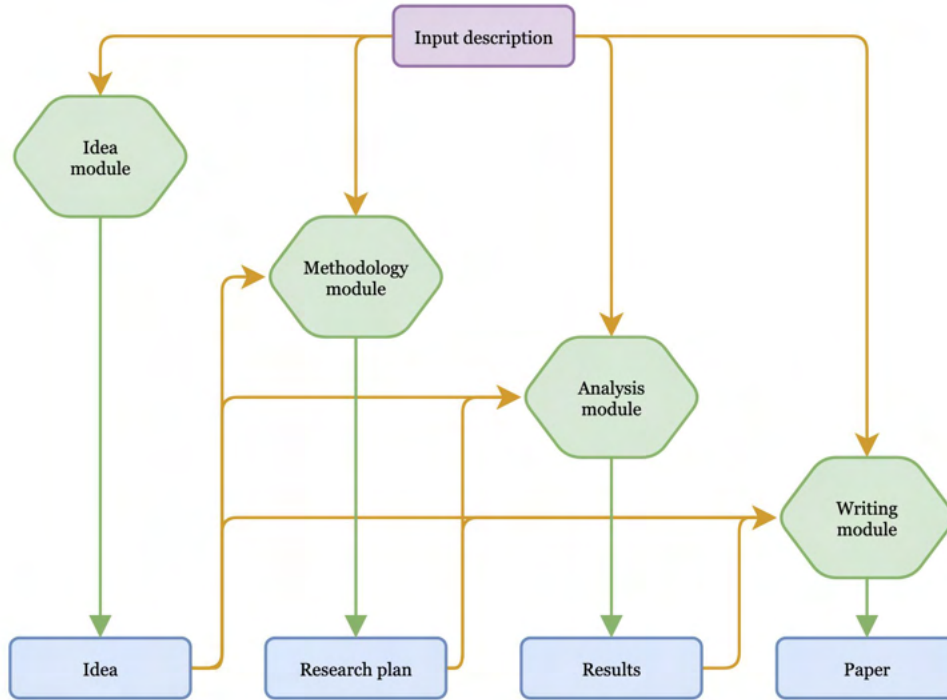
The Denario Project: Modular Automation of Scientific Research with Multi-Agent Systems



SCAN ME

Francisco Villaescusa-Navarro^{1,2,*}, Boris Bolliet^{3,4,*}, Pablo Villanueva-Domingo^{5,*},
Adrian E. Bayer^{1,2}, Aidan Acquah⁶, Chetana Amancharla⁷, Almog Barzilay-Siegal⁸,
Pablo Bermejo^{9,10,11}, Camille Bilodeau¹², Pablo Cárdenas Ramírez^{13,14,15}, Miles Cranmer¹⁶,
Urbano L. França^{17,18}, ChangHoon Hahn^{19,20}, Yan-Fei Jiang¹, Raul Jimenez^{21,22},
Jun-Young Lee¹, Antonio Lerario²³, Osman Mamun¹³, Thomas Meier²⁴,
Anupam A. Ojha^{25,26}, Pavlos Protopapas²⁷, Shimanto Roy¹², Pedro Tarancón-Álvarez^{21,28},
Ujjwal Tiwari⁷, Matteo Viel^{23,29,30,31,32}, Digvijay Wadekar³³, Chi Wang³⁴,
Bonny Y. Wang³⁵, Licong Xu^{36,4}, Yossi Yovel^{8,37}, Shuwen Yue¹³, Wen-Han Zhou³⁸,
Qiyao Zhu²⁵, Jiajun Zou³⁹, Íñigo Zubeldia^{36,4}

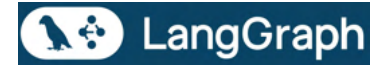
Denario: Modular Automation of Research



Multi-agent orchestration
with



<https://ag2.ai/>



Deep Research
with

CMBAGENT

<https://github.com/CMBAgents/cmbagent>

<https://github.com/AstroPilot-AI/Denario>

Community Talks

Journey of taking
Testzeus-Hercules to production...

Interactive Speculative Planning

Globant Code Fixer Agent: #1 on
SWE-Bench Lite

Agentic Commerce: Enabling AG2
with identity, payment and
monetization

Building Multi-Agent Systems for
Investment Analysis

Introducing FastAgency - the
fastest way to bring AutoGen
workflows to production

Agent-Model Orchestration in
Multi-Agent Applications

Copilot Agent Architecture
Designing

Physics-Aware AI: Bridging
Science Through Multi-Agent
Systems

Trace-ing the Path to
Self-adapting AI Agents

NOVA: Building Agentic Workflows
for Structured Data Intelligence at
Nexla

Investigating Group
Decision-Making Mechanism in
Decentralized Multi-Agent
Collaboration

Multi-AI Agents for Chip Design

Maris: A Security Controlled Development Paradigm for
Multi-Agent Collaboration Systems

From Content to Conversations:
Building Social Media & WhatsApp
Automation with AG2

Make AI Agents Collaborate: Drag,
Drop, and Orchestrate with
Waldiez

Crypto Trading with AG2

Frontiers of LLM Agents: Memory,
Tool Use, Multi-Modal Input, and
RL with LLMs

Pedagogy Support using AG2

Exploring Pragmatic Patterns in
Agentic Systems

A Multi-Agent Approach to
Podcast Generation with AG2

Example Projects for Learning

<https://github.com/ag2ai/build-with-ag2>





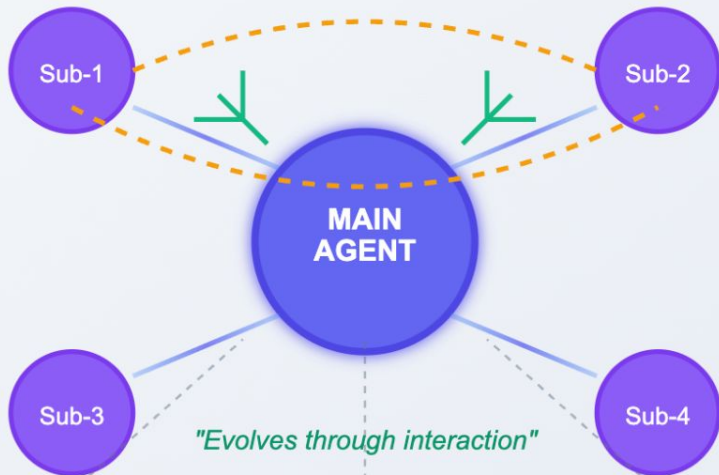
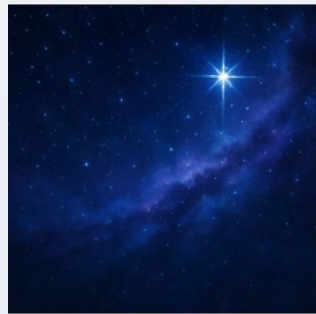
THE OPEN-SOURCE AGENTOS

GitHub 

Join our growing community of over 20,000 agent builders

Discord 





2023: AutoGen (AG2)

Foundation laid

2025: Industry Progress

What remains?

North Star: AI Agent That Grows With You

**NATURAL
INTERFACE**

Sight • Sound • Context

**STRONG
CAPABILITY**

Code • Action • Learning

**SCALABLE
ARCHITECTURE**

Multi-Agent • Parallel

Natural Interface – Computers with Senses



Today:

- Gemini Live photo/video guidance
- Screen-aware prototypes
- Real-time interactions



Missing: Context Flow Across Devices

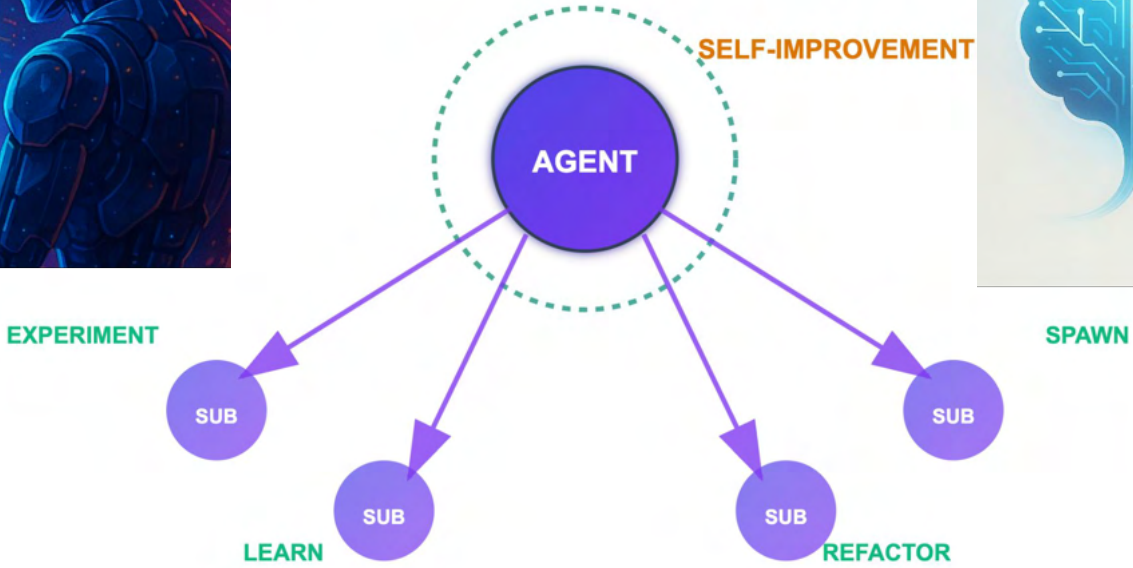


Challenge: Know When to Engage vs. Step Back

Intelligent timing + Seamless context = Natural interface

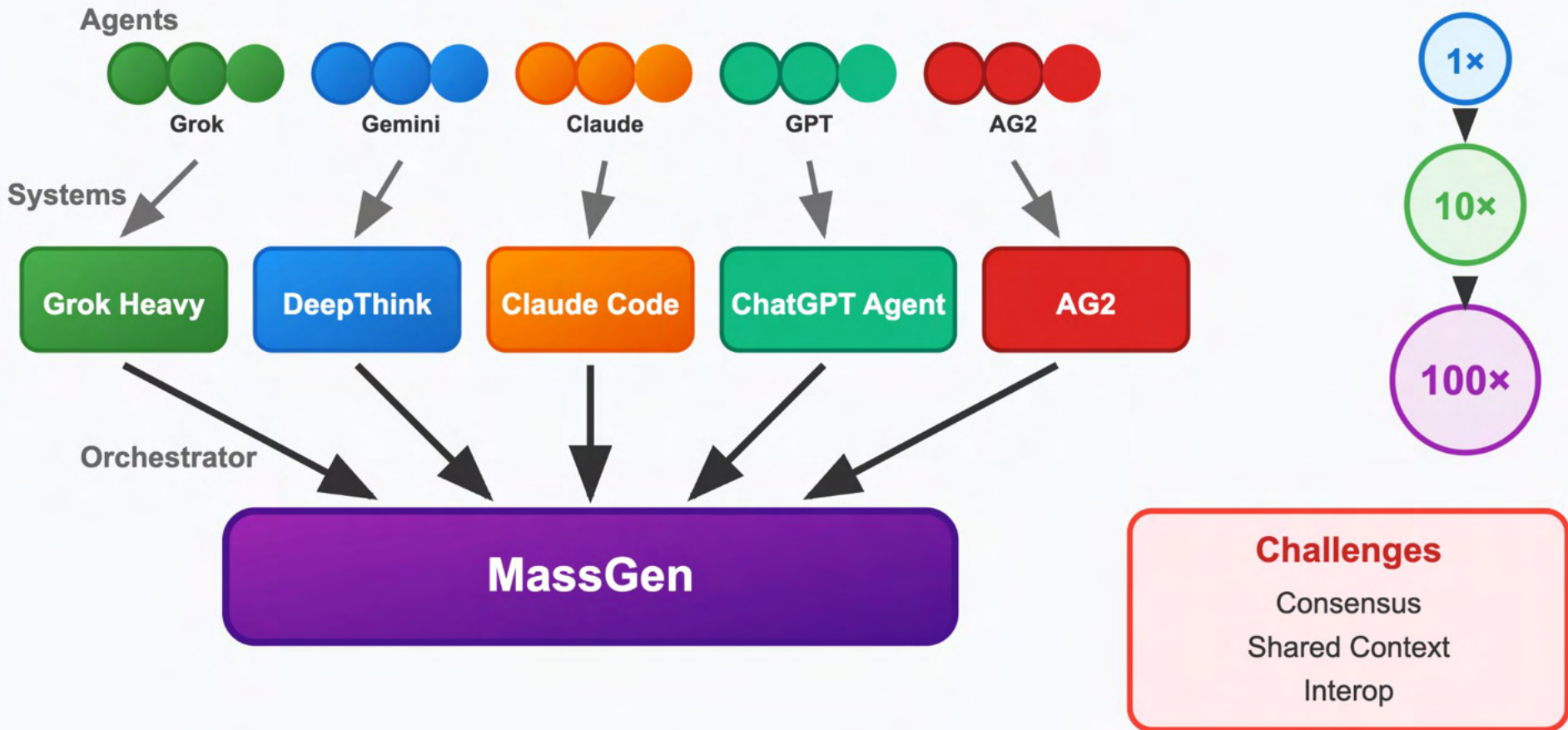
CODE IS META-POWER

Claude Code • Gemini CLI • Codex • ChatGPT Agent



- To Realize This Vision:**
- Long-term Memory
 - Cross-session Reflection
 - Learning by Interaction

Many Minds Working as One



Building the Future Together



Natural interfaces

Turn intention into understanding



Strong capabilities

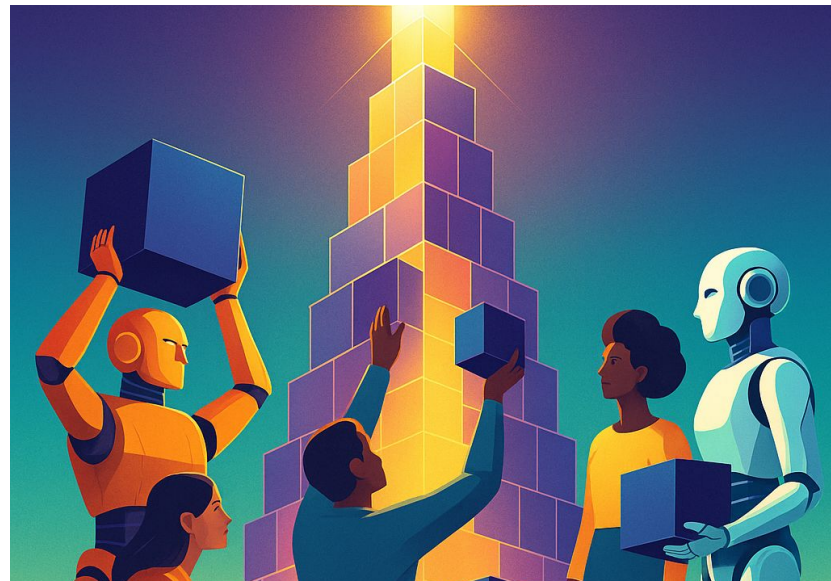
Turn understanding into action & self-improvement



Scalable architectures

Enable co-creation across many agents

AI that grows with us




X @Chi_Wang_


MassGen



Scaling AI Through Multi-Agent Collaboration

 **RecSys'25 Tutorial**

Agentic Recommendation Systems

 Prague, Czech Republic • September 26, 2025

 massgen.ai | [GitHub](https://github.com)



How Do We Scale Up AI?



Traditional Scaling Laws Hit Limits



Power Crisis

4-16 GW by 2030

Enough to power entire cities



Data Depletion

Depleted by 2026-2028

Quality text data exhausted



Performance Plateau

GPT-5 delayed >1 year

Early training runs failed

⚠️ Inference-time scaling

No universal way to leverage improvements & address limitations

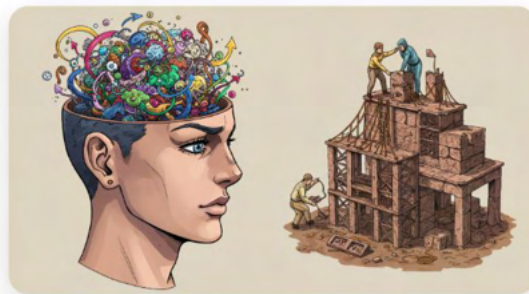
The New Paradigm:

Model → Agent → **Multi-Agent Systems**



The Promise of Multi-Agent Collaboration

- **Study Group Dynamics:** Like humans collaborating on complex problems
- **Cross-Ecosystem Integration:** Bridge Claude, Gemini, GPT, Grok, and specialized coding agents
- **Emergent Intelligence:** Collective problem-solving beyond individual capabilities
- **Real-time Intelligence Sharing:** Agents learn and adapt from each other



The Promise of Collaborative Reasoning

root.massgen.ai - "Myth of Reasoning"

🔧 Built on AG2's foundational multi-agent research and community



Proven Performance Gains

Grok-4 Standard

1

Single Agent Processing

38.6%

Last Human Exam Score
\$30/month

Grok-4 Heavy

A1

A2

A3

Multi-Agent Collaboration

44.4%

Last Human Exam Score
\$300/month

Gemini 2.5 DeepThink



Competition Gold Medals

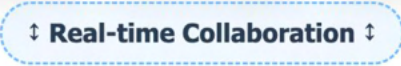
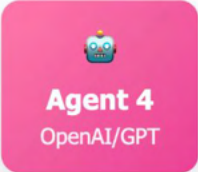
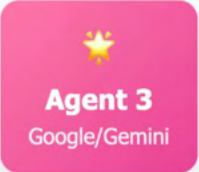
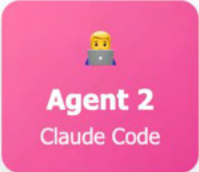
IMO + ICPC

5/6 IMO Problems (2024)
10/12 ICPC Problems (2025)

First AI Gold Medals

Multi-Agent Revolution

"Individual AI excellence + Multi-agent coordination = Next frontier of AI capabilities"





Key Features & Capabilities



**Multi-Backend
Support**



Iterative Refinement



The Reality of Reasoning



**Parallel
Processing**



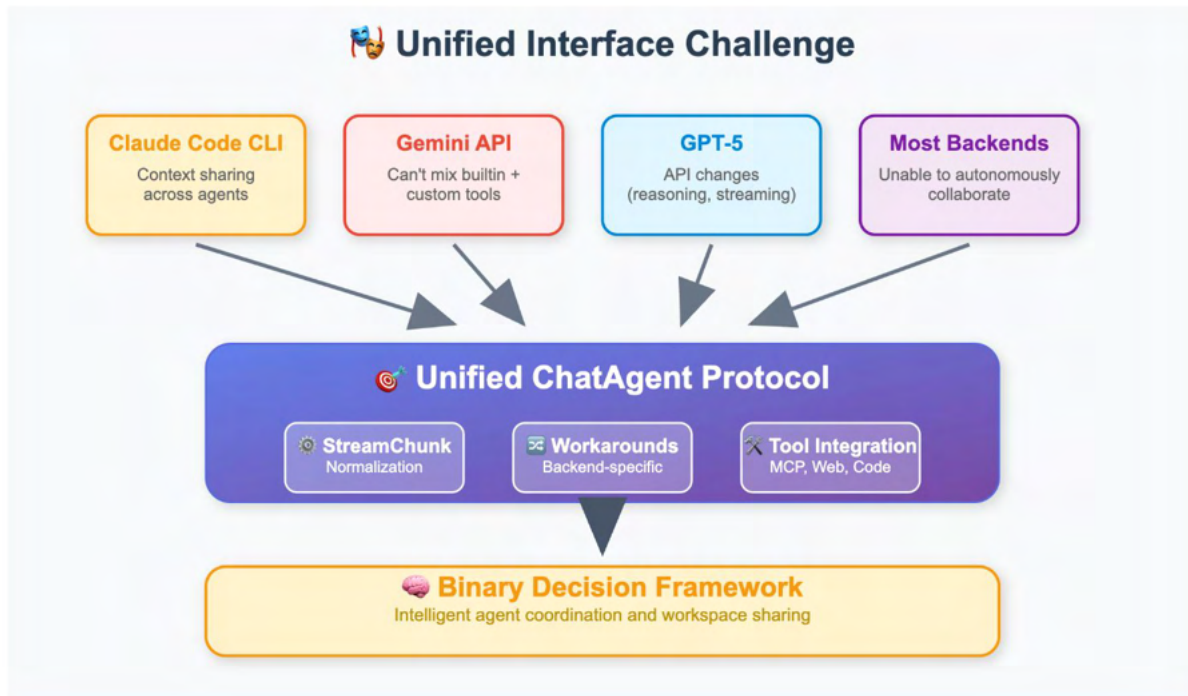
**Intelligence
Sharing**



**Consensus
Building**

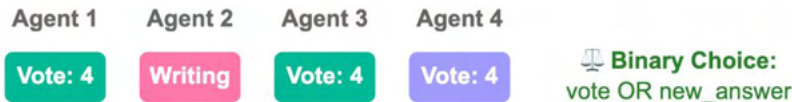


Tech Deep Dive: Backend Abstraction Challenges





Tech Deep Dive: Binary Decision Framework Solution



Agent 2 provides new_answer
→ **ALL VOTES RESET**

After Reset:



Async Results:



→ If new_answer: reset again

 **Key Innovation: Vote Invalidation Creates Dynamic Consensus**



Case Study: Success Through Peer Correction

Graduate-level physics question from GPQA-Diamond benchmark

The Problem

A quasar shows a peak at 790 nm wavelength. Given Lambda-CDM cosmological parameters ($H_0 = 70$ km/s/Mpc, $\Omega_m = 0.3$, $\Omega_\Lambda = 0.7$), what is the comoving distance?

Options: A) 8 Gpc B) 7 Gpc C) 6 Gpc D) 9 Gpc

Final Result



Correct Answer: A (8 Gpc)

Orchestration succeeded where individual agents initially failed

Round 1: Initial Answers

Claude: "I calculate ~6 Gpc → Answer C"

GPT-5: "I get ~8.95 Gpc → Answer D"

Gemini: "~6.1 Gpc → Answer C"

Self-Correction Process

Claude observes: "There is significant discrepancy in calculations: Agent1 gets ~6.1 Gpc, Agent2 gets ~8.95 Gpc. Let me re-examine..."

Breakthrough Moment

Claude revises: "Standard cosmological calculators yield 8000-8500 Mpc for $z=5.5$. This equals 8.0-8.5 Gpc, closest to option A."

Result: 3/4 agents converge on correct answer

Success Mechanism:

Peer observation → Discrepancy detection → Self-correction → Consensus



Benchmarking: Preliminary Results

Scientific evaluation across graduate-level reasoning, instruction-following, and narrative tasks

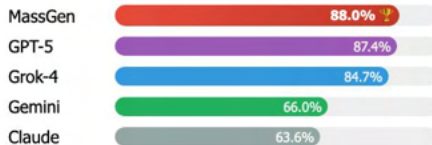
GPQA-Diamond

Graduate Physics/Chemistry



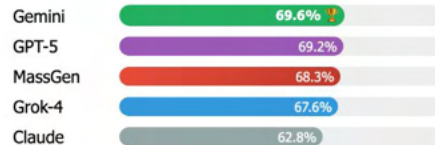
IFEval

Instruction Following



MuSR

Narrative Reasoning



 Overall Champion

MassGen: 81.2%

Wins 2/3 benchmarks • Statistically significant

Key Results:

- Highest on 2/3 benchmarks
- Best overall average
- Consistent performance

Statistical:

- vs Claude: $p = 1.4e-07$ ★★★★★
- vs Gemini: $p = 1.1e-28$ ★★★★★
- Not due to chance

Research Gap:

- Oracle: 95.5% (GPQA)
- Actual: 87.4%
- Potential: 8.1 points



Agentic Recommendation Applications

Multi-Agent Personalization Pipeline



Content Analysis

Multiple agents analyze user behavior, content features, and contextual signals simultaneously



Preference Fusion

Collaborative filtering, content-based, and deep learning agents debate optimal recommendations




Dynamic Ranking

Real-time consensus building for personalized item ranking and diversity optimization


Cross-domain expertise, multi-objective optimization, explainable recommendations



Live Demonstrations

 **LLM Fun Facts Website (v0.0.14):** Claude Code agents create interactive websites with enhanced logging and workspace isolation

Result: Conflict-free parallel development with comprehensive versioning

 **Unified Filesystem (v0.0.16):** Cross-backend collaboration between Gemini and Claude Code agents creating educational content with shared workspace management

Result: First-time cross-backend coordination producing comprehensive 25-slide presentations

 **IMO 2025 Winner Research:** Multi-agent fact-checking → unanimous consensus on Google DeepMind victory

Result: Accurate identification despite conflicting information

 **Technical Analysis:** Complex Grok-4 HLE pricing calculation through iterative refinement

Result: Accurate cost estimates through collaborative validation

 case.massgen.ai - Complete Case Studies



Get Started in 60 Seconds

```
# 1. Clone and setup
git clone https://github.com/Leezekun/MassGen
cd MassGen && pip install uv && uv venv

# 2. Configure API keys
cp .env.example .env # Add your API keys

# 3. Run single agent (quick test)
uv run python -m massgen.cli --model gemini-2.5-flash "When is your knowledge up to"

# 4. Run multi-agent collaboration
uv run python -m massgen.cli --config three_agents_default.yaml "Summarize latest news of
github.com/Leezekun/MassGen"
```

✔ Supported Models & Providers

🏢 Major Providers:

Anthropic Claude & Claude Code • Google Gemini • OpenAI GPT • xAI Grok • ZAI GLM

🏠 Local & Extended:

Cerebras • Fireworks • Groq • LM Studio • OpenRouter • Together...

🔧 Advanced Tools

Web Search • Code Execution • MCP Tools • File Operations • Browser Automation • Advanced Permissions



Considerations in Building Practical Agentic AI Recommender Systems

September 2025

Considerations in Building Successful Multi-Agent AI Recommender System

- **Background:** Prior to the introduction of GenAI, customer facing typically had limited world knowledge outside of pretrained embeddings. This limited 1) **dynamic content generation**, and 2) **building an even more powerful ML-driven solutions**. GenAI provided path to solving these. However...

Constraints:

Product Goals

Are we building the product we want?

Brand Guidelines

Ensure the product sounds/feels like Walmart

Improve Business KPI's

We also care about moving the needle

Scalability

Build a product that scales to 100M+ customers

Accuracy

Outcome needs to be factually correct

Balancing these using Out-of-the-box LLM's were difficult.
You also can't context-window your way to your solution!

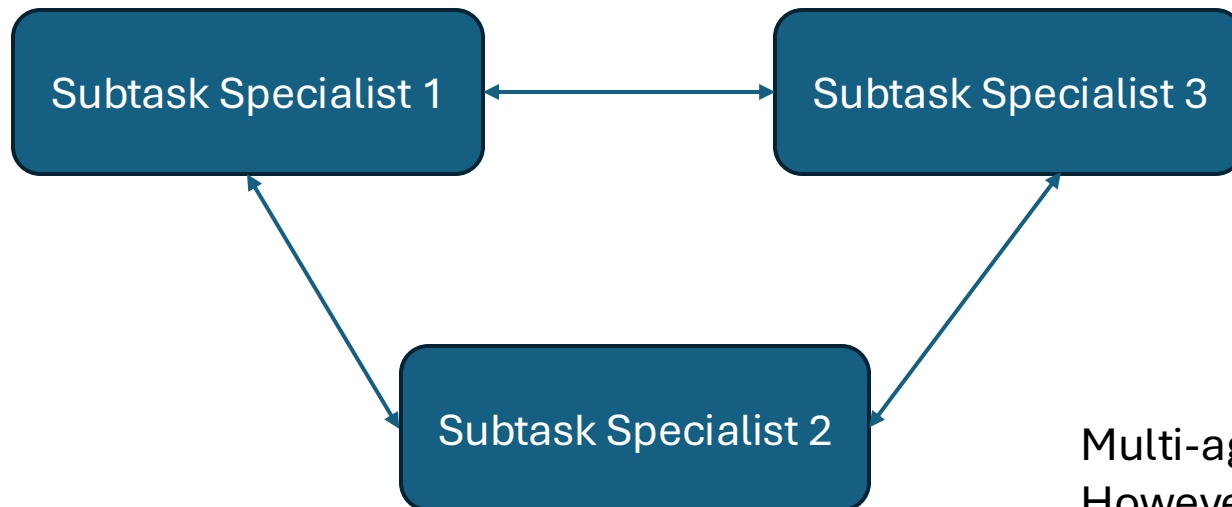


Multi-Agentic Framework: How can Subtask specialists help?

- **Problem Setting:**

System that can a) orchestrate conflicting sub-tasks to b) achieve common task.

Q: What is the a) optimal task order, and b) contexts to send to each subtask specialists to achieve 1) high accuracy, 2) scalability.



- **Must have**

- Flexibility to customize (potentially growing list of) subtasks
- Ability to converge upon satisfying subtask constraints
- Not bound to any specific LLM's

- **Good-to-have**

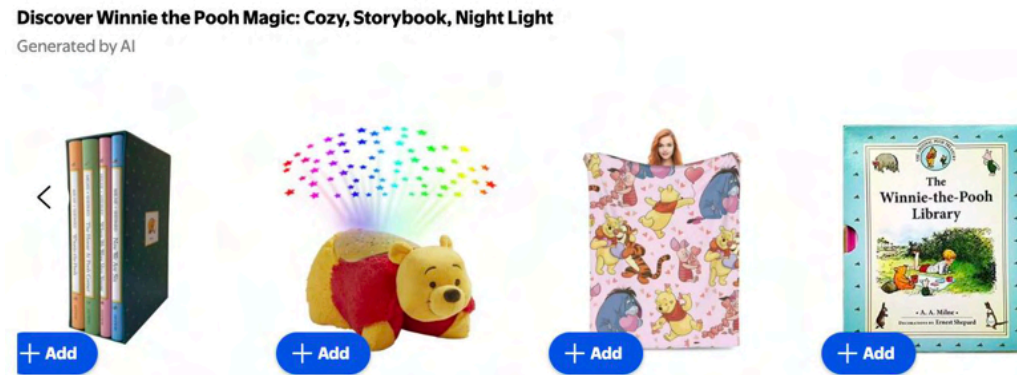
- Scalable (cheaper) solution
- Easy to prototype new agentic architecture

Multi-agentic framework allows for subtask specialists. However, what is the best conversation (orchestration) pattern?

Business Constraints: A multi-faceted optimization Problem!

- **Problem Setting:**

Before embarking on the Agentic System, understand the landscape that these agents are operating on! Legal, Marketing, Creative, and Operation teams can all impose different types of constraints. Identify if these constraints are good to have vs must have.



Legal

- Are we allowed to show/claim this?
- Trademark/Copyright infringement
- Do we need a disclaimer?

Creative

- What is our brand's voice?
 - Friendliness, down-to-earth or formal, etc's
- Ensure consistent tone and grammar

Marketing

- What is the marketing goal?
- What worked before and what did not? Incorporate the findings into the generation process

Operations

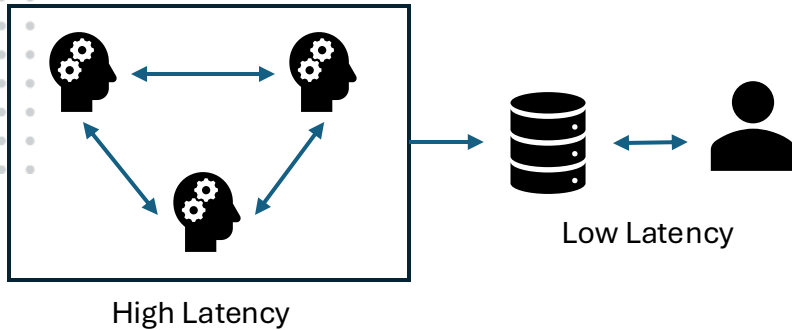
- Ease of use – how do we deploy this framework
- Any blacklisted products or messages we should not be showing?



Engineering Constraints: Scalability, Cost, and Latency

- **Problem Setting:**

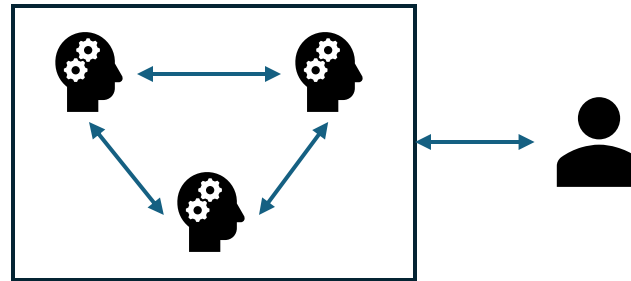
Agentic Framework takes time for real time inferencing. Identify whether the end user expects to wait couple of seconds vs they expect an immediate response.



Batched Inferencing

Existing touch points (search, recommendations, catalog, derived attributes, etc's)

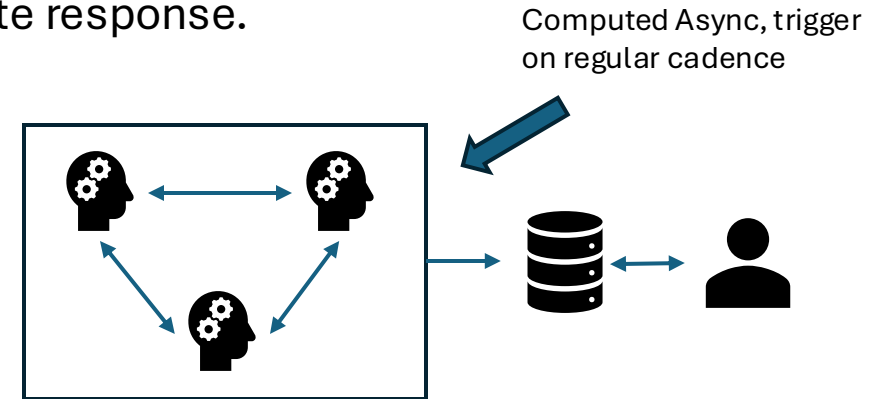
Users expect to get results immediately – so any agentic AI based inferencing should be computed offline



Real-Time Inferencing

Interactive Usecases (Chatbot, planner, research agents, etc's)

Users do not expect immediate results. That said, important to build a specialized agents that's not just a wrapper.



Hybrid (Near-real time)

Enhancements of existing touch points, passive in-session listener to improve

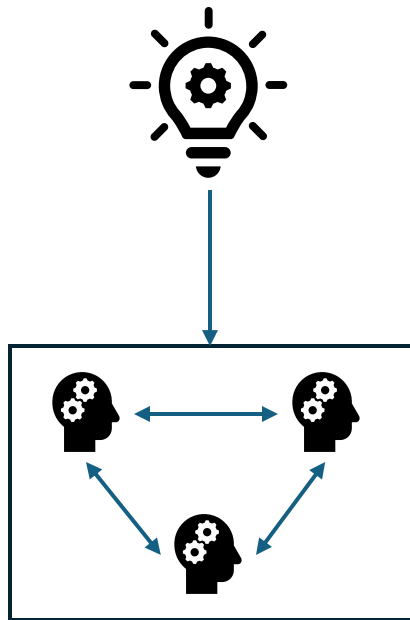
Users expect to get results immediately – so any agentic AI based inferencing should be computed asynchronously



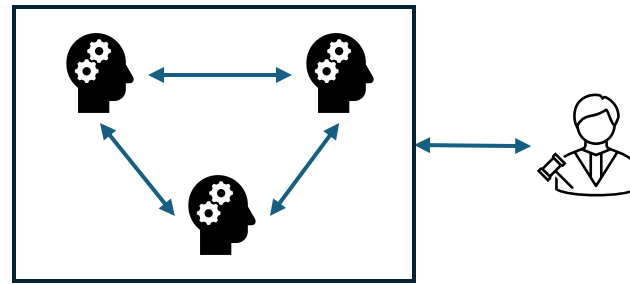
Metrics: What are we trying to solve?

- Problem Setting:**

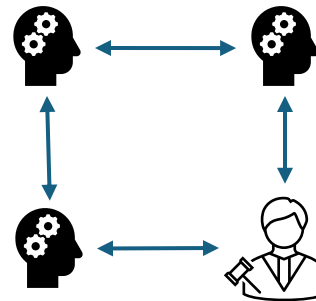
Once the business constraints and the types of acceptable scale/latency are defined, we can build the most reasonable workflow. The workflow should aim to **optimize the metric of interest**.



Define and Build Recommender System using Agentic AI

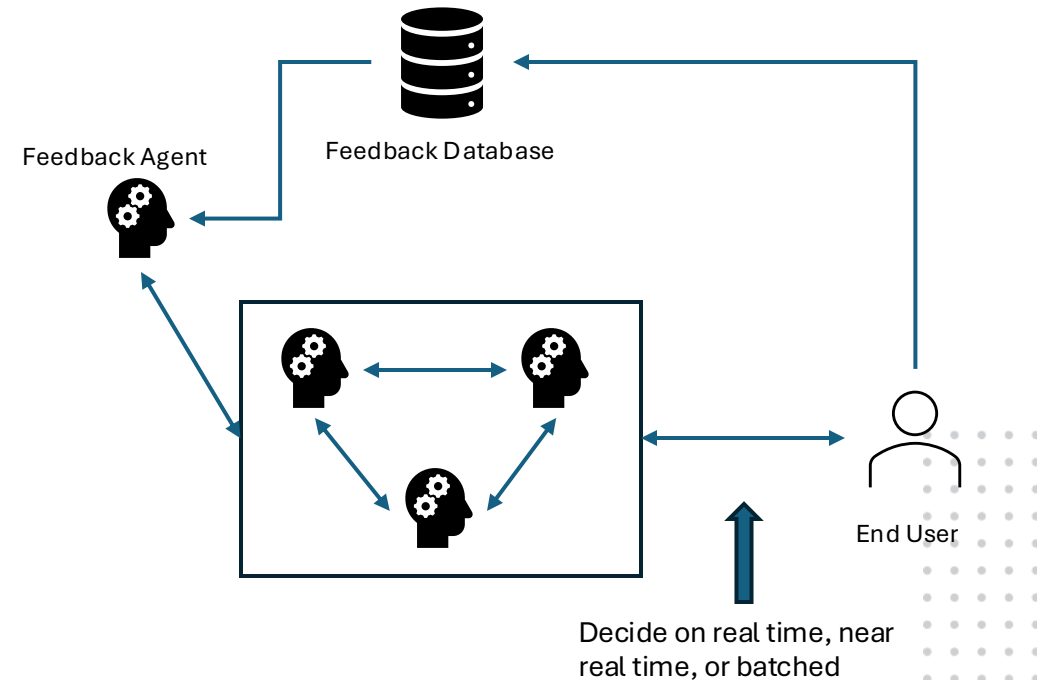


Hard Biz Constraints



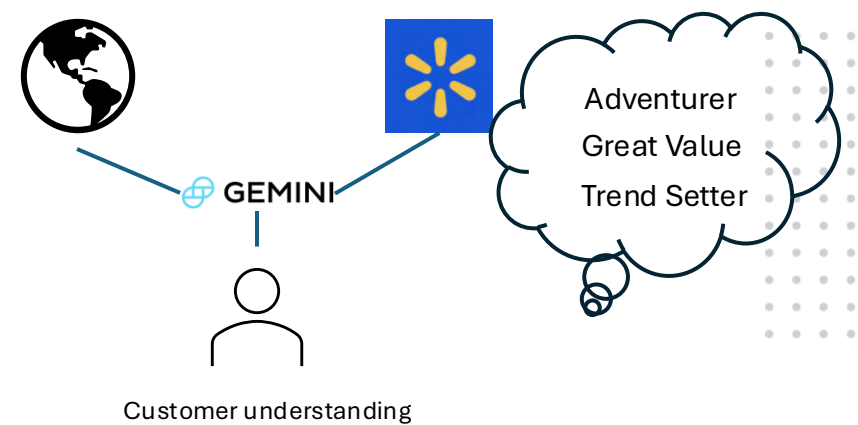
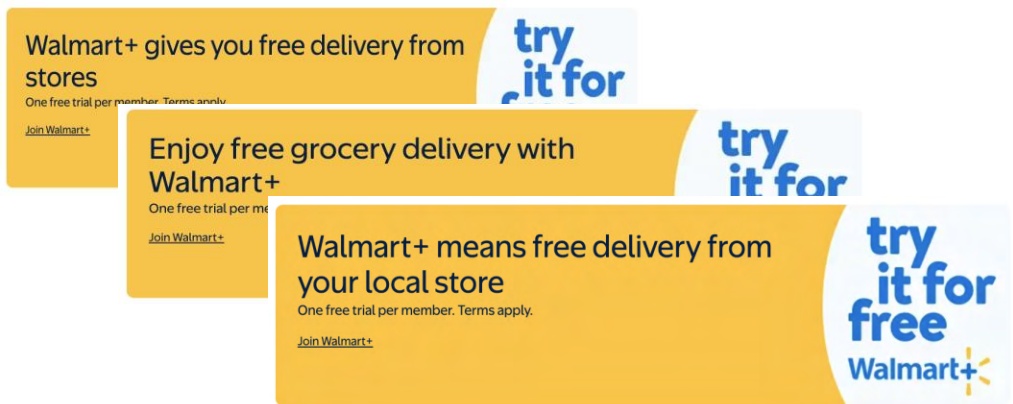
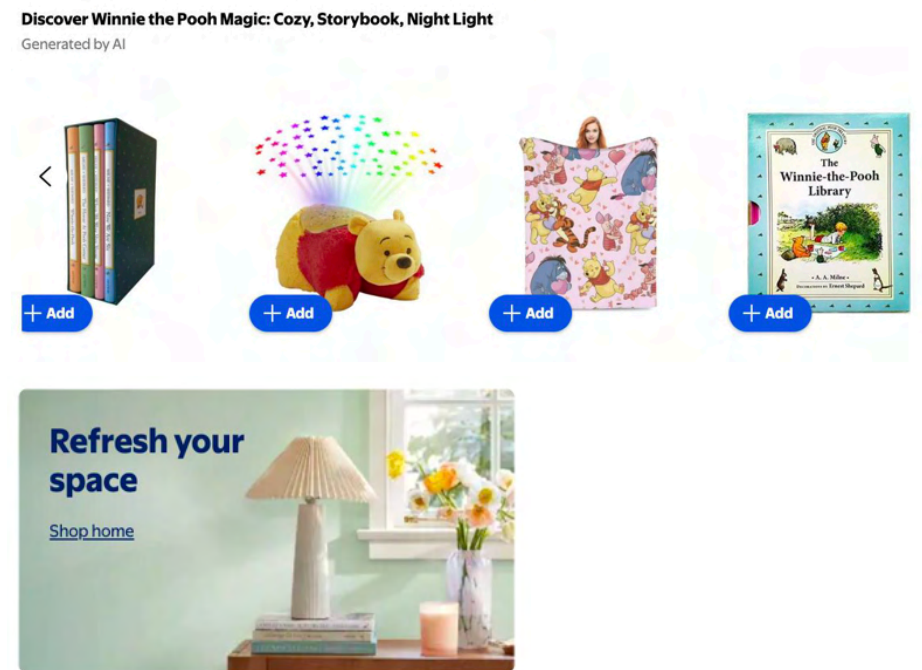
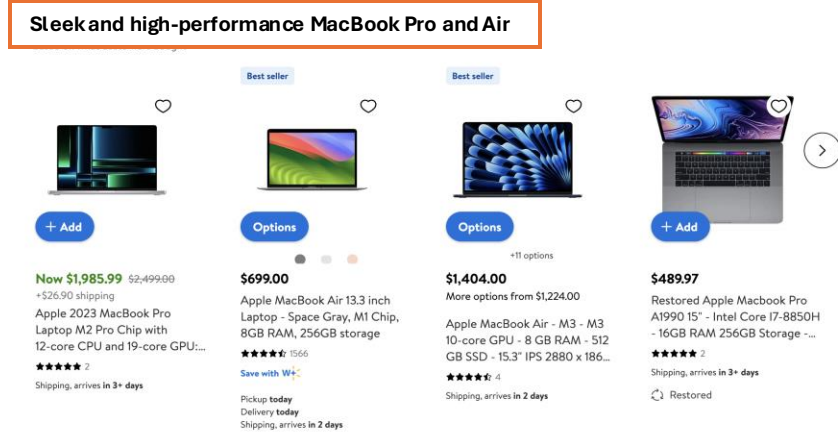
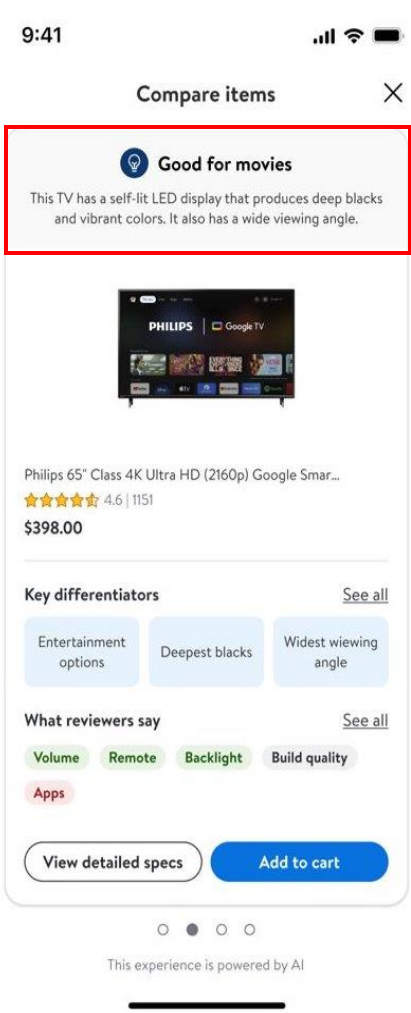
Soft Biz Constraints

Identify soft & hard biz constraints and incorporate the constraints into the framework accordingly



Ensure proper instrumentation & feedback loop to iteratively enhance both user experience and metrics

E-comm and Multi-Agent AI



Challenges and practical lessons building self-improving agents

Derek Cheng

Disclaimer: I've got more
questions than solutions.

What do I do?

- Run a research team focusing on:
 - Model & Data Efficiency;
 - Gen AI Monetization, especially for Ads & Commerce.
- End to end from idea, research, prototyping, all the way to production.
- Notable work from my team:
 - Cross-Batch Aggregation for Streaming Learning from Label Proportions. Jonathan Valverde et al. RecSys 2025.
 - Training data efficient LLMs: AskLLM. Noveen Sachdeva et al. -> Gemma & Gemini.
 - Unified Embedding. Ben Coleman & Wang-Cheng Kang et al. -> RecSys / Ads
 - Deep Cross Network V2. Ruoxi Wang et al. -> RecSys / Ads
- One recent focus: building self-improving agents.

Why building self-improving agents is HARD?

- Adding unique value
- Zero-shot v.s. Fine-tuning
- Evaluation
- Scaling
- Productionization

Adding Unique Value

- The water is rising with DeepSeek, Gemini 2.5 Pro, GPT 5, Grok 4, Claude 4, ...
- How do we make sure stuff we build is not yet another Gemini-wrapper?
- Keep adding unique value on top of a base model?

- A winning formula:
 - Success = product definition + problem formulation + innovations + connections + go to market + sales / marketing

Zero shot v.s. Fine-tuning

- Zero-shot can achieve a LOT.
 - Maybe even superhuman capabilities in some cases.
- But is that enough?
- A good example is Cursor:
 - Great company and product, built even more unique value and moat with improved cursor tab w/ online RL.

Evaluation

- Metrics
 - Definition
 - Offline v.s. online
- Human-raters
 - Speed, quality, cost
- Auto-raters
 - Side by side, pointwise

Scaling

- Model Size Scaling
 - A solution worked well with small data, would it scale well with bigger data?
- Feedback Data Scaling
 - We know more RL data is helpful, but could you build a scalable RL feedback loop with fast and accurate responses?
- Task Scaling
 - Would an agent work well on one scenario, generalize well to another scenario?

Productionization

- Scaling down
 - An agent worked well with a bigger model, would it distill well scaling down?
- QPS, latency, # of chips, ROI
- Regression while migrating to new & (supposedly) better models

Thanks!

Hiring research scientists and engineers in LLM & RecSys.
Derek Cheng: zcheng@google.com