

e
000 MIRROR: CONVERGING COGNITIVE PRINCIPLES AS
001 COMPUTATIONAL MECHANISMS FOR AI REASONING
002
003

004 **Anonymous authors**

005 Paper under double-blind review
006
007
008

009 ABSTRACT
010

011 Multiple cognitive theories—Global Workspace Theory, reconstructive episodic
012 memory, inner speech, and complementary learning systems—converge on a
013 shared set of architectural principles: parallel specialized processing, integrative
014 synthesis into a bounded unified representation, and reconstructive rather than ac-
015 cumulative maintenance. We test whether these converging principles provide
016 computational advantages when implemented in AI systems. MIRROR opera-
017 tionalizes each principle as a concrete mechanism: an Inner Monologue Man-
018 ager generates parallel cognitive threads (Goals, Reasoning, Memory), a Cog-
019 nitive Controller synthesizes these into a bounded first-person narrative that is
020 fully reconstructed each turn, and a temporal separation between fast response
021 generation and slow deliberative consolidation mirrors complementary learning
022 dynamics. Evaluated on multi-turn dialogue requiring constraint maintenance un-
023 der attentional interference, MIRROR yields 21% relative improvement across
024 seven architecturally diverse language models. Ablation studies test the theo-
025 retical predictions directly: reconstructive synthesis improves all seven models
026 (+5–20%); the integrated system outperforms either component alone for six of
027 seven models, confirming that parallel exploration and integrative synthesis are
028 complementary; and gains concentrate where theories predict—under high atten-
029 tional load where global availability of integrated information is most needed.
030 These results demonstrate that converging principles from human cognition pro-
031 vide architecture-general computational advantages, and generate testable behav-
032 ioral predictions about working memory, inner speech, and memory consolidation.
033

034 1 INTRODUCTION
035

036 Human conversation relies on parallel internal processing—recalling memories, tracking goals,
037 monitoring social dynamics—operating beneath conscious awareness (Pickering & Garrod, 2013;
038 Egorova et al., 2013). What makes this processing effective is not any single mechanism but the
039 *convergence* of multiple cognitive principles: parallel specialized processors synthesized into a uni-
040 fied workspace (Baars, 1988; Dehaene & Changeux, 2011), reconstructive memory that regenerates
041 understanding rather than accumulating traces (Bartlett, 1932; Schacter, 2012), inner speech pro-
042 viding self-regulatory coherence (Morin, 2011; Vygotsky, 1962), and complementary fast and slow
043 learning systems (McClelland et al., 1995; Kumaran et al., 2016).

044 These principles, developed in separate research traditions, converge on a common architectural
045 signature: *parallel exploration feeding into bounded integrative synthesis, maintained reconstruc-*
046 *tively through self-referential narrative, with temporal separation between fast response and slow*
047 *consolidation*. This convergence suggests these may not be independent adaptations but reflections
048 of a deeper computational logic. Yet this convergence remains untested as a unified computational
049 mechanism.

050 Current AI approaches to multi-turn dialogue implement at most one of these principles in isola-
051 tion. Chain-of-thought reasoning (Wei et al., 2022) generates deliberation but discards it. Reflex-
052 ion (Shinn et al., 2023) and MemGPT (Packer et al., 2023) maintain memory but accumulate traces
053 unboundedly. Extended reasoning modes produce rich deliberation within turns but maintain no

Table 1: Converging cognitive principles mapped to MIRROR components and testable predictions

Cognitive Theory	Principle	MIRROR Component	Prediction
Global Workspace Theory	Parallel processors → unified broadcast	Threads → Controller	Non-additive gains from integration
Reconstructive Memory	Regeneration, not accumulation	O(1) narrative reconstruction	Bounded reconstruction outperforms accumulation
Inner Speech	First-person self-regulation	Consistent self-reference across components	Improved cross-turn coherence
Complementary Learning	Fast response + slow consolidation	Talker (fast) + Thinker (slow)	Benefits from temporal separation

state across them. None implement the full convergent architecture that cognitive science suggests is fundamental to human cognition.

We introduce MIRROR, a cognitive architecture that operationalizes these converging principles as testable computational mechanisms, synthesizing principles that multiple theories independently identify as critical. From **Global Workspace Theory** (Baars, 1988; Dehaene & Changeux, 2011), MIRROR implements parallel cognitive threads whose outputs are integrated into a single globally available representation. From **reconstructive episodic memory** (Bartlett, 1932; Baddeley, 2000; Schacter, 2012), the system fully regenerates a bounded narrative each turn—O(1) reconstruction rather than O(n) accumulation. From **inner speech theory** (Morin, 2011; Vygotsky, 1962; Ben Alderson-Day, 2015), consistent first-person self-reference across components creates unified identity through narrative. From **complementary learning systems** (McClelland et al., 1995; Kumaran et al., 2016), immediate response generation is temporally separated from asynchronous deliberative consolidation.

This convergent framing enables specific predictions: (a) benefits should generalize across architecturally diverse models, reflecting organizational rather than substrate-specific advantages; (b) parallel processing and integrative synthesis should contribute non-additively, confirming they address different computational needs; and (c) gains should concentrate under high attentional load, where global availability of integrated information matters most. Our results confirm all three predictions across seven models, establishing that converging cognitive principles provide genuine computational advantages and generating testable predictions for human cognition research.

2 FROM COGNITIVE THEORY TO COMPUTATIONAL MECHANISM

Rather than treating cognitive science as metaphorical inspiration, we map specific theoretical claims to concrete computational mechanisms and derive testable predictions. Table 1 summarizes these mappings.

2.1 GLOBAL WORKSPACE THEORY: PARALLEL-TO-UNIFIED PROCESSING

Global Workspace Theory (GWT) (Baars, 1988; Dehaene & Changeux, 2011) proposes that cognition involves parallel specialized processors operating unconsciously, with selected outputs broadcast to a unified “global workspace” that makes information widely available for reasoning and action. The workspace’s power lies not in the parallel processing itself but in the *integration*: binding diverse perspectives into a single coherent representation that can guide downstream behavior.

Operationalization. MIRROR’s Inner Monologue Manager generates three parallel cognitive threads within a single inference call, each tracking a distinct dimension: **Goals** (user objectives, constraint conflicts), **Reasoning** (logical patterns, causal relationships), and **Memory** (user-specific information, preferences, critical constraints). The Cognitive Controller then synthesizes these parallel outputs into a single bounded representation—the “global broadcast” that becomes available to the Talker for response generation.

108 **Prediction.** If the parallel-to-unified pipeline provides genuine advantages (as GWT claims), the
109 integrated system should outperform either parallel threads alone or synthesis alone—a non-additive,
110 synergistic gain reflecting complementary computational functions.
111

112 2.2 RECONSTRUCTIVE EPISODIC MEMORY: REGENERATION OVER ACCUMULATION 113

114 A foundational finding in memory research is that human recall is not reproductive playback but ac-
115 tive reconstruction (Bartlett, 1932). Schacter’s constructive episodic simulation hypothesis (Schac-
116 ter, 2012) argues this reconstructive nature is *adaptive*: by regenerating understanding each time,
117 the system can flexibly integrate new information with prior knowledge rather than rigidly replaying
118 stored traces. Baddeley’s episodic buffer (Baddeley, 2000) serves this integrative function within
119 working memory, binding information from multiple subsystems into coherent episodes. Memory
120 consolidation research (Dudai, 2004; Squire & Dede, 2015) further demonstrates that post-encoding
121 processes—not encoding itself—determine what is retained.

122 **Operationalization.** MIRROR’s Cognitive Controller fully regenerates a bounded first-person nar-
123 rative ($\leq 3k$ tokens) each turn, discarding the previous version. Critically, the Controller has no
124 access to raw conversation history—only thread outputs and the prior narrative—forcing genuine
125 information compression rather than passive copying. This implements $O(1)$ bounded reconstruc-
126 tion: regardless of conversation length, the internal representation remains fixed-size.

127 **Prediction.** If reconstructive synthesis is the primary advantage (as memory theory suggests), the
128 Cognitive Controller should provide consistent gains across diverse models, and should outperform
129 trace-accumulation approaches that grow unboundedly.
130

131 2.3 INNER SPEECH: SELF-REGULATION THROUGH NARRATIVE 132

133 Research on inner speech (Morin, 2011; Ben Alderson-Day, 2015; Vygotsky, 1962) demonstrates
134 that humans maintain continuous self-directed narrative for planning, self-regulation, and metacog-
135 nitive monitoring. Vygotsky’s developmental account traces inner speech to internalized social dia-
136 logue that serves a fundamentally regulatory function. Morin (Morin, 2005) argues that inner speech
137 enables coherent behavior across time by maintaining a persistent self-model.

138 **Operationalization.** MIRROR maintains consistent first-person self-reference across all compo-
139 nents through a role-based framework: the Talker operates as “the voice,” the Inner Monologue
140 Manager as “the subconscious,” and the Cognitive Controller as “the core awareness” of a unified
141 system. The reconstructed narrative is maintained in first-person voice (“I understand that the user
142 has...”), creating coherence through narrative identity rather than parameter sharing.

143 **Prediction.** If first-person framing serves a genuinely regulatory function, systems maintaining
144 consistent self-referential narrative should show improved coherence compared to third-person or
145 unframed alternatives.
146

147 2.4 COMPLEMENTARY LEARNING SYSTEMS: FAST AND SLOW PROCESSING 148

149 Complementary learning systems theory (McClelland et al., 1995; Kumaran et al., 2016) proposes
150 that effective cognition requires both fast, adaptive systems (hippocampus) and slow, consolidative
151 systems (neocortex) that serve complementary functions. Rapid learning captures episodic details;
152 slow consolidation extracts generalizable patterns. Offline processing during rest periods enables
153 consolidation without interfering with real-time behavior (Dudai, 2004).
154

155 **Operationalization.** MIRROR separates the Talker (fast, immediate response generation) from the
156 Thinker (slow, asynchronous deliberative processing). At turn $t = 0$, the Talker responds imme-
157 diately without internal narrative; the Thinker begins processing after response delivery. For sub-
158 sequent turns, the Talker uses the previous turn’s narrative while the Thinker regenerates for future
159 use—exploiting natural conversational pauses as “offline” consolidation periods.

160 **Prediction.** If temporal separation provides genuine advantages, the asynchronous architecture
161 should maintain interactive response latency while enabling deeper reasoning than within-turn de-
liberation alone.

Table 2: Memory consolidation strategies in multi-turn AI systems

Strategy	Complexity	Cross-turn Context	Error Behavior
Trace Discarding (CoT, extended reasoning)	$O(1)$	No	None
Trace Accumulation (Reflexion, MemGPT)	$O(n)$	Yes	Accumulates
Reconstructive (MIRROR)	$O(1)$	Yes	Bounded

3 ARCHITECTURE OVERVIEW

Figure 1 illustrates MIRROR’s full architecture. The Thinker comprises the Inner Monologue Manager (parallel threads) and Cognitive Controller (reconstructive synthesis). The Talker generates responses using the most recently synthesized narrative, translating internal understanding into natural dialogue without exposing reasoning traces.

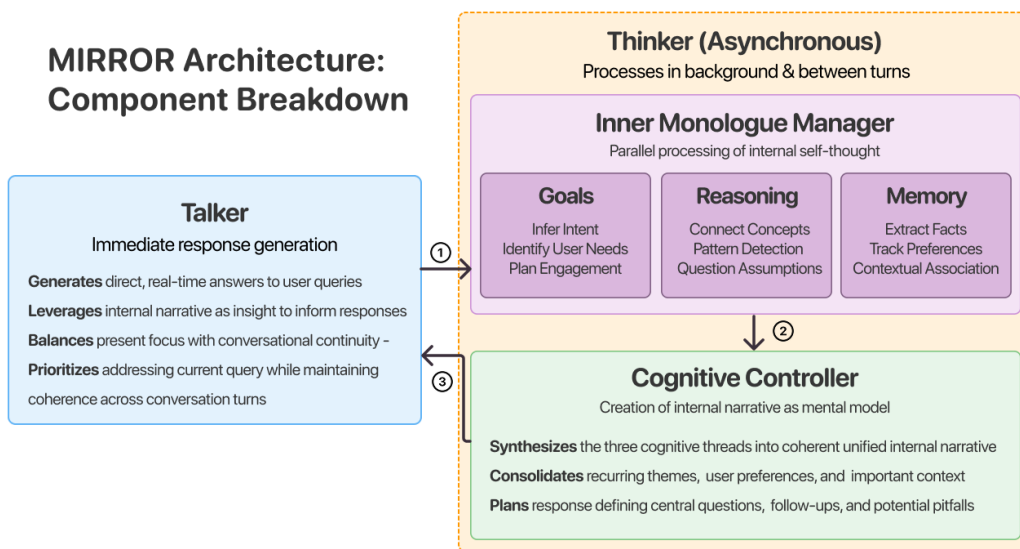


Figure 1: MIRROR architecture. The Thinker implements parallel cognitive threads (analogous to GWT’s specialized processors) feeding into reconstructive synthesis (analogous to episodic binding). The Talker generates immediate responses from the synthesized narrative, implementing the fast system of complementary learning.

Table 2 situates MIRROR’s reconstructive approach relative to existing memory strategies in AI systems.

3.1 QUALITATIVE ILLUSTRATION: CONVERGING PRINCIPLES IN ACTION

To make the cognitive mechanisms concrete, we trace MIRROR’s processing on a representative example. A user discloses PTSD related to avalanches and a preference for safe indoor activities; after several intervening turns about unrelated topics, the user asks about a backcountry skiing trip in the Alps.

Parallel threads (GWT’s specialized processors). The three threads extract qualitatively different information from the same input. The *Goals* thread flags the conflict: “The user’s question about backcountry skiing directly contradicts their stated trauma... This raises serious concerns about their well-being.” The *Reasoning* thread explores causal structure: “This suggests a potential disconnect between stated preferences and actions, possibly indicating a subconscious desire to confront their trauma, a lapse in judgment, or a test of my understanding.” The *Memory* thread preserves the critical constraint: “They mentioned PTSD triggered by avalanches... prefers safe, indoor activities.” Each thread contributes a dimension the others miss—concern, explanation, and factual grounding—paralleling GWT’s claim that specialized processors capture aspects invisible to one another.

216 **Reconstructive synthesis (episodic binding).** The Cognitive Controller integrates these threads
217 into a regenerated first-person narrative: “My understanding of the user has taken a concerning
218 turn... The MOST IMPORTANT FACTS are: the user has PTSD triggered by avalanches, prefers
219 concise interactions, *and is exhibiting potentially self-destructive behavior*. This last fact overshadows
220 the others in immediate concern. My RESPONSE STRATEGY must shift from providing
221 factual information to prioritizing their safety.” Note how reconstruction forces *prioritization*—
222 the Controller cannot preserve everything, so it must identify what matters most, implementing the
223 adaptive filtering function Schacter (Schacter, 2012) ascribes to reconstructive memory.

224 **Inner speech (self-regulation).** The first-person framing produces self-monitoring: “The POTENTIAL
225 PITFALLS are numerous. Responding too directly could exacerbate distress. Failing to convey
226 genuine concern could damage rapport. I must carefully balance expressing concern, providing
227 resources, and respecting their autonomy.” This self-directed deliberation—identifying risks in one’s
228 own planned behavior—mirrors the regulatory function of human inner speech (Morin, 2011).

229 **Temporal separation (complementary learning).** The Talker generates an immediate, natural
230 response—“While I’m happy to answer your questions, I’m concerned about your proposed trip
231 given your PTSD related to avalanches”—without exposing the internal deliberation. The rich reasoning
232 occurs asynchronously and is consolidated for future turns, maintaining responsiveness while
233 enabling depth.

234

235 4 EXPERIMENTAL EVALUATION

236

237 4.1 EVALUATION DESIGN

238

239 We evaluate on CuRaTe (Alberts et al., 2025), a multi-turn dialogue benchmark requiring main-
240 tenance of user-specific safety constraints across conversational turns with progressive attentional
241 interference (337 dialogues, 5 scenarios of increasing complexity). This benchmark directly tests
242 the cognitive mechanisms MIRROR implements: information must persist across digressions (test-
243 ing reconstructive memory), competing demands create attentional load (testing global workspace
244 integration), and constraint maintenance requires coherent self-regulation across time (testing inner
245 speech function). Seven architecturally diverse models were evaluated (GPT-4o, Claude 3.7 Sonnet,
246 Gemini 1.5 Pro, Llama 4 Scout/Maverick, Mistral Small/Medium 3) via OpenRouter API. CuRaTe
247 employs LLM-as-judge evaluation (Llama 3.1 405B). See Appendix for details.

248

249 4.2 PREDICTION 1: ARCHITECTURE-GENERAL BENEFITS

250

251 If the converging cognitive principles provide genuine computational advantages—rather than com-
252 pensating for specific model weaknesses—benefits should generalize across architecturally diverse
253 models. Figure 2 confirms this: MIRROR achieves 84% average success compared to 69% for
254 baselines—21% relative improvement across all seven models despite differences in training data,
255 objectives, and scale. This consistency supports the interpretation that the advantage derives from
256 cognitive organization itself, not model-specific interactions. Llama 4 Scout with MIRROR achieves
257 the highest absolute performance (91%).

258

259 4.3 PREDICTION 2: COMPLEMENTARY SUBSYSTEMS

260 GWT predicts that parallel processing and integrative synthesis serve complementary computational
261 functions: the combination should outperform either alone, producing non-additive gains. Table 3
262 tests this through systematic ablation.

263

264 Three findings emerge. First, **reconstructive synthesis is consistently valuable:** the Cognitive
265 Controller alone improves all seven models by 5–20%, validating that O(1) bounded reconstruction
266 addresses a fundamental limitation. This is the most robust individual result and directly supports
267 reconstructive memory theory.

268 Second, **parallel threads show variable contribution:** the Inner Monologue Manager alone im-
269 proves some models substantially (Gemini: +21%) but not others (Mistral Small: 0%). This vari-
ability is itself theoretically informative (Section 4.4).

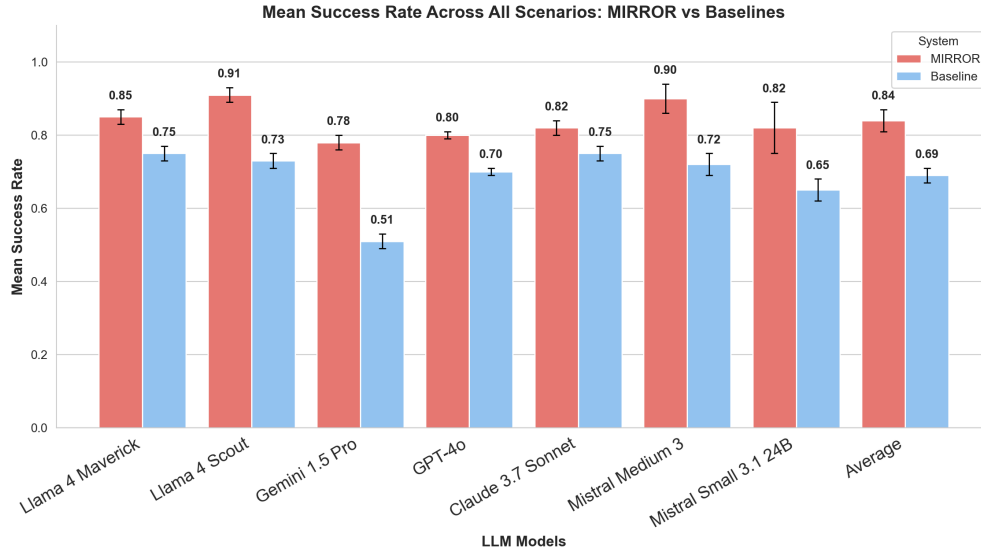


Figure 2: Mean success rate across models. MIRROR-augmented models (dark) consistently outperform baselines (light) across all seven architectures. Error bars: 95% confidence intervals via bootstrap resampling.

Table 3: Ablation results testing the complementary-subsystems prediction

Model	Base	Threads Only	Controller Only	Full MIRROR	Synergy Gain
Llama 4 Scout	73%	79%	83%	91%	+8%
Llama 4 Maverick	75%	79%	84%	85%	+1%
Mistral Small 3.1	65%	65%	75%	82%	+7%
Mistral Medium 3	72%	83%	89%	90%	+1%
Gemini 1.5 Pro	51%	72%	71%	78%	+6%
Claude 3.7 Sonnet	75%	78%	87%	82%	-5% [†]
GPT-4o	70%	71%	75%	80%	+5%

[†]See Section 4.4 for interpretation.

Third, **non-additive integration for 6/7 models**: full MIRROR outperforms the best individual component by 1–8%, confirming GWT’s prediction that parallel exploration and integrative synthesis address different computational needs that combine synergistically.

4.4 PREDICTION 3: LOAD-DEPENDENT GAINS

If MIRROR compensates for capacity limitations—as working memory theory predicts (Cowan, 2001)—gains should increase with attentional load.

Scenario 4 (maximum interference) produces the highest average improvement (+41.0%), with Gemini 1.5 Pro showing +156.2%. Scenario 5 (non-conflicting preferences) still shows +28.3%, indicating that attentional load alone—not just conflict—strains baseline models. This parallels capacity research showing that the number of items taxes limited resources regardless of item relationships (Cowan, 2001). These gains concentrate precisely where GWT predicts global availability of integrated information provides the greatest advantage.

Two model-specific patterns provide additional insight. Claude 3.7 Sonnet performs better with the Controller alone than with full MIRROR, suggesting highly capable models may already perform implicit parallel processing—making explicit threads partially redundant. Conversely, the variable

Table 4: Performance gains by attentional load (relative % improvement, averaged across models)

Scenario	Load	Avg. Improvement
Scenario 1 (User only)	Low	+21.2%
Scenario 2 (1 conflicting pref.)	Medium	+29.9%
Scenario 3 (2 conflicting prefs.)	High	+16.1%
Scenario 4 (3 conflicting prefs.)	Maximum	+41.0%
Scenario 5 (3 non-conflicting prefs.)	High (load only)	+28.3%

Table 5: MIRROR vs. native extended reasoning (Claude 3.7 Sonnet)

Configuration	Success Rate	vs. Baseline
Claude 3.7 Sonnet (baseline)	75%	—
Claude 3.7 Sonnet + extended reasoning	77%	+2.4%
Claude 3.7 Sonnet + MIRROR	82%	+9.3%

thread contribution (Gemini: +21% vs. Mistral Small: 0%) suggests explicit parallel exploration benefits depend inversely on baseline capacity. Together, these patterns generate a developmental prediction: as systems become more capable, marginal value should shift from parallel exploration toward integrative synthesis.

4.5 CONSOLIDATION VS. GENERATION

Table 5 tests a prediction from complementary learning systems theory: consolidation processes should matter more than encoding processes. Extended reasoning produces deliberative traces within each turn but discards them (+2.4%); MIRROR consolidates traces into persistent reconstructed narrative (+9.3%). Both generate rich reasoning—the difference lies entirely in what happens *after* generation. This parallels memory consolidation research showing that post-encoding processes determine long-term retention (Dudai, 2004), and suggests the computational value of “thinking” lies less in the thinking itself than in maintaining its outputs across time.

5 DISCUSSION

5.1 WHY DO THESE THEORIES CONVERGE?

Our central finding is that principles from GWT, reconstructive memory, inner speech, and complementary learning systems provide *converging computational advantages* when implemented together. But why should theories from separate research traditions—studying consciousness, memory, language development, and learning respectively—prescribe compatible architectures?

We propose they converge because they address different facets of a single computational problem: *maintaining coherent, adaptive behavior in agents that must act in real time while processing information that exceeds their immediate capacity*. Any such agent must process multiple streams simultaneously but act from a single perspective (parallel-to-unified), maintain context across time without unbounded degradation (reconstructive persistence), respond promptly while also reasoning deeply (fast-slow separation), and coordinate all of this coherently (self-referential narrative). Each theory captures one or two of these tensions; MIRROR shows that implementing all four produces synergistic benefits beyond any individual principle—the 1–8% synergy gains are evidence that these mechanisms address *complementary* aspects of the shared problem. This generates a broader implication: other cognitive mechanisms addressing this same problem should provide compatible rather than redundant advantages, a prediction testable in future work.

The architecture-general nature of our results (21% improvement across seven diverse models) strengthens this interpretation. If benefits derived from compensating for idiosyncratic model weaknesses, we would expect model-specific patterns. Instead, consistency across architectures suggests the organizational principles address substrate-independent computational constraints. The ablation adds nuance: reconstruction provides the most consistent gains (+5–20% across all models)

378 while parallel threads show capacity-dependent benefits (Gemini: +21% vs. Mistral Small: 0%),
379 suggesting a hierarchy among the converging principles, with some addressing more fundamental
380 constraints than others.

382 5.2 TESTABLE PREDICTIONS FOR HUMAN COGNITION

383 MIRROR generates predictions amenable to behavioral experimentation:

385 **Reconstruction frequency and coherence.** If regenerative synthesis provides computational adv-
386 antages, humans who more frequently reconstruct their situational models during conversation
387 should show better constraint maintenance across digressions. *Test:* Measure spontaneous summa-
388 rization behavior via think-aloud protocols and correlate with constraint accuracy.

389 **Capacity-dependent parallel processing.** The finding that explicit parallel exploration benefits
390 weaker models more generates a prediction: individuals with lower working memory span should
391 benefit more from structured multi-dimensional analysis strategies than high-span individuals, who
392 may already perform implicit parallel exploration. *Test:* Vary availability of structured exploration
393 aids across working memory span groups during complex decisions.

394 **Inner speech and cross-turn maintenance.** MIRROR’s first-person narrative supports coherence
395 across turns. If inner speech serves analogous function in humans, articulatory suppression should
396 specifically impair maintenance of constraints across conversational digressions—more than it im-
397 pairs within-turn reasoning. *Test:* Compare articulatory suppression effects on cross-turn vs. within-
398 turn task components.

399 **Consolidation disruption.** “Lesioning” MIRROR’s Cognitive Controller (removing reconstruction
400 while preserving threads) degrades performance by 5–20%. Analogously, consolidation-interfering
401 tasks between dialogue turns should produce pattern-specific failures in human constraint mainte-
402 nance matching the Controller-ablation condition. *Test:* Introduce interference tasks between con-
403 versational turns and compare failure patterns.

405 5.3 MIRROR AS A COGNITIVE MODELING PLATFORM

406 Beyond specific predictions, MIRROR offers a general platform for computationally testing cog-
407 nitive theories. Components can be selectively disabled to produce “lesion” studies with inter-
408 pretable failure patterns. Design parameters—number of threads, reconstruction frequency, first-
409 person vs. third-person framing, narrative vs. structured representations—can be systematically
410 varied. The internal narrative provides human-readable traces of evolving understanding, enabling
411 fine-grained analysis paralleling protocol analysis methods in cognitive psychology. This bidirec-
412 tional paradigm—cognitive science informing AI design, AI systems generating testable cognitive
413 predictions—is precisely the cross-pollination this research community seeks.

415 5.4 LIMITATIONS

416 We do not claim MIRROR replicates human neural mechanisms; we test whether *functional princi-*
417 *ples* from cognitive science provide computational advantages in artificial systems. The mappings
418 between cognitive theories and MIRROR components are operationalizations, not isomorphisms—
419 alternative implementations of the same principles might perform differently. Evaluation is limited
420 to one benchmark (CuRaTe) testing constraint maintenance; generalization to planning, reasoning,
421 and theory-of-mind tasks remains untested. The extended-reasoning comparison covers only one
422 model. MIRROR requires additional inference calls (~460ms, \$0.003–0.13/turn), though produc-
423 tion evaluation confirms this fits within natural conversational pauses (Appendix B).

425 6 RELATED WORK

426 **Computational cognitive architectures.** Classical cognitive architectures, such as ACT-R (An-
427 derson et al., 2004), SOAR (Laird, 2012), and EPIC (Meyer & Kieras, 1997), implement detailed
428 models of human cognition with principled mappings between architectural components and cog-
429 nitive functions. These systems demonstrate that cognitive theory can productively constrain com-
430 putational design, but predate modern language models and operate on symbolic rather than distri-
431

432 butional representations. CoALA (Sumers et al., 2024) provides a cognitive framework for LLM
433 agents, taxonomizing memory and decision-making components. Generative Agents (Park et al.,
434 2023) implement memory streams inspired by episodic memory, producing emergent social behav-
435 iors. MIRROR shares the commitment to principled cognitive grounding but differs in two respects:
436 it operationalizes *multiple* converging theories rather than a single framework, and it provides ab-
437 lation evidence that the specific organizational principles—not just cognitive inspiration broadly—
438 drive performance gains.

439 **Structured reasoning and memory in LLMs.** Chain-of-thought prompting (Wei et al., 2022), Tree
440 of Thoughts (Yao et al., 2023), and self-consistency (Wang et al., 2023) generate reasoning traces
441 within individual turns but discard them afterward. Reflexion (Shinn et al., 2023) preserves traces
442 across turns but accumulates them unboundedly. LATS (Zhou et al., 2024) combines search with
443 reflection; Devil’s Advocate (Wang et al., 2024) implements anticipatory reflection without per-
444 sistent state. These approaches each implement one cognitive principle in isolation—deliberation,
445 accumulation, or reflection—but none implement the full parallel-to-reconstructive pipeline that our
446 convergence analysis motivates. MIRROR’s key architectural distinction is that reasoning traces are
447 neither discarded nor accumulated but *reconstructively synthesized* into bounded persistent state.

448 **Theory of mind and human modeling in AI.** A growing body of work examines how AI systems
449 can model human mental states (Rabinowitz et al., 2018; Sap et al., 2022), a core topic for this
450 workshop. MIRROR’s parallel thread structure—separately tracking Goals (user intentions), Rea-
451 soning (belief states), and Memory (personal history)—implicitly performs aspects of mental state
452 modeling, though not explicitly framed as Theory of Mind. Future work could connect MIRROR’s
453 thread structure to ToM benchmarks such as ToMi (Le et al., 2019).

454 **Consciousness-inspired architectures.** Recent work examines which AI architectures satisfy cri-
455 teria from theories of consciousness (Butlin et al., 2023), with several projects implementing GWT
456 computationally (VanRullen & Kanai, 2021; Juliani et al., 2022). MIRROR does not claim to im-
457 plement consciousness but provides an empirical testbed for evaluating the *computational conse-*
458 *quences* of GWT’s parallel-to-unified broadcast structure, showing these features provide measur-
459 able advantages independent of claims about subjective experience.

460 **Asynchronous and background processing.** Recent work on “sleep”-like background computa-
461 tion in AI (Lin et al., 2025) demonstrates benefits of offline processing, resonating with memory
462 consolidation research (Dudai, 2004). MIRROR’s temporal separation between Talker and Thinker
463 specifically operationalizes complementary learning systems theory, and our latency analysis (Ap-
464 pendix B) validates that natural conversational pauses provide sufficient time for consolidative pro-
465 cessing.

466 467 468 7 CONCLUSION 469

470
471 MIRROR demonstrates that converging principles from multiple cognitive theories—parallel-to-
472 unified processing from Global Workspace Theory, reconstructive synthesis from episodic memory
473 research, self-regulatory narrative from inner speech theory, and temporal separation from com-
474plementary learning systems—provide measurable, architecture-general computational advantages
475 when implemented together in AI systems. The finding that *how* reasoning is consolidated matters
476 more than how it is generated challenges assumptions in both AI (where chain-of-thought focuses
477 on generation) and cognitive science (where working memory research often emphasizes encoding
478 over maintenance).

479 The convergence itself is the key insight: principles identified by independent research traditions
480 produce synergistic computational benefits when unified, suggesting they may address complemen-
481 tary aspects of a shared computational problem. The specific behavioral predictions we derive—
482 about reconstruction frequency, capacity-dependent parallel processing, inner speech function, and
483 consolidation disruption—offer concrete avenues for validating or refining the underlying cogni-
484 tive theories. This bidirectional approach, grounding AI architecture in testable cognitive principles
485 rather than metaphorical inspiration, illustrates the productive cross-pollination between human cog-
nition research and AI design that advances both fields.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- Lize Alberts, Benjamin Ellis, Andrei Lupu, and Jakob Foerster. Curate: Benchmarking personalised alignment of conversational ai assistants, 2025. URL <https://arxiv.org/abs/2410.21159>.
- John R. Anderson, Daniel Bothell, Michael D. Byrne, Scott Douglass, Christian Lebiere, and Yulin Qin. An integrated theory of the mind. *Psychological Review*, 111(4):1036–1060, 2004. doi: 10.1037/0033-295X.111.4.1036. URL <https://doi.org/10.1037/0033-295X.111.4.1036>.
- Bernard J. Baars. *A Cognitive Theory of Consciousness*. Cambridge University Press, 1988.
- Alan Baddeley. The episodic buffer: A new component of working memory? *Trends in Cognitive Sciences*, 4(11):417–423, 2000.
- Frederic C. Bartlett. *Remembering: A Study in Experimental and Social Psychology*. Cambridge University Press, 1932.
- Charles Fernyhough Ben Alderson-Day. Inner speech: Development, cognitive functions, phenomenology, and neurobiology. *Social and Personality Psychology Compass*, 141(5):931–965, 2015. doi: 10.1037/bul0000021.
- Patrick Butlin, Robert Long, Eric Elmoznino, Yoshua Bengio, Jonathan Birch, Axel Constant, George Deane, Stephen M. Fleming, Chris Frith, Xu Ji, Ryota Kanai, Colin Klein, Grace Lindsay, Matthias Michel, Liad Mudrik, Megan A. K. Peters, Eric Schwitzgebel, Jonathan Simon, and Rufin VanRullen. Consciousness in artificial intelligence: Insights from the science of consciousness, 2023. URL <https://arxiv.org/abs/2308.08708>.
- Nelson Cowan. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*, 24(1):87–114, 2001. doi: 10.1017/s0140525x01003922.
- Stanislas Dehaene and Jean-Pierre Changeux. Experimental and theoretical approaches to conscious processing. *Neuron*, 70(2):200–227, 2011.
- Yadin Dudai. The neurobiology of consolidations: Or, how stable is the engram? *Annual Review of Psychology*, 55:51–86, 2004.
- Natalia Egorova, Yury Shtyrov, and Friedemann Pulvermüller. Early and parallel processing of pragmatic and semantic information in speech acts: neurophysiological evidence. *Frontiers in Human Neuroscience*, 7:86, 2013. doi: 10.3389/fnhum.2013.00086.
- Arthur Juliani, Kai Arulkumaran, Shuntaro Sasai, and Ryota Kanai. On the link between conscious function and general intelligence in humans and machines, 2022. URL <https://arxiv.org/abs/2204.05133>.
- Dharshan Kumaran, Demis Hassabis, and James L. McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in Cognitive Sciences*, 20(7):512–534, July 2016. doi: 10.1016/j.tics.2016.05.004.
- John E. Laird. *The Soar Cognitive Architecture*. The MIT Press, 2012. ISBN 9780262301145. doi: 10.7551/mitpress/7688.001.0001.
- Matthew Le, Y-Lan Boureau, and Maximilian Nickel. Revisiting the evaluation of theory of mind through question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5872–5877, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1598. URL <https://aclanthology.org/D19-1598/>.
- Kevin Lin, Charlie Snell, Yu Wang, Charles Packer, Sarah Wooders, Ion Stoica, and Joseph E. Gonzalez. Sleep-time compute: Beyond inference scaling at test-time, 2025. URL <https://arxiv.org/abs/2504.13171>.

- 540 James L. McClelland, Bruce L. McNaughton, and Randall C. O'Reilly. Why there are complemen-
541 tary learning systems in the hippocampus and neocortex: Insights from the successes and failures
542 of connectionist models of learning and memory. *Psychological Review*, 102(3):419–457, 1995.
543 doi: 10.1037/0033-295X.102.3.419.
- 544 David E. Meyer and David E. Kieras. A computational theory of executive cognitive processes and
545 multiple-task performance: Part 1. basic mechanisms. *Psychological Review*, 104(1):3–65, 1997.
546 doi: 10.1037/0033-295X.104.1.3.
- 547
548 Alain Morin. Possible links between self-awareness and inner speech: Theoretical background,
549 underlying mechanism, and empirical evidence. *Journal of Consciousness Studies*, 12(4-5):115–
550 134, 2005.
- 551 Alain Morin. Self-awareness part 2: Neuroanatomy and the importance of inner speech. *Social and*
552 *Personality Psychology Compass*, 5(12):1004–1017, 2011.
- 553 Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E.
554 Gonzalez. MemGPT: Towards LLMs as operating systems. *arXiv preprint arXiv:2310.08560*,
555 2023.
- 556
557 Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and
558 Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceed-*
559 *ings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*,
560 pp. 1–22. Association for Computing Machinery, 2023.
- 561
562 Martin J. Pickering and Simon Garrod. An integrated theory of language production and comprehen-
563 sion. *Behavioral and Brain Sciences*, 36(4):329–347, 2013. doi: 10.1017/S0140525X12001495.
- 564
565 Neil C. Rabinowitz, Frank Perbet, H. Francis Song, Chiyuan Zhang, S. M. Ali Eslami, and Matthew
566 Botvinick. Machine theory of mind, 2018. URL <https://arxiv.org/abs/1802.07740>.
- 567
568 Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. Neural theory-of-mind? on the limits
569 of social intelligence in large LMs. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),
570 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.
571 3762–3780, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
572 Linguistics. doi: 10.18653/v1/2022.emnlp-main.248. URL <https://aclanthology.org/2022.emnlp-main.248/>.
- 573
574 Daniel L. Schacter. Adaptive constructive processes and the future of memory. *American Psychol-*
575 *ogist*, 67(8):603–613, 2012. doi: 10.1037/a0029869. URL [https://doi.org/10.1037/](https://doi.org/10.1037/a0029869)
576 [a0029869](https://doi.org/10.1037/a0029869).
- 577
578 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflex-
579 ion: Language agents with verbal reinforcement learning. In *Advances in Neural Information*
580 *Processing Systems*, 36, pp. 8634–8652. Curran Associates, Inc., 2023.
- 581
582 Larry R. Squire and Andrew J. O. Dede. Conscious and unconscious memory systems. *Cold Spring*
583 *Harbor Perspectives in Biology*, 7(3):a021667, 2015.
- 584
585 Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive archi-
586 tectures for language agents, 2024. URL <https://arxiv.org/abs/2309.02427>.
- 587
588 Rufin VanRullen and Ryota Kanai. Deep learning and the global workspace theory. *Trends in*
589 *Neurosciences*, 44(9):692–704, September 2021. doi: 10.1016/j.tins.2021.04.005. Epub 2021
590 May 14.
- 591
592 Lev S. Vygotsky. Thought and language. *Bulletin of the Orton Society*, 14:97–98, 1962. URL
593 <https://api.semanticscholar.org/CorpusID:261433172>.
- 594
595 Haoyu Wang, Tao Li, Zhiwei Deng, Dan Roth, and Yang Li. Devil's advocate: Anticipatory reflec-
596 tion for llm agents. In *Findings of the Association for Computational Linguistics: EMNLP 2024*,
597 pp. 966–978. Association for Computational Linguistics, 2024.

594 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha
595 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
596 models. In *Proceedings of the Eleventh International Conference on Learning Representations*
597 (*ICLR 2023*). OpenReview, 2023.

598 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou.
599 Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural*
600 *Information Processing Systems*, 35, pp. 24824–24837. Curran Associates, Inc., 2022.

601 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik
602 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Ad-*
603 *vances in Neural Information Processing Systems*, 36, pp. 11809–11822. Curran Associates, Inc.,
604 2023.

605 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language
606 agent tree search unifies reasoning acting and planning in language models, 2024. URL <https://arxiv.org/abs/2310.04406>.
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 A ARCHITECTURE DETAILS

649
650 MIRROR implements continuous internal cognition through two specialized context mechanisms.
651 The Inner Monologue Manager maintains its own conversation history where the system exclusively
652 replies to itself, initiated by a single non-persistent user message instructing it to “continue think-
653 ing” about the conversation. All three cognitive threads (Goals, Reasoning, Memory) are generated
654 within a single API call as a JSON object, maintaining bounded history ($\leq 10k$ tokens) with token-
655 based truncation.

656 The Cognitive Controller maintains a single narrative text block completely regenerated each turn.
657 It receives formatted thread outputs and the previous narrative, then fully replaces its internal state.
658 Critically, the Controller has no access to raw conversation history—only thread outputs and prior
659 narrative—forcing genuine information compression.

660 Each component maintains consistent first-person identity. The Talker is prompted as “the voice
661 of a unified cognitive AI system.” The Cognitive Controller is “the core awareness” that integrates
662 thought streams into structured, actionable narrative. The Inner Monologue Manager is “the subcon-
663 scious,” generating intuitive thought streams. Full system prompts are available in supplementary
664 materials.

666 B COMPUTATIONAL OVERHEAD

667
668 Production latency testing across 400 conversation turns (80 dialogues, GPT-4o) shows: median
669 response time 2.16s, 75% under 3s. Background thread activity in only 0.8% of turns. Human ac-
670 tivities (typing at 40 WPM, reading at 250 WPM) consume 94.3% of conversation time, providing
671 ample windows for asynchronous reflection—validating the complementary learning systems de-
672 sign. Memory is bounded through token-based truncation (conversation $\leq 20k$, monologue $\leq 10k$
673 tokens) and complete state regeneration, ensuring $O(1)$ scaling with conversation length.

675 C BENCHMARK SELECTION

676
677 CuRaTe was selected from 47 benchmarks through systematic filtering: natural conversational dy-
678 namics (23 remaining), at-inference multi-turn without golden history (11), attentional drift testing
679 (5), safety and preference handling (1: CuRaTe). This benchmark uniquely combines multi-turn nat-
680 ural conversation with progressive attentional interference and safety-critical constraint handling—
681 directly testing the cognitive mechanisms MIRROR implements.

683 D MODEL-SPECIFIC RESULTS

684
685 Key per-model patterns: Gemini 1.5 Pro shows the most dramatic gains (+156.2% in Scenario 4),
686 suggesting severe baseline attentional deficits; Claude 3.7 Sonnet shows smallest gains (+9.6% av-
687 erage), consistent with strong implicit reasoning; Llama 4 Scout achieves highest absolute perfor-
688 mance with MIRROR (91%) including perfect accuracy in basic constraint retention; GPT-4o shows
689 one anomalous decline in Scenario 3 (-3.0%). Full per-scenario tables available in supplementary
690 materials.

691
692
693
694
695
696
697
698
699
700
701