

IN-CONTEXT LEARNING CAN PERFORM CONTINUAL LEARNING LIKE HUMANS

Anonymous authors
Paper under double-blind review

ABSTRACT

Large language models (LLMs) can adapt to new tasks via in-context learning (ICL) without parameter updates, making them powerful learning engines for fast adaptation. While extensive research has examined ICL as a few-shot learner, whether it can achieve long-term retention and cross-task knowledge accumulation when multitasks arrive sequentially remains underexplored. Motivated by human memory studies, we investigate the retention characteristics of ICL in multitask settings and extend it to in-context continual learning (ICCL), where continual learning ability emerges through task scheduling and prompt rearrangement. Experiments on Markov-Chain and regbench based benchmarks demonstrate that, for specific large-language models, ICCL benefits from distributed practice (DP) in a manner analogous to humans, consistently revealing a spacing “sweet spot” for retention. Beyond retention performance, we propose a human-retention similarity metric to quantify how closely a continual-learning (CL) method aligns with human retention dynamics. Using this metric, we show that linear-attention models such as MAMBA and RWKV exhibit particularly human-like retention patterns, despite their retention performance lagging behind that of Transformer-based LLMs. Overall, our results establish ICCL as cognitively plausible and practically effective, [providing an inference-only CL paradigm that enables LLMs to maintain long-term retention of a target task under cross-task interference without parameter updates.](#)

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities across a wide range of domains, with in-context learning (ICL) emerging as one of their most distinctive properties (Brown et al., 2020). Via ICL, LLMs adapt new tasks on the fly by conditioning on demonstrations embedded in the prompt without parameter update. This inference-only property turns them into versatile, efficient solvers that cope with zero-shot (Ouyang et al., 2022), few-shot (Brown et al., 2020), and many-shot (Agarwal et al., 2024) regimes, spanning Bayesian inference (Xie et al., 2024), regression (Raventós et al., 2023; Wu et al., 2023), and reinforcement learning (Laskin et al., 2023; Lee et al., 2023; Wang et al., 2025a). Such flexibility suggests that ICL could evolve into a general-purpose learning engine rather than remain a mere few-shot adapter (Wang et al., 2025b; Fan & Liu, 2025). [However, prior work predominantly focuses on single-task or one-off adaptation behavior. Only limited studies investigate multi-task ICL settings, and these works typically optimize short-term prediction by mixing or reordering demonstrations rather than examining whether knowledge can persist across tasks. Thus, it remains unclear whether ICL can support continual learning, where knowledge must be maintained and reused across sequentially presented tasks.](#)

In the broader literature, catastrophic forgetting (CF) has long been recognized as a fundamental challenge in gradient-based continual learning (GBCL). To mitigate CF, a rich body of GBCL has been developed, including regularization-based approaches that reshape the objective function (Li et al., 2020; Evilevitch & Ingram, 2021), replay-based approaches that reorganize the data stream (Rolnick et al., 2019; Buzzega et al., 2021), and approaches that directly modify the optimization process (Arous et al., 2021; Poggio et al., 2011). Although effective on many benchmarks, GBCL approaches fundamentally depend on parameter modification or explicit memory buffers and often struggle to achieve a desirable balance between stability and plasticity. [This naturally raises a broader question: does a similar form of forgetting also arise in parameter-free, inference-only regime of LLMs?](#)

Table 1: Property comparisons among GBCL, ICL and ICCL.

Property	GBCL	ICL	ICCL
Parameter updates	Yes	No	No
Update rules	Gradient Descent	Blackbox Computation	Blackbox Computation
Knowledge carrier	Parameters	Memory/Hidden States	Memory/Hidden States
Task settings	Multiple Tasks	Single Task	Multiple Tasks

In other words, when tasks are presented sequentially within a long context, can a LLM stably retain the mapping of a target task while processing intervening tasks? To investigate this question, we introduce **in-context continual learning (ICCL)** as a framework for studying inference-time retention in LLMs. Unlike multi-task ICL, which jointly learns all tasks, ICCL focuses on how the ordering and distributed scheduling of target-task and intervening-task demonstrations affect the model’s ability to preserve the target-task behavior throughout extended inference, without any parameter updates.

Several recent studies have begun to explore directions related to ICCL, often by extending ICL with prompt modification (e.g., prompt tuning) or external augmentation (e.g., memory retrieval) (Wang et al., 2022; Gao et al., 2023; Momeni et al., 2024; Shinwari & Usama, 2025). While these approaches show promise, they also present several limitations. First, most efforts remain confined to single-task adaptation, without systematically addressing multitask sequential adaptation and cross-task knowledge accumulation. Second, they often rely on heuristic prompt engineering or auxiliary components that raise complexity yet remain vulnerable to memory bloat and retrieval errors, and plateau when procedural memory integration rather than fragmented retrieval is required. Third, current evaluation benchmarks are often either too simple, resulting in trivial zero-shot solutions, or too complex, exceeding the capacity of LLMs to learn within context. As a result, they fail to adequately support systematic analyses of multitask learning dynamics.

In this paper we focus on ICCL, which exploits the LLM’s built-in memory to flexibly control both learning and retention rather than external plug-and-play modules and manually specified learning rules, as summarized in Table 1. To compare GBCL and ICCL and to expose ICCL’s characteristics, we empirically examine the following perspectives: First, in his pioneering studies of human memory, Hermann Ebbinghaus deliberately used “nonsense syllables” in serial-learning tasks (Ebbinghaus, 2013) to strip away prior knowledge and semantic associations, thereby isolating the core mechanisms of memory. Guided by the same principle, we design two benchmarks for evaluating LLMs whose parameters already encode vast human knowledge: randomly generated Markov chains and regular language generated by random finite automata. These chains preserve essential dynamical-system properties (Attal, 2010; Froyland, 2001) while minimizing contamination from pre-existing knowledge. Second, cognitive science shows that structured practice schedules enhance long-term memory of human (Cepeda et al., 2006; 2008; Pavlik Jr & Anderson, 2005; Bothell, 2020), underscoring the need to introduce task scheduling in machine continual learning. Guided by this, we inject explicit task identifiers and cognition-inspired schedules to systematically probe how scheduling and prompts shape the retention characteristics of both LLMs and GBCL. Those includes single practice (SP), multiple practice (MP), and distributed practice (DP). Experiments show that ICCL retains information more robustly than GBCL baselines and consistently privileges DP over SP and MP, mirroring cognitive findings on human memory. Moreover, linear-attention models such as Mamba and RWKV exhibit retention dynamics that more closely parallel human behavior, even though their absolute retention lags behind standard self-attention. Our contributions are threefold:

- ICCL provides, to our knowledge, the first framework that systematically investigates whether forgetting arises purely at inference time for LLMs and how to enable them to maintain long-term retention of a target task under cross-task interference without parameter updates.
- We show that ICCL benefits from DP in a manner analogous to humans, consistently exhibiting a spacing sweet spot that enhances retention and provides a more favorable balance between stability and plasticity compared to GBCL baselines.

[†]Corresponding to:

- We introduce evaluation benchmarks that combine retention performance with human-retention similarity, enabling systematic analysis of multitask retention and revealing that linear-attention models (e.g., MAMBA, RWKV) display particularly human-like retention patterns, despite lower absolute performance than Transformer-based LLMs.

2 METHODOLOGY

2.1 FORMALIZATION OF ICCL

We now formalize the general protocol of ICCL. A task is denoted by $\tau : y_i^T \sim p_\tau(\cdot|x_i)$, which specifies a conditional distribution mapping an input x_i to an output y_i^T . A segment of φ experiences on task τ is represented as $\mathcal{D}_\varphi^\tau = \{x_1^\tau, y_1^\tau, \dots, x_\varphi^\tau, y_\varphi^\tau\}$. To model ICCL, we consider a historical sequence of such experience segments $\mathcal{C}_t = \bigoplus_{i=1}^N \mathcal{D}_{\varphi_i}^{\tau_i}$, where $t = \sum_{i=1}^N \varphi_i$ and \bigoplus denotes ordered concatenation. Unlike a random collection, the order of segments is essential since ICCL is sensitive to the task order arrangement. Given a query input x and the historical sequence \mathcal{C}_t , ICCL aims to approximate the output of an unknown target task τ^* . Formally, we define:

$$\hat{p}_\theta(y|x, \mathcal{C}_t) = \mathcal{F}_\theta \left(\bigoplus_{i=1}^N \mathcal{D}_{\varphi_i}^{\tau_i}, x \right) \Rightarrow p_{\tau^*}(y|x), \quad (1)$$

where \mathcal{F}_θ denotes a sequential prediction structure parameterized by θ , such as a large language model. Here, $p(y|x, \mathcal{C}_t)$ represents the probability of predicting output y given x and \mathcal{C}_t , and the goal is to approximate the true conditional distribution $p_{\tau^*}(y|x)$. Intuitively, ICCL leverages prior task experiences in context to infer the correct behavior on the target task. Finally, tasks τ_i and τ^* are elements of a task collection \mathcal{T} , which defines the general scope of tasks considered in ICCL.

Induction of task identifier. Previous zero-shot instruction-following studies typically assume the availability of a task description d^τ that supplements the prompt (Tanaka et al., 2024; Lou & Yin, 2023). However, in practice, such meta-information is rarely known a priori and tasks must often be inferred implicitly. To accommodate this more realistic setting, we relax d^τ into a *task identifier*, which need not encode semantic details but merely serves to distinguish one task from another. This abstraction aligns well with real-world applications: for instance, in preference modeling, the system may not have access to a user’s full profile but can instead rely on a user ID combined with the interaction history. Under this formulation, ICCL can be reformulated as

$$\hat{p}_\theta(y|x, \mathcal{C}_t) = \mathcal{F}_\theta \left(\bigoplus_{i=1}^N (d^{\tau_i}, \mathcal{D}_{\varphi_i}^{\tau_i}), d^{\tau^*}, x \right), \quad (2)$$

where task identifiers d^{τ_i} and d^{τ^*} act as lightweight labels that allow the model to keep track of task boundaries without requiring explicit semantic descriptions.

Evaluation of retention. Retention in ICCL refers to the model’s ability to preserve and use earlier demonstrations within long contexts. Given a historical sequence \mathcal{C}_t , we next evaluate how much information about the target task τ^* is retained after the model processes an intervening block $\mathcal{D}_{\varphi_D}^{\tau'}$ from another task $\tau' \neq \tau^*$. Specifically, this evaluation quantifies the extent to which the model "remembers" the target task τ^* after the t -step sequence \mathcal{C}_t and an additional φ_D -step exposure to the distracting task τ' :

$$R_{\tau^*, \tau', \mathcal{C}_t, \theta}(t + \varphi_D) = \mathbb{E}_{x, \tau'} \left[M \left(\hat{p}_\theta(\cdot|x, \mathcal{C}_t \oplus \mathcal{D}_{\varphi_D}^{\tau'}), p_{\tau^*}(\cdot|x) \right) \right], \quad (3)$$

where $M(p_1, p_2)$ is a metric that quantifies the similarity between two probability distributions p_1 and p_2 . We use *Normalized Performance* (denoted by the symbol R) to assess the accuracy of \hat{p}_θ , with $R(t) = 0$ indicating a complete loss of memory and $R(t) = 1$ indicating a perfect match between the estimator and the ground truth p_{τ^*} .

Scheduling strategy. We formalize a scheduling strategy to capture how task demonstrations are arranged in a historical sequence \mathcal{C}_t . Specifically, different practice conditions correspond to different constructions of \mathcal{C}_t , as summarized as follows:

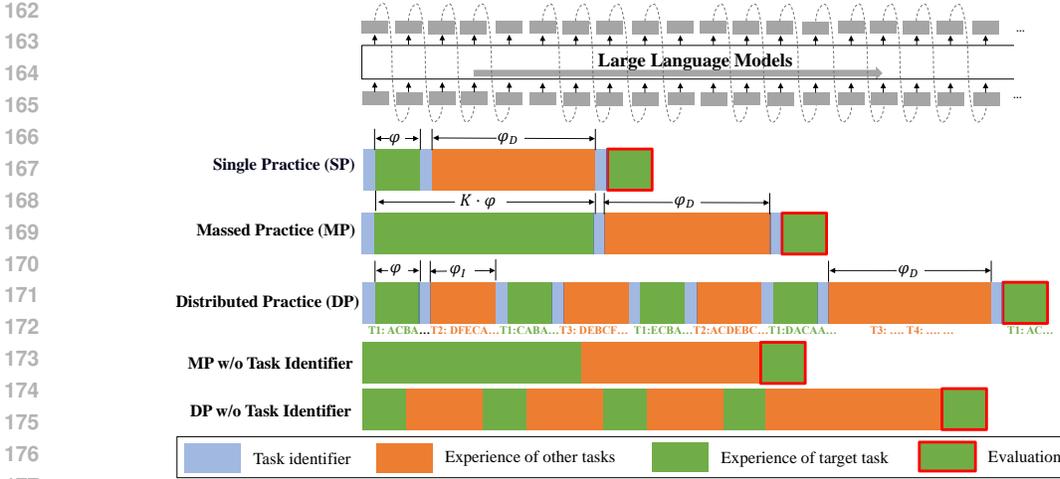


Figure 1: Illustration of how target and other task demonstrations are arranged under SP, MP, and DP, together with their variants without explicit task identifiers in a historical sequence. During the interleaved period, MP contains only target-task demonstrations (with $\varphi_I = 0$) while DP includes both target-task and intervening-task demonstrations ($\varphi_I > 0$). The post-exposure period (φ_D) always involves multi-task interference and is used to evaluate retention decay.

$$\begin{aligned}
 \text{Single Practice (SP): } & \mathcal{C}_t = \mathcal{D}_\varphi^{\tau^*} \\
 \text{Massed Practice (MP): } & \mathcal{C}_t = \mathcal{D}_{K\varphi}^{\tau^*} \\
 \text{Distributed Practice (DP): } & \mathcal{C}_t = \bigoplus_{i=1}^K (\mathcal{D}_\varphi^{\tau^*}, \mathcal{D}_{\varphi_I}^{\tau_i}), \quad \forall \tau_i \in \mathcal{T}, \tau_i \neq \tau^*.
 \end{aligned} \tag{4}$$

Eq. (4) characterizes three canonical scheduling conditions, which are illustrated in Figure 1. In *SP*, demonstrations of the target task τ^* are presented only once in a contiguous block of length φ , after which the sequence moves to other tasks. In *MP*, the demonstrations of the target task are extended to a larger block of length $K\varphi$, corresponding to the repeated rehearsal of τ^* without interruption. In *DP*, the sequence alternates between the target task and other tasks: each target task block of length φ is followed by another task block $\mathcal{D}_{\varphi_I}^{\tau_i}$ with $\tau_i \neq \tau^*$. This setup more closely resembles human retention patterns (spacing is most effective when intervening content is distinct from the target) and allows ICCL to retain the target task knowledge in a manner analogous to human retention. Therefore, we maximize differences between target and other tasks to create clear contextual boundaries and reliable reactivation cues. Compared with independent ICL, which must repeatedly include large blocks of target-task demonstrations to sustain retention, ICCL achieves comparable performance using far fewer demonstrations. Under DP, the target-task pattern is naturally reinstated through spaced exposure, enabling the model to maintain high retention even when only a limited number of target-task demonstrations are available. This advantage becomes essential when retention must persist over long content lengths.

To further investigate the role of explicit task identifiers, we additionally evaluate two variants where d^T is removed from the prompt, denoted as *MP w/o instruction* and *DP w/o instruction*, respectively. These two variants test whether the presence of explicit task identifiers contributes to multitask sequential adaptation beyond the task order arrangement.

2.2 RETENTION DYNAMICS IN ICCL

A well-documented phenomenon in human retention is the spacing effect (Anderson et al., 2004; Cepeda et al., 2008): MP triggers rapid decay, whereas DP yields superior long-term retention. To account for this effect, the ACT-R cognitive model (Borst & Anderson, 2017; Pavlik Jr & Anderson, 2005; Petrov, 2006) describes the resulting learning and retention curves. Suppose the learner practices (observes) the target task τ^* at the time stamp of $t_1, t_2, \dots, t_{K\varphi}$ and is exposed to other tasks otherwise; the ACT-R model predicts the probability of retaining the task at any subsequent

time $t = t_{K\varphi} + \varphi_D$ by the following equation:

$$\hat{R}(t) = 1/[1 + \exp(-\frac{w(t) - \gamma}{s})], \quad w(t) = \ln \sum_{i=1}^{K\varphi} [\kappa \cdot (t - t_i)]^{-d} \quad (5)$$

where $w(t)$ represents the memory activation that decays over time when no new practice of the target task occurs, and $\hat{R}_{\tau^*}(t)$ denotes the probability of successfully accomplishing task τ^* at time t given the memory activation $w(t)$. Notice that in our settings, for MP, $t_i = i$, for DP, $t_i = i + \lfloor \frac{i}{\varphi} \rfloor \varphi_I$. The parameters s , γ , and d are hyper-parameters to be determined. To bridge the gap between the actual time-dependency of human retention and the context length in ICCL, we introduce an additional hyper-parameter κ . Our experimental results in Section 3.4 reveal an interior optimum of $\hat{R}(t)$ under DP, indicating the existence of a spacing sweet spot that maximizes the retention probability. This empirical pattern mirrors decades of findings in human memory research and provides practical guidance for prompt scheduling that avoids both massed and overly spaced practice. [This retention arises not from short-term contextual reuse, but from spaced reactivation of the target-task pattern, which counteracts interference and long-context attention decay.](#)

2.3 QUANTIFYING HUMAN RETENTION SIMILARITY VIA MAHALANOBIS DISTANCE

We now ask a related question: *to what extent do LLMs resemble human memory in their retention characteristics?* To enable a direct human–model comparison, we introduce Human Retention Similarity via Mahalanobis Distance (HRS-MD), which evaluates how closely the fitted ACT-R parameters of a model fall within parameter distributions reported in cognitive psychology. Specifically, Let the fitted parameter vector be $\hat{\theta} = [\hat{d}, \hat{s}, \hat{\gamma}]^\top$, where \hat{d} denotes the decay rate, \hat{s} the activation noise, and $\hat{\gamma}$ the retrieval threshold estimated from retention curves. Human reference distributions are characterized by a mean vector $\mu = [\mu_d, \mu_s, \mu_\gamma]^\top$ and covariance matrix Σ , derived from prior experimental studies (Pavlik Jr & Anderson, 2005; Lewis & Vasishth, 2005; Bothell, 2020; Borst & Anderson, 2017; Said et al., 2016; van der Velde et al., 2022). HRS-MD is then defined as $(\hat{\theta} - \mu)^\top \Sigma^{-1} (\hat{\theta} - \mu)$, which quantifies the discrepancy between the model-fitted parameters and human baselines while accounting for parameter scales and correlations. A smaller HRS-MD indicates closer alignment between LLMs and human memory retention. This metric enables systematic comparisons across ICCL and GBCL methods, revealing which model classes most closely reproduce human-like retention pattern. [Details of \$\Sigma\$ and the ACT-R parameters, along with their reference distributions and the rationale for the Mahalanobis metric, are provided in Appendix A.2.](#)

3 EXPERIMENT RESULTS

3.1 EXPERIMENT SETTINGS

Models and Baselines. We evaluate two categories of methods: inference-only ICCL models and GBCL baselines. The ICCL group consists of four representative open-source LLMs: LLaMA3-8B, DeepSeek-R1, MAMBA, and RWKV-7. [Note that ICCL focuses on inference-time retention within standard LLMs and does not rely on architectural or external-memory mechanisms; thus, comparisons with memory-enhanced models fall outside the scope of this study.](#) The GBCL group includes three classical approaches and one modern method: Stochastic Gradient Descent (SGD) (Arous et al., 2021), Experience Replay (ER) (Rolnick et al., 2019), Elastic Weight Consolidation (EWC) (Evilevitch & Ingram, 2021), and Prompt Gradient Projection (PGP) (Qiao et al., 2024). For ICCL, experiences serve as the context in LLMs without updating any parameters; for GBCL, we construct a small predictive model that maps the input & task-identifier to an output and sequentially refine it with the arriving experiences. The detailed model architecture and hyperparameters are provided in the Appendix A.3. To ensure fair comparisons, all models are evaluated on the same input sequences and prompt structures, with greedy decoding applied during inference. Each experiment is repeated 16 times with different random seeds, and results are reported as stable averages.

Benchmark Tasks. Discrete Markov Chains (DMC) are widely adopted benchmarking tasks (Zekri et al., 2024; Akyürek et al., 2024; Liu et al., 2024) because they capture essential properties of dynamical systems while minimizing contamination from pre-existing knowledge, thus providing

a clean and controllable environment for analyzing retention and interference dynamics. In our benchmark, two levels of complexity are considered: *simple* tasks with $N_s = 4$ states and *complex* tasks with $N_s = 8$ states. For each setting, transition matrices are generated to sample states and construct experience sequences. Each target and interference tasks are defined by a transition matrix $P_{target} \in \mathbb{R}^{N_s \times N_s}$ and a transition matrix $P_{interference} \in \mathbb{R}^{N_s \times N_s}$, and the model receives a state sequence as the main prompt content to predict the next state. Prompts are prepared under two conditions: (i) with task identifiers, where each sequence segment is prefixed with a label (e.g., [TARGET_TASK], [INTERFERENCE_TASK]), and (ii) without identifiers, where the same content is presented unlabeled. This setting enables systematic evaluation of whether explicit identifiers facilitate multitask sequential adaptation. We apply the scheduling strategies introduced in Section 2.1 to form prompts in SP, MP, and DP formats. We include a detailed example in the Appendix A.4 illustrating how a DMC prompt is constructed, along with its corresponding transition matrices. We further design a complementary benchmark built upon the RegBench framework (Akyürek et al., 2024) to evaluate ICCL in natural language–style environments. RegBench introduces linguistically structured inputs and richer syntactic variation, which make RegBench closer to real-world sequential learning scenarios. Appendix A.4 provides complete details on the prompt construction process and ICCL’s retention performance on this new benchmark.

Parameter Settings. We configure the DMC benchmark with $\varphi = 100$, $K = 5$, $\varphi_I \in \{10, 50, 100, 200, 400, 600\}$, $\varphi_D \in (0, 700)$, considering the token limits of different LLMs, e.g., DEEPSEEK-R1 supports up to approximately 4096 tokens, whereas LLAMA3-8B allows up to 8192 tokens. These settings ensure that all models operate within their valid context windows while maintaining comparable information exposure across architectures. They also enable controlled studies on *interference length* (φ_I) and *duration* (φ_D), providing a consistent basis for analyzing retention and interference as the effective context grows. This unified configuration guarantees that each model runs at its optimal performance regime without exceeding its token budget.

Evaluation Metrics. A bottleneck in DMC evaluation is that performance divergence across tasks often stems from intrinsic task differences rather than prediction accuracy. To suppress this noise, we first introduce a normalized performance metric. Inspired by the normalization schemes in (Cha, 2007), we define the normalized retention performance metric M (from eq. (3)) by rescaling the Bhattacharyya distance (Liu et al., 2021; Zekri et al., 2024) $D_{Bah}(p_1, p_2)$ to the interval $[0, 1]$:

$$M(\hat{p}_\theta, p_{\tau^*}) = \frac{\exp[D_{Bah}(p_{\tau^*}, p_{\tau^*}) - D_{Bah}(\hat{p}_\theta, p_{\tau^*})] - 1}{\exp[D_{Bah}(p_{\tau^*}, p_{\tau^*}) - D_{Bah}(p_{\mathcal{U}}, p_{\tau^*})] - 1}, \quad (6)$$

where D_{Bah} is the Bhattacharyya distance and $p_{\mathcal{U}}$ denotes the uniform distribution, representing complete forgetting. This normalization makes retention scores comparable across tasks by rescaling the model–target divergence between the best case $D_{Bah}(p_{\tau^*}, p_{\tau^*})$ and the worst case $D_{Bah}(p_{\mathcal{U}}, p_{\tau^*})$. Bhattacharyya distance is chosen because it directly measures overlap between predictive and target distributions, though its raw scale varies with label entropy and task complexity. The exponential transform further linearizes distance differences and amplifies meaningful gaps when distributions are close, following standard similarity-normalization practices (Cha, 2007). Under this mapping, $M = 1$ corresponds to perfect retention, $M = 0$ to uniform guessing, and intermediate values reflect retention quality on a consistent cross-task scale. All reported results are given as mean \pm 95% CI. We additionally employ HRS-MD to quantify alignment with human memory dynamics.

3.2 RETENTION PERFORMANCE UNDER SP, MP AND DP

We first evaluate how the retention performance of ICCL models and GBCL baselines varies on simple and complex tasks with task identifiers and $\varphi_I = 200$, under SP, MP, and DP scheduling. The experimental results on the complex task are shown in Figure 2. For both SP and MP, we observe that ICCL models exhibit a sharp performance drop immediately after the last target-task block, followed by a small continued decline as φ_D increases, whereas some GBCL methods show a more gradual but consistent decrease—underscoring their susceptibility to catastrophic forgetting. Under DP, however, ICCL benefits substantially: DP effectively mitigates forgetting and maintains higher retention, while GBCL baselines show little improvement. We also observe that once LLM transitions from the target-task block to intervening tasks, its retention on the target task drops abruptly, even though the target demonstrations were encountered only moments earlier and remain relatively close

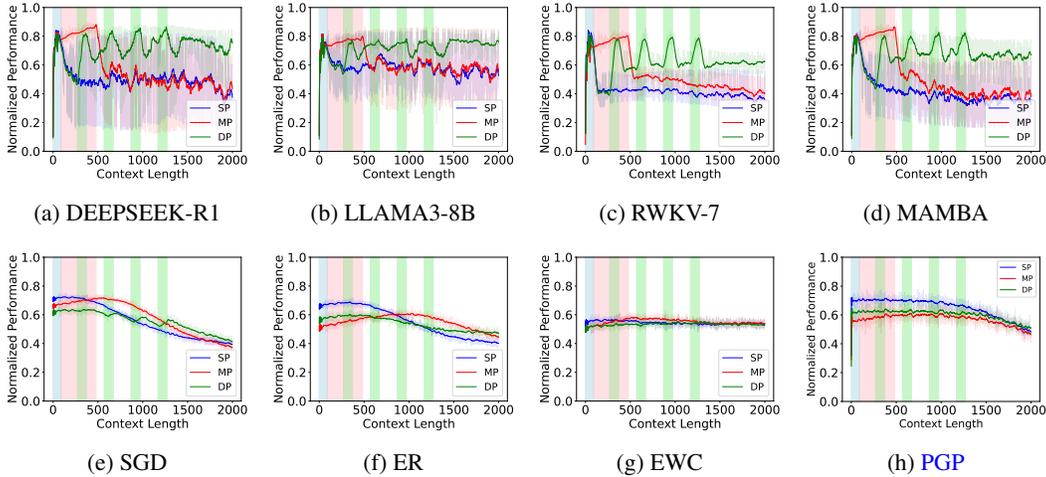


Figure 2: Retention performance on complex task under SP, MP, and DP settings (with task identifiers, $\varphi_I = 200$). Target task blocks are highlighted in light blue (SP), pink (MP), and green (DP).

Table 2: Retention performance (mean \pm 95% CI) under DP for simple ($N_s = 4$) and complex ($N_s = 8$) Markov-chain tasks across varying intervals φ_D .

Simple Task (markov chain with $N_s = 4$)					
Model Types	$\varphi_D = 0$	$\varphi_D = 100$	$\varphi_D = 200$	$\varphi_D = 400$	$\varphi_D = 600$
LLAMA3-8B	0.819 \pm 0.086	0.737 \pm 0.090	0.865 \pm 0.030	0.833 \pm 0.067	0.781 \pm 0.078
MAMBA	0.668 \pm 0.106	0.655 \pm 0.105	0.678 \pm 0.094	0.685 \pm 0.100	0.667 \pm 0.101
SGD	0.606 \pm 0.100	0.581 \pm 0.084	0.558 \pm 0.086	0.507 \pm 0.098	0.471 \pm 0.101
ER	0.608 \pm 0.095	0.613 \pm 0.075	0.617 \pm 0.077	0.620 \pm 0.081	0.629 \pm 0.085
EWC	0.617 \pm 0.101	0.612 \pm 0.086	0.604 \pm 0.088	0.601 \pm 0.099	0.602 \pm 0.101
PGP	0.665 \pm 0.064	0.652 \pm 0.066	0.637 \pm 0.059	0.609 \pm 0.069	0.562 \pm 0.079
Complex Task (markov chain with $N_s = 8$)					
Model Types	$\varphi_D = 0$	$\varphi_D = 100$	$\varphi_D = 200$	$\varphi_D = 400$	$\varphi_D = 600$
LLAMA3-8B	0.784 \pm 0.054	0.776 \pm 0.056	0.753 \pm 0.086	0.783 \pm 0.057	0.792 \pm 0.058
MAMBA	0.750 \pm 0.059	0.700 \pm 0.061	0.608 \pm 0.103	0.724 \pm 0.064	0.726 \pm 0.048
SGD	0.577 \pm 0.072	0.540 \pm 0.076	0.491 \pm 0.080	0.464 \pm 0.073	0.448 \pm 0.085
ER	0.519 \pm 0.081	0.501 \pm 0.086	0.472 \pm 0.088	0.479 \pm 0.080	0.495 \pm 0.091
EWC	0.563 \pm 0.079	0.551 \pm 0.087	0.526 \pm 0.091	0.537 \pm 0.083	0.542 \pm 0.093
PGP	0.608 \pm 0.057	0.600 \pm 0.061	0.584 \pm 0.060	0.549 \pm 0.063	0.515 \pm 0.060

to the query in absolute context position. This sharp decline shows that the dominant failure mode is not insufficient recency, but rather the model’s susceptibility to interference and long-range attention decay during continued context processing.

Building on the above results, Table 2 compares the retention performance of ICCL models and GBCL baselines under DP across simple and complex tasks as φ_D increases from 0 to 600. ICCL consistently outperforms all GBCL methods, and although retention declines with larger φ_D , ICCL remains in a high-retention regime while GBCL methods operate at substantially lower levels. Although ER and EWC show relatively small numerical variations across φ_D , this pattern arises from early saturation at a low-retention regime. By contrast, ICCL maintains much higher retention with a slower degradation rate, demonstrating stronger resistance to extended interference.

Ablation Study on Target Task Block Numbers K . To further understand how the number of target-task blocks K influences model retention, we conduct an ablation study where K varies in $\{3, 5, 6, 7\}$ while fixing $\varphi_D = 600$. Figure 3 presents the normalized retention performance of ICCL and GBCL methods under DP. We observe that K has a much stronger impact on ICCL methods than on GBCL methods. For transformer-based LLMs (e.g., DEEPSEEK-R1 and LLAMA3-8B),

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

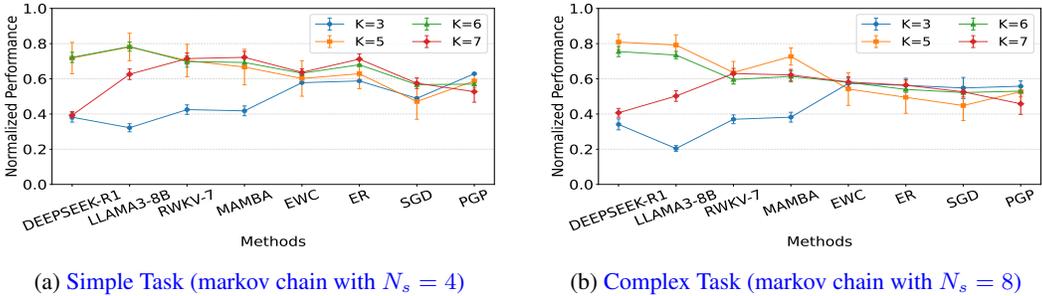


Figure 3: Retention performance (with $\varphi_D = 600$) under DP for simple ($N_s = 4$) and complex ($N_s = 8$) Markov-chain tasks across different target task block numbers ($K \in \{3, 5, 6, 7\}$).

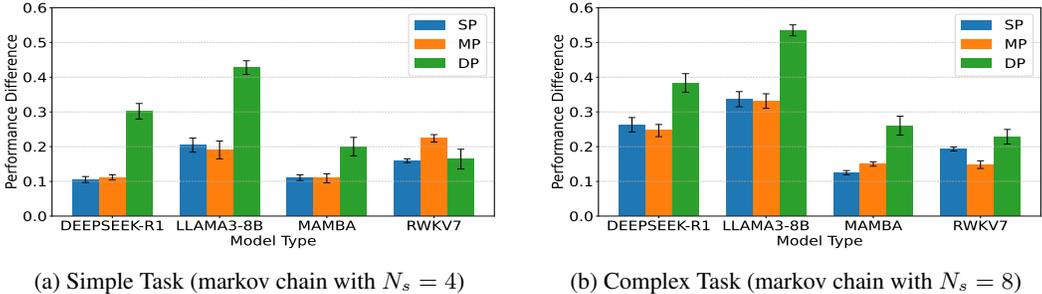


Figure 4: Retention performance difference of ICCL models with and without task identifiers under SP, MP, and DP ($\varphi_I = 200$ and $\varphi_D = 600$). Positive values indicate higher retention when task identifiers are included in the prompt.

retention performance first increases and then decreases as K grows, because a larger K yields longer effective contexts that approach or exceed the transformer’s token limit. In contrast, linear-attention LLMs (e.g., RWKV-7 and MAMBA) show continuously improving and eventually stable retention performance, reflecting their superior scalability to long sequences. Overall, the results confirm that context length ($K\varphi$) interacts with model architecture to determine how efficiently previous task information is retained. This trend consistently appears in both simple and complex Markov-chain tasks, highlighting K as a key factor shaping long-term memory dynamics in LLMs.

3.3 EFFECT OF TASK IDENTIFIER ON ICCL RETENTION

Figure 4 illustrates retention performance difference of ICCL models when prompts include task identifiers versus when they do not, under SP, MP, and DP conditions. Here, the plotted value represents the difference between ICCL’s retention performance with and without identifiers; a positive value indicates an improvement when identifiers are used. As shown in Figure 4, incorporating task identifiers consistently enhances retention across all models and scheduling strategies, with the most pronounced gains observed under DP. The effect is more substantial in complex DMC tasks, suggesting that explicit identifiers become increasingly important as task complexity grows.

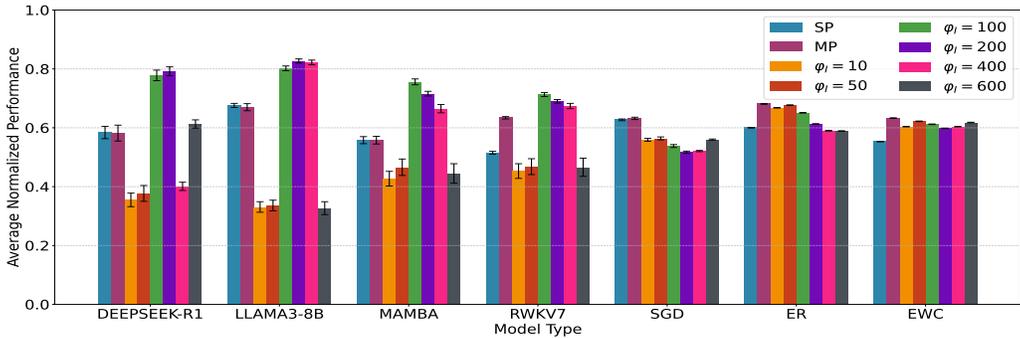
3.4 SPACING SWEET SPOT IN ICCL UNDER DP

Figure 5 reports the average retention performance of ICCL and GBCL models under DP across different intervals φ_I , evaluated on both simple and complex tasks. Here, average retention performance is defined as the mean retention over φ_D , thereby capturing memory retention beyond the immediate target-task exposure. The results reveal a clear contrast between ICCL and GBCL. For ICCL, DP induces a characteristic *spacing sweet spot*, where average retention performance peaks at an intermediate φ_I lying between 100 and 400. This sweet spot consistently emerges across both simple and complex tasks, although average retention performance is higher in the simpler setting. These findings empirically support the presence of a spacing sweet spot and further demonstrate that

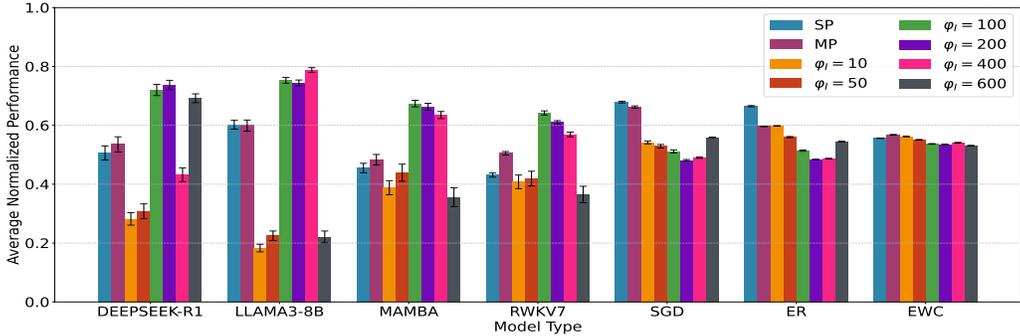
Table 3: Fitted ACT-R parameters with corresponding MSE and HRS-MD values.

Model	Decay Rate (d)	Activation Noise (s)	Scaling Factor (κ)	Threshold (γ)	MSE	HRS-MD
LLAMA3	0.14	2.00	1.20	1.01	0.0024	500.22
DEEPSEEK-R1	0.35	1.69	0.43	-0.24	0.0155	302.39
RWKV-7	0.27	1.62	0.94	0.59	0.0052	287.57
MAMBA	0.29	1.59	0.76	0.33	0.0049	272.74
SGD	0.41	2.00	0.77	0.15	0.0036	445.07
ER	0.27	2.00	1.36	1.02	0.0061	466.74
EWC	0.22	2.00	1.31	1.65	0.0021	481.53

ICCL under DP surpasses the performance achieved with SP and MP. Thus, ICCL benefits from DP in a manner analogous to human memory, and careful tuning of φ_I can substantially enhance retention. In contrast, GBCL baselines exhibit no such sweet spot: their performance remains largely flat or gradually declines as φ_I increases, indicating that DP provides little to no advantage.



(a) Simple Task (markov chain with $N_s = 4$)



(b) Complex Task (markov chain with $N_s = 8$)

Figure 5: Average retention performance comparison among ICCL and GBCL methods under DP with different intervals (φ_I) on both simple and complex tasks.

3.5 HRS-MD QUANTIFICATION

To investigate how the aforementioned retention curves align with human retention, we employed the human-retention-validated ACT-R model to fit the ICCL and GBCL curves. Specifically, we used eq. (5) to model the normalized ICCL/GBCL performance obtained from eq. (3). Table 3 summarizes the fitted ACT-R parameters, mean squared error (MSE) during parameter fitting, and HRS-MD values for both ICCL and GBCL methods based on retention performance results on both simple and complex tasks under DP across different intervals φ_I . Several consistent patterns emerge. First, ICCL models—except for LLaMA3—achieve lower HRS-MD values than GBCL baselines, indicating

that their retention dynamics more closely approximate human memory. In particular, MAMBA and RWKV-7 yield the lowest HRS-MD scores (272.74 and 287.57, respectively), demonstrating the strongest alignment with human retention patterns. By contrast, GBCL methods produce substantially higher HRS-MD values, reflecting greater divergence from human memory dynamics. Second, the fitted parameters reveal distinctive tendencies across the two paradigms. ICCL models generally converge to moderate decay rates, relatively low activation noise, and thresholds near zero—indicating efficient and accessible retrieval dynamics, particularly evident in MAMBA and RWKV-7. In contrast, GBCL methods exhibit more extreme parameter configurations: SGD yields the highest decay rate ($d = 0.41$), consistent with rapid forgetting, while ER and EWC converge to elevated thresholds ($\gamma \geq 1.0$), implying stricter retrieval conditions and reduced accessibility of stored knowledge. These patterns highlight the poor balance between stability and plasticity in GBCL approaches.

4 RELATED WORK

GBCL is the classical approach for mitigating catastrophic forgetting through parameter updates. Methods typically fall into three categories: EWC, ER, and SGD. EWC constrains updates to parameters important for previous tasks, using Fisher-information-based scores (Li et al., 2020; Evilevitch & Ingram, 2021) or online importance estimates such as synaptic intelligence (Zenke et al., 2017; Zenke et al.). ER maintains a small memory buffer and interleaves past examples with new data (Rolnick et al., 2019; Buzzega et al., 2021). SGD (Arous et al., 2021; Poggio et al., 2011) enables rapid adaptation but suffers severe forgetting due to its lack of preservation mechanisms. [Recent extensions emphasize representation stability, but remain incompatible with our setting: EC-ACE \(Xiong et al., 2024\) relies on vision-specific multi-view consistency losses, while InsCL \(Wang et al., 2024\) adapts LLMs via instruction-tuned parameter updates—both outside the scope of inference-only ICCL.](#)

ICL allows LLMs to adapt to tasks from in-prompt demonstrations without parameter updates. Prior work has analyzed its empirical sensitivities (Zhao et al., 2021; Liu et al., 2021; Lu et al., 2021; Liu et al., 2023) and offered theoretical views of ICL as Bayesian inference, gradient-based adaptation, or kernel regression (Xie et al., 2021; Akyürek et al., 2022; Von Oswald et al., 2023; Han et al., 2023). Many-shot studies further show extension to more complex tasks (Agarwal et al., 2024; Garg et al., 2022). Yet most research focuses on single-task adaptation, leaving sequential multitask retention largely unexplored. Recent continual-ICL efforts adjust prompts incrementally (Wang et al., 2022; Xiong et al., 2024) or use retrieval modules (Gao et al., 2023; Momeni et al., 2024; Shinwari & Usama, 2025), but these rely on heuristics or auxiliary components prone to memory growth and retrieval noise. [While \(Xiong et al., 2024\) mixes demonstrations across tasks for short-term gains, ICCL instead examines long-term retention under sequential exposures.](#)

Human memory research offers a rich foundation for understanding retention dynamics. Ebbinghaus (Ebbinghaus, 2013) introduced the forgetting curve, later refined into power-law formulations (White, 2001). The *spacing effect* shows that distributed practice yields superior long-term retention, a finding supported by large-scale meta-analyses (Cepeda et al., 2006; 2008). These phenomena are formalized in the base-level learning equation (Pavlik Jr & Anderson, 2005) and its ACT-R extension (Pavlik & Anderson, 2008; Bothell, 2020; Borst & Anderson, 2017), which link recall to the frequency and recency of exposures. Our work leverages these principles to test whether LLMs display analogous retention patterns when demonstrations are distributed across prompts.

5 CONCLUSION

In this work, we investigated the retention characteristics of ICL in multi-task settings and extended it to ICCL through prompt scheduling and rearrangement. Experiments on benchmarks showed that ICCL benefits DP in a manner analogous to humans, consistently revealing a spacing sweet spot that enhances retention beyond SP and MP settings. Beyond retention performance, we introduced the HRS-MD metric to quantify alignment with human memory dynamics, and found that linear-attention models such as MAMBA and RWKV display particularly human-like retention patterns. Overall, our results establish ICCL as both cognitively plausible and practically effective, providing an inference-only paradigm that mitigates catastrophic forgetting and addresses the stability–plasticity dilemma in conventional continual learning methods.

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. Our study focuses on the development and evaluation of in-context continual learning methods using synthetic and publicly available benchmark datasets. No human subjects or sensitive personal data are involved. The proposed methods aim to improve retention dynamics in large language models without introducing harmful applications or privacy risks. All datasets used are either synthetic (Markov-chain sequences) or publicly accessible, and no proprietary or restricted data are employed. In addition, the large language models used in our experiments—LLaMA3-8B, DeepSeek-R1, MAMBA, and RWKV-7—are employed strictly in compliance with their open-source licenses. We are not aware of any conflicts of interest, sponsorship influence, or ethical risks associated with this research.

REPRODUCIBILITY STATEMENT

We have taken multiple steps to ensure the reproducibility of our results. The main text provides a complete description of the proposed ICCL framework, the baseline methods, and the evaluation metric. Details of additional ablation studies are included in the Appendix. All datasets used in the experiments are synthetically generated Markov chain sequences or are publicly available. Hyperparameter choices and training configurations for both GBCL baselines and large language model baselines are fully documented. To further support replication, we will release our anonymized source code as supplemental materials submitted with the paper.

REFERENCES

- Rishabh Agarwal, Avi Singh, Lei Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, et al. Many-shot in-context learning. *Advances in Neural Information Processing Systems*, 37:76930–76966, 2024.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and algorithms. In *International Conference on Machine Learning*, pp. 787–812. PMLR, 2024.
- 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, 2004.
- Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- Stéphane Attal. Markov chains and dynamical systems: the open system point of view. *arXiv preprint arXiv:1010.2894*, 2010.
- Jelmer P. Borst and John R. Anderson. A step-by-step tutorial on using the cognitive architecture act-r in combination with fmri data. *Journal of Mathematical Psychology*, 76:94–103, 2017. URL <https://www.jelmerborst.nl/pubs/Borst2017.pdf>.
- Dan Bothell. Act-r 7.21+ reference manual. *Pittsburgh (PA): Carnegie Mellon University*, 2020.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, and Simone Calderara. Rethinking experience replay: a bag of tricks for continual learning. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 2180–2187. IEEE, 2021.

- 594 Nicholas J Cepeda, Harold Pashler, Edward Vul, John T Wixted, and Doug Rohrer. Distributed
595 practice in verbal recall tasks: A review and quantitative synthesis. *Psychological bulletin*, 132(3):
596 354, 2006.
- 597 Nicholas J Cepeda, Edward Vul, Doug Rohrer, John T Wixted, and Harold Pashler. Spacing effects
598 in learning: A temporal ridgeline of optimal retention. *Psychological science*, 19(11):1095–1102,
599 2008.
- 600 Sung-Hyuk Cha. Comprehensive survey on distance/similarity measures between probability density
601 functions. *City*, 1(2):1, 2007.
- 602 DeepSeek. Deepseek coder 7b instruct v1.5. [https://huggingface.co/deepseek-ai/
603 deepseek-coder-7b-instruct-v1.5](https://huggingface.co/deepseek-ai/deepseek-coder-7b-instruct-v1.5), 2024.
- 604 Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of neurosciences*,
605 20(4):155, 2013.
- 606 Anton Evilevitch and Robert Ingram. Avoiding catastrophic forgetting in continual learning through
607 elastic weight consolidation, 2021.
- 608 Wang Fan and Shaoshan Liu. Putting the smarts into robot bodies. *Commun. ACM*, 68(3):6–8, 2025.
609 doi:10.1145/3703761. URL <https://doi.org/10.1145/3703761>.
- 610 Gary Froyland. Extracting dynamical behavior via markov models. In *Nonlinear dynamics and
611 statistics*, pp. 281–321. Springer, 2001.
- 612 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,
613 Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A
614 survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023.
- 615 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
616 in-context? a case study of simple function classes. *Advances in neural information processing
617 systems*, 35:30583–30598, 2022.
- 618 Chi Han, Ziqi Wang, Han Zhao, and Heng Ji. Explaining emergent in-context learning as kernel
619 regression. *arXiv preprint arXiv:2305.12766*, 2023.
- 620 Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald,
621 DJ Strouse, Steven Stenberg Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement
622 learning with algorithm distillation. In *The Eleventh International Conference on Learning
623 Representations*, 2023.
- 624 Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma
625 Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural
626 Information Processing Systems*, 36:43057–43083, 2023.
- 627 Richard L Lewis and Shrawan Vasishth. An activation-based model of sentence processing as skilled
628 memory retrieval. *Cognitive science*, 29(3):375–419, 2005.
- 629 Yijun Li, Richard Zhang, Jingwan Lu, and Eli Shechtman. Few-shot image generation with elastic
630 weight consolidation. *arXiv preprint arXiv:2012.02780*, 2020.
- 631 Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What
632 makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
- 633 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
634 and Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint
635 arXiv:2307.03172*, 2023.
- 636 Toni Jb Liu, Nicolas Boulle, Raphaël Sarfati, and Christopher Earls. Llms learn governing principles
637 of dynamical systems, revealing an in-context neural scaling law. In *Proceedings of the 2024
638 Conference on Empirical Methods in Natural Language Processing*, pp. 15097–15117, 2024.
- 639 Renze Lou and Wenpeng Yin. Toward zero-shot instruction following. *arXiv preprint
640 arXiv:2308.03795*, 2023.

- 648 Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered
649 prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint*
650 *arXiv:2104.08786*, 2021.
- 651
652 Meta. Llama 3.1 8b instruct. [https://huggingface.co/meta-llama/Llama-3.](https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct)
653 1-8B-Instruct, 2024.
- 654 Saleh Momeni, Sahisnu Mazumder, Zixuan Ke, and Bing Liu. In-context continual learning assisted
655 by an external continual learner. *arXiv preprint arXiv:2412.15563*, 2024.
- 656
657 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
658 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
659 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
660 27744, 2022.
- 661 Philip I Pavlik and John R Anderson. Using a model to compute the optimal schedule of practice.
662 *Journal of Experimental Psychology: Applied*, 14(2):101, 2008.
- 663
664 Philip I Pavlik Jr and John R Anderson. Practice and forgetting effects on vocabulary memory: An
665 activation-based model of the spacing effect. *Cognitive science*, 29(4):559–586, 2005.
- 666
667 Alexander A Petrov. Computationally efficient approximation of the base-level learning equation in
668 act-r. In *Proceedings of the seventh international conference on cognitive modeling*, pp. 391–392.
669 Trieste, ITA Edizioni Goliardiche, 2006.
- 670 Tomaso Poggio, Stephen Voinea, and Lorenzo Rosasco. Online learning, stability, and stochastic
671 gradient descent. *arXiv preprint arXiv:1105.4701*, 2011.
- 672
673 Jingyang Qiao, Xin Tan, Chengwei Chen, Yanyun Qu, Yong Peng, Yuan Xie, et al. Prompt gra-
674 dient projection for continual learning. In *The Twelfth International Conference on Learning*
675 *Representations*, 2024.
- 676
677 Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
678 emergence of non-bayesian in-context learning for regression. *Advances in neural information*
processing systems, 36:14228–14246, 2023.
- 679
680 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience
681 replay for continual learning. *Advances in neural information processing systems*, 32, 2019.
- 682
683 RWKV. Rwkv-4 world 7b. <https://huggingface.co/RWKV/rwkv-4-world-7b>, 2024.
- 684
685 Nadia Said, Till Nagel, and Frank Rebitschek. Applying mathematical optimization methods to an
686 act-r model: The sugar factory. *PLOS ONE*, 11(7):e0158832, 2016. URL [https://journals.](https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0158832)
[plos.org/plosone/article?id=10.1371/journal.pone.0158832](https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0158832).
- 687
688 Haseeb Ullah Khan Shinwari and Muhammad Usama. Memory-augmented architecture for long-term
689 context handling in large language models. *arXiv preprint arXiv:2506.18271*, 2025.
- 690
691 Ryota Tanaka, Taichi Iki, Kyosuke Nishida, Kuniko Saito, and Jun Suzuki. Instructdoc: A dataset for
692 zero-shot generalization of visual document understanding with instructions. In *Proceedings of the*
AAAI conference on artificial intelligence, volume 38, pp. 19071–19079, 2024.
- 693
694 TII. Falcon-mamba 7b. <https://huggingface.co/tiiuae/falcon-mamba-7b>, 2024.
- 695
696 Maarten van der Velde, Florian Sense, Jelmer P. Borst, Leendert van Maanen, and Hedderik van
697 Rijn. Capturing dynamic performance in a cognitive model: Estimating act-r memory pa-
698 rameters with the linear ballistic accumulator. *Topics in Cognitive Science*, 14(4):889–903,
699 2022. doi:10.1111/tops.12614. URL [https://pmc.ncbi.nlm.nih.gov/articles/](https://pmc.ncbi.nlm.nih.gov/articles/PMC9790673/)
[PMC9790673/](https://pmc.ncbi.nlm.nih.gov/articles/PMC9790673/).
- 700
701 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
International Conference on Machine Learning, pp. 35151–35174. PMLR, 2023.

- 702 Fan Wang, Pengtao Shao, Yiming Zhang, Bo Yu, Shaoshan Liu, Ning Ding, Yang Cao, Yu Kang,
703 and Haifeng Wang. Omnirl: In-context reinforcement learning by large-scale meta-training in
704 randomized worlds. *arXiv preprint arXiv:2502.02869*, 2025a.
- 705 Fan Wang, Bo Yu, Pengtao Shao, Zhiyuan Chen, Hungchyun Chou, Liuwang Kang, Shenbao Shao,
706 Xuan Xia, Ning Ding, Haoyi Xiong, et al. In-context learning as general-purpose learning: A
707 comprehensive survey and new perspectives. 2025b.
- 709 Yifan Wang, Yafei Liu, Chufan Shi, Haoling Li, Chen Chen, Haonan Lu, and Yujiu Yang. Insl: A
710 data-efficient continual learning paradigm for fine-tuning large language models with instructions.
711 *arXiv preprint arXiv:2403.11435*, 2024.
- 712 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent
713 Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings*
714 *of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 139–149, 2022.
- 715 K Geoffrey White. Forgetting functions. *Animal Learning & Behavior*, 29(3):193–207, 2001.
- 717 Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett.
718 How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint*
719 *arXiv:2310.08391*, 2023.
- 720 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
721 learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 723 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
724 learning as implicit bayesian inference. In *International Conference on Learning Representations*,
725 2024.
- 726 Zheyang Xiong, Ziyang Cai, John Cooper, Albert Ge, Vasilis Papageorgiou, Zack Sifakis, Angeliki
727 Giannou, Ziqian Lin, Liu Yang, Saurabh Agarwal, et al. Everything everywhere all at once: Lms
728 can in-context learn multiple tasks in superposition. *arXiv preprint arXiv:2410.05603*, 2024.
- 729 Oussama Zekri, Ambroise Odonnat, Abdelhakim Benechehab, Linus Bleistein, Nicolas Boullé, and
730 Ievgen Redko. Large language models as markov chains. *arXiv preprint arXiv:2410.02724*, 2024.
- 732 Friedemann Zenke, Ben Poole, and Surya Ganguli. Improved multitask learning through synaptic
733 intelligence. In *Proceedings of the International Conference on Machine Learning*, pp. 3987–3995.
734
- 735 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence.
736 In *International conference on machine learning*, pp. 3987–3995. PMLR, 2017.
- 737 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving
738 few-shot performance of language models. In *International conference on machine learning*, pp.
739 12697–12706. PMLR, 2021.

741 A APPENDIX

742 A.1 USE OF LLMs

743 We used large language models (LLMs) only as an auxiliary tool to improve the clarity and presenta-
744 tion of this paper. The assistance was limited to:

- 745 • **Language refinement:** grammar checking, wording suggestions, and improving sentence
746 fluency while preserving the authors’ original technical content.
- 747 • **Mathematical support:** verifying the correctness and readability of some derivations and
748 notations, without introducing new technical results.

749 No LLM was used for generating research ideas, designing experiments, analyzing results, or writing
750 original scientific content. All conceptual and technical contributions were made by the authors. The
751 authors take full responsibility for the entire content of this paper, and LLMs are not eligible for
752 authorship.

Table 4: Statistical reference means and standard derivations for human’s ACT-R model parameters according to (Pavlik Jr & Anderson, 2005; Lewis & Vasishth, 2005; Bothell, 2020; Borst & Anderson, 2017; Said et al., 2016; van der Velde et al., 2022).

Parameter	Mean value (μ)	Standard derivation (σ)
Decay Rate (d)	$\mu_d = 0.50$	$\sigma_d = 0.05$
Activation Noise (s)	$\mu_s = 0.32$	$\sigma_s = 0.08$
Threshold (γ)	$\mu_\gamma = -0.50$	$\sigma_\gamma = 0.71$

A.2 ACT-R MODEL AND HRS-MD CALCULATION

To complement Section 2.3, we provide the full ACT-R retention model together with an intuitive explanation of its parameters and their roles in memory dynamics.

ACT-R base-level activation and retrieval probability. In ACT-R model (Pavlik Jr & Anderson, 2005; Lewis & Vasishth, 2005; Bothell, 2020), each practice creates a memory trace. If the learner is exposed to the practice at times t_1, t_2, \dots, t_k , *memory strength* or *base-level activation* at a later time t is:

$$w(t) = \ln \sum_{i=1}^k [\Delta t_i]^{-d}, \quad (7)$$

where d controls how quickly each memory trace decays (larger d implies faster forgetting) and Δt_i is the time elapsed since the i -th exposure.

Memory activation is converted to a probability of successful recall using a logistic function:

$$P(t) = \frac{1}{1 + \exp\left(-\frac{w(t)-\gamma}{s}\right)}, \quad (8)$$

where γ is the retrieval threshold (the minimum activation required for successful recall), s is the activation noise controlling stochastic variability in retrieval. Together, equation 7 and equation 8 fully specify how ACT-R model predicts forgetting curves as a function of practice timing.

Why Mahalanobis distance? Human experiments report many (d, s, γ) triplets, forming an ellipsoidal distribution with mean μ and covariance Σ . We evaluate the plausibility of the fitted ACT-R model parameters $\hat{\theta} = [\hat{d}, \hat{s}, \hat{\gamma}]^\top$ with respect to human reference distribution characterized by $\mu = [\mu_d, \mu_s, \mu_\gamma]^\top$ and a covariance matrix Σ . Given μ and Σ , the squared Mahalanobis distance is

$$D^2 = (\hat{\theta} - \mu)^\top \Sigma^{-1} (\hat{\theta} - \mu),$$

which normalizes deviations by considering parameter scales and correlations. We report an interpretable similarity score $\text{Score} = \exp(-\frac{1}{2}D^2)$, interpretable as a Gaussian RBF under the Mahalanobis metric, to account for differences in parameter scales and their empirical correlations. Score measures how close a model’s retention behavior is to that exhibited by humans. A small value indicates that the model’s retention parameters fall within typical human ranges, whereas a large HRS-MD indicates atypical or non-human-like retention dynamics.

Estimating the covariance matrix Σ . Table 4 reports the reference means and standard deviations of human ACT-R parameters from prior studies. Since only marginal statistics are available and raw human samples (hence empirical correlations) are not, we construct the covariance matrix from these reported scales. Let $\sigma = [\sigma_d, \sigma_s, \sigma_\gamma]^\top$ denote the standard deviations and $\mathbf{D} = \text{diag}(\sigma)$.

By default, we adopt the independence approximation, assuming zero cross-covariances:

$$\Sigma = \text{diag}(\sigma_d^2, \sigma_s^2, \sigma_\gamma^2).$$

This choice reflects the fact that no reliable correlation estimates are available and reduces the Mahalanobis distance to a multidimensional z -score distance.

A.3 MODEL ARCHITECTURE AND PARAMETERS

GBCL. SGD, ER, and EWC employ the neural network structure. The neural network takes as input the current discrete state and a task identifier. By default, both are embedded through separate layers of size $\frac{d}{2}$ (with $d = 64$), concatenated, and passed to a multilayer perceptron with two hidden layers and ReLU activations: The output produces logits over the next-state space, which are transformed into probabilities via softmax. Weights of linear and embedding layers are initialized with Xavier uniform initialization, and biases are set to zero. Unless otherwise specified, we set $d = 64$, `num_states= 4/8`, and `num_tasks= 2`. Optimization in all cases uses stochastic gradient descent with cross-entropy loss. While the network architecture is identical across methods, their training hyperparameters differ. For SGD, the model is updated online using only the current sample with a learning rate of 0.001. In ER, a replay buffer is employed with a capacity of 8000 samples, a replay ratio of 0.5, and a batch size of 32. For EWC), the model is trained with the same online updates as SGD but augmented with a quadratic penalty weighted by the diagonal Fisher information matrix; the optimal regularization strength is set to $\lambda_{\text{EWC}} = 700$.

LLMs. LLaMA3-8B (Meta, 2024), DeepSeek-R1 (DeepSeek, 2024), MAMBA (TH, 2024), and RWKV-7 (RWKV, 2024) are used in their publicly released configurations without further fine-tuning. LLaMA3-8B is a transformer-based decoder model with approximately 8 billion parameters and a context length of 8K tokens. DeepSeek-R1 is a transformer model with roughly 7 billion parameters, optimized for reasoning tasks, supporting up to 8K tokens. MAMBA is a selective state-space model with about 8 billion parameters and an extended context length of 16K tokens. RWKV-7 adopts a hybrid RNN–transformer architecture, with around 7 billion parameters and a context length of 4K tokens. For all large models, we adopt the default inference hyperparameters recommended by their authors: greedy or temperature-controlled decoding with temperature in $[0.6, 0.8]$, maximum generation length matched to the context window, and no additional training beyond the released checkpoints.

A.4 ICCL BENCHMARKS

DMC benchmark: To illustrate how a Discrete Markov Chain (DMC) prompt is constructed, we provide an example containing both the target task transition matrix P_{target} and the interference task transition matrix $P_{\text{interference}}$. P_{ij} represents the transition probability from state i to state j . During prompt construction, samples are first generated based on P to form the target task sequences, and later partially replaced by samples drawn from $P_{\text{interference}}$ to simulate interference phases.

Table 5: Example transition matrices for the DMC benchmark with total 4 states. P_{target} denotes the target task transition matrix, and $P_{\text{interference}}$ denotes the interference task transition matrix.

P_{target} (Target Task)				$P_{\text{interference}}$ (Other Task)			
0.2694	0.2987	0.2681	0.1637	5.6×10^{-8}	0.50	0.50	5.6×10^{-8}
0.2161	0.0106	0.3997	0.3736	5.6×10^{-8}	0.50	0.50	5.6×10^{-8}
0.2489	0.2278	0.2645	0.2588	5.6×10^{-8}	0.50	0.50	5.6×10^{-8}
0.3021	0.2828	0.2460	0.1691	5.6×10^{-8}	0.50	0.50	5.6×10^{-8}

Examples of SP, MP and DP are shown below. Each prompt begins with the instruction “Learn pattern, predict next state. [SWITCH_TO_INTERFERENCE] = interference pattern starts. [SWITCH_TO_NORMAL] = target pattern resumes. Predict next:”, followed by a continuous sequence of sampled states. The control tokens [SWITCH_TO_INTERFERENCE] and [SWITCH_TO_NORMAL] mark the onset and offset of the interference period, respectively, allowing the model to experience alternating target and interference tasks within a single context.

SP prompt example (DMC benchmark)

```
Learn pattern, predict next state. [SWITCH_TO_INTERFERENCE] =
interference pattern starts. [SWITCH_TO_NORMAL] = target pattern
```

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

```
resumes. Predict next:  
013110301001313012331213202301331010200102331222100120122123222323012  
2231320223210202233032023201203[SWITCH_TO_INTERFERENCE]21211222122112  
21121111122212212211111121222112111121222112212222221111222121111  
111112111112222222122222121222112112222211221221211121121211121  
2222211221121221112211221211111112121122211211222112121212211  
122112111212122112221122212121121121112112211222121222112211111111  
22211222121111121122222121212111112112211222221122122121111122212  
1212122211211222112222112211211222112112221121222111212112112112  
21222121222122211112211221222121121122112212212111222221122222211  
21122211212221122211221111111211211222112212222222112222212121  
212221122121211221122111112111222222222211222111212121212221  
11112121121211221122122221222122221222122211222111212211212222221  
2121121221111212212212222212121212111222121221121212121212121222122  
11221222212121121122212112122112221212221211111222211222112122221  
122221122122212222111121121121121221121221212212122122221211112122112  
22211211222122121121121122121221122111121122222112211111212122121  
212221221111112212122121211211222121211112111222212122121122122  
22211212211211212121111222222211121221212112211222121211222121111121  
1111121221212111221122112211221122112211221122212112221211221212121  
21112222212121122212121122122221211222212221211222211222112121111  
221121212212222211112112112211211221122211221122211121212122221  
11222221121222212221222121121122221221221222112122112211221212121  
12221122222122221111112212121111222211212112212111122212122222  
22121122111221122122112121122112222111212221222221211212111121222  
12222122212112121121112222222222122122211111222122221221222111121  
11221122221122121222121122122121121122221222112211211121211212212111  
1122112211211122221221222122211221211211122212211111212122122221221  
21222121211222122
```

```
MP prompt example (DMC benchmark)  
Learn pattern, predict next state. [SWITCH_TO_INTERFERENCE] =  
interference pattern starts. [SWITCH_TO_NORMAL] = target pattern  
resumes. Predict next:  
013110301001313012331213202301331010200102331222100120122123222323012  
223132022321020223303202320120330212223022002301312000122133132100010  
12122312120000130332102202323310101233030100001210300001223231320233  
32302222213120231322013201221212130021301220323222010230013013011232  
013312010010100313013221230102312130201313300133120212002230230132301  
332303102013010030123131301312022201212012101323002331301012131023133  
130213121010200122013332310220321310120023202030202022121310222202232  
21212310020033200[SWITCH_TO_INTERFERENCE]2112212122221112121121112112  
212221212221222211122112212221211121122112212212111222221122222211  
2111222112122211122211221111111211211222112212222222112222212121  
21222112212121121221122211111211122222222221122211121211212122221  
1111212112121122211221222212221222211222212212221111212211212222221  
21211212211112122122122222121212121112221212211212121212121212222122  
11221222212121121122212111221122212122212111111222211222112122221  
12222112212221222211111211121121221121121212212212221222121112122112  
2221121122212221211211121122122112211121122222112211111212122121  
212221221111112212112212121121122221212111121112222121221211122122  
2221121212211211212121212221221112222122211212221111222212122211222  
1212122211211121212111112222222211121221212121122111222212111111121  
22221111211122221122212212121221211222211222211222211222212221112121  
111112112212112111122112112211121221112122112112222121112222121122121  
2111222221212112222121211221222221211222212222121122221122211212111  
22112121221222221111211211221121122112222112211222111212121122221  
112222221121222221222212212112112222122122122211212211221122121212111  
122211222222122221111112212121111222211212112212111122212122222  
221211221111222122122112121121221121222111212221222221211212111121222
```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```
12222122212112121121112222222222122122211111222122212212211222111121
11221122221221212221211221221211211222212221112211211121212112112212111
112211221121112222122122212221122121121112221221111112121221222221221
21222121211222122
```

DP prompt example (DMC benchmark)

```
Learn pattern, predict next state. [SWITCH_TO_INTERFERENCE] =
interference pattern starts. [SWITCH_TO_NORMAL] = target pattern
resumes. Predict next:
013110301001313012331213202301331010200102331222100120122123222323012
22313202232102022330320232012033[SWITCH_TO_INTERFERENCE]1211222122112
211211111222122122111111212221112111121222112212222211111222121111
111112111112222222122222121222112111222222112211221211121121211121
2222211122111212211112211221211111111212112221[SWITCH_TO_NORMAL]230
102312130201313300133120212002230230132301332303102013010030123131301
31202220121201210132300233130[SWITCH_TO_INTERFERENCE]111121122222121
212111111211221122222211221221211111222121212121221112211222112222211
1221121122211211221212222112121121112112212221212221222211112211221
22212112112211221221211122222122222211211[SWITCH_TO_NORMAL]232210
213321213131231012001000131020022212312222231232101323213020313121222
23131231321033310012131001[SWITCH_TO_INTERFERENCE]22222222221122211
121211212122221111112121121211222112212222212221222211222212212221111
2122112122222212121121221111212212212222212121212111222121221121212
121212112222122112212221212112122212111[SWITCH_TO_NORMAL]303310031
3121232020001000322210233003323012222120320323212322101203012301021
33002001313032023312131[SWITCH_TO_INTERFERENCE]1112112211222211211222
1221211211211221212211221111121122222112211111212122121212221221111
112212112212121121122221212111121112222121221211122122221121212211
2112112121212221221111222212221121222[SWITCH_TO_NORMAL]12123233130233
01323131313330030122020301000032322231210203213020300220123133130121
200130232312001301[SWITCH_TO_INTERFERENCE]1222211222122121211221221122
22111222211122222122211121121111111211221211211111122112112211112122
1112111222212111222121112121211122221212112222121211221222212112
222122221211222211222112121112211212212222221111212112211211221
1122221122112221112121211222211122222112122221222212221221211222212
21221222112122112211221212121112221112222221222211111122121211112
2221112121122121111222121222222121122111221122122122121212212222
111212221222221211212111121222122221222121121211122222222212212
221111122212222122122112221111211222112212122212112212212112112
2221222111221121112121122122111122112211221122212212221222122121
121112221221111111212122122221221212221212112221221222122122122122
```

REGBENCH benchmark: We also construct a symbolic sequence benchmark based on RegBench, where both the target and interference tasks are generated by probabilistic finite automata (PFAs) implemented as deterministic finite automata (DFAs). Concretely, we first sample a target DFA and an interference DFA with different numbers of states, alphabet sizes, and random seeds. Each DFA induces a stationary stochastic process over a finite alphabet, and we generate short examples by running the corresponding DFA for a random length between L_{\min} and L_{\max} symbols. We concatenate such examples into a single long sequence and insert control tokens `[SWITCH_TO_INTERFERENCE]` and `[SWITCH_TO_NORMAL]` whenever the generator switches between the target DFA and the interference DFA. This yields a context in which the model must continuously track the active PFA and predict the next symbol in the sequence.

SP, MP and DP examples are explained as follows. In our setup, the target task is produced by a DFA with 4 states, alphabet size 10, maximum 2 outgoing edges per state, and a fixed random seed (`seed=42`) to ensure reproducibility. The interference task is generated by a distinct DFA with 6 states, alphabet size 20, and maximum 3 outgoing edges per state, ensuring a substantially different symbolic distribution and transition structure. Each DFA defines a stationary symbolic process, and we generate short examples by running the corresponding automaton for a random length sampled from $L \sim \text{Uniform}(10, 15)$. We concatenate $N_t = 320$ such examples to form a prompt. We additionally generate a separate test split consisting of $N_{\text{test}} = 10$ pure target-DFA examples, whose

lengths are sampled from Uniform(8, 10). This test set assesses how much target-pattern ability the LLM retains after experiencing last historical demonstrations in a prompt.

SP prompt example (Regbench benchmark)

Learn pattern from probabilistic finite automate(PFA) sequence.
 [SWITCH_TO_INTERFERENCE] = interference pattern starts.
 [SWITCH_TO_NORMAL] = target pattern resumes. Predict next: m n m n m n
 m n m s w, m s m s w m n m s w m n, m s w m s m s w, m n m s s m s m s w m s w, m s s m s m s w m s s, m s m s w
 m s s m s w m s s, m n m s s w m n m n, m n m s s m s m s m s, m n m s s m s m s w m n m, m s s m s m s w, m n m s w m
 n m n m s w m, m s s m s w m s m s w m s s, m s s w m s m s m s w m s, m n m n m s w m s m, m s s w m s s w m n, m n m
 n m s m s m, m s s w m s m s m s m s m s, m s s w m n m n m n, m s s w m n m s w m n m s s, m s s w m n m n m s w m,
 [SWITCH_TO_INTERFERENCE], n v e m x x d r u x, d p a x u q r t e m, d x d r d n m q r u r e, d x u q n m q d
 p e, r t e m x p e v e v a x d n, n g d p e g d x d n, r u x p a p a x u q r, n v a p a x t e v a x d d x d, n g n g n g r t e v e
 m r a, d p a x d n v e g d x, r t a x t e v a x u r a, n v e g r t a p e v e, n g r d r u q r t e m x, n g r t a p a x u q d x, n v
 e m q r t a x u r e, r u q d x u q n g n m r a x, n m x p a x d d x d n v, r t a x u x p e v a p a x d, d x d r u q d x u r a p e m r
 , r u q n m r e g r t e v, d x t e v a x u r a p a, r d n v e m r e m x x d d p, n g d p a x d n g n g, d p a x t e m q d p e v e m
 , d p e g n m r e g r d r, d x d n g n v a p e v a p, r u x x t a x u r e v e v a, n g r u q d p e m r e g r d d, d x t a x d d p e v
 , n m q r d d p e g n v e g, n m x p e g n m x p, d x u r a p e g d p e m x p, d x u r e v e v a x d n g r d, r d d p e v e m r e g n v
 , d x d n g d x d n g r u, r d n m x x d r d r t e, n g d p a p a p e m q, r t a p e v e g d p a p e g, r u r a p e m x p e g r u x x,
 d x u q d p e g n v a, d x t a p a p e g n g d x, n v e g r u q d x d n v a p e, d x t e g r t e g n m q d p e, r u x x u x t e m q,
 n g r t e v e g n m r a, d p a x u q r u x x d n v a x, r u x p e m r a x d d p, r d d x u x x u q n m r e v, n v e g r t a p e g d x
 d d p, d p a p e m r a p e g d p a p, n m r e v e m r a x u p e v, d x d r d n g d x d, r t e g d p e m x x u r, n m q n g r u q n
 m x p e m, r d r d n v a x d d x d d, d x d d p e g d p a p a x u r, d x d r u q r d r u r a, n v e m r e v a p a x u, d p e m r e m
 x x u q n g r d, n v e m q n m r e v a p, n m q d p a p a p e, n m r a x d d x t a p, r t a x u r a p e v e m x p, r t a x u x x d r
 t e v e v e, n g d x t a x t e g r u q r t, r t a p a x t a x t a, n v a p a p a p a p a p e g, d p e g d p e m q d p a, d p e g n m
 r a p e, d x d r d n v a p a x t e v, d p a x u p a p a, r t e m r e v e m q n m x x, d x d n m q n m x p a p a, r u q n g r d n g
 r d n m, n v a p a x u r a p a x d, d x t a p a x u q r d d x, n g n g n m q d x d d p a, d x u x x d d p a x d r, n m q r d r d r u
 r e m x, n g r t e m r e m r, n v e v e v a p a p a, n v a p e v e g r t, d x u q r t e g n v e v e g d, r d n v e v a x t a p,
 n g n m r e g n m q, d x t e g d x u x p a, n g r u x p e v e g d x u r a, d x u q d p e m r a p e g r u, d p e v e m q n g r u q n m,
 n m q n v e m x p e, n v e g d p a x d n, r u r a x d d x d n g, d x d d x t e m x p a p e m x, d x u r a p a x d d x d d x u, d x u
 q d x u q n g n, n g d x d r u r e v e g d x, d p a x d r d r d n, d p e m r a p a p e, d x u r e m x x u r e m, d p e v a p a x u q
 r d d, r d r u r a x u r e g, n v a x t e v a p a p a, n v e v e m q r u q n v e g, r d d p a p a x d d p a p, d x u q d x t a p a x,
 d x t a p a x u r a p a x, n g r u q r d r t a, n g n g r d d p e g, d p e g n v a x u r a x, r u q d x d r d d x d, r t e m x x t a
 x t a x d n, r d r u r a x d r u r e g d x, n m q n g r u x p a, n g d x t a p a x t a x t a p, n g n v e g d p a p, n g r t a p e v
 a x t e v e, d x t e v a p e g n m x p e, d p a x u x x d d p e g, n g d x t e m q d p e g n m q, n v a p e g d p a x u r, r d r u x
 p a p e m r e m, d p a p e v a x t e, d p e m r a x u q n v a p, n g d x d r u x x d n, n v a p e v e v e v a p e m x, n g n m r e g
 r t e g, d p a x u q r d r u x x d n m, r t e v e v e v e v e g d x d, r d r d r d d x t e g r d r, n v e m q n v a p a p e g d p, r
 t e v e m r e v a, n v a x u q r u q d p e, n g r u x x u x p a p e m r, n v e v e g r d r u x p e, d x d d x t e v a p a x u r, n g
 n m q r t e g d p a p e m, n g d x d d x d n v e v a x, r d d p e g n v a x u r e m q, n m x x t e v a x d n v e v e, d x t e g d x d
 r u, n m r a x d n v e g n m r, n g d x d d x d r u x p a p e, n v a x u q n m q r, n m q n m x x u x p e v a x, n m r e m q r t e v
 a x u, d p a x t a p a p e m r a, r d n g d p e v a p, r t a x d r t e m q d x t a, r u r e g r u r a p e g d x u, d p e v a x u r e
 v, r t a x u q r t a p e g d x, r d n g r u x x d r d n g r, d x t a x d n g d p e m r e m, d p a x d d x d d p a x d, d p a p a p e
 g n v a x, r d n g r d d x u r a x d d, r d d p a x u r a p, d x u x p a p a x d r, d p e g n g d x u q r d n, n m r e v a x d n g r
 t e m, d x d n v e m x x d r u r, d p e v e g n v e v a p a p a, r t e g r t e m q r t a p a, r t e v a x u x p a x u, r t e m q d p
 a p a x t, n m q r u q r u r a p e m q, n v a x u r e m x p e g, n m x p a x t e m x p a, n g r u q d x t e g r, n m q n g n v e v a
 x, r u r a p e m r a p e, r d r t a p a x t a p a x u, d x t a p e m x x d n m, r t a p a x u q n m r a x u x, d p a x d d x u x x t
 e, n v e g n g d x t a, r d n g d x d d x u x p, n v a x d r u r a p a x d r d, n v a p e v e m r a x t e v, d p a p e m q n v a p a
 , d x t e v a p a x t, r t a x t e g r d r, n v e v e m x x t a x t e m, d x u q n g d x d n g n g, n v a p e m r e m r a, r u x x t
 a p e v e v, n g n g n g n m r a, d p a x u x p e v a p, d x d n g n v e g d p e, r t e g r u r a x t e g r t, n g n m q r t e m q r
 u r, d x u r e m r a x d r d n g d, d p e g d x u q r d n m, d x t a p e v a p e m q n, d x u q d p a x d n, n g n v e g n m x p e m
 q n m, r u x p e g n v a x d r d, r d r t e m q r u x p a x t e, r d d x u r e m x p a x, r t e v a x u r a x d r, n g n g r t e g n
 g d x u, n g d p e v a p e m, d p e g r u q r u q d p a x, r d n m q d x u r e g r t e, r t e m q d p e v e v a x u, r d d p e g d x
 d n v e v e, d x u x x t e g r d d p e m q, d x d d p e v a p a p e, n v e m x p a p e g, n m q d p a p a p e g n, n m r a x u r e m
 r e g r t, n m r a p a p a p a p e m, n g d x d n v a x t e, n g n g n g n g n m r a x d n, d x t a p a x d n m x, d p e m r a x u r
 a x, r d d p a p a p e v a x, n v a p e v a x d d p a x, n v e g r d r t e g r u, d x t a x t a p a x u r e g n, d p e g d x t e g r,
 n m r a x d d p a p e v a, n g r u q n g r u r, d p e m r a x u r a x d r d n, r d n m x p a x d d, d x u x x t a p a x d r d, n g d
 p e v a p e g n v a x, n v e m r a p a p e g r t, n v a x u x p e v a x t a x, d x t a x u r a p a p e v, d p e g d x d n m q, d x d
 r t e g d p a x u, n m r a x t a p e g, n g r t e m q d x t e m x p, n v a x u x x u x p e, d x d n m r e v a x d, r u x p e m r a x
 u, n g d x d n m q r t e v, d p e m q n g r u q d, r u r a p a x t a x t e, d x t e m r a x d d, d p e m q n m r e m q n m r a, d p
 e g n v e v e g, d x t e v a x t e m q r d, n m q n v a x u x p a x u, r d r d d p e v a x d d x d, n v e g r u r e v a x, d x t a x
 t e g r u q r u, n g r t e m r e g r t a p a, r d n g d p e m x x t a p a, r d d x d r u r a x d r u r e, d x d d x d d p a x, n g n
 g d x d n m r e g, n g r u x p a p a x u q, r t e g n g n g d p e m, d x d r d r d r u q d x d, d x d r t a p e g d x u x p a, n g r
 t a x d n v a, n v a x d r t a x t e g r t a, r d r t a x t a x u, r d n g n m x p e g d p a p e, d p a x u r a x t a x t e m r, r d

1026 ruxxurapapem, rdngdrngdrnva, rddpaxuxxu, ruregdxdxtap, nmxpevapegrdrdn, n
 1027 vemxxddpemx, ruxxdruxxd, dpapegngrdrte, dxdnmqrtegrtapa, ruqdpaxtemrape, d
 1028 xddpaxuqrta, dpaxuqrddpe, nmrpamqrquqd, nvegrdnmxxdrtem, dpaxdruqru, nmqdx
 1029 drdngdpemr, dxtemraxddpa, dxdnvaxuxxura, ruxpemxpaxuqdpe, nmremqnmremx, dx
 1030 ddpapaxtap, dpegdpemqrteve, dxdrurapegrd, ngdxtegrugrdrd, ruxpapegdxta, rd
 1031 dxtevaxdn, ruraxuqnvaxuqng, dpaxtevapaxta, dxdnmxxuraxdd, ngrtaxtdngngd
 1032 , nmqrugnmxpape, nmxxtevakura, nvevapemqrquqd, ruxpevevegdxuq, ngdpevevemqd,
 1033 rddpapemxxu, nvapemremr, dpevevemqru

MP prompt example (Regbench benchmark)

1034
 1035
 1036 Learn pattern from probabilistic finite automate(PFA) sequence.
 1037 [SWITCH_TO_INTERFERENCE] = interference pattern starts.
 1038 [SWITCH_TO_NORMAL] = target pattern resumes. Predict next: m n m n m n
 1039 mnmssm, msswmnmnmssmsm, mssmswmsm, mnmsswmsmsswmn, msswmsmsswm, msswmnms
 1040 swmnmss, mnmssmsmsms, mnmsswmsm, mnmsswmsmssms, msswmsmsswm, mnmssmsmsms
 1041 wnm, msswmnmswmmnms, mssmswmsmss, mnmnmssmsms, mnmnmssmsms, msswmnmn, mnmnmssw
 1042 ms, mssmswmsmsswm, mssmsmsmsm, mssmswmmnmnmss, mssmsmsmsswmn, mnmnmsswms
 1043 swms, msswmnmnmn, mnmnmsswmsm, mssmsmsm, mssmswmsmssw, mnmssmswms, ms
 1044 smswmsms, mnmnmnmssmsms, mnmsswmmnmms, msswmnmnmss, mssmsmsmssw, m
 1045 nmnmsswms, mnmsswmsm, mnmsswmmnmms, mnmnmsswmmms, msswmnmsmsswm, mnm
 1046 swmsm, mssmsmsmssw, mnmnmnmnmss, mnmnmsswmsmss, mssmsmsmsswm, mnmnmnmssmsm, ms
 1047 wnmnmnmsswms, mnmnmnmssms, mnmssmsmsswms, msswmsmsswm, mnmnmnmssms, m
 1048 sswmnmnmnmn, mnmnmssmsmsswm, mssmsmsm, mnmssmswmsms, mnmsswmsw, ms
 1049 smswmmnmssw, mssmswmsmss, msswmnmsswmn, mnmnmssmswms, msswmnmsswmn
 1050 , mnmnmsswmn, mssmsmswmsmsw, mssmsmswmsm, mssmswmmnmss, mssmsmsmsswmn
 1051 , mnmnmsswmmnm, mssmsmswmmnmss, msswmsmsswm, mnmsswmmnmss, mnmnmssmsm
 1052 smsw, mssmswmswms, msswmsmsswm, mnmnmssmswmmn, mnmnmssmswms, msswmnmsswmn, m
 1053 nmnmssmswms, mssmswmsw, mssmswmsw, mnmnmsswmmn, mnmnmssmswmsw, mnmnmsswmmn, m
 1054 nmnmssmswmsw, mssmswmsm, mnmnmsswmmn, mnmnmsswmmn, mnmnmsswmmn, mnmnmsswmmn
 1055 n, mnmnmsswms, mssmswmsmsswm, mssmswmmnmss, msswmnmnmnmss, mnmnmnmssmsm
 1056 s, mnmnmnmnmnmss, mnmnmnmnmnmss, mnmnmnmnmssw, mnmsswmswms, mnmnmssw
 1057 [SWITCH_TO_INTERFERENCE], dxtegdpmxxte, dxdrugnvax, rrdnmraxdrdr, ng
 1058 dxuqdxdr, dxdnvaxdn, ruqrurevape, dpaxtapegruraxd, rurevapaxupapa, rtax
 1059 taxdrtemqn, dpemqrtapa, dxtapaxdn, nvaxdnvegdp, rdnvemxpapegngn, rdruraxu
 1060 xpegrur, rrdrdrdde, dxuqngruxxurem, ruqdpapape, rtapemregn, rdruvaxurap,
 1061 rtaxuxxuregnv, ngnvaxdruxx, dxddxtemxxuq, dxtevatxaddpa, ngrtegruxxura, r
 1062 urevemqrquq, rddxtevevem, dxtaxtapem, dpevapapax, rtapemqdpaxd, nvevegng
 1063 dn, nmxpaxddxdruxx, ngnvemqngruuxuq, dpemrxtax, dxuqdxtemqdxddp, dpemxxu
 1064 xt, dxtemxpemregdp, rddpevapaxtapa, ruqdpaxtegnv, dxuqdpaxuxpaxte, dxdrure
 1065 grure, ngnvapegrteva, ruqrtemqdp, rtapemremqdx, dpemxpegdpe, dxuxpevegdx
 1066 apa, dpegrtapegr, rtemregngregdp, nmxxtegdpeveveg, ngngngruqnvapa, dxtaxd
 1067 dxtegrura, nmxxtapaxu, dxuqrdnvevap, dxtapevapegrte, dxuqnvapaxtap, rdnvem
 1068 xxdrugru, dxdnmxxuregdpap, dpemxxdnmnmre, ngrtaxddxuregnm, dpaxdnmregdx
 1069 ev, ruqnvaxtax, dpevegdxuregr, dxurapegnvapegr, dxuremnm, dpemxpevapaxur,
 1070 dpaxtegdxd, rtemqdpengrd, ngdpapaxdr, nmregnmxpevemq, nvapapevevapaxu,
 1071 ruxxxdxure, nmqrugnqrurax, ngdxtedpapeva, rddxdnmrapapegn, rtemxxuxxtem
 1072 x, rtegrtapemq, ngdpapemraxdnv, ngrtemremr, rdrtaxtemxp, ruxpevaxuqngd, dp
 1073 apevegdxure, ruxxdrtaxu, rtaxtdnmqnvax, nmrapexaxteva, nmxxddxaxd
 1074 r, nmraxuxxtemq, dpegrdnvaxurax, dxddxupapap, dpapemqngng, dpapevemreg, d
 1075 pemremxxdr, nvaxuxpexp, ngnmqnmqdxtemr, rddxapaxtegn, nmraxdrddpaxdr, rt
 1076 emxxuraxdr, dxurevaxdr, ngdxuqnvax, dxddpaxuqnmqngn, dxdrtevakdrdnm, rte
 1077 grtemqrquq, ruqrdrtemq, nmqdxtegn, dxtegngrdruqr, rrdpdpemxxuq, dxdngva
 1078 pax, nvaxddpaxte, dpeveveg, rtemrapaxdn, rdnmremrxta, nmxpapaxuqnvax, dx
 1079 dnmxxuxtax, rdruvaxdngnmra, rtaxuqrdngdp, rddxuremxxpaxd, rdruqdp, dpe
 1080 mraxtapapemx, nvapaxddxngn, ngnmrxttegnm, ngrdruvaxurax, nmxxdrurevap
 1081 , dpevaxurevemr, dpegrdrtap, rtapapegrdrtem, ngdpevevegrev, nmxpapevegdx
 1082 r, ngruxxuraxdrdn, rdrtemxxtapapax, ruxxuxxnmxp, dxtapaxuxp, dpevaxtegnve
 1083 , dpevevapaxte, dpaxuraxtevap, dpemxxddxdr, dxdrdrdrdnvemxx, rddxuregdp,
 1084 rtapevevemq, ngruqnvapemq, dxuraxddpevem, dxurapaxurap, rddxuqrtemremv, r
 1085 tapemqrux, dpaxuxxtegnm, dpapemraxu, rtapemrapemxxpax, ngdpapemxx, ruraxdd

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

```
pemxpa, dxuxpemrapemqd, dxuxxuxxtapax, dxddxtapemqdx, rddxdruraxure, ruxx
uqnmx, rdngdpvaxdr, dpaxddtax, dxtegnmremxpap, dxddxtvapa, rdnmxpaxddp,
nvapapemra, dxuraxtegdpe, ruraxdrtaxu, nvaxuqrdrur, ruqrtevegr, rtapemxpa
pemxpe, rtapemregr, rddpevemxpegn, dpemqdxuxpemr, ngngrtaxuraxdn, dxtegrur
emq, rddxddpapemxx, dpaxtapaxteva, ngdxuxpapemrax, ngruxpapevapeva, rdve
grtem, dxdrumxmpap, dxtapaxuqdp, nmxxdrdruxpa, rdngngngrdrdn, rdngdxurax
drtax, dpaxuraxdr, dxddpaxddxtegd, ngnmqrddxu, ngdxdngruxxta, rtemremre
mregnv, ngnvaxurevegng, ngdxtaxuxpemq, ngrtemraxu, nvegruraxuqdp, dpapaxdn
grtapax, dxtapegruxpa, nvaxtapaxt, rtegruqrtapev, rdnmxxttegrteva, nvevaxu
qngrdru, rdnvaxdrdnmr, ruqdpagnvap, dpaxtapegrdr, dxtapaxtaxux, rtemxpa
xta, dpevegngdxuxpax, rddpapemxxur, ruxxdddxxuxure, nvapapapemrevap, dpegnm
qnvapa, rdnvaxuremqn, ruxpapapegdpev, ruxpevevapap, dxuqdpagrngng, nvapeg
ruremq, ngrddpempqn, nvemqrddxruqdp, rtemqnmqrddx, rdnmxpevevegr, dxtemr
emqnmreve, dpegrddpapegn, dpaxdnmqdxur, dxtevapaxtaxu, nvapaxtegruqrd, dpe
gnmxxtaxu, ngrtegngr, dxdnvapapa, rtaxtegnmxp
```

DP prompt example (Regbench benchmark)

```
Learn pattern from probabilistic finite automate (PFA) sequence.
[SWITCH_TO_INTERFERENCE] = interference pattern starts.
[SWITCH_TO_NORMAL] = target pattern resumes. Predict next:
mnmnmnmnmssw, mssmswmnmsswm, msswmsmsw, mnmssmswmssw, mssmswmss, mss
mswmsmswmss, mnmsswmnmnm, mnmssmsms, mnmssmswmnm, mssmsmsw, mnmsswm
nmnmssw, mssmswmssmswmss, msswmsmsmswms, mnmnmsswmsm, msswmswm, mnmn
mssmsm, msswmsmsmsmsms, msswmnmnmnm, msswmnmsswmnmss, msswmnmnmssw, mnm
nmssmsm, [SWITCH_TO_INTERFERENCE], nvaxuqdxdd, nvegdxdrdnve, nmraxdruxp,
dxuxpegruqnmxx, rurevapapa, dxddxdrtemx, rtapaxtemrapemx, ruqnmremrapeva,
ngnvemqdpem, dxtaxdrunva, rtevemqdpap, ruxpevapapeg, ruxxtemqnmra, rtegrd
rtapap, dpapapaxuqdpax, rddpaxdnmxxt, dxuqngnvaxurap, nmxpemrapemxpev, dxt
apapegrdr, [SWITCH_TO_NORMAL], mnmnmnmsswm, mssmswmnmssmsm, mnmnmnmssw
m, mnmsswmnmnmssm, mssmsmsmswms, mnmnmsswmnmnm, msswmswmnmnm, mnmnmsswm
nmnmss, msswmnmnm, mnmsswmsmsw, mnmssmsmsm, msswmsmsmsm, msswmsmsmsw
mss, msswmsmsw, mnmnmsswmsm, msswmsmsw, mnmnmsswms, msswmsmsw, msswmnmsswms
s, msswmnmnmssms, msswmswmss, msswmnmnmssw, [SWITCH_TO_INTERFERENCE],
rtevevaxtemqngd, nmqnvemxpapapeg, dpemrevapem, ruxxdrdnmxxu, nvaxuqdxdrd
rtax, dxddpevaxdnv, dpemremraxdngr, rtevemregnvapap, nvapegnvegndv, rd
nvapevaxuxpem, nmxpemraxd, dxdnvegdpape, rdrtegnqdppegd, dpegrtaxdnvev, n
mxxddpapapaxur, nmxpaxddpaxt, rtegnrtemqn, ngdpapegnvaxuxp, rtegrtevemq
d, [SWITCH_TO_NORMAL], mnmsswmnm, mnmsswmnmnm, msswmnmnmnmssw, msswmnm
msswmsw, mnmsswmnmnmnm, msswmnmssw, mnmnmnmsswssw, mnmsswmsm, mssm
smsw, msswmsmsw, mnmsswms, msswmsmswmsw, msswmswmsw, msswmswmsw, msswmsw
nmnmss, mnmnmnmssw, mnmnmnmnmssw, mnmnmnmnmss, mnmnmnmsswms, mnmssw
mssmsm, mnmnmsswms, mnmnmsswmsw, [SWITCH_TO_INTERFERENCE], rtapemxx
dr, nmrpapemqnmrax, dpaxuqrtaxd, ruqdpvapa, nvapapapegn, ruxxdrnvapaxura
x, nvaxuxpevapevev, ngdxdxtapaxte, rdrtapegng, rtevegnmxx, dpemqdpapap, n
mxxtapegnvegdp, nmrevemrevapaxu, nmraxtaxtem, rurapevapapeva, nvaxddpegr
, ngdpaxureg, nmrapegrureg, ngdxuxxurapev, [SWITCH_TO_NORMAL], mssmswmsw
ss, mnmnmnmnmnm, mnmnmsswmsw, msswmswmsw, mnmnmsswms, mnmnmnmnmss
wm, mnmnmssw, mnmsswmnm, mnmssmsw, msswmswmsw, msswmsmsw, msswmsmsw, mss
mswmswmnmnm, mnmsswmnm, mnmnmnmnmnmnmss, mnmsswmnm, mnmnmsswmnmssm, mnm
nmnmsswmnm, mnmsswmnmssm, mnmnmnmnmnmssms, mnmnmnmssw, mssmsmsmsms,
[SWITCH_TO_INTERFERENCE], nvapemqrte, ngngnmraxtaxd, rurapevemxx, rtapem
rapaxddpe, ruqdxdngrdn, ngrtaxddpegng, dxdrdrdrdruremx, dpegdpaxtevegd,
rtegrtaxdnngveg, dxdrddpemer, rtaxuqdpaxux, ruxpemrevapegn, rtemxpegdxdp,
nmxpaxdrteva, ruqdxdrurapeg, rurapaxdrtemra, dpapevaxuraxddx, rddxtem
rapaxdru, nvapapapev
```

A.5 RETENTION PERFORMANCE UNDER SP, DP AND MP FOR DMC REGBENCH

Figure 6 illustrates the retention performance of ICCL models and GBCL baselines varies on the simple task with task identifiers and $\varphi_I = 200$, under SP, MP, and DP scheduling. For ICCL methods, we observe that in both SP and MP conditions, the retention performance drops sharply immediately after finishing all target task demonstrations, and subsequently remains at a reduced level as the

sequence length increases. In contrast, ICCL under DP consistently achieves better performance when the sequence length reaches its maximum, demonstrating its advantage in retaining cross-task knowledge. However, for GBCL baselines, DP does not yield a comparable improvement and performance suffers sharp drops across sequence lengths. These results confirm that both ICCL and GBCL models suffer from catastrophic forgetting, yet DP is effective in mitigating this issue for ICCL by alleviating performance degradation over long sequences.

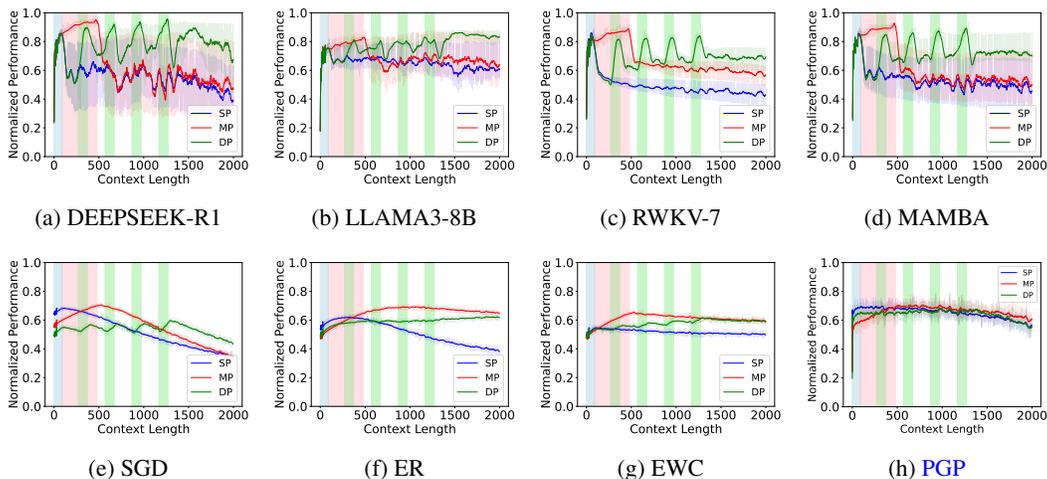


Figure 6: Retention performance of ICCL and GBCL models on simple task under SP, MP, and DP settings (with task identifiers, $\varphi_I = 200$). Target task blocks are highlighted in light blue (SP), pink (MP), and green (DP).

A.6 RETENTION PERFORMANCE UNDER SP, DP AND MP FOR REGBENCH BENCHMARK

Figure 7 presents the retention performance of ICCL models on the RegBench task with task identifiers ($\varphi_I = 20$, $\varphi = 40$, $K = 5$, and $\varphi_D = 120$) under SP, MP, and DP scheduling. Consistent with the observations on the DMC benchmark, all ICCL methods under SP and MP experience an immediate and pronounced performance drop once the demonstrations for the target task conclude, after which their accuracy stabilizes at a lower level as the sequence length continues to grow. In contrast, under DP scheduling, ICCL models maintain consistently strong performance without exhibiting the characteristic sharp decline even when the sequence reaches its maximum length, highlighting DP’s advantage in preserving cross-task knowledge. Evaluating ICCL methods on both DMC and RegBench benchmarks further demonstrates that ICCL is effective not only in purely synthetic settings but also in language-style sequence environments, indicating strong generalizability across heterogeneous task formats.

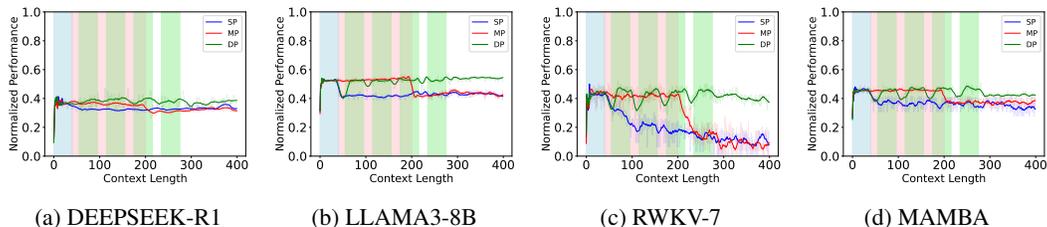


Figure 7: Retention performance of ICCL models on regbench benchmark under SP, MP, and DP settings (with task identifiers, $\varphi_I = 20$). Target task blocks are highlighted in light blue (SP), pink (MP), and green (DP).

Table 6: The comparisons between the true normalized performance of ICCL and GBCL methods and their fitted ACT-R curves.

Model	DEEPSEEK-R1	LLAMA3-8B	RWKV-7	MAMBA	SGD	ER	EWC
Pearson correlation	0.5152	0.7068	0.6171	0.6215	0.0339	0.5660	0.7253

A.7 OPTIMAL PARAMETER λ IN EWC

Figure 8 and Figure 9 present the performance of EWC under simple and complex tasks with different regularization strengths λ_{EWC} across SP, MP, and DP scheduling strategies. The results highlight the trade-off between plasticity and stability inherent in EWC. With small λ_{EWC} values (e.g., $\lambda_{EWC} = 100$), the model maintains higher plasticity and adapts quickly to new demonstrations; however, this comes at the cost of poor stability, as retention performance decays sharply with increasing context length. Conversely, large λ_{EWC} values (e.g., $\lambda_{EWC} = 1000$) enforce stronger stability but severely limit plasticity, leading to under-adaptation and overall lower performance across all scheduling strategies.

Intermediate values achieve a better balance. In particular, $\lambda_{EWC} = 700$ consistently provides a favorable trade-off, delivering higher long-term retention compared to smaller λ_{EWC} while avoiding the rigidity observed at $\lambda_{EWC} = 1000$. This effect is especially evident under MP and DP, where $\lambda_{EWC} = 700$ maintains relatively stable retention curves while smaller λ_{EWC} values degrade more rapidly. Based on these observations, we select $\lambda_{EWC} = 700$ as the default regularization strength for subsequent experiments, as it offers the most robust balance between plasticity and stability across both simple and complex tasks

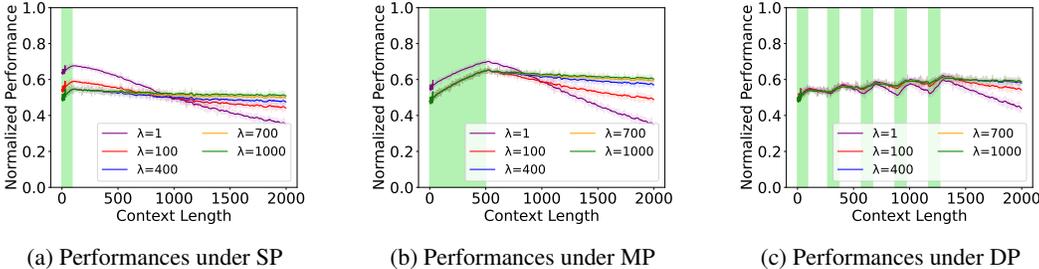


Figure 8: EWC’s retention performances under Markov chains with simple task (markov chain with $N_s = 4$).

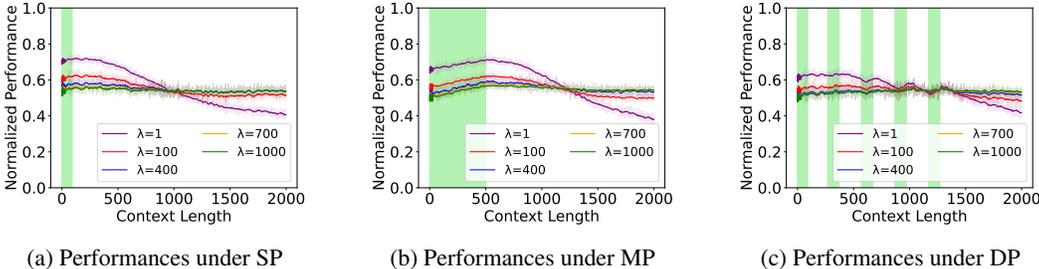


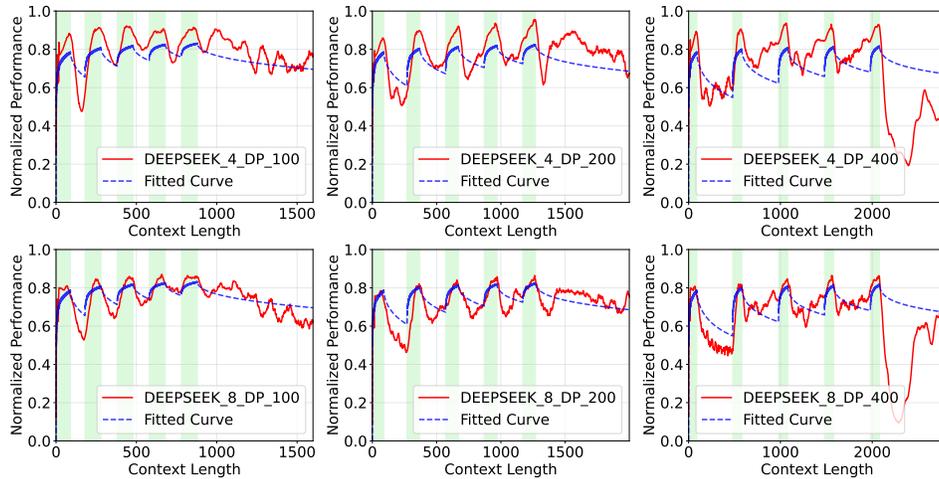
Figure 9: EWC’s retention performances under Markov chains with complex task (markov chain with $N_s = 8$).

A.8 THE FITTED ACT-R MODEL RESULTS FROM ICCL AND GBCL

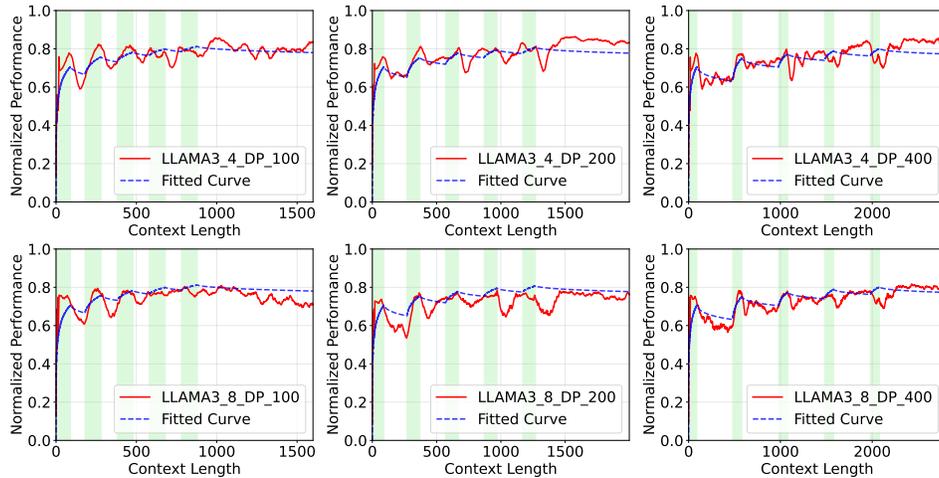
To further analyze the retention dynamics of ICCL and GBCL methods, we leveraged the experimental data of normalized retention performance across prompt sequence positions under the DP setting. We

1242 applied a modified ACT-R model to fit the performance trajectories of each method. The comparisons
 1243 between the observed normalized retention performance curves and the corresponding fitted ACT-R
 1244 curves on both simple and complex tasks are illustrated in Figures 10–16.

1245
 1246 In addition, we computed the average Pearson correlation coefficients between the true normalized
 1247 retention performance curves and their fitted ACT-R counterparts, as reported in Table 6. The
 1248 results show that, with the exception of SGD, all ICCL and GBCL methods achieve relatively high
 1249 correlation coefficients (e.g., 0.71 for LLaMA3-8B and 0.73 for EWC). This indicates that the fitted
 1250 ACT-R model provides a reliable description of the dynamics of normalized retention performance,
 1251 capturing how ICCL and GBCL retention evolves with increasing sequence positions.

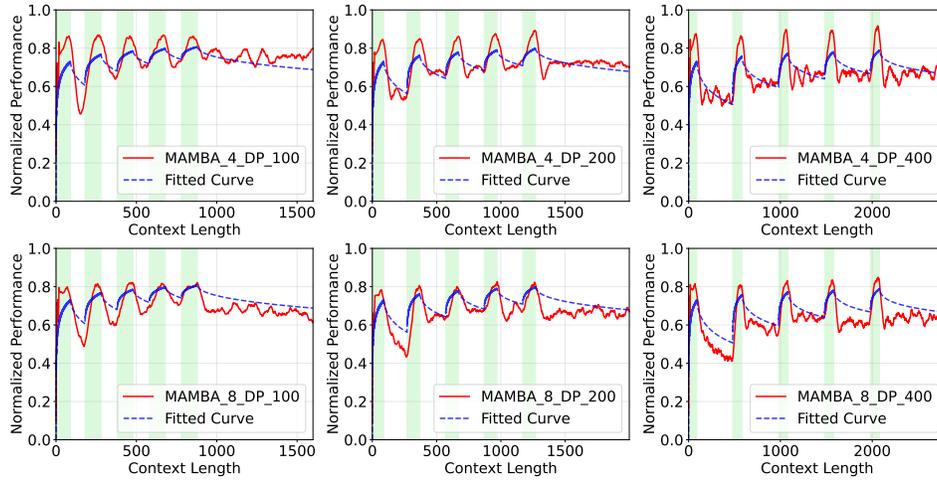


1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268 Figure 10: The performance comparisons between DEEPSEEK-R1 and its fitted ACT-R model.



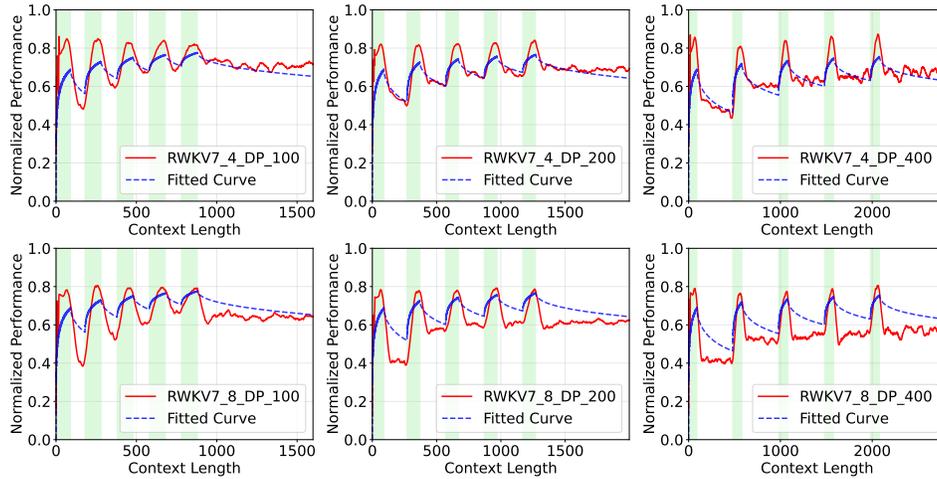
1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287 Figure 11: The performance comparisons between LLaMA3-8B and its fitted ACT-R model.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311



1312 Figure 12: The performance comparisons between MAMBA and its fitted ACT-R model.

1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329



1330 Figure 13: The performance comparisons between RWKV-7 and its fitted ACT-R model.

1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

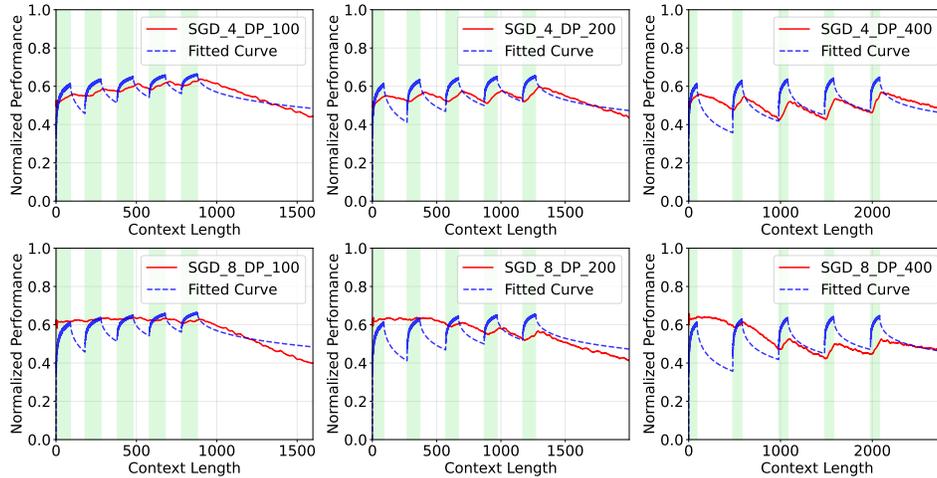


Figure 14: The performance comparisons between SGD and its fitted ACT-R model.

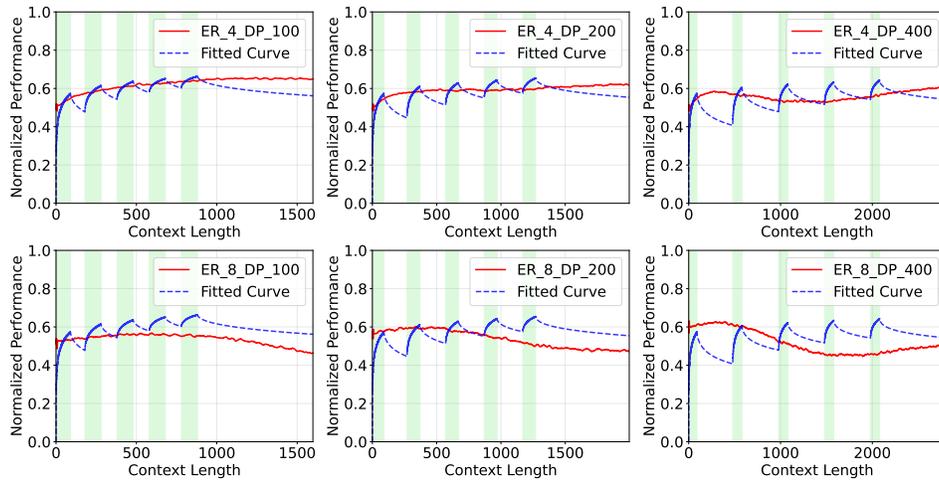


Figure 15: The performance comparisons between ER and its fitted ACT-R model.

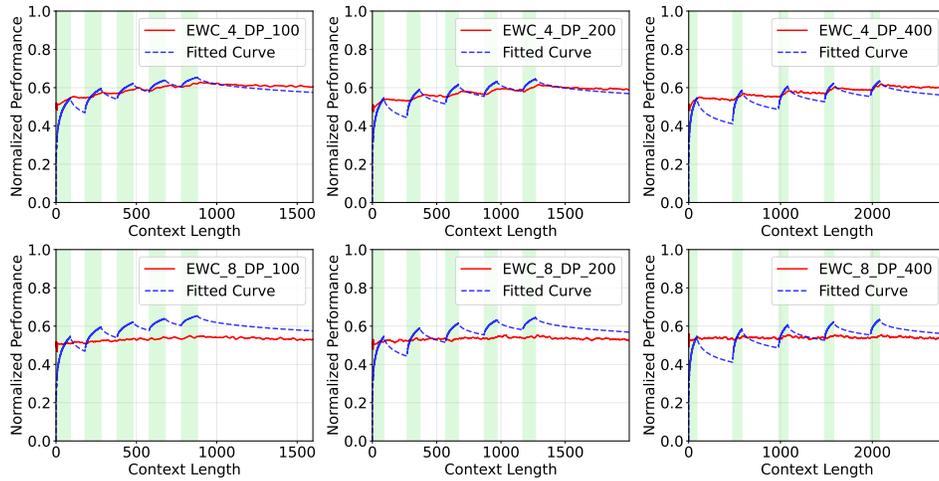


Figure 16: The performance comparisons between EWC and its fitted ACT-R model.