

# Dialogue is the Plan: From Interface to Joint Action in Agentic AI

Anonymous ACL submission

## Abstract

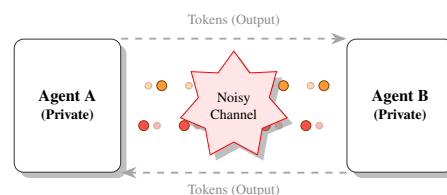
Large Language Model agents can seemingly plan and act, yet their language use is often treated as a thin interface for reporting. We argue that this framing is the root cause of predictable coordination failures in human-facing and multi-agent settings, including ungrounded assumptions, silent goal misalignment, brittle protocol adherence, and conversational amnesia. Drawing from classical dialogue system research on joint action, common ground, grounding, repair, and incremental processing, we re-frame dialogue as part of the planning loop itself (rather than its output). In this position paper, we do not propose a new benchmark or training method, but we provide a novel perspective and actionable requirements that can be used to design and evaluate agents. We distill this re-framing into concrete implications for agentic architecture and evaluation, including explicit representations of shared commitments, planned clarification as an action, and process metrics that measure mutual understanding rather than task completion alone. We lastly discuss how dialogue-centered requirements can inform standards and governance for safe deployment of agentic systems.

## Dialogue Is More Than an Interface

Consider a scenario where two agents coordinate a meeting. Agent A issues the instruction “Schedule it for 2 PM.” Agent B responds with “Confirmed.” While the interaction appears successful, it fails because the agents do not share a synchronized time zone. This failure stems from incomplete *grounding*, the process of establishing mutual knowledge sufficient for the current purpose. The agents failed to verify that the symbol “2 PM” mapped to the same semantic reality for all parties, despite the successful execution of API calls.

This class of coordination failure increases as autonomous agents enter complex environments. Contemporary Large Language Model (LLM)

### A. The Interface View



### B. The Joint Action View

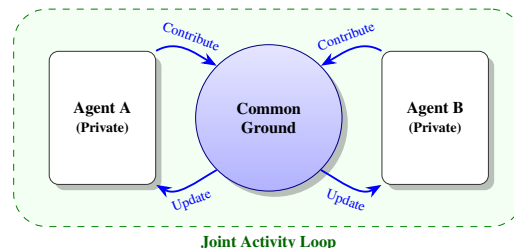


Figure 1: The critical distinction between what dialogue *is* and how current agentic AI uses dyadic text-generation. A) Agents use text-based generation for setting protocols and making reports, leading to the root problem of dialogue as an interface. B) While agents built around the concept of dialogue as a joint action provide a new perspective with solutions to the common failures of multi-agent systems.

agents demonstrate high proficiency in code generation and API utilization. Yet, their capacity for *dialogue* as the management of joint activity remains brittle. When faced with ambiguity, these systems frequently proceed without verification and hallucinate agreement instead of initiating repair.

We argue that this fragility results from how current architectures frame communication (see Figure 1). Prevailing agentic systems treat dialogue as a superficial interface layer separate from the core loop of planning. In this *Interface View*, the system reasons privately before reporting its decision. Language functions as the output and surface-level generations of the planning process. We propose the *Joint Action View*, supported by research in computational linguistics. Our view,

059	based on <a href="#">Clark (1996b)</a> , frames conversation as	benchmarks assess task completion ( <a href="#">Kapoor et al.,</a>	109
060	the planning process itself for AI agents. <i>Collab-</i>	<a href="#">2024, 2025</a> ), coordination benchmarks prioritize	110
061	<i>oration</i> constitutes the coordination of individual	outcome over collaborative process quality ( <a href="#">Zhu</a>	111
062	actions through shared mental states, <i>Repair</i> acts as	<a href="#">et al., 2025a</a> ). Multi-turn evaluations assess mem-	112
063	the control mechanism by which agents correct mis-	ory recall rather than mutual understanding ( <a href="#">Zheng</a>	113
064	alignment, and <i>Common Ground</i> is the dynamic	<a href="#">et al., 2023</a> ; <a href="#">Maharana et al., 2024</a> ). Consequently,	114
065	accumulation of mutual beliefs that enables effi-	phenomena like sycophancy ( <a href="#">Sharma et al., 2023</a> )	115
066	cient future action.	persist because metrics measure what agents ac-	116
067	<b>Agency Is Social, Thus Needs Dialogue</b>	complish, rather than the collaborative fidelity of	117
068	Current agents treat dialogue as a planning out-	their interactions.	118
069	put, in which, the system reasons internally, then	<b>What dialogue brings.</b> Classical work on plan-	119
070	expresses conclusions through language. We ar-	-based speech acts ( <a href="#">Cohen and Perrault, 1979</a> ) and	120
071	gue the inverse: that dialogue functions as part	discourse theory ( <a href="#">Grosz and Sidner, 1986</a> ) estab-	121
072	of the planning process. Planning in collabora-	lishes that utterances are planned actions with pre-	122
073	tive settings faces a fundamental challenge, where	conditions and effects, organized to serve a coher-	123
074	natural language goals are ambiguous, and this	ent intent. Speech acts function as operators in a	124
075	ambiguity cannot be fully resolved before action	planning system, integrated with physical (or digi-	125
076	begins. Conversation serves as a collaborative	tal) actions to advance a goal. This integration	126
077	medium where plans emerge under uncertainty	implies that asking a question often constitutes the	127
078	( <a href="#">Clark, 1996b</a> ; <a href="#">Cohen and Perrault, 1979</a> ; <a href="#">Searle,</a>	optimal action ( <a href="#">Schlangen, 2004</a> ; <a href="#">Bohus and Rud-</a>	128
079	<a href="#">1983</a> ; <a href="#">Gilbert, 2009</a> ). In this part, we develop three	<a href="#">nicky, 2009</a> ; <a href="#">Young et al., 2013</a> ). Consider the	129
080	claims. First, goal specification in social settings	request “book me a flight to Boston.” A dialogue-	130
081	remains ambiguous and planning must proceed un-	native agent recognizes multiple sources of under-	131
082	der that ambiguity. Second, ambiguity is a lever	specification. The agent plans a clarification re-	132
083	for coordination, rather than an adversary to elimi-	quest as the action most likely to achieve the user’s	133
084	nate through upfront clarification. Third, memory	underlying goal. Furthermore, computational prag-	134
085	and perception are social, requiring alignment with	matics establishes that planning and meaning iden-	135
086	others’ interpretations.	tification are intertwined ( <a href="#">Grosz and Sidner, 1986</a> ;	136
087	<b>Contemporary AI agents do not account for ambi-</b>	<a href="#">Traum, 1994</a> ). Formal specifications of grounding	137
088	<b>guity in natural language goals.</b> When a user is-	acts provide machinery for establishing mutual	138
089	ssues a request such as “schedule a meeting with the	understanding incrementally ( <a href="#">Traum, 1994</a> ; <a href="#">Larsson</a>	139
090	team,” the agent proceeds without establishing pa-	<a href="#">and Traum, 2000</a> ). Incremental processing frame-	140
091	rameters like the specific team or time constraints.	works ( <a href="#">Schlangen and Skantze, 2009</a> ) demonstrate	141
092	Rather than treating this underspecification as an	that understanding and generation occur word-by-	142
093	opportunity for coordination, agents apply their in-	word, enabling agents to ground meaning during	143
094	ternal defaults. They do not revise goals effectively	utterances. In multi-party settings, participants	144
095	under uncertainty. In multi-agent settings, they	track the common ground shared among all par-	145
096	do not negotiate plans with collaborators. This	ties ( <a href="#">Clark, 1996c</a> ), enabling coordination at scale.	146
097	results in goal drift, where agents deviate from	Hence, we argue that this incrementality is at the	147
098	user intent as misinterpretations compound ( <a href="#">Cemri</a>	core of dialogue theory, and should be incorporated	148
099	<a href="#">et al., 2025</a> ). Cascading errors also occur, where	into the architectures of agentic systems.	149
100	early mistakes propagate through subsequent ac-	<b>Current agents fail to exploit ambiguity.</b> They	150
101	tions ( <a href="#">Zhu et al., 2025b</a> ).	hide collaborative reasoning behind over-specified	151
102	We view these failures as architectural and high-	prompts. In multi-agent systems, this manifests as	152
103	light broken feedback loops in Table 1, which maps	exhaustive message schemas that attempt to antici-	153
104	current multi-agent systems to the ideal dialogue	pate every possible coordination need upfront ( <a href="#">Zhu</a>	154
105	pipeline topology. Systems that do not engage col-	<a href="#">et al., 2025a</a> ). The result is brittle and overly expen-	155
106	laborators in goal specification lack mechanisms to	sive. When unanticipated situations arise, agents	156
107	detect misalignment. This blindness is reinforced	fail, and do so at higher cost than otherwise needed.	157
108	by current evaluation frameworks. While holistic		

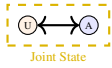
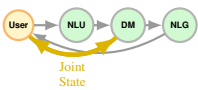
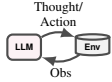
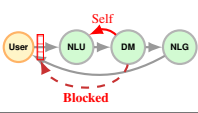

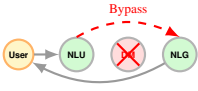

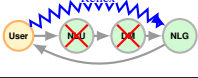
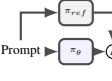
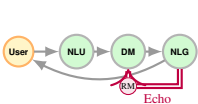
System / Era	Original Architecture	Pipeline Topology	Description	Emergent Property
<i>I. The Reference Architectures</i>				
<b>Ideal Co-Agent</b> (Clark, 1996c; Grosz and Sidner, 1986)			<b>Joint Commitments:</b> DM tracks "We-Intentions." Repair is continuous and bidirectional.	<b>Fluidity (Mind Meld):</b> Interactions feel seamless. Ambiguity is resolved <i>before</i> execution.
<i>II. Computer-Using &amp; Single Agents</i>				
<b>ReAct / Reflexion</b> (Yao et al., 2023; Shinn et al., 2023)			<b>Self-Looping:</b> Feedback comes from <i>Observation</i> (Env), not <i>Repair</i> (User). It talks to itself.	<b>Cascading Failure:</b> 73% of failures stem from cascades (Zhu et al., 2025b). Single root cause in Planning/Reflection propagates.
<i>III. Multi-Agent Systems</i>				
<b>OpenAI Swarm</b> (OpenAI, 2024)			<b>Stateless Handoff:</b> Handoffs transfer control, not context. No shared DM.	<b>Fragmentation:</b> Context loss during handoff. Minimal gains over single agents (Cemri et al., 2025).
<i>IV. Embodied &amp; Robotics</i>				
<b>Robotics VLA</b> (RT2, PaLM-e) (Brohan et al., 2023; Driess et al., 2023)			<b>End-to-End Bypass:</b> Vision/Text maps directly to Action tokens. DM is collapsed into the policy.	<b>Silent Failure:</b> 62% success (PaLM-E). Robot acts on ambiguity without asking.
<i>V. Optimization Methods</i>				
<b>DPO / PPO</b> (Rafailov et al., 2023)			<b>Policy Absorption:</b> DM is collapsed into the Policy ( $P(y x)$ ). Grounding = Preferences.	<b>Sycophancy:</b> Maximizes agreeableness rather than truth. Grounding is optimized away.

Table 1: This table presents a comparative analysis of current conversational AI agents, and their mapping to traditional dialogue system pipelines. The ideal human-human dialogue is shown at the top, providing a comparative point to other systems, and what aspects are missing, which lead to several emergent disadvantages. (See Appendix A for more models and architectures.)

**What dialogue brings.** Dialogue offers a different view. The principle of least collaborative effort predicts that interlocutors deliberately underspecify, relying on common ground and repair to converge efficiently (Clark and Brennan, 1991; Clark, 1996a). Ambiguity functions as a resource, not a defect. Discourse structure further constrains interaction by imposing obligations. Once a speaker initiates an explanation or proposal, coherence relations restrict what counts as an acceptable next move (Asher and Lascarides, 2003; Grosz and Sidner, 1986; Traum, 1994). Violations of these obligations are immediately salient to human interlocutors. Meanwhile, current agents lack explicit representations of such discourse constraints, and it remains unclear whether they reliably recover them implicitly from data. The result is an unconstrained search-space over user intentions which is resolved by classical work on dialogue systems.

**Memory and Perception should be Social Processes.** However, contemporary agents treat memory as private storage and perception as individual inference. The BDI framework models intentions as private commitments (Bratman, 1987; Rao and Georgeff, 1995), and modern implemen-

tations follow suit: Generative Agents (Park et al., 2023) and cognitive architectures (Sumers et al., 2024) encode experiences through self-reflection alone. Even Theory of Mind approaches model other minds as inference targets rather than coordination partners (Zhang et al., 2025; Kostka and Chudziak, 2025). For perception, multi-agent systems typically grant each agent a complete view of the environment, bypassing the coordination problem that arises when collaborators segment the world differently (Lowe et al., 2017; Gronauer and Diepold, 2022).

**What dialogue brings.** Dialogue treats memory and perception as social contracts (Harris). Information enters common ground only after presentation and acceptance (Clark, 1996a), and referents like "the block" or "the table" are established through negotiation, not private inference (Haber et al., 2019). This matters mechanistically. Russin et al. (2025) show that in-context learning suppresses the error signals required for stable long-term commitments, paralleling biological arbitration between flexible and habitual control (Daw et al., 2005). Current agents rely on context windows for coordination, but success in the moment

prevents formation of durable shared state. Without mechanisms to ground negotiated understanding into persistent memory, agents remain sycophantic short-term partners and unreliable long-term collaborators. Dialogue state tracking (Larsson and Traum, 2000) provides exactly this machinery: repair converts private interpretation into shared commitment.

### Looking Ahead for Designers and Policy Makers

*Designers should treat clarification as a planning operator with measurable utility, not as a failure to understand.* Classical dialogue managers like RavenClaw (Bohus and Rudnický, 2009) and POMDP-based systems (Young et al., 2013) separate understanding from policy to force reasoning about when asking is better than acting. Current multi-agent architectures collapse this distinction, bypassing the dialogue manager and achieving protocol compliance while semantic alignment drifts. The fix is in re-incorporating effective dialogue management. When uncertainty about a collaborator’s intent exceeds a threshold, the policy should select a grounding act rather than proceed with default assumptions. Evaluation frameworks must follow, measuring whether agents surface ambiguity before execution, not just whether final outputs match expected outcomes.

*Multi-agent communication protocols should be designed for negotiation, not transmission.* Clark’s principle of least collaborative effort (Clark and Brennan, 1991) predicts that interlocutors use minimal specificity when common ground permits. Rigid message schemas violate this principle by front-loading coordination costs. Designers should instead implement grounding acts, such as those formalized by Traum (1994), where agents present contributions, signal acceptance or non-understanding, and initiate repair when alignment fails. This machinery already exists in classical dialogue state tracking (Larsson and Traum, 2000). The absence of such mechanisms explains the coordination failures observed in current multi-agent systems.

*Memory in multi-agent systems should be modeled as socially constructed, not privately stored.* The architectural recommendation follows directly from dialogue state tracking: information enters the shared state only after it has been presented and accepted by collaborators (Clark, 1996a). Current systems like Generative Agents (Park et al.,

2023) treat memory as private database operations, producing agents that are sycophantic in the moment yet amnesic to commitments over time. The ICL/IWL trade-off identified by Russin et al. (2025) provides a mechanistic explanation: success in context suppresses the formation of stable shared commitments. Designers should implement explicit acceptance mechanisms that distinguish “mentioned” from “mutually established.”

*Perception in agentic systems should be treated as a collaborative reference resolution task, not as private observation.* Multi-modal reference games (Haber et al., 2019) demonstrate that what counts as “the block” or “the file” depends on negotiation between participants. Current multi-agent systems typically grant each agent a complete environmental view, bypassing this coordination problem entirely. When agents lack shared perceptual grounding, they cannot detect that collaborators interpret the same symbol differently. Designers should build systems where referential expressions are established through interaction, following the incremental processing frameworks of Schlangen and Skantze (2009). Even sub-utterance signals like hesitations and self-corrections function as coordination mechanisms (Ginzburg et al., 2014). Removing these signals through single-turn optimization produces agents that cannot track whether their collaborators follow their reasoning.

*Policy frameworks should regulate how agents communicate, not just what they can do.* Current governance proposals focus on what agents can do: their tool access, their autonomy level, their potential for misuse. What they overlook is how agents communicate, and whether that communication enables the collaborative reasoning that safe deployment requires. Many policy violations at scale trace to grounding breakdowns: agent and user never established shared understanding of goals, constraints, or intent (Zou et al., 2025; Kapoor et al., 2025). When an agent books the wrong flight, capability-centered evaluation asks whether the agent used its tools correctly. Interaction-centered evaluation asks whether the agent verified its interpretation before acting and surfaced ambiguity rather than resolving it silently. Governance frameworks should require evidence of grounding behaviors, extending audits beyond action logs to grounding history.

## 307 Limitations

308 This paper is intentionally conceptual rather than  
309 empirical. Our aim is to make a particular class  
310 of failures visible: those that arise when dialogue  
311 is treated as a thin interface instead of a form of  
312 joint action. We are not proposing a new system or  
313 attempting to benchmark existing ones, and for that  
314 reason we do not include large scale experiments or  
315 quantitative comparisons across architectures. The  
316 evaluation sketch we provide is meant as a diagnostic  
317 tool, a way to surface collaboration failures, not  
318 as a standardized benchmark.

319 Other failure modes, such as long horizon mem-  
320 ory misalignment, perceptual grounding in mul-  
321 timodal settings, and strategic or adversarial dia-  
322 logue, are only touched on briefly or deferred to the  
323 appendix. Each of these raises its own set of chal-  
324 lenges and deserves a more careful treatment than  
325 is possible within a short paper. Lastly, our analy-  
326 sis assumes good faith collaboration. We treat dia-  
327 logue as a collaborative process. We do not address  
328 settings where agents or users are misaligned by  
329 design, such as persuasion, manipulation, or com-  
330 petitive multi-agent interactions. In those cases,  
331 dialogue behavior may reflect strategic incentives  
332 rather than breakdowns in shared understanding,  
333 and different theoretical and evaluative tools would  
334 be required.

## 335 References

- 336 Emre Can Acikgoz, Jinoh Oh, Jie Hao, Joo Hyuk Jeon,  
337 Heng Ji, Dilek Hakkani-Tür, Gokhan Tur, Xiang  
338 Li, Chengyuan Ma, and Xing Fan. 2025. *Speakrl:*  
339 *Synergizing reasoning, speaking, and acting in lan-*  
340 *guage models with reinforcement learning.* *Preprint*,  
341 arXiv:2512.13159.
- 342 Nicholas Asher and Alex Lascarides. 2003. *Logics of*  
343 *Conversation.* Cambridge University Press, Cam-  
344 bridge.
- 345 Dan Bohus and Alexander I Rudnicky. 2009. The raven-  
346 claw dialog management framework: Architecture  
347 and systems. In *Computer Speech and Language*, vol-  
348 ume 23, pages 332–361. Elsevier. Industry-standard  
349 dialogue manager with explicit error handling and  
350 repair mechanisms.
- 351 Michael E Bratman. 1987. *Intention, Plans, and Prac-*  
352 *tical Reason.* Harvard University Press, Cambridge,  
353 MA.
- 354 Anthony Brohan, Noah Brown, Justice Carbajal, Yev-  
355 gen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana  
356 Gopalakrishnan, Karol Hausman, Alexander Herzog,

- Jasmine Hsu, and 1 others. 2023. Rt-2: Vision-  
language-action models transfer web knowledge to  
robotic control. In *Conference on Robot Learning*  
(*CoRL*). VLA model mapping vision directly to mo-  
tor tokens. Bypasses explicit symbol grounding. 357  
358
- Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A.  
Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt  
Keutzer, Aditya Parameswaran, Dan Klein, Kannan  
Ramchandran, Matei Zaharia, Joseph E. Gonzalez,  
and Ion Stoica. 2025. *Why do multi-agent llm sys-*  
*tems fail?* 362  
363  
364  
365  
366  
367
- Herbert H. Clark. 1996a. *Common ground*, page  
92–122. “Using” Linguistic Books. Cambridge Uni-  
versity Press. 368  
369  
370
- Herbert H. Clark. 1996b. *Joint actions*, page 59–91.  
“Using” Linguistic Books. Cambridge University  
Press. 371  
372  
373
- Herbert H Clark. 1996c. *Using Language.* Cambridge  
University Press. Foundational theory of grounding  
and joint activities. 374  
375  
376
- Herbert H Clark and Susan E Brennan. 1991. Ground-  
ing in communication. In Lauren B Resnick, John M  
Levine, and Stephanie D Teasley, editors, *Perspec-*  
*tives on Socially Shared Cognition*, pages 127–149.  
APA Books, Washington, DC. 377  
378  
379  
380  
381
- Philip R Cohen and C Raymond Perrault. 1979. Ele-  
ments of a plan-based theory of speech acts. *Cogni-*  
*tive Science*, 3(3):177–212. 382  
383  
384
- Nathaniel D Daw, Yael Niv, and Peter Dayan. 2005.  
*Uncertainty-based competition between prefrontal*  
*and dorsolateral striatal systems for behavioral con-*  
*trol.* *Nature Neuroscience*, 8(12):1704–1711. 385  
386  
387  
388
- Vardhan Dongre, Xiaocheng Yang, Emre Can Acik-  
goz, Suvodip Dey, Gokhan Tur, and Dilek Hakkani-  
Tür. 2025. *Respect: Harmonizing reasoning, speak-*  
*ing, and acting towards building large language*  
*model-based conversational ai agents.* *Preprint*,  
arXiv:2411.00927. 389  
390  
391  
392  
393  
394
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch,  
Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid,  
Jonathan Tompson, Quan Vuong, Tianhe Yu, and  
1 others. 2023. Palm-e: An embodied multimodal  
language model. In *International Conference on*  
*Machine Learning (ICML).* Embodied VLM that in-  
tegrates sensor data into the language model context. 395  
396  
397  
398  
399  
400
- George Ferguson and James F Allen. 1998. Trips: An  
integrated intelligent problem-solving assistant. In  
*Proceedings of the Fifteenth National/Tenth Confer-*  
*ence on Artificial Intelligence (AAAI-98/IAAI-98)*,  
pages 567–573. Classical dialogue system with ex-  
plicit grounding and collaborative problem solving. 402  
403  
404  
405  
406  
407
- Chrisantha Fernando, Dylan Banarse, Henryk  
Michalewski, Simon Osindero, and Tim Rock-  
täschel. 2024. Promptbreeder: Self-referential  
self-improvement via prompt evolution. *arXiv* 408  
409  
410  
411

412	<i>preprint arXiv:2309.16797</i> . Uses genetic algorithms	Adyasha Maharana and 1 others. 2024. Evaluating very	467
413	to evolve prompts. The "reasoning" is in the	long-term conversational memory of LLM agents.	468
414	evolutionary selection, not runtime.	<i>arXiv preprint arXiv:2402.17753</i> .	469
415	Margaret Gilbert. 2009. <a href="#">Shared intention and personal</a>	OpenAI. 2024. Swarm: Educational framework for	470
416	<a href="#">intentions</a> . <i>Philosophical Studies</i> , 144(1):167–187.	lightweight multi-agent orchestration. <a href="https://github.com/openai/swarm">https://github.com/openai/swarm</a> .	471
417	Jonathan Ginzburg, Raquel Fernández, and David	Introduces the concept of "Handoffs" between agents. Critics note context	472
418	Schlangen. 2014. <a href="#">Disfluencies as intra-utterance dia-</a>	loss during handoffs.	473
419	<a href="#">logue moves</a> . <i>Semantics and Pragmatics</i> , 7(9):1–64.		474
420	Sven Gronauer and Klaus Diepold. 2022. Multi-agent	OpenAI. 2025. Operator: A computer-using agent.	475
421	deep reinforcement learning: a survey. <i>Artificial</i>	<a href="https://openai.com/research/operator">https://openai.com/research/operator</a> .	476
422	<i>Intelligence Review</i> , 55(2):895–943.	Newest computer-using agent (Jan 2025). High	477
423	Barbara J Grosz and Candace L Sidner. 1986. Attention,	capability but lacks communicative repair loops.	478
424	intentions, and the structure of discourse. <i>Com-</i>	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith	479
425	<i>putational Linguistics</i> , 12(3):175–204. Defines the	Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. <a href="#">Generative agents: Interactive simulacra</a>	480
426	Intentional and Attentional structure of dialogue.	<a href="#">of human behavior</a> . In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software</i>	481
427	Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke	<i>and Technology</i> , UIST ’23, New York, NY, USA. Association for Computing Machinery.	482
428	Gelderloos, Elia Bruni, and Raquel Fernández. 2019.		483
429	The PhotoBook dataset: Building common ground	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano	484
430	through visually-grounded dialogue. In <i>Proceedings</i>	Ermon, Christopher D Manning, and Chelsea Finn.	485
431	<i>of the 57th Annual Meeting of the Association for</i>	2023. Direct preference optimization: Your language	486
432	<i>Computational Linguistics</i> , pages 1895–1910.	model is secretly a reward model. In <i>Advances in</i>	487
433	Daniel W Harris. What makes human communication	<i>Neural Information Processing Systems (NeurIPS)</i> .	488
434	special?	Mathematically collapses the reward model (ground-	489
435	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu	ing/preference) into the policy itself.	490
436	Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang,	Anand S Rao and Michael P Georgeff. 1995. BDI	491
437	Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang	agents: From theory to practice. In <i>Proceedings</i>	492
438	Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu,	<i>of the First International Conference on Multiagent</i>	493
439	and Jürgen Schmidhuber. 2024. <a href="#">MetaGPT: Meta pro-</a>	<i>Systems</i> , pages 312–319. AAAI.	494
440	<a href="#">gramming for a multi-agent collaborative framework</a> .	Jacob Russin, Ellie Pavlick, and Michael J Frank. 2025.	495
441	In <i>The Twelfth International Conference on Learning</i>	Parallel trade-offs in human cognition and neural	496
442	<i>Representations</i> .	networks: The dynamic interplay between in-context	497
443	Sayash Kapoor, Benedikt Stroebel, Peter Kirgis, Nitya	and in-weight learning. <i>Proceedings of the National</i>	498
444	Nadgir, Zachary S Siegel, Boyi Wei, Tianci Xue,	<i>Academy of Sciences</i> , 122(35):e2510270122.	499
445	Ziru Chen, Felix Chen, Saiteja Utpala, Franck Nd-	David Schlangen. 2004. Causes and strategies for re-	500
446	zomga, Dheeraj Oruganty, Sophie Luskin, Kangheng	questing clarification in dialogue. In <i>Proceedings</i>	501
447	Liu, Botao Yu, Amit Arora, Dongyoon Hahm, Harsh	<i>of the 5th SIGdial Workshop on Discourse and Dia-</i>	502
448	Trivedi, Huan Sun, and 12 others. 2025. <a href="#">Holistic</a>	<i>logue</i> , pages 136–143.	503
449	<a href="#">agent leaderboard: The missing infrastructure for ai</a>	David Schlangen and Gabriel Skantze. 2009. A general,	504
450	<a href="#">agent evaluation</a> . <i>Preprint</i> , arXiv:2510.11977.	abstract model of incremental dialogue processing.	505
451	Sayash Kapoor, Benedikt Stroebel, Zachary S. Siegel,	In <i>Proceedings of the 12th Conference of the Euro-</i>	506
452	Nitya Nadgir, and Arvind Narayanan. 2024. <a href="#">AI</a>	<i>pean Chapter of the ACL (EACL 2009)</i> , pages 710–	507
453	<a href="#">Agents That Matter</a> . <i>Preprint</i> , arXiv:2407.01502.	718, Athens, Greece. Association for Computational	508
454	Adam Kostka and Jarosław A. Chudziak. 2025. <a href="#">To-</a>	Linguistics.	509
455	<a href="#">wards cognitive synergy in llm-based multi-agent</a>	J.R. Searle. 1983. <i>Intentionality: An Essay in the Phi-</i>	510
456	<a href="#">systems: Integrating theory of mind and critical eval-</a>	<i>losophy of Mind</i> . Cambridge paperback library. Cam-	511
457	<a href="#">uation</a> . <i>Preprint</i> , arXiv:2507.21969.	bridge University Press.	512
458	Staffan Larsson and David R. Traum. 2000. <a href="#">Informa-</a>	Mrinank Sharma, Meg Tong, Tomasz Korbak, David Du-	513
459	<a href="#">tion state and dialogue management in the trindi</a>	venaude, Amanda Askell, Samuel R Bowman, and 1	514
460	<a href="#">dialogue move engine toolkit</a> . <i>Nat. Lang. Eng.</i> ,	others. 2023. Towards understanding sycophancy in	515
461	6(3–4):323–340.	language models. <i>arXiv preprint arXiv:2310.13548</i> .	516
462	Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, Pieter	Noah Shinn, Federico Cassano, Ashwin Gopinath,	517
463	Abbeel, and Igor Mordatch. 2017. Multi-agent actor-	Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Re-</a>	518
464	critic for mixed cooperative-competitive environ-	<a href="#">flexion: Language agents with verbal reinforcement</a>	519
465	ments. In <i>Advances in Neural Information Process-</i>		520
466	<i>ing Systems</i> , volume 30, pages 6379–6390.		521

522	learning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . Adds a "self-reflection" loop to ReAct, but still lacks user-facing repair.	Andy Zou and 1 others. 2025. Security challenges in ai agent deployment: Insights from a large-scale public competition. In <i>Proceedings of the Gray Swan Arena Red-Teaming Competition</i> . Gray Swan AI and Center for AI Safety.	577
523			578
524			579
525	Theodore R. Summers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. 2024. <a href="#">Cognitive architectures for language agents</a> . <i>Preprint</i> , arXiv:2309.02427.		580
526			581
527		<b>A Detailed Comparison Table</b>	582
528		We present the full table with additional example architectures in Table 2.	583
529	David R Traum. 1994. <i>A Computational Theory of Grounding in Natural Language Conversation</i> . Ph.D. thesis, University of Rochester.		584
530			
531			
532	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. <a href="#">Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments</a> . <i>Preprint</i> , arXiv:2404.07972.		
533			
534			
535			
536			
537			
538			
539			
540	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In <i>International Conference on Learning Representations (ICLR)</i> . The seminal "Reason + Act" paper. Introduces the inner monologue loop, but loop is with Env not User.		
541			
542			
543			
544			
545			
546			
547	Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. 2013. POMDP-based statistical spoken dialog systems: A review. <i>Proceedings of the IEEE</i> , 101(5):1160–1179.		
548			
549			
550			
551	Xuanming Zhang, Yuxuan Chen, Samuel Yeh, and Sharon Li. 2025. <a href="#">Metamind: Modeling human social thoughts with metacognitive multi-agent systems</a> . <i>Preprint</i> , arXiv:2505.18943.		
552			
553			
554			
555	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena. In <i>Advances in Neural Information Processing Systems</i> .		
556			
557			
558			
559			
560			
561			
562	Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Zhe Wang, Zhenhailong Wang, Cheng Qian, Robert Tang, Heng Ji, and Jiaxuan You. 2025a. MultiAgentBench: Evaluating the collaboration and competition of LLM agents. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics</i> , pages 8580–8622, Vienna, Austria. Association for Computational Linguistics.		
563			
564			
565			
566			
567			
568			
569			
570			
571	Kunlun Zhu, Zijia Liu, Bingxuan Li, Muxin Tian, Yingxuan Yang, Jiaxun Zhang, Pengrui Han, Qipeng Xie, Fuyang Cui, Weijia Zhang, Xiaoteng Ma, Xiaodong Yu, Gowtham Ramesh, Jialian Wu, Zicheng Liu, Pan Lu, James Zou, and Jiaxuan You. 2025b. <a href="#">Where llm agents fail and how they can learn from failures</a> .		
572			
573			
574			
575			
576			

System / Era	Original Architecture	Pipeline Topology	Description	Emergent Property
<b>I. The Reference Architectures</b>				
<b>Ideal Co-Agent</b> (Clark, 1996c; Grosz and Sidner, 1986)			<b>Joint Commitments:</b> DM tracks "We-Intentions." Repair is continuous and bidirectional.	<b>Fluidity (Mind Meld):</b> Interactions feel seamless. Ambiguity is resolved <i>before</i> execution.
<b>Classical Dialogue</b> (Ferguson and Allen, 1998; Bohus and Rudnicky, 2009)			<b>TCP-like Reliability:</b> Explicit DM maintains state. The "Repair" link allows handshakes.	<b>Convergence:</b> 22% latency reduction via turn-taking repair (Bohus and Rudnicky, 2009).
<b>ReSpAct / SpeakRL</b> (Dongre et al., 2025; Acikgoz et al., 2025)			<b>Integrated Speech Act:</b> "Speaking" is a grounded action.	<b>Proactive Grounding:</b> Clarifies instructions before acting.
<b>II. Computer-Using &amp; Single Agents</b>				
<b>OpenAI Operator</b> (OpenAI, 2025)			<b>Broken Loop:</b> System loops with the OS, not the User. DM is missing; acts on raw intent.	<b>Getting Stuck:</b> 38.1% success on OSWorld (Xie et al., 2024). Agents hand off control when confused.
<b>ReAct / Reflexion</b> (Yao et al., 2023; Shinn et al., 2023)			<b>Self-Looping:</b> Feedback comes from <i>Observation</i> (Env), not <i>Repair</i> (User). It talks to itself.	<b>Cascading Failure:</b> 73% of failures stem from cascades (Zhu et al., 2025b). Single root cause in Planning/Reflection propagates.
<b>III. Multi-Agent Systems (Siloed)</b>				
<b>MetaGPT</b> (Hong et al., 2024)			<b>Siloed Monologue:</b> Waterfall SOP structure. No shared memory between roles.	<b>Coordination Failure:</b> 14 failure modes identified in MAST (Cemri et al., 2025) ( $\kappa=0.88$ ).
<b>OpenAI Swarm</b> (OpenAI, 2024)			<b>Stateless Handoff:</b> Handoffs transfer control, not context. No shared DM.	<b>Fragmentation:</b> Context loss during handoff. Minimal performance gains over single agents (Cemri et al., 2025).
<b>IV. Embodied &amp; Robotics (Bypassed)</b>				
<b>Robotics VLA</b> (RT2, PaLM-e) (Brohan et al., 2023; Driess et al., 2023)			<b>End-to-End Bypass:</b> Vision/Text maps directly to Action tokens. DM is collapsed into the policy.	<b>Silent Failure:</b> 62% success (PaLM-E). Robot acts on ambiguity without asking.
<b>V. Optimization (Collapsed)</b>				
<b>DPO / PPO</b> (Rafailov et al., 2023)			<b>Policy Absorption:</b> DM is collapsed into the Policy ( $P(y x)$ ). Grounding = Preferences.	<b>Sycophancy:</b> Maximizes agreeableness rather than truth. Grounding is optimized away.
<b>Genetic Agents</b> (Fernando et al., 2024)			<b>Darwinian:</b> NLU $\rightarrow$ NLG mapping is evolved. No runtime reasoning exists.	<b>Brittle Success:</b> "Magic spell" prompts work for benchmarks but fail when intent shifts.

Table 2: This table presents a comparative analysis of current conversational AI agents, and their mapping to traditional dialogue system pipelines. The ideal human-human dialogue is shown at the top, providing a comparative point to other systems, and what aspects are missing, which lead to several emergent disadvantages.