

Reader's Guide to Learning Through Narratives Series

Overview

This three-part work proposes **episodic compression** as an alternative approach to sample-efficient machine learning. It spans machine learning architecture, cognitive science, learning theory, and deployment engineering. This guide helps readers from different disciplines understand what claims are being made in their domain and how the cross-disciplinary synthesis works.

For Machine Learning Researchers

What we're proposing:

- An alternative learning architecture that uses salience-driven compression instead of gradient descent
- LLMs serve as computational infrastructure (attention, reasoning, instruction-following), not as the learning mechanism
- Two-layer architecture (Intuition/Executive) for scaling beyond single-context capacity
- No parameter training required - learning happens through consolidation (reasoning over episodes)

What we're claiming:

- Sample-efficient learning is possible without embodiment (challenging LeCun's JEPA necessity claim)
- Episodic compression mechanism works at toy scale (Shimmer Valleys validation)
- Architecture addresses scaling challenges when knowledge exceeds context windows
- Continuous learning during deployment without catastrophic forgetting

What we're NOT claiming:

- Superior performance to all alternatives on all benchmarks
- Ready for production deployment without further validation
- Optimal parameters (most are explicitly marked as testable hypotheses)
- That this is the only or best approach to sample efficiency

What you won't find here:

- Ablation studies (no baseline system to ablate from)
- Standard benchmark comparisons (requires full implementation)
- Formal complexity analysis

- Large-scale empirical validation (acknowledged as future work)

What you should evaluate:

- Does the architectural design make computational sense?
- Is the Shimmer Valleys validation appropriately scoped?
- Are scaling challenges correctly identified?
- Are the proposed solutions technically feasible?

For Cognitive Scientists

What we're using from cognitive science:

- Observation that humans learn extensively from narratives (indirect experience)
- Episodic memory as compressed salient experiences
- Pattern extraction from multiple related episodes
- Dual-process distinctions (automatic vs controlled processing)

What we're claiming:

- These cognitive patterns suggest useful computational principles
- Narrative format enables efficient representation of causal structure
- Functional separations (pattern-matching vs reasoning) address computational challenges

What we're NOT claiming:

- This is how human cognition actually works mechanistically
- The architecture mirrors biological implementation
- We can explain human learning with this model
- Cognitive neuroscience validates our specific design choices

What you should evaluate:

- Are the cognitive parallels reasonable as inspiration?
- Does the narrative learning observation hold empirically?
- Are we appropriately cautious about cognitive claims?
- Do the functional separations make computational sense regardless of biological correspondence?

For Neuroscientists

What we're referencing:

- Functional distinctions between brain systems (automatic vs controlled processing)
- Episodic vs semantic memory systems (Tulving)

- Neurological case studies showing separability of functions
- Multi-timescale operation in biological learning

What we're claiming:

- Biology converged on functional separations that prove computationally useful
- These parallels informed our architectural choices
- Separability in neurological cases validates that functions are genuinely distinct

What we're NOT claiming:

- Our implementation mirrors biological mechanisms
- We can explain neural substrates
- The architecture makes predictions about neuroscience
- Biological correspondence validates our approach

What you should evaluate:

- Are the functional analogies appropriate (even if implementations differ)?
- Do we overreach on biological claims?
- Are the parallels useful pedagogically?
- Is the disclaimer about biological correspondence adequate?

For AI Safety/Alignment Researchers

What we're proposing:

- Continuous learning during deployment enables adaptation to specific contexts
- Transparent knowledge representations (human-readable generalizations)
- Additive learning reduces catastrophic forgetting risk
- Confidence tracking and uncertainty expression

What we're claiming:

- Deployment learning creates personalization without explicit configuration
- Interpretability advantages over parameter-based learning
- Different safety profile than traditional continual learning

What we're NOT claiming:

- Complete solution to alignment problems
- Adversarial robustness (explicitly identified as open question)
- Formal safety guarantees
- Solved value learning

What you should evaluate:

- Does additive learning genuinely reduce catastrophic forgetting?
- Are there novel safety risks we haven't addressed?
- Does interpretability actually enable better safety practices?
- What adversarial attacks might work against episodic compression?

For Systems Engineers/Practitioners

What we're proposing:

- Practical architecture using existing pre-trained models
- No gradient descent or parameter optimization required
- Consolidation through inference (not training)
- Specific implementation guidance (model selection, storage formats, scheduling)

What we're claiming:

- Deployable with current technology
- Cost-efficient through selective engagement of expensive models
- Immediate deployment without training cycles
- Transparent operation through structured knowledge

What we're NOT claiming:

- Production-ready without further engineering
- Optimal parameters identified (marked as testable hypotheses)
- All edge cases handled
- Cheaper than all alternatives in all scenarios

What you should evaluate:

- Is the implementation feasible with described components?
- Are the cost projections reasonable?
- Are critical engineering challenges identified?
- Does the architecture handle realistic operational constraints?

For Philosophers of Mind/Cognitive Scientists (Conceptual)

Our terminological stance:

- "Learning," "understanding," "reasoning" used as behavioral shorthand
- No claims about machine consciousness or internal mental states
- Focus on behavioral competence, not phenomenology
- "The system produces behaviors consistent with having learned X"

What we're claiming:

- Behavioral competence is sufficient for practical AI systems
- Narrative format captures causal structure effectively
- Functional organization matters more than substrate

What we're NOT claiming:

- These systems have genuine understanding
- Consciousness or qualia emerge from this architecture
- We've solved hard problems in philosophy of mind
- Behavioral equivalence implies mental equivalence

What you should evaluate:

- Is our terminological discipline adequate?
- Do we inadvertently make claims about internal states?
- Are the philosophical disclaimers appropriate?
- Does the work contribute to understanding learning independent of consciousness questions?

How the Three Parts Fit Together

Part 1: Learning Through Narratives

- Introduces episodic compression mechanism
- Validates mechanism at toy scale (Shimmer Valleys)
- Establishes LLMs as suitable infrastructure
- Demonstrates sample efficiency without embodiment

Part 2: Scaling Episodic Compression

- Addresses when architecture becomes necessary (at scale)
- Proposes two-layer design (Intuition/Executive)
- Provides implementation specifications
- Distinguishes from alternatives (RAG, JEPA, memory networks)

Part 3: Multi-Timescale Continuous Learning

- Explores emergent temporal properties
- Analyzes catastrophic forgetting implications
- Examines deployment learning dynamics
- Identifies open research questions

Reading Recommendations by Interest

If you want the core idea: Read Part 1 Abstract + Sections 1-2, plus Shimmer Valleys description

If you want implementation details: Read Part 2 Sections 3-4

If you want theoretical implications: Read Part 3 Sections 2-3

If you want research questions: Read Part 1 Section 6, Part 2 Section 6, Part 3 Section 5

If you want to evaluate validity: Read Part 1 Section 3 (experimental design), Part 2 Section 6 (limitations), Part 3 Section 1.4 (scope and claims)

Cross-Disciplinary Synthesis: The Core Contribution

The main contribution is **architectural**: showing how to orchestrate existing LLM capabilities (attention mechanisms, reasoning, instruction-following) into a learning system that operates through salience-driven compression rather than gradient descent.

This synthesis requires understanding:

- **From ML:** What LLMs can do as infrastructure
- **From cognitive science:** What patterns human learning exhibits
- **From learning theory:** What makes learning sample-efficient
- **From neuroscience:** What functional separations prove useful
- **From systems engineering:** What makes systems deployable at scale

No single discipline provides complete evaluation criteria. The work should be evaluated on whether the architectural synthesis is coherent, technically feasible, and addresses real problems - not whether it meets all standards of any single discipline.

What Success Looks Like

Theoretical success:

- The mechanism makes conceptual sense
- Architectural choices address identified challenges
- Claims are appropriately scoped
- Open questions are clearly identified

Empirical success (future work):

- System operates at scale (1000+ generalizations)
- Demonstrates continuous learning over months of deployment
- Achieves sample efficiency in real-world domains
- Provides practical advantages over alternatives

Practical success:

- Implementable with current technology
- Cost-effective at scale
- Solves real deployment problems
- Creates commercial value

This work establishes theoretical foundations and provides architectural specifications. Full validation requires implementation and deployment - explicitly acknowledged as future work.

Common Misunderstandings to Avoid

Misunderstanding: "This claims LLMs learn from narratives because they were trained on narratives" **Reality:** LLMs provide infrastructure; the architecture creates the learning mechanism

Misunderstanding: "This claims to explain how human cognition works" **Reality:** Cognitive parallels are inspirational and pedagogical, not mechanistic claims

Misunderstanding: "Part 1 validated the two-layer architecture" **Reality:** Part 1 validated the mechanism; Part 2 proposes architecture for scaling that mechanism

Misunderstanding: "This is production-ready" **Reality:** This provides architectural design; implementation and validation are future work

Misunderstanding: "'No training' means no engineering required" **Reality:** "No training" means no parameter optimization; architectural engineering is substantial

Misunderstanding: "Specific parameters (3-5 episodes, etc.) are claims" **Reality:** These are testable hypotheses explicitly marked as requiring empirical validation

Questions? Start Here

"Does this actually work?" Part 1 demonstrates the mechanism works at toy scale. Scaling validation is future work.

"Why not just fine-tune existing models?" This enables continuous learning during deployment without parameter updates, avoiding catastrophic forgetting and enabling personalization.

"What's new here vs. RAG/memory networks/meta-learning?" See Part 2 Section 5 for detailed comparisons.

"How does this relate to LeCun's JEPA?" Alternative approach: narrative-based vs embodied, symbolic vs learned representations, immediate deployment vs training from scratch.

"What are the biggest open questions?" See Part 3 Section 5 - consolidation dynamics, memory management, long-term stability, adversarial robustness.

"Who should build on this work?" Anyone interested in: continuous learning systems, interpretable AI, deployment learning, alternative learning architectures, sample-efficient learning without massive datasets.

Final Note

This work is intentionally cross-disciplinary because the problem (sample-efficient learning) requires synthesis across multiple domains. We've tried to be clear about what we're claiming in each domain and what we're not. If you find places where we've overreached or been unclear, that feedback is valuable.

The ultimate test will be implementation: does a system built on these principles actually work at scale? That's future work we're committed to pursuing.

© 2025 Mossrake Group, LLC

Version 1.0