



A Multi-Domain Evaluation of Scaling in a General Episodic Memory

Nate Derbinsky, Justin Li, John E. Laird

University of Michigan



supported by

RESEARCH QUESTION

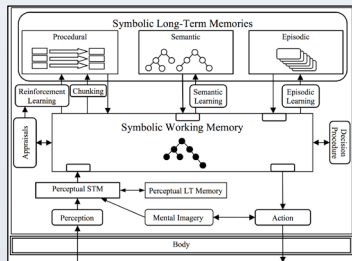
To what extent is Soar's episodic memory effective and efficient for real-time agents that persist for long periods of time across a variety of tasks?

Approach: **Multi-Domain Evaluation**

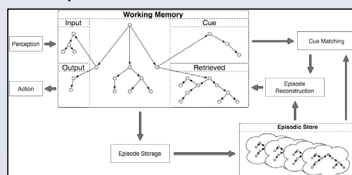
- Existing agents, diverse tasks (49)
- Long runs (hours-days; 10^5 - 10^8 episodes)
- Evaluate at every X episodes: memory, max. cue-matching time (>100 task-relevant cues, 7 general capabilities)

AGENT INTEGRATION

The Soar cognitive architecture



Episodic operations:



Representation

- Episode: connected di-graph
- Store: temporal sequence

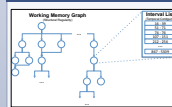
Storage

- Automatic, no dynamics (e.g. forgetting)

Cue Matching

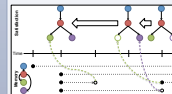
- Cue: acyclic di-graph
- Find the most recent episode that shares the most leaf nodes in common with cue

ALGORITHMIC OVERVIEW



Storage

Encode WM-changes (Δ 's) as temporal intervals in novel dynamic-graph index



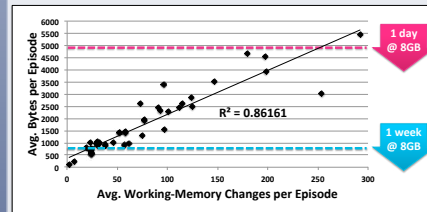
Cue Matching

Reverse temporal walk of cue-relevant Δ 's

- 2-phase search: only graph-match episodes that have *all* cue features independently
- Only evaluate episodes that have changes relevant to cue features via priority queue of b+-tree pointers
- Incrementally re-score episodes via novel dynamic discrimination network

STORAGE CHARACTERIZATION

- Memory scales linearly with Δ 's



RETRIEVAL CHARACTERIZATION

Assumptions

- Temporal Contiguity**
Few changes per episode
- Structural Regularity**
Representational re-use
- Small cues (relative to state size)

Scaling Parameters (w.r.t. cue features)

- Search distance
 - Temporal Selectivity**: Δ frequency
 - Co-Occurrence**: related to |state space|
- Episode scoring
 - Structural Selectivity**: how many ways can cue unify with episode

WORD SENSE DISAMBIGUATION

Experimental Setup

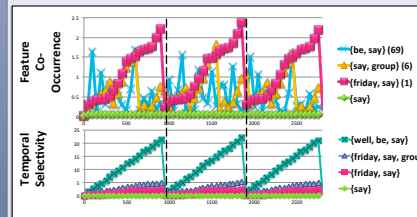
- Input: <"word", POS> + result
 - Corpus: SemCor
- Output: sense
- Agent: maintain context as n-gram
 - on input, query EpMem for context
 - > if success, output next result

Accuracy	Trial #1	Trial #2
2-gram	14.57%	92.82%
3-gram	2.32%	99.47%

Results

- Storage: 234 bytes/ep. (avg)
- Cue-Matching
 - All 1-, 2-, 3-grams <50 msec.
 - 0.2% of 4-grams exceed 50 msec.

Retrieval Time (msec) vs. Episodes (x1000)



PLANNING

Experimental Setup

- 12 domains converted from PDDL
 - Logistics, Blocksworld, Grid, ...
 - 44 problem instances (e.g. # blocks)
- Agent: randomly explore state space
 - 50K episodes, measure every 1K

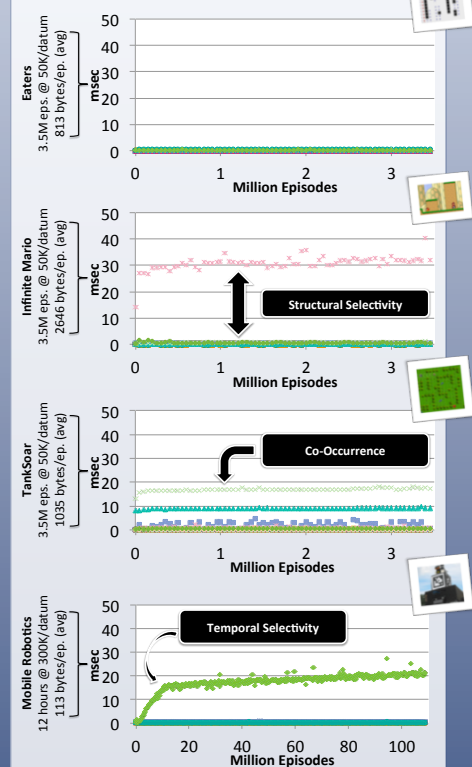
Results

- Storage:
 - Time: <12.04 msec./ep.
 - Memory: 562-5454 bytes/ep.
- Cue-Matching (time <50 msec.)
 - Full State: 12 smallest-state domains
 - Relational: no domains
 - Schema: all domains (max 0.08 msec.)

VIDEO GAMES & MOBILE ROBOTICS

- Hand-coded cues (per domain)

Cue-Matching Time (msec) vs. Episodes



SUMMARY OF RESULTS

Storage

- <50 msec./episode (except in Mario 🎮 due to temporal discontinuity)
- 0.18 – 4 kb/episode (days – months)

Retrieval

- <50 msec. cue matching for many cues
- No cue-matching time growth