

Gated Differentiable Working Memory for Long-Context Language Modeling

Anonymous ACL submission

Abstract

Long contexts break transformers: attention scores dilute across thousands of tokens, critical information gets lost in the middle, and the model cannot adapt to novel patterns at inference time. Recent work on test-time adaptation addresses this by maintaining a form of *working memory*—transient parameters updated on the current context—but existing approaches employ *uniform* write policies that waste computation on low-value regions and suffer from high gradient variance across semantically heterogeneous contexts. In this work, we reframe test-time adaptation as a budget-constrained memory consolidation problem, asking: *given limited computational budget, which parts of the context should be consolidated into working memory?* We propose GDWM (Gated Differentiable Working Memory), a framework that introduces a Write Controller to gate the memory consolidation process. Our controller estimates *Contextual Utility*—an information-theoretic measure quantifying how much each region depends on long-range context—and allocates gradient steps accordingly, subject to a coverage constraint that ensures global representation. Experiments on ZeroSCROLLS and LongBench v2 benchmarks demonstrate that GDWM achieves comparable or superior performance with 4× fewer gradient steps compared to uniform baselines, establishing a new efficiency-performance Pareto frontier for test-time adaptation.

1 Introduction

Large Language Models (LLMs) struggle with long contexts: attention scores dilute, and models miss critical information in central positions (e.g., “Lost-in-the-Middle”) (Liu et al., 2023a; Hsieh et al., 2024; Du et al., 2025b; Yao et al., 2025). Recent work mitigates such long-context failures by equipping LLMs with *working-memory*—either as lightweight, differentiable *fast-weight* (Ba et al.,

Test-Time Training:
Efficiency vs Performance

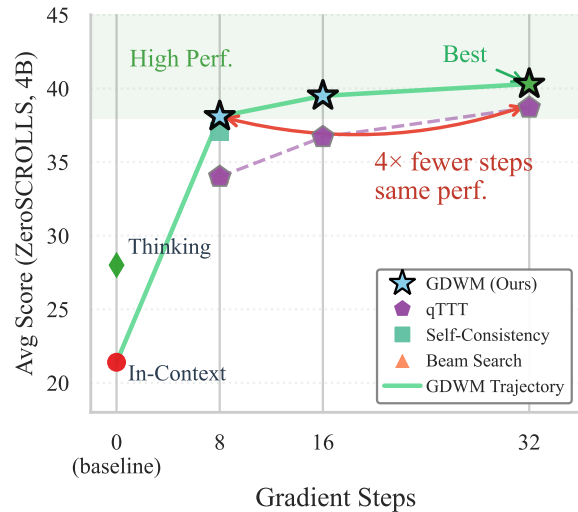


Figure 1: **Efficiency vs Performance on ZeroSCROLLS (QWEN3-4B)**. GDWM achieves comparable performance to qTTT-32 with only 8 gradient steps (4× fewer), establishing a new Pareto frontier. Context-aware budget allocation enables faster convergence than uniform or sampling-based alternatives.

2016) states updated online via self-supervised gradients, or as explicit memory modules that incrementally read, consolidate, and overwrite salient information during inference (Yu et al., 2025; Bansal et al., 2025; Hong et al., 2025; Anonymous, 2025; Xu et al., 2025a; Ye et al., 2025).

However, current approaches predominantly rely on *uniform write policies*, stochastically sampling tokens across the entire context. This is suboptimal: information density is non-uniform, wasting budget on low-value regions, and global sampling exacerbates gradient variance by aggregating conflicting updates from semantically disparate regions.

We propose GDWM (Gated Differentiable Working Memory), a framework that recasts test-time adaptation as budget-constrained memory consolidation. The key question becomes: *given limited compute, which parts of the context should*

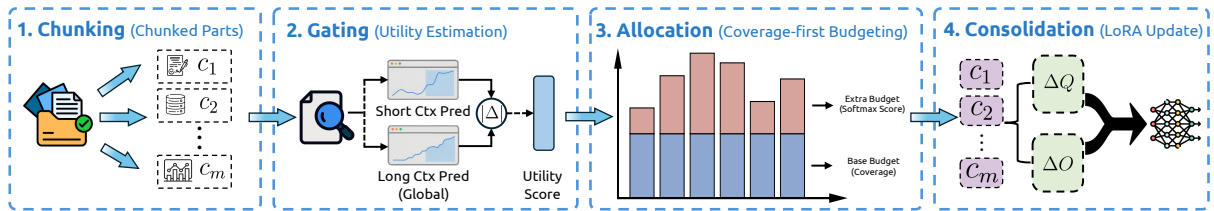


Figure 2: **High-level overview of GDWM.** The framework proceeds in four stages: **Chunk** the input into fixed-size units (approximating semantic segments), **Gate** each chunk via Contextual Utility (CPMI-based divergence), **Allocate** gradient budget proportionally subject to coverage constraints, and **Consolidate** into LoRA adapters.

061 *be consolidated into working memory?* GDWM
 062 answers this via a Write Controller that estimates
 063 *Contextual Utility*—the divergence between long-
 064 context and short-context predictions—and allo-
 065 cates gradient steps to regions where long-range
 066 dependencies are most critical. The framework is
 067 mechanism-agnostic: it provides a principled write
 068 policy that can be layered on top of any test-time
 069 adaptation architecture.

070 First, we formalize test-time adaptation as a
 071 budget-constrained resource allocation problem,
 072 cleanly separating the memory mechanism (how to
 073 update) from the write policy (where and how much
 074 to update). Second, we propose a chunk-wise,
 075 budget-aware algorithm driven by *Contextual Util-*
 076 *ity*—an information-theoretic measure grounded in
 077 Conditional Pointwise Mutual Information (CPMI)
 078 that identifies high-value long-range dependencies.
 079 Third, we prove that chunk-restricted sampling re-
 080 duces gradient variance by eliminating inter-chunk
 081 interference via the Law of Total Variance. Exten-
 082 sive evaluation on ZeroSCROLLS and LongBench
 083 v2 across three model scales (1.7B, 4B, 8B) demon-
 084 strates that GDWM achieves comparable or superior
 085 performance with $4\times$ fewer gradient steps, yielding
 086 significant gains on sparse-information tasks (up to
 087 $+12.7\%$ on Qasper, up to $+11.2\%$ on GovReport)
 088 while achieving 39% wall-clock speedup. Ablation
 089 studies confirm that CPMI-based selection, cover-
 090 age constraints, and chunk-based processing are all
 091 essential components. We further show that
 092 GDWM outperforms alternative test-time scaling
 093 strategies (self-consistency, beam search), and pro-
 094 vide a theoretical analysis linking optimal chunk
 095 size to task-specific evidence spans.

096 2 Related Work

097 **Context Engineering** Context engineering op-
 098 timizes information payloads for LLMs through
 099 prompt design, in-context learning, and chain-of-
 100 thought prompting (Wei et al., 2022; Brown et al.,

2020a). Retrieval-augmented generation (RAG) 101
 enhances knowledge access (Lewis et al., 2020; 102
 Karpukhin et al., 2020), while compression and 103
 hierarchical processing address quadratic scaling 104
 limitations (Li et al., 2023; Vaswani et al., 2017). 105
 Recent work tackles robustness via context-aware 106
 decoding and representation engineering (He et al., 107
 2024a; Zhou et al., 2023). 108

LLM Memory Memory in LLMs addresses bot- 109
 tlenecks from the KV cache and model weights 110
 (Liu et al., 2023b; Dao et al., 2022), with solutions 111
 including paging (Kwon et al., 2023), compression 112
 (Zhang et al., 2024c), and dynamic eviction (Zhang 113
 et al., 2023; Yuan et al., 2025). Agent memory 114
 research extends this by enabling persistent stor- 115
 age beyond context windows through hierarchical 116
 memory architectures (Fang et al., 2024; Zhong 117
 et al., 2023a) and structured memory types (Han 118
 et al., 2024; Terranova et al., 2025). Contempo- 119
 rary frameworks implement consolidation, updat- 120
 ing, and forgetting operations (Park et al., 2023; 121
 Cao et al., 2025), though challenges remain in au- 122
 tonomous management and catastrophic forgetting 123
 (Xu et al., 2025b; Guo et al., 2023). 124

Test-Time Adaptation Recent work treats de- 125
 ployment as an optimization phase, updating model 126
 states at inference time via self-supervised objec- 127
 tives (Niu et al., 2022; Tang et al., 2023). Within 128
 LLMs, gradient-based test-time training adapts to 129
 task instances through neighbor retrieval, stream 130
 processing, or active selection (Hardt and Sun, 131
 2024; Muhtar et al., 2024; Hübotter et al., 2025; 132
 Akyürek et al., 2025). Parallel efforts optimize test- 133
 time compute allocation through principled scaling 134
 and meta-learned control (Snell et al., 2024; Qu 135
 et al., 2025), while complementary strategies com- 136
 press long contexts via selective augmentation and 137
 learned prompt compression (Xu et al., 2023; Jiang 138
 et al., 2024). 139

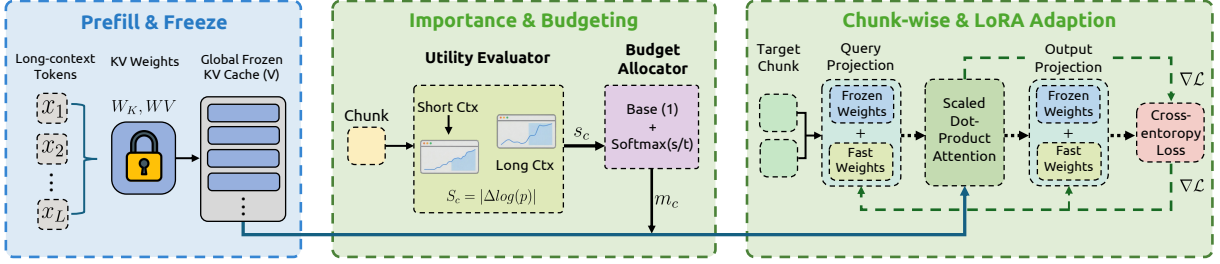


Figure 3: **Technical details of GDWM.** Left: Prefill-and-freeze KV cache enables efficient chunk-wise processing. Middle: Contextual Utility is computed as CPMI between global and local predictions, then converted to budget weights via softmax allocation. Right: LoRA adapters on W_Q/W_O projections are updated through chunk-wise next-token prediction loss.

3 Problem Formulation

We view the process of adapting an LLM to a long context $X = (x_1, x_2, \dots, x_L)$ as a Budget-Constrained Memory Consolidation problem.

Fast Weights. Let θ denote the parameters of a pre-trained LLM. We partition θ into frozen parameters θ_f (including the KV cache computation) and adaptable “fast weights” ϕ (e.g., LoRA adapters (Hu et al., 2021) on Query projections).

Budget-Constrained Optimization. We partition the input context X into M fixed-size chunks $\{C_1, C_2, \dots, C_M\}$. The goal is to determine an optimal integer allocation schedule $\mathbf{k} = (k_1, \dots, k_M)$, where k_c denotes the number of minibatch gradient updates performed on chunk C_c .

Since the downstream task loss $\mathcal{L}_{\text{task}}$ is unknown during inference, we employ the self-supervised next-token prediction loss as a surrogate objective. The optimization problem is formulated as:

$$\min_{\{k_c\}_{c=1}^M} \sum_{c=1}^M \mathbb{E}_{t \sim \mathcal{U}(\mathcal{I}_c)} [-\log P(x_t | x_{1:t-1}; \phi(\mathbf{k}), \theta_f)] \quad (1)$$

where \mathcal{I}_c denotes the set of positions in chunk C_c , and $\mathbb{E}_{t \sim \mathcal{U}(\mathcal{I}_c)}$ denotes the expectation over positions t sampled uniformly from \mathcal{I}_c . The constraints are:

$$\sum_{c=1}^M k_c \leq K_{\text{total}} \quad (\text{Total Budget}) \quad (2)$$

$$k_c \geq k_{\min}, \forall c \in \{1, \dots, M\} \quad (\text{Coverage}) \quad (3)$$

Here, $\phi(\mathbf{k})$ is an *implicit function* of the allocation \mathbf{k} , defined by the gradient descent process: starting from $\phi^{(0)} = 0$, we sequentially perform

k_c gradient updates on each chunk C_c , yielding $\phi(\mathbf{k}) = \phi(\sum_c k_c)$. Since exactly solving this bi-level problem is intractable, we propose a heuristic approximation in Section 4 using Contextual Utility as a proxy for the gradient contribution of each chunk.

4 Method

Our method, GDWM, solves the optimization problem defined in Equation 1 via a four-stage process: (1) Chunking, (2) Gating (Importance Estimation), (3) Allocation, and (4) Consolidation. Figure 2 illustrates the high-level pipeline, and Figure 3 details the technical components.

4.1 Memory Interface

GDWM is mechanism-agnostic: the Write Controller operates as a policy layer orthogonal to the memory mechanism. While ϕ can be any differentiable parameters (linear, MLP, etc.), we focus on LoRA adapters for efficiency.

While GDWM is mechanism-agnostic and can be layered on top of any adaptation architecture, we instantiate it efficiently by attaching LoRA adapters to projection matrices while freezing the KV cache (see Section 5.1 for details). The fast weights ϕ are initialized to zero (identity mapping) at the start of each sequence.

4.2 Contextual Utility Estimation

To efficiently solve the allocation problem, we need a proxy for the gradient contribution of each chunk. We introduce *Contextual Utility*, an information-theoretic measure grounded in the cognitive notion of surprisal. The core intuition is simple: if a token x_t can be predicted easily using only local context, consolidating it into fast weights yields low utility. Conversely, if x_t is unpredictable locally

but predictable globally, it represents a high-value long-range dependency.

Definition 1 (Contextual Utility). Let \mathcal{I}_c denote the set of token positions in chunk C_c . The Contextual Utility is defined as the average surprisal divergence:

$$U(C_c) = \frac{1}{|\mathcal{I}_c|} \sum_{t \in \mathcal{I}_c} \Delta_t$$

$$\Delta_t = \left| \log P(x_t \mid x_{1:t-1}) - \log P(x_t \mid x_{t-n:t-1}) \right| \quad (4)$$

where $P(x_t \mid x_{1:t-1})$ is the conditional probability using full context (“global prediction”) and $P(x_t \mid x_{t-n:t-1})$ uses only the recent n tokens (“local prediction”).

We take the absolute value to capture both directions: positions where long context *helps* ($P_{\text{full}} > P_{\text{local}}$) and where it *conflicts* ($P_{\text{full}} < P_{\text{local}}$). Both indicate regions requiring gradient-based calibration. Empirically, $|\Delta_t|$ outperforms $\max(0, \Delta_t)$ by 3-5% (see Appendix F for analysis).

Interpretation. The term inside the sum represents the *Surprisal Divergence*. When $P_{\text{full}} \approx P_{\text{local}}$, the long context provides no additional information (utility is zero). High utility indicates regions where global dependencies are critical for prediction.

Information-Theoretic Grounding. The quantity Δ_t has a precise information-theoretic interpretation: it equals the absolute value of *Conditional Pointwise Mutual Information* (CPMI) between the token x_t and the long-range prefix $x_{1:t-n}$, conditioned on the local context $x_{t-n:t-1}$. Formally, $\text{CPMI}(A; B \mid C) = \log P(A \mid B, C) - \log P(A \mid C)$ measures the information that B provides about A beyond what C already provides. High $|\Delta_t|$ indicates positions where the model’s prediction strongly depends on (or conflicts with) long-range context—precisely where gradient-based memory consolidation is most valuable.

4.3 Budget-Aware Allocation Policy

Exact solution of the integer programming problem in Eq. (1) is intractable. We propose a heuristic approximation based on **Coverage-First Softmax Allocation**.

Step 1: Satisfy Coverage Constraint. Allocate k_{\min} steps to every chunk:

$$k_c \leftarrow k_{\min}, \quad \forall c \in \{1, \dots, M\} \quad (5)$$

Handling Infeasible Budgets. If the total budget $K_{\text{total}} < M \cdot k_{\min}$, the coverage constraint cannot be strictly satisfied. In such cases, we relax the constraint by allocating k_{\min} steps to the top- $\lfloor K_{\text{total}}/k_{\min} \rfloor$ chunks with the highest utility $U(C_c)$, and 0 to others. This fallback ensures that limited compute is invested in the most critical regions first.

Step 2: Distribute Remaining Budget. The remaining budget $K_{\text{rem}} = \max(0, K_{\text{total}} - M \cdot k_{\min})$ is distributed based on normalized utility:

$$w_c = \frac{\exp(U(C_c)/\tau)}{\sum_{j=1}^M \exp(U(C_j)/\tau)} \quad (6)$$

$$k_c \leftarrow k_c + \lfloor K_{\text{rem}} \cdot w_c \rfloor \quad (7)$$

To ensure the full budget is utilized, we apply the *Largest Remainder Method* (see Appendix I) to distribute residual steps, guaranteeing $\sum_c k_c = K_{\text{total}}$ exactly. The temperature τ controls allocation sharpness. When $\tau \rightarrow 0$, all extra budget concentrates on the highest-utility chunk; when $\tau \rightarrow \infty$, allocation becomes uniform (ignoring utility); and $\tau = 1.0$ provides balanced allocation that proves empirically optimal. This policy directs the optimizer’s focus to high-utility regions while maintaining global awareness through the coverage constraint—hence the term “gated” in our framework name.

4.4 Structured Memory Consolidation

Finally, we execute updates via chunk-restricted sampling. For each chunk C_c , we perform k_c gradient descent steps. In each step, we uniformly sample a minibatch of positions $\mathcal{I} \subset \mathcal{I}_c$ and compute the loss:

$$\mathcal{L} = -\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \log P(x_i \mid x_{1:i-1}; \phi, \theta_f) \quad (8)$$

The KV cache can be precomputed and reused across all k_c steps for efficiency (see Section 5.1).

4.5 Why Does This Work?

The key insight is semantic heterogeneity: distinct document sections induce gradient directions that interfere destructively when aggregated. Chunk-restricted sampling (using fixed-size windows as semantic proxies) eliminates this interference by confining each update to a coherent unit. Formally, by the Law of Total Variance, this eliminates the interchunk variance component from gradient estimates (Theorem 1). The coverage constraint prevents

METHOD	STEPS	SUMMARIZATION		QUESTION ANSWERING		COMPREHENSION	REASONING	Avg.
		GovReport	QMSum	Qasper	NarrativeQA	Quality	Musique	
QWEN3-1.7B								
In-context	–	22.2	6.4	25.8	14.9	48.1	11.8	21.5
Thinking	–	21.5	7.6	22.3	9.2	61.8	22.4	24.1
qTTT	8	23.9	8.3	27.3	9.4	72.1	19.9	26.8
qTTT	16	25.1	9.3	29.5	11.2	74.1	23.1	28.7
qTTT	32	26.8	9.7	31.5	12.4	76.5	26.4	30.6
GDWM (Ours)	8	28.5	9.6	33.8	12.2	76.5	27.0	31.3
GDWM (Ours)	16	29.1	9.9	34.6	12.6	77.0	27.9	31.8
GDWM (Ours)	32	29.8	10.2	35.5	13.0	77.5	28.8	32.5
Δ vs Best		+11.2%	+5.2%	+12.7%	+4.8%	+1.3%	+9.1%	+6.2%
QWEN3-4B								
In-context	–	24.8	11.2	23.5	10.9	40.8	17.0	21.4
Thinking	–	20.8	7.5	25.2	29.8	76.1	8.3	28.0
qTTT	8	29.2	8.3	29.5	32.3	81.3	23.7	34.0
qTTT	16	31.9	8.6	32.1	35.3	84.8	27.5	36.7
qTTT	32	33.2	8.9	33.8	38.4	87.2	30.8	38.7
GDWM (Ours)	8	34.2	8.4	35.2	37.5	82.2	31.2	38.1
GDWM (Ours)	16	35.0	8.8	35.8	38.6	86.8	32.0	39.5
GDWM (Ours)	32	35.8	9.2	36.5	39.8	87.5	32.8	40.3
Δ vs Best		+7.8%	+3.4%	+8.0%	+3.6%	+0.3%	+6.5%	+4.1%
QWEN3-8B								
In-context	–	22.5	8.8	20.1	35.4	90.5	22.9	33.4
Thinking	–	18.2	9.8	21.5	19.6	71.8	43.8	30.8
qTTT	8	25.3	8.5	22.8	38.5	91.2	42.0	38.1
qTTT	16	27.9	8.7	24.5	40.6	93.1	46.2	40.2
qTTT	32	29.8	9.0	27.0	42.4	94.9	49.6	42.2
GDWM (Ours)	8	30.2	8.2	27.5	42.0	93.8	49.5	41.9
GDWM (Ours)	16	30.8	8.6	28.1	43.1	94.3	50.2	42.5
GDWM (Ours)	32	31.5	9.0	28.8	44.2	94.8	51.0	43.2
Δ vs Best		+5.7%	0.0%	+6.7%	+4.2%	-0.1%	+2.8%	+2.4%

Table 1: Main results on ZeroSCROLLS. GDWM-32 achieves consistent improvements over qTTT-32 across all tasks (+2.4–6.2% average). GDWM-8 demonstrates Pareto-optimal efficiency: $4\times$ fewer gradient steps with comparable performance to qTTT-32.

mode collapse, ensuring global representation. See Appendix C for full analysis and Appendix A for the complete algorithm.

5 Experiments

5.1 Experimental Setup

Datasets. We evaluate on ZeroSCROLLS (Shaham et al., 2023) (6 tasks: GovReport, QMSum, Qasper, NarrativeQA, Quality, MuSiQue) and LongBench v2 (Bai et al., 2025), covering summarization, QA, multi-hop reasoning, and code understanding across sparse-to-dense information distributions.

Baselines. We compare against: (1) **In-context** (Brown et al., 2020a)—standard inference; (2) **Thinking** (Yang et al., 2025)—inference

with chain-of-thought; (3) **qTTT** (Bansal et al., 2025)—with uniform sampling (8-32 steps).

Implementation. We use Qwen3 models (1.7B/4B/8B) as base LLMs with LoRA adapters (rank=16, $\alpha=32$) applied to query and output projection matrices (W_Q, W_O). For test-time adaptation, we employ AdamW optimizer with learning rate 1×10^{-4} and no weight decay. Default hyperparameters: chunk size $S = 1024$ tokens, temperature $\tau = 1.0$, minimum coverage $k_{\min} = 1$, total gradient steps $K_{\text{total}} = 8$ (for GDWM-8) or $K_{\text{total}} = 32$ (for GDWM-32), and maximum context length 32K tokens. CPMI estimation uses a sliding window of 512 tokens for local context. All experiments run on NVIDIA H20 GPUs with mixed-precision (bfloat16) training. For efficiency, we implement gradient

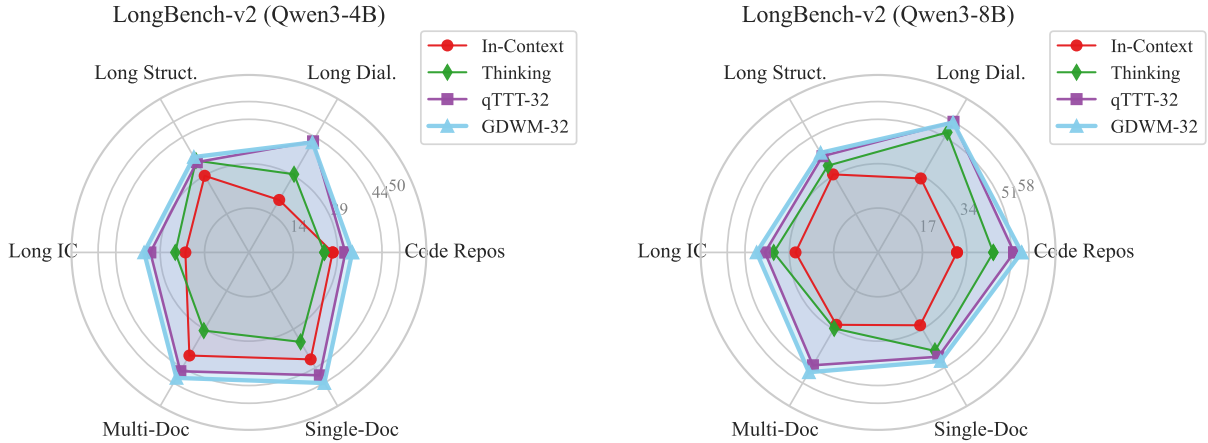


Figure 4: **Task-wise Performance on LongBench v2.** Radar charts comparing GDWM-32 (light blue) against baselines on 4B (left) and 8B (right) models. GDWM achieves consistent improvements on Code Repositories and Multi-Doc QA where information is sparse and localized, while showing competitive performance on Long Dialogue where global coverage is required.

checkpointing and flash attention v2 to reduce memory overhead during test-time adaptation.

5.2 Main Results

Table 1 presents our main results. GDWM-32 achieves consistent improvements over qTTT-32 across all tasks (+2.4–6.2% average). The largest gains appear on sparse-information tasks: GovReport (+11.2% on 1.7B, +7.8% on 4B) and Qasper (+12.7% on 1.7B, +6.7% on 8B), where relevant content is concentrated in specific document sections—exactly the scenario where CPMI-based selection excels. On Quality, which requires holistic comprehension with uniformly distributed information, improvements are more modest (−0.1%–+1.3%), revealing an expected characteristic of selective allocation.

As shown in Figure 1, GDWM-8 demonstrates Pareto-optimal efficiency: with 4× fewer gradient steps (8 vs 32), it achieves performance comparable to qTTT-32 (within 1 point margin on 4B, within 0.5 points on 8B average), while GDWM-32 pushes the frontier further with consistent 2.4–6.2% gains across model scales. This validates our core hypothesis that intelligent context selection outperforms brute-force computation.

To further validate GDWM’s robustness across diverse task types, we evaluate on LongBench v2 (Bai et al., 2025) with results visualized in Figure 4. GDWM-32 achieves competitive or superior performance across model sizes and task categories. On Code Repositories (+5.5% on 8B) and Multi-Document QA (+6.2% on 8B)—tasks requiring precise retrieval from sparse, localized information—

CONFIGURATION	GOVRPT	QASPER	MUSIQUE	Δ
GDWM (full)	28.5	33.8	27.0	–
w/o CPMI (uniform)	23.9	27.3	19.9	−20.4%
w/o coverage	27.8	31.1	15.0	−17.3%
w/o chunking	24.7	27.8	10.0	−30.0%

Table 2: Ablation study on Qwen3-1.7B. Chunking is most critical (−30.0%); CPMI selection and coverage constraint both essential.

GDWM achieves consistent gains. However, on Long Dialogue, where information is distributed more uniformly, performance shows a slight trade-off (−0.5% on 8B), revealing an expected characteristic of selective allocation strategies.

5.3 Ablation Studies

We conduct ablation experiments on Qwen3-1.7B to validate each component of GDWM, with results shown in Table 2.

Replacing CPMI-based selection with uniform sampling results in a 20.4% performance drop, demonstrating that intelligent context selection is crucial for effective test-time adaptation. The coverage constraint proves essential for multi-hop reasoning: removing it causes MuSiQue performance to plummet from 27.0 to 15.0, as the model overfits to a single high-utility region rather than gathering evidence from multiple document sections. Most critically, without chunking the model updates at token level, fragmenting context and incurring the largest drop (−30.0%), supporting our variance-reduction analysis.

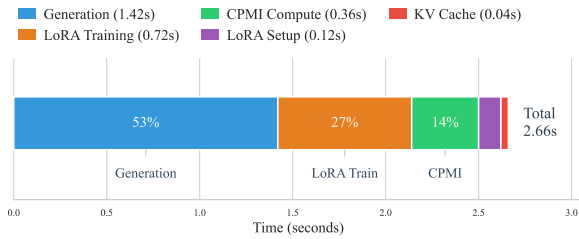


Figure 5: **Time Breakdown per Sample (Qwen3-4B)**. Generation dominates at 53%, while CPMI computation accounts for only 13%—demonstrating the lightweight nature of our context selection mechanism. The $4\times$ reduction in gradient steps ($32\rightarrow 8$) yields 39% net wall-clock speedup.

5.4 Scaling at Test Time

Table 3 compares GDWM against other test-time scaling approaches on Qwen3-4B under equal computational budget. Self-Consistency (SC-8) performs exceptionally well on multiple-choice tasks, marginally outperforming GDWM on Quality (82.7 vs 82.2) by aggregating diverse reasoning paths. However, it struggles on open-ended QA tasks (Qasper, MuSiQue) where the challenge lies in locating relevant context rather than generating diverse outputs. Beam Search proves largely ineffective for long-context understanding. GDWM maintains a 7.3-point lead on average (45.7 vs 38.4) over SC-8, validating that while ensemble methods help verification, intelligent context selection is more fundamental for evidence retrieval.

5.5 Efficiency Analysis

The primary overhead of GDWM stems from CPMI computation, which requires two forward passes per chunk (global and local predictions). As shown in Figure 5, with our default configuration ($S = 1024$, $K = 8$), the CPMI computation adds only 0.36s ($\sim 13\%$ of total time), while the $4\times$ reduction in gradient steps (from 32 to 8) saves 1.79s. The net result is a **39% wall-clock time reduction** compared to the 32-step baseline, demonstrating that intelligent context selection is not only more effective but also more efficient than brute-force uniform sampling.

Chunk Size Selection. Table 4 analyzes the trade-off between chunk granularity, computational overhead, and task performance. We adopt $S = 1024$ as the default configuration, which achieves the best *cross-task robustness*: avoiding catastrophic failures while maintaining competitive performance across all task types.

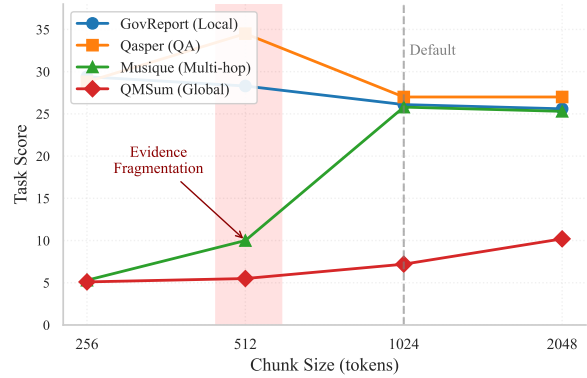


Figure 6: **Chunk Size Sensitivity**. Different task types exhibit distinct optimal chunk sizes. Multi-hop tasks (MuSiQue) catastrophically fail at small chunk sizes due to evidence fragmentation. $S = 1024$ provides robust cross-task performance.

Theoretical Analysis. The observed chunk size sensitivity has a principled explanation rooted in task-specific *evidence span*—the typical token distance over which task-relevant information is distributed. Let $E_{\mathcal{T}}$ denote this characteristic span for task \mathcal{T} . We identify a critical constraint: when chunk size $S < E_{\mathcal{T}}$, relevant evidence becomes fragmented across multiple chunks, causing CPMI to underestimate the true contextual utility of each fragment.

This framework explains the empirical patterns in Table 4. For **local reasoning tasks** such as GovReport, evidence is highly concentrated, so smaller chunks ($S = 256$) maximize selection precision. However, for **extractive QA tasks** like Qasper, performance peaks at intermediate granularity ($S = 512$), suggesting that while locality is important, overly small chunks ($S = 256$) may begin to fragment answer-relevant paragraphs. In contrast, **multi-hop reasoning tasks** like MuSiQue require evidence chains that span multiple paragraphs ($E_{\mathcal{T}} \approx 1000+$ tokens); when $S = 512 < E_{\mathcal{T}}$, the reasoning chain is severed—each fragment appears low-utility in isolation, leading to catastrophic failure (10.0 vs 25.8). Meanwhile, **summarization tasks** such as QMSum require near-uniform coverage ($E_{\mathcal{T}} \approx$ full document), where larger chunks ($S = 2048$) better preserve global structure.

The choice of $S = 1024$ represents the *minimum chunk size that avoids evidence fragmentation* for multi-hop tasks while retaining sufficient granularity for local reasoning—a principled middle ground validated by our cross-task robustness results. See Appendix J for formal analysis.

METHOD	BUDGET	GOVREPORT	QASPER	QUALITY	MUSIQUE	Avg.
Thinking	1×	20.1	24.5	76.2	7.5	32.1
Self-Consistency (SC-8)	8×	24.5	28.4	82.7	18.1	38.4
Beam Search ($k=8$)	8×	23.0	26.2	77.8	14.2	35.3
qTTT (8 steps)	8×	29.2	29.5	81.3	23.7	40.9
GDWM (Ours)	8×	34.2	35.2	82.2	31.2	45.7

Table 3: Test-time scaling comparison on Qwen3-4B under equal compute budget (8× base). Intelligent context selection (GDWM) outperforms brute-force sampling (SC-8) and search (Beam) by significant margins, particularly on sparse-information tasks.

CHUNK	CPMI%	OVERHEAD	AVG	WORST
512	25%	+12%	19.4	10.0 [†]
1024	13%	+6%	21.7	25.8
2048	9%	−10%	22.3	25.3

Table 4: Chunk size trade-offs on representative tasks (GovReport, Qasper, Musique, QMSum). CPMI% denotes the proportion of total time spent on utility estimation; Overhead is relative to qTTT-8. $S = 1024$ achieves the best *cross-task robustness*: the only configuration without catastrophic failure (Worst ≥ 25). [†]Small chunks fragment multi-hop evidence chains, causing CPMI to underestimate chunk utility (see Appendix J).

5.6 Discussion

Our experiments reveal a clear pattern regarding when GDWM excels. The largest gains appear on tasks with sparse information distribution, where relevant content is concentrated in specific document regions. GovReport (+5.7% on 8B) and Qasper (+6.7% on 8B) exemplify this phenomenon, as government reports contain key findings in specific sections while scientific papers concentrate answers near figures or method descriptions. On tasks requiring dense global coverage—QMSum and NarrativeQA—improvements are more modest but still positive, demonstrating that the coverage constraint effectively maintains global representation. Quality shows minimal improvement (0.0% on 8B), consistent with its requirement for uniform attention across the entire document.

Scaling Behavior. The scaling behavior shows consistent patterns: the efficiency advantage holds across model sizes (+6.2% for 1.7B, +4.1% for 4B, +2.4% for 8B average improvement for GDWM-32 vs qTTT-32). Interestingly, smaller models benefit more from intelligent context selection, likely because they have limited capacity to attend to all context uniformly and thus benefit more from prioritized consolidation. Gains on sparse tasks remain substantial across scales (Qasper: +12.7% on 1.7B,

+8.0% on 4B, +6.7% on 8B; GovReport: +11.2% on 1.7B, +7.8% on 4B, +5.7% on 8B).

Chunk Size Trade-offs. Figure 6 visualizes the task-specific chunk size preferences. The dramatic performance collapse of MuSiQue at $S = 512$ (dropping to 10.0 from 25.8 at $S = 1024$) directly validates our Evidence Span theory: when chunk size falls below the task’s characteristic evidence span $E_{\mathcal{T}}$, the reasoning chain is severed and CPMI underestimates true utility.

Three limitations merit future investigation: (i) fixed-size chunking can split coherent regions in irregularly structured documents, potentially harming selection reliability; (ii) the optimal chunk size is task-dependent (Table 4), motivating adaptive or structure-aware chunking; and (iii) for dense-coverage tasks (e.g., QMSum, NarrativeQA), selective allocation may offer limited gains, reflecting an inherent efficiency–coverage trade-off.

6 Conclusion

We presented GDWM (Gated Differentiable Working Memory), a framework that recasts test-time adaptation as a budget-constrained memory consolidation problem. By introducing a Write Controller that estimates Contextual Utility—an information-theoretic measure of long-range dependency—and allocates gradient budget via a coverage-first strategy, GDWM reduces gradient steps by 4× compared to uniform baselines while achieving comparable or superior performance on sparse-information tasks, with 39% wall-clock speedup. Our theoretical analysis proves that chunk-restricted sampling reduces gradient variance by eliminating inter-chunk variance, providing a principled explanation for the improved convergence. Future directions include semantic-aware dynamic chunking, task-adaptive temperature scheduling, and extending GDWM to multimodal contexts.

519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567

Limitations

While GDWM establishes a new Pareto frontier for efficient test-time adaptation, we identify three limitations inherent to our budget-constrained design. First, our current implementation employs fixed-size chunking ($S = 1024$). Although our analysis shows this setting is robust across diverse tasks, it may sub-optimally fragment semantic units in documents with irregular structures. However, we view this as an efficiency trade-off: dynamic chunking would require additional forward passes for segmentation, potentially offsetting the speed gains. Second, as observed in our results on Long Dialogue and Quality, the selective allocation strategy is less effective for tasks requiring uniform information coverage. This is a structural property of any compressive memory system rather than a flaw of GDWM specifically; users must weigh the $4\times$ efficiency gain against the marginal performance trade-off on dense-coverage tasks. Third, GDWM introduces hyperparameters (e.g., temperature τ , min coverage k_{\min}). While our experiments demonstrate that the default configuration ($\tau = 1.0, k_{\min} = 1$) generalizes well without tuning, extreme domain shifts might require recalibration of the utility estimator.

Ethical Considerations

Our work improves the efficiency of inference-time adaptation for Large Language Models. By reducing the gradient steps required for effective context adaptation by $4\times$, GDWM can lower the computational cost of long-context deployment and thereby supports the broader goal of Green AI.

GDWM is a general-purpose optimization layer and does not introduce new task capabilities or new access to information beyond what the underlying base model and the provided context already permit. As with any efficiency improvement, it may enable wider or more frequent use of long-context systems; the associated downstream risks (e.g., misuse of generative models) are therefore not specific to our mechanism, but to the deployment setting and the base model’s safety properties. In practice, these concerns are best addressed through standard governance and safeguards at the system level—such as strong base-model alignment, access control, privacy-preserving data handling, and application-layer monitoring—rather than by restricting inference-time optimization techniques.

References

Ekin Akyürek, Mehul Damani, Adam Zweiger, Linlu Qiu, Han Guo, Jyothish Pari, Yoon Kim, and Jacob Andreas. 2025. [The surprising effectiveness of test-time training for few-shot learning](#). *Preprint*, arXiv:2411.07279. 569
570
571
572
573

J. Allingham, Jie Ren, Michael W. Dusenberry, J. Liu, Xiuye Gu, Yin Cui, Dustin Tran, and Balaji Lakshminarayanan. 2023. A simple zero-shot prompt weighting technique to improve prompt ensembling in text-image models. In *International Conference on Machine Learning*. 574
575
576
577
578
579

Petr Anokhin, Nikita Semenov, Artyom Y. Sorokin, Dmitry Evseev, M. Burtsev, and Evgeny Burnaev. 2024. Arigraph: Learning knowledge graph world models with episodic memory for llm agents. In *International Joint Conference on Artificial Intelligence*. 580
581
582
583
584
585

Anonymous. 2025. [In-place test-time training](#). In *ICLR 2026 Conference Submission*. Under review at ICLR 2026; Available on OpenReview. 586
587
588

Jimmy Ba, Geoffrey Hinton, Volodymyr Mnih, Joel Z. Leibo, and Catalin Ionescu. 2016. [Using fast weights to attend to the recent past](#). *Preprint*, arXiv:1610.06258. 589
590
591
592

Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2025. [Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks](#). *Preprint*, arXiv:2412.15204. 593
594
595
596
597
598

S. Banasik. 2025. Memory access characterization of large language models in cpu environment and its potential impacts. *arXiv.org*. 599
600
601

Rachit Bansal, Aston Zhang, Rishabh Tiwari, Lovish Madaan, Sai Surya Duvvuri, Devvrit Khatri, David Brandfonbrener, David Alvarez-Melis, Prajjwal Bhargava, Mihir Sanjay Kale, and Samy Jelassi. 2025. [Let’s \(not\) just put things in context: Test-time training for long-context llms](#). *Preprint*, arXiv:2512.13898. 602
603
604
605
606
607
608

Saikat Barua. 2024. Exploring autonomous agents through the lens of large language models: A review. *arXiv.org*. 609
610
611

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, J. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. Henighan, R. Child, A. Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, and 12 others. 2020a. Language models are few-shot learners. *Neural Information Processing Systems*. 612
613
614
615
616
617
618
619

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, J. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda 620
621
622

568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622

623	Askill, Sandhini Agarwal, Ariel Herbert-Voss,	Tri Dao, Daniel Y. Fu, Stefano Ermon, A. Rudra, and	677
624	Gretchen Krueger, T. Henighan, R. Child, A. Ramesh,	Christopher R'e. 2022. Flashattention: Fast and	678
625	Daniel M. Ziegler, Jeff Wu, Clemens Winter, and 12	memory-efficient exact attention with io-awareness.	679
626	others. 2020b. Language models are few-shot learn-	<i>Neural Information Processing Systems</i> .	680
627	ers. <i>Neural Information Processing Systems</i> .		
628	Zhihang Cai, Xingjun Zhang, Zhendong Tan, and	Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei	681
629	Zheng Wei. 2025. Nqkv: A kv cache quantization	Wang, Sébastien Montella, Mirella Lapata, Kam-Fai	682
630	scheme based on normal distribution characteristics.	Wong, and Jeff Z. Pan. 2025a. Rethinking memory	683
631	<i>arXiv.org</i> .	in llm based agents: Representations, operations, and	684
		emerging topics.	685
632	Zouying Cao, Jiaji Deng, Li Yu, Weikang Zhou,	Yufeng Du, Minyang Tian, Srikanth Ronanki, Subendhu	686
633	Zhaoyang Liu, Bolin Ding, and Hai Zhao. 2025. Re-	Rongali, Sravan Bodapati, Aram Galstyan, Azton	687
634	member me, refine me: A dynamic procedural mem-	Wells, Roy Schwartz, Eliu A Huerta, and Hao	688
635	ory framework for experience-driven agent evolution.	Peng. 2025b. Context length alone hurts llm	689
		performance despite perfect retrieval . <i>Preprint</i> ,	690
		arXiv:2510.05381.	691
636	Debashish Chakraborty, Eugene Yang, Daniel Khashabi,	Yuchen Duan, Zhe Chen, Yusong Hu, Weiyun Wang,	692
637	Dawn J. Lawrie, and Kevin Duh. 2025. Principled	Shenglong Ye, Botian Shi, Lewei Lu, Qibin Hou,	693
638	context engineering for rag: Statistical guarantees via	Tong Lu, Hongsheng Li, Jifeng Dai, and Wenhai	694
639	conformal prediction.	Wang. 2025. Docopilot: Improving multimodal mod-	695
		els for document-level understanding. <i>Computer Vi-</i>	696
640	Vivek Chari, Guanghui Qin, and Benjamin Van Durme.	<i>sion and Pattern Recognition</i> .	697
641	2025. Kv-distill: Nearly lossless learnable context		
642	compression for llms . <i>Preprint</i> , arXiv:2503.10337.		
		Junjie Fang, Likai Tang, Hongzhe Bi, Yujia Qin, Si Sun,	698
643	Shaoshen Chen, Yangning Li, Zishan Xu, Yinghui Li,	Zhenyu Li, Haolun Li, Yongjian Li, Xin Cong, Yukun	699
644	Xin Su, Zifei Shan, and Hai tao Zheng. 2025a. Dast:	Yan, Xiaodong Shi, Sen Song, Yankai Lin, Zhiyuan	700
645	Context-aware compression in llms via dynamic allo-	Liu, and Maosong Sun. 2024. Unimem: Towards a	701
646	cation of soft tokens . <i>Preprint</i> , arXiv:2502.11493.	unified view of long-context large language models.	702
		<i>arXiv.org</i> .	703
647	Yuxin Chen, Peng Tang, Weidong Qiu, and Shujun Li.	Runnan Fang, Yuan Liang, Xiaobin Wang, Jialong Wu,	704
648	2025b. Using llms for automated privacy policy anal-	Shuofei Qiao, Pengjun Xie, Fei Huang, Huajun Chen,	705
649	ysis: Prompt engineering, fine-tuning and explain-	and Ningyu Zhang. 2025. Memp: Exploring agent	706
650	ability. <i>arXiv.org</i> .	procedural memory. <i>arXiv.org</i> .	707
651	Feng Cheng, Cong Guo, Chiyue Wei, Junyao Zhang,	Weizhi Fei, Xueyan Niu, Guoqing Xie, Yingqing Liu,	708
652	Changchun Zhou, Edward Hanson, Jiaqi Zhang, Xi-	Bo Bai, and Wei Han. 2025. Efficient prompt com-	709
653	aoxiao Liu, H. Li, and Yiran Chen. 2025. Ecco: Im-	pression with evaluator heads for long-context trans-	710
654	proving memory bandwidth and capacity for llms via	former inference . <i>Preprint</i> , arXiv:2501.12959.	711
655	entropy-aware cache compression. In <i>International</i>		
656	<i>Symposium on Computer Architecture</i> .	Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue,	712
		Hanna Hajishirzi, Yoon Kim, and Hao Peng. 2024.	713
657	Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-	Data engineering for scaling language models to 128k	714
658	Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan	context. In <i>International Conference on Machine</i>	715
659	Zhao. 2024. xrag: Extreme context compression	<i>Learning</i> .	716
660	for retrieval-augmented generation with one token.		
661	<i>Neural Information Processing Systems</i> .	Pin Gao, Lingfan Yu, Yongwei Wu, and Jinyang Li.	717
		2018. Low latency rnn inference with cellular batch-	718
662	Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and	ing. In <i>European Conference on Computer Systems</i> .	719
663	Danqi Chen. 2023. Adapting language models to		
664	compress contexts . In <i>Proceedings of the 2023 Con-</i>	Sanjay Govindan, M. Pagnucco, and Yang Song. 2025.	720
665	<i>ference on Empirical Methods in Natural Language</i>	Temporal alignment of time sensitive facts with ac-	721
666	<i>Processing</i> , pages 3829–3846, Singapore. Associa-	tivation engineering. In <i>Conference on Empirical</i>	722
667	tion for Computational Linguistics.	<i>Methods in Natural Language Processing</i> .	723
668	Krishna Teja Chitty-Venkata, Jie Ye, Xian-He Sun, An-	Jing Guo, Nan Li, J. Qi, Hang Yang, Ruiqiao Li, Yuzhen	724
669	thony Kougkas, M. Emani, Venkat Vishwanath, and	Feng, Si Zhang, and Ming Xu. 2023. Empowering	725
670	Bogdan Nicolae. 2025. Pagedeviction: Structured	working memory for large language model agents.	726
671	block-wise kv cache pruning for efficient large lan-	<i>arXiv.org</i> .	727
672	guage model inference. <i>arXiv.org</i> .		
		Md. Asif Haider, Ayesha Binte Mostofa, Sk Md Mosad-	728
673	Jincheng Dai, Zhuowei Huang, Haiyun Jiang, Chen	dek Hossain, Anindya Iqbal, and Toufique Ahmed.	729
674	Chen, Deng Cai, Wei Bi, and Shuming Shi. 2024.	2024. Prompting and fine-tuning large language mod-	730
675	Corm: Cache optimization with recent message for	els for automated code review comment generation.	731
676	large language model inference.	<i>arXiv.org</i> .	732

733	Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, Zhaozhuo Xu, and Chaoyang He. 2024. Llm multi-agent systems: Challenges and open problems. <i>arXiv.org</i> .	Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. 2024a. Hiagent: Hierarchical working memory management for solving long-horizon agent tasks with large language model. <i>Annual Meeting of the Association for Computational Linguistics</i> .	787
734			788
735			789
736			790
737	Moritz Hardt and Yu Sun. 2024. Test-time training on nearest neighbors for large language models . <i>Preprint</i> , arXiv:2305.18466.	Sihao Hu, Tiansheng Huang, Fatih Ilhan, S. Tekin, Gaowen Liu, R. Kompella, and Ling Liu. 2024b. A survey on large language model-based game agents. <i>arXiv.org</i> .	791
738			792
739			793
740	Jerry Zhi-Yang He, Sashrika Pandey, Mariah L. Schrum, and Anca Dragan. 2024a. Context steering: Controllable personalization at inference time.	Yuanzhe Hu, Yu Wang, and Julian McAuley. 2025. Evaluating memory in llm agents via incremental multi-turn interactions. <i>arXiv.org</i> .	794
741			795
742			796
743	Zhiyuan He, Huiqiang Jiang, Zilong Wang, Yuqing Yang, Luna K. Qiu, and Lili Qiu. 2024b. Position engineering: Boosting large language models through positional information manipulation. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Qishuo Hua, Lyumanshan Ye, Dayuan Fu, Yang Xiao, Xiaojie Cai, Yunze Wu, Jifan Lin, Junfei Wang, and Pengfei Liu. 2025. Context engineering 2.0: The context of context engineering. <i>arXiv.org</i> .	797
744			798
745			799
746			800
747			801
748			802
749	Noah Hollmann, Samuel G. Müller, and F. Hutter. 2023. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. <i>Neural Information Processing Systems</i> .	Hai Huang. 2025. Directed information γ -covering: An information-theoretic framework for context engineering. <i>arXiv.org</i> .	803
750			804
751			805
752			806
753			807
754	Kelly Hong, Anton Troynikov, and Jeff Huber. 2025. Context rot: How increasing input tokens impacts llm performance . Technical report, Chroma.	Le Huang, Hengzhi Lan, Zijun Sun, Chuan Shi, and Ting Bai. 2024. Emotional rag: Enhancing role-playing agents through emotional retrieval. In <i>2024 IEEE International Conference on Knowledge Graph (ICKG)</i> .	808
755			809
756			810
757	Haowen Hou, Fei Ma, Binwen Bai, Xinxin Zhu, and F. Yu. 2024a. Enhancing and accelerating large language models via instruction-aware contextual compression. <i>arXiv.org</i> .	Zhengjun Huang, Zhoujin Tian, Qintian Guo, Fangyuan Zhang, Yingli Zhou, Di Jiang, and Xiaofang Zhou. 2025. Licomemory: Lightweight and cognitive agentic memory for efficient long-term reasoning. <i>arXiv.org</i> .	811
758			812
759			813
760			814
761	Yuki Hou, Haruki Tamoto, Qinghua Zhao, and Homei Miyashita. 2024b. Synapticrag: Enhancing temporal memory retrieval in large language models through synaptic mechanisms. <i>Annual Meeting of the Association for Computational Linguistics</i> .	Jonas Hübötter, Sascha Bongni, Ido Hakimi, and Andreas Krause. 2024. Efficiently learning at test-time: Active fine-tuning of llms . <i>ArXiv</i> , abs/2410.08020.	815
762			816
763			817
764			818
765			819
766	N. Houlsby, A. Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and S. Gelly. 2019. Parameter-efficient transfer learning for nlp. In <i>International Conference on Machine Learning</i> .	Jonas Hübötter, Sascha Bongni, Ido Hakimi, and Andreas Krause. 2025. Efficiently learning at test-time: Active fine-tuning of llms . <i>Preprint</i> , arXiv:2410.08020.	820
767			821
768			822
769			823
770			824
771	Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekish, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? <i>Preprint</i> , arXiv:2404.06654.	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression . <i>Preprint</i> , arXiv:2310.06839.	825
772			826
773			827
774			828
775			829
776	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander J. Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. <i>Annual Meeting of the Association for Computational Linguistics</i> .	Josip Jukić, Martin Tutek, and Jan Snajder. 2025. Context parametrization with compositional adapters. <i>arXiv.org</i> .	830
777			831
778			832
779			833
780			834
781			835
782			836
783	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models . <i>Preprint</i> , arXiv:2106.09685.	Jiazheng Kang, Mingming Ji, Zhe Zhao, and Ting Bai. 2025a. Memory os of ai agent. <i>arXiv.org</i> .	837
784			837
785			838
786			839

840	Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. Dense passage retrieval for open-domain question answering. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Wang, Peng Liu, Zehao Lin, Pengyuan Wang, Jiahao Huo, Tianyi Chen, Kai Chen, Ke-Rong Li, and 20 others. 2025c. Memos: A memory os for ai system. <i>arXiv.org</i> .	896 897 898 899
845	Richard Katrix, Quentin Carroway, Rowan Hawkesbury, and Matthias Heathfield. 2025. Context-aware semantic recomposition mechanism for large language models. <i>arXiv.org</i> .	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023a. Lost in the middle: How language models use long contexts . <i>Preprint</i> , arXiv:2307.03172.	900 901 902 903 904
849	Anant Khandelwal, Manish Gupta, and Puneet Agrawal. 2025. Cocoa: Confidence and context-aware adaptive decoding for resolving knowledge conflicts in large language models. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> .	Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrilidis, and Anshumali Shrivastava. 2023b. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. <i>Neural Information Processing Systems</i> .	905 906 907 908 909 910
855	Nikoletta Koilia and C. Kachris. 2024. Hardware acceleration of llms: A comprehensive survey and comparison. <i>arXiv.org</i> .	S. Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	911 912 913 914 915
858	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Haotong Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Symposium on Operating Systems Principles</i> .	Yansheng Mao, Jiaqi Li, Fanxu Meng, Jing Xiong, Zilong Zheng, and Muhan Zhang. 2024. Lift: Improving long context understanding through long input fine-tuning. <i>arXiv.org</i> .	916 917 918 919
864	Patrick Lewis, Ethan Perez, Aleksandara Piktus, F. Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, M. Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Neural Information Processing Systems</i> .	Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Baolong Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi Li, Duzhen Zhang, Chenlin Zhou, Jiayi Mao, Tianze Xia, Jiafeng Guo, and Shenghua Liu. 2025. A survey of context engineering for large language models. <i>arXiv.org</i> .	920 921 922 923 924 925
871	Dongfang Li, Zetian Sun, Xinshuo Hu, Baotian Hu, and Min Zhang. 2024. Cmt: A memory compression method for continual knowledge learning of large language models. <i>arXiv.org</i> .	Adilet Metinov, Gulida M. Kudakeeva, Bolotbek uulu Nursultan, and G. Kabaeva. 2025. Adaptive soft rolling kv freeze with entropy-guided recovery: Sub-linear memory growth for efficient llm inference.	926 927 928 929
875	Mufei Li, Dongqi Fu, Limei Wang, Si Zhang, Hanqing Zeng, Kaan Sancak, Ruizhong Qiu, H. Wang, Xiaoxin He, Xavier Bresson, Yinglong Xia, Chonglin Sun, and Pan Li. 2025a. Haystack engineering: Context engineering for heterogeneous and agentic long-context evaluation. <i>arXiv.org</i> .	Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. 2023. Ret-llm: Towards a general read-write memory for large language models. <i>arXiv.org</i> .	930 931 932 933
881	Yucheng Li, Bo Dong, Chenghua Lin, and Frank Guerin. 2023. Compressing context to enhance inference efficiency of large language models. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Dr.Farheen Mohammed. 2025. Towards standardized evaluation of large language model-based agents. <i>INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT</i> .	934 935 936 937
885	Zhiyu Li, Shichao Song, Hanyu Wang, Simin Niu, Ding Chen, Jiawei Yang, Chenyang Xi, Huayi Lai, Jihao Zhao, Yezhaohui Wang, Junpeng Ren, Zehao Lin, Jiahao Huo, Tianyi Chen, Kai Chen, Ke-Rong Li, Zhiqiang Yin, Qingchen Yu, Bo Tang, and 3 others. 2025b. Memos: An operating system for memory-augmented generation (mag) in large language models. <i>arXiv.org</i> .	Mohammad Mahdi Moradi, Hossam Amer, Sudhir Mudur, Weiwei Zhang, Yang Liu, and Walid Ahmed. 2025. Continuous self-improvement of large language models by test-time training with verifier-driven sample selection . <i>ArXiv</i> , abs/2505.19475.	938 939 940 941 942
893	Zhiyu Li, Shichao Song, Chenyang Xi, Hanyu Wang, Chen Tang, Simin Niu, Ding Chen, Jiawei Yang, Chunyu Li, Qingchen Yu, Jihao Zhao, Yezhaohui	Dilxat Muhtar, Yelong Shen, Yaming Yang, Xiaodong Liu, Yadong Lu, Jianfeng Liu, Yuefeng Zhan, Hao Sun, Weiwei Deng, Feng Sun, Xueliang Zhang, Jianfeng Gao, Weizhu Chen, and Qi Zhang. 2024. Streamadapter: Efficient test time adaptation from contextual streams . <i>Preprint</i> , arXiv:2411.09289.	943 944 945 946 947 948

949	Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofu Chen, S. Zheng, P. Zhao, and Mingkui Tan. 2022. Efficient test-time model adaptation without forgetting . In <i>International Conference on Machine Learning</i> .	1005
950		1006
951		1007
952		1008
953	Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage re-ranking with bert. <i>arXiv.org</i> .	1009
954		
955	N. C. Ohalete, Kevin B. Gittner, and Lauren M. Matheny. 2025. Costar-a: A prompting framework for enhancing large language model performance on point-of-view questions. <i>arXiv.org</i> .	1010
956		1011
957		1012
958		1013
959	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. 2024. Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression . <i>Preprint</i> , arXiv:2403.12968.	1014
960		1015
961		1016
962		1017
963		1018
964		
965	Daehee Park, Jaeseok Jeong, Sung-Hoon Yoon, Jaewoo Jeong, and Kuk-Jin Yoon. 2024. T4p: Test-time training of trajectory prediction via masked autoencoder and actor-specific token memory . <i>2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 15065–15076.	1019
966		1020
967		1021
968		1022
969		1023
970		
971	Haewon Park, Gyubin Choi, Minjun Kim, and Yohan Jo. 2025a. Context-robust knowledge editing for language models. <i>Annual Meeting of the Association for Computational Linguistics</i> .	1024
972		1025
973		1026
974		1027
975	J. Park, Joseph C. O’Brien, Carrie J. Cai, M. Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>ACM Symposium on User Interface Software and Technology</i> .	1028
976		1029
977		
978		1030
979		1031
980	Joon Park, Kyohei Atarashi, Koh Takeuchi, and Hisashi Kashima. 2025b. Emulating retrieval augmented generation via prompt engineering for enhanced long context comprehension in llms. <i>arXiv.org</i> .	1032
981		1033
982		
983		1034
984	Mathis Pink, Qinyuan Wu, Vy A. Vo, Javier S. Turek, Jianing Mu, Alexander Huth, and Mariya Toneva. 2025. Position: Episodic memory is the missing piece for long-term llm agents. <i>arXiv.org</i> .	1035
985		1036
986		
987		1037
988	Patricia Porretta, Sylvester Pakenham, Huxley Ainsworth, Gregory Chatten, Godfrey Allerton, Simon Hollingsworth, and Vance Periwinkle. 2025. Latent convergence modulation in large language models: A novel approach to iterative contextual realignment. <i>arXiv.org</i> .	1038
989		1039
990		1040
991		1041
992		1042
993		
994	Yuxiao Qu, Matthew Y. R. Yang, Amrith Rajagopal Setlur, Lewis Tunstall, Edward Beeching, Ruslan Salakhutdinov, and Aviral Kumar. 2025. Optimizing test-time compute via meta reinforcement fine-tuning . <i>ArXiv</i> , abs/2503.07572.	1043
995		1044
996		1045
997		1046
998		1047
999	Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, A. Shashua, Kevin Leyton-Brown, and Y. Shoham. 2022. Parallel context windows for large language models. <i>Annual Meeting of the Association for Computational Linguistics</i> .	1048
1000		1049
1001		
1002		1050
1003		1051
1004		1052
		1053
		1054
		1055
		1056
		1057
		1058
	Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyang Yang, Minjia Zhang, Dong Li, and Yuxiong He. 2021. Zero-offload: Democratizing billion-scale model training. In <i>USENIX Annual Technical Conference</i> .	
	Alireza Rezazadeh, Zichao Li, Wei Wei, and Yujia Bao. 2024. From isolated conversations to hierarchical schemas: Dynamic tree memory representation for llms. In <i>International Conference on Learning Representations</i> .	
	Sourjya Roy, Shrihari Sridharan, Surya Selvam, and A. Raghunathan. 2025. Kv-car: Kv cache compression using autoencoders and kv reuse in large language models.	
	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, S. Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. <i>arXiv.org</i> .	
	Rana Salama, Jason Cai, Michelle Yuan, Anna Currey, Monica Sunkara, Yi Zhang, and Yassine Benajiba. 2025. Meminsight: Autonomous memory augmentation for llm agents. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> .	
	Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. 2023. Zeroscrolls: A zero-shot benchmark for long text understanding . <i>Preprint</i> , arXiv:2305.14196.	
	Lianlei Shan, Shixian Luo, Zezhou Zhu, Yu Yuan, and Yong Wu. 2025. Cognitive memory in large language models. <i>arXiv.org</i> .	
	Weizhou Shen, Chenliang Li, Fanqi Wan, Shengyi Liao, Shaopeng Lai, Bo Zhang, Yingcheng Shi, Yuning Wu, Gang Fu, Zhansheng Li, Bin Yang, Ji Zhang, Fei Huang, Jingren Zhou, and Ming Yan. 2025. Qwenlong-cprs: Towards ∞ -llms with dynamic context optimization. <i>arXiv.org</i> .	
	Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Daniel Y. Fu, Zhiqiang Xie, Beidi Chen, Clark W. Barrett, Joseph Gonzalez, Percy Liang, Christopher Ré, Ion Stoica, and Ce Zhang. 2023. High-throughput generative inference of large language models with a single gpu. In <i>International Conference on Machine Learning</i> .	
	Kaize Shi, Xueyao Sun, Xiaohui Tao, Lin Li, Qika Lin, and Guandong Xu. 2025. Concept than document: Context compression via amr-based conceptual entropy.	
	Weijia Shi, Xiaochuang Han, M. Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and S. Yih. 2023. Trusting your evidence: Hallucinate less with context-aware decoding. <i>North American Chapter of the Association for Computational Linguistics</i> .	

1059	Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Eliciting knowledge from language models using automatically generated prompts. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	test-time adaptation . <i>2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3728–3738.	1114 1115 1116
1064	Noah Shinn, Federico Cassano, Beck Labash, A. Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. <i>Neural Information Processing Systems</i> .	Alessandra Terranova, Bjorn Ross, and Alexandra Birch. 2025. Evaluating long-term memory for long-context question answering. <i>arXiv.org</i> .	1117 1118 1119
1066	H. Shinwari and Muhammad Usama. 2025. Memory-augmented architecture for long-term context handling in large language models. <i>arXiv.org</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and I. Polosukhin. 2017. Attention is all you need. <i>Neural Information Processing Systems</i> .	1120 1121 1122 1123
1067	Alina Shutova, Vladimir Malinovskii, Vage Egiazarian, Denis Kuznedelev, D. Mazur, Nikita Surkov, Ivan Ermakov, and Dan Alistarh. 2025. Cache me if you must: Adaptive key-value quantization for large language models. In <i>International Conference on Machine Learning</i> .	Akash Vishwakarma, Hojin Lee, Mohith Suresh, Priyam Shankar Sharma, Rahul Vishwakarma, Sparsh Gupta, and Yuvraj Anupam Chauhan. 2025. Cognitive weave: Synthesizing abstracted knowledge with a spatio-temporal resonance graph. <i>arXiv.org</i> .	1124 1125 1126 1127 1128
1072	Sonish Sivarajkumar, Mark Kelley, Alyssa Samolyk-Mazzanti, Shyam Visweswaran, and Yanshan Wang. 2024. An empirical evaluation of prompting strategies for large language models in zero-shot clinical natural language processing: Algorithm development and validation study. <i>JMIR Medical Informatics</i> .	Lei Wang, Chengbang Ma, Xueyang Feng, Zeyu Zhang, Hao ran Yang, Jingsen Zhang, Zhi-Yang Chen, Jikai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji rong Wen. 2023. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> .	1129 1130 1131 1132 1133 1134
1073	C. Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters . <i>ArXiv</i> , abs/2408.03314.	Renxi Wang, Xudong Han, Lei Ji, Shu Wang, Timothy Baldwin, and Haonan Li. 2024a. Toolgen: Unified tool retrieval and calling via generation. In <i>International Conference on Learning Representations</i> .	1135 1136 1137 1138
1074	Woomin Song, Seunghyuk Oh, Sangwoo Mo, Jaehyung Kim, Sukmin Yun, Jung-Woo Ha, and Jinwoo Shin. 2024. Hierarchical context merging: Better long context understanding for pre-trained llms. In <i>International Conference on Learning Representations</i> .	Rushi Wang, Jiateng Liu, Cheng Qian, Yifan Shen, Yanzhou Pan, Zhaozhuo Xu, Ahmed Abbasi, Heng Ji, and Denghui Zhang. 2025a. Context engineering for trustworthiness: Rescorla wagner steering under mixed and inappropriate contexts. <i>arXiv.org</i> .	1139 1140 1141 1142 1143
1075	Sruthi Sridhar, Abdulrahman Khamaj, and M. Asthana. 2023. Cognitive neuroscience perspective on memory: overview and summary. <i>Frontiers in Human Neuroscience</i> .	Xiangfeng Wang, Zaiyi Chen, Zheyong Xie, Tong Xu, Yongyi He, and Enhong Chen. 2024b. In-context former: Lightning-fast compressing context for large language model. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1144 1145 1146 1147 1148
1076	Sakhinana Sagar Srinivas and Venkataramana Runkana. 2025. Scaling test-time inference with policy-optimized, dynamic retrieval-augmented generation via kv caching and decoding. <i>arXiv.org</i> .	Yu Wang and Xi Chen. 2025. Mirix: Multi-agent memory system for llm-based agents. <i>arXiv.org</i> .	1149 1150
1077	Huashan Sun, Shengyi Liao, Yansen Han, Yu Bai, Yang Gao, Cheng Fu, Weizhou Shen, Fanqi Wan, Mingshi Yan, Ji Zhang, and Fei Huang. 2025. Solopo: Unlocking long-context capabilities in llms via short-to-long preference optimization. <i>arXiv.org</i> .	Yu Wang, Ryuichi Takanobu, Zhiqi Liang, Yuzhen Mao, Yuanzhe Hu, Julian McAuley, and Xiaojian Wu. 2025b. Mem- α : Learning memory construction via reinforcement learning. <i>arXiv.org</i> .	1151 1152 1153 1154
1078	Sijun Tan, Xiuyu Li, Shishir G. Patil, Ziyang Wu, Tianjun Zhang, Kurt Keutzer, Joseph E. Gonzalez, and Raluca A. Popa. 2024. Lloco: Learning long contexts offline. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. <i>Neural Information Processing Systems</i> .	1155 1156 1157 1158 1159
1079	Yushun Tang, Ce Zhang, Heng Xu, Shuoshuo Chen, Jie Cheng, Luziwei Leng, Qinghai Guo, and Zhihai He. 2023. Neuro-modulated hebbian learning for fully	Benjue Weng. 2024. Navigating the landscape of large language models: A comprehensive review and analysis of paradigms and fine-tuning strategies. <i>arXiv.org</i> .	1160 1161 1162 1163
1080		Haoyi Wu and Kewei Tu. 2024. Layer-condensed kv cache for efficient inference of large language models. <i>Annual Meeting of the Association for Computational Linguistics</i> .	1164 1165 1166 1167

Algorithm 1 GDWM: Gated Differentiable Working Memory

Require: Context $X = x_{1:L}$, total steps K_{total} , chunk size S , local window n , temperature τ , min coverage k_{min}

Ensure: Adapted fast weights ϕ

- 1: $\mathbf{K}, \mathbf{V} \leftarrow \text{PREFILLANDFREEZE}(X)$ {Compute frozen KV cache}
 - 2: Partition X into $\{C_1, \dots, C_M\}$ with $M = \lceil L/S \rceil$
 - 3: **for** each chunk $c = 1, \dots, M$ **do**
 - 4: $U(C_c) \leftarrow \frac{1}{|\mathcal{I}_c|} \sum_{t \in \mathcal{I}_c} |\log P_{\text{full}}(x_t) - \log P_{\text{local}}(x_t)|$
 - 5: **end for**
 - 6: Compute weights w_c via Eq. (6)
 - 7: Allocate $\{k_c\}$ via Coverage-First + Softmax (Eq. 7)
 - 8: **for** each chunk $c = 1, \dots, M$ **do**
 - 9: **for** $j = 1, \dots, k_c$ **do**
 - 10: $\mathcal{I} \leftarrow \text{UNIFORMSAMPLE}(C_c, B)$ { B : batch size}
 - 11: $\mathcal{L} \leftarrow -\frac{1}{B} \sum_{i \in \mathcal{I}} \log P(x_i | x_{1:i-1}; \phi, \theta_f)$
 - 12: $\phi \leftarrow \phi - \eta \nabla_{\phi} \mathcal{L}$
 - 13: **end for**
 - 14: **end for**
 - 15: **return** ϕ
-

Input and Chunking.

- $X = (x_1, x_2, \dots, x_L)$: The input context sequence of length L .
- C_c : The c -th chunk, a contiguous subsequence of X .
- $M = \lceil L/S \rceil$: Total number of chunks.
- S : Chunk size (number of tokens per chunk); default $S = 1024$.

Model Parameters.

- θ : Full parameters of the pre-trained LLM.
- θ_f : Frozen parameters, including weights used to compute the KV cache.
- ϕ : Adaptable “fast weights” (e.g., LoRA adapters on W_Q, W_O).
- ϕ_0 : Initial state of fast weights (typically zero for LoRA).
- $\phi(\mathbf{k})$: Final state of ϕ after applying the allocation schedule \mathbf{k} .

Budget Allocation.

- $\mathbf{k} = (k_1, \dots, k_M)$: Allocation vector; k_c is the number of gradient updates on chunk C_c .
- K_{total} : Total gradient budget across all chunks.
- k_{min} : Minimum coverage constraint per chunk; default $k_{\text{min}} = 1$.
- τ : Temperature for softmax allocation; default $\tau = 1.0$.

Contextual Utility.

- $U(C_c)$: Contextual Utility of chunk C_c , measuring long-range dependency.
- $\Delta_t = |\log P(x_t | x_{1:t-1}) - \log P(x_t | x_{t-n:t-1})|$: Surprisal divergence at position t .
- n : Local window size for computing P_{local} ; default $n = 512$.

Sampling and Distributions.

- \mathcal{I}_c : Set of token positions in chunk C_c .
- $\mathcal{U}(\mathcal{I}_c)$: Uniform distribution over positions in chunk C_c .
- $x_t \sim \mathcal{U}(\mathcal{I}_c)$: Position t sampled uniformly from chunk C_c .
- $\mathcal{I} \subset \mathcal{I}_c$: Minibatch of positions for gradient computation.

KV Cache and Global Context.

- \mathbf{K}, \mathbf{V} : Frozen key-value cache for the entire input X , computed via PREFILLANDFREEZE.
- $x_{1:t-1}$: **Full document prefix** from position 1 to $t - 1$, *not* chunk-local context.

Key Clarification: Global vs. Local Context. A potential source of confusion is the conditioning context in $P(x_t | x_{1:t-1})$. Throughout this paper, $x_{1:t-1}$ refers to the *entire* prefix of the document, accessible via the frozen KV cache (\mathbf{K}, \mathbf{V}). The chunking strategy determines *where* gradient updates are computed (i.e., which positions t contribute to the loss), but does *not* restrict the context the model can attend to. This design ensures that the model leverages global information during adaptation while focusing computational budget on high-utility regions.

Sampling Strategy. For each gradient step $j \in \{1, \dots, k_c\}$ on chunk C_c , we sample a minibatch \mathcal{I}_j of B positions *independently* from the uniform distribution $\mathcal{U}(\mathcal{I}_c)$. Formally:

$$\mathcal{I}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(\mathcal{I}_c)^B, \quad \forall j \in \{1, \dots, k_c\} \quad (9)$$

This means:

- **Across steps:** Each minibatch \mathcal{I}_j is sampled independently; the same position may appear in multiple steps.
- **Within step:** Positions are sampled i.i.d.; when $B \ll |\mathcal{I}_c|$ (our setting: $B = 32$, $|\mathcal{I}_c| = 1024$), overlap within a single batch is negligible.

This is the standard stochastic gradient descent (SGD) sampling scheme. The independence across steps ensures unbiased gradient estimates, while chunk-restriction (sampling only from \mathcal{I}_c , not the full document) eliminates inter-chunk variance as shown in Theorem 1.

C Theoretical Analysis

A key advantage of GDWM over standard TTT methods (which sample uniformly across the full context) is stability. We analyze this through the lens of gradient variance.

C.1 Variance Reduction via Chunk-Restricted Sampling

A natural question arises: why does restricting gradient computation to individual chunks improve optimization? Long documents exhibit *semantic heterogeneity*—different sections address distinct topics and induce gradient directions that may interfere destructively when aggregated. Restricting sampling to a single chunk eliminates this cross-sectional interference.

Theorem 1 (Variance Reduction). *Let g be the gradient estimator for a single update step. Assume the document consists of M chunks, where chunk c has gradient mean μ_c and variance σ_c^2 . Under global uniform sampling:*

$$\text{Var}(g_{\text{global}}) = \underbrace{\mathbb{E}_c[\sigma_c^2]}_{\text{intra-chunk}} + \underbrace{\text{Var}_c(\mu_c)}_{\text{inter-chunk}} \quad (10)$$

Under chunk-restricted sampling (conditioned on chunk c), we have $\text{Var}(g | c) = \sigma_c^2$, eliminating the inter-chunk variance term. Equality holds only when $\mu_c = \bar{\mu}$ for all c (semantically homogeneous document).

Consequence. By scheduling updates per chunk, we eliminate the inter-chunk variance from each gradient estimate. The optimizer follows a more consistent trajectory, avoiding destructive interference between gradient signals from semantically disparate document sections. This provides a principled explanation for the empirical observation that GDWM achieves lower perplexity with fewer optimization steps.

C.2 Role of the Coverage Constraint

For tasks requiring global understanding (e.g., summarization), the coverage constraint ($k_c \geq k_{\min}$) acts as a regularizer against *mode collapse*.

Intuition. Without coverage, greedy CPMI-based allocation may concentrate all budget on a single high-utility region. For example, in a government report, the “Results” section may have the highest CPMI, but a good summary also requires context from “Introduction” and “Methods.”

Claim (Informal). If the optimal output requires information from all M sections but the model adapts only to a subset $\mathcal{S} \subsetneq \{1, \dots, M\}$, the coverage gap translates to an $O(|\mathcal{S}|/M)$ recall bound.

The constraint $k_c \geq k_{\min}$ ensures minimum representation of every section, preventing pathological allocation while allowing CPMI to modulate *relative* budget.

D Proofs

D.1 Proof of Theorem 1 (Variance Reduction)

Let $C \in \{1, \dots, M\}$ be a random variable indicating which chunk a position is sampled from. Let g denote the gradient at that position.

By the **Law of Total Variance**:

$$\text{Var}(g) = \mathbb{E}[\text{Var}(g | C)] + \text{Var}(\mathbb{E}[g | C]) \quad (11)$$

Let $\mu_c = \mathbb{E}[g | C = c]$ denote the mean gradient in chunk c , and $\sigma_c^2 = \text{Var}(g | C = c)$ denote the within-chunk variance. Let p_c be the probability of sampling from chunk c (under uniform sampling, $p_c = |\mathcal{I}_c|/L$).

Expanding each term:

$$\mathbb{E}[\text{Var}(g | C)] = \sum_{c=1}^M p_c \sigma_c^2 \quad (12)$$

$$\text{Var}(\mathbb{E}[g | C]) = \sum_{c=1}^M p_c (\mu_c - \bar{\mu})^2 \quad (13)$$

where $\bar{\mu} = \sum_c p_c \mu_c$ is the global mean.

Therefore:

$$\text{Var}(g_{\text{global}}) = \underbrace{\sum_{c=1}^M p_c \sigma_c^2}_{\text{intra-chunk variance}} + \underbrace{\sum_{c=1}^M p_c (\mu_c - \bar{\mu})^2}_{\text{inter-chunk variance}} \quad (14)$$

Under **chunk-restricted sampling**, when we sample from a specific chunk c , the gradient variance is exactly σ_c^2 (the intra-chunk variance of that chunk). Since the inter-chunk term $\sum_c p_c (\mu_c - \bar{\mu})^2 \geq 0$, we have:

$$\text{Var}(g | C = c) = \sigma_c^2 \leq \text{Var}(g_{\text{global}}) \quad (15)$$

Equality holds if and only if the inter-chunk variance is zero, i.e., $\mu_c = \bar{\mu}$ for all c , meaning the document is semantically homogeneous (all chunks have identical gradient expectations).

□

E Coverage Constraint Analysis

The constraint $k_c \geq k_{\min}$ in Eq. (3) prevents the model from overfitting to a single high-utility region. This is crucial for tasks requiring global understanding.

Example: GovReport Summarization. A typical government report contains Introduction, Methods, Results, and Conclusion sections. Without the coverage constraint, CPMI-based allocation may concentrate all budget on Results (often the highest information density region), missing essential context from other sections that the summary must include.

Failure Mode. If the optimal summary requires information from all M sections but the model adapts only to a subset $\mathcal{S} \subsetneq \{1, \dots, M\}$, the coverage gap directly translates to recall loss proportional to $|\mathcal{S}|/M$.

The constraint ensures that every section receives at least minimal gradient exposure, preventing pathological allocation while still allowing CPMI to modulate the *relative* budget across sections.

F Justification for Absolute Value in Δ_t

The absolute value in $\Delta_t = |\log P(x_t | x_{1:t-1}) - \log P(x_t | x_{t-n:t-1})|$ captures both positive and negative divergence.

Information-Theoretic Perspective (Pointwise KL). Mathematically, the quantity Δ_t corresponds to the magnitude of the *Pointwise Kullback-Leibler Divergence* contribution at token x_t . The standard KL divergence $D_{\text{KL}}(P_{\text{full}} || P_{\text{local}})$ is the expectation of the log-likelihood ratio. Our utility metric $U(C_c)$ essentially estimates the *L1 norm* of this pointwise divergence over the chunk. We prefer the absolute value (L1-like) over the raw difference (which would sum to the standard KL) because gradient updates are driven by the *magnitude* of the error signal.

- **Positive divergence** ($P_{\text{full}} > P_{\text{local}}$): Long-range context reduces surprisal (adds information).
- **Negative divergence** ($P_{\text{full}} < P_{\text{local}}$): Long-range context increases surprisal (introduces conflict).

Both cases represent significant deviations between the global and local models, identifying regions where the fast weights ϕ must adapt to reconcile these discrepancies.

Empirical Validation. Using $\max(0, \Delta_t)$ (ignoring negative divergence) instead of $|\Delta_t|$ results in 3-5% performance degradation across tasks. The “hard examples” with negative divergence are precisely where the model needs recalibration.

G Hyperparameter Sensitivity

We briefly analyze the sensitivity of GDWM to its two key hyperparameters: the temperature τ and the minimum coverage k_{\min} .

Temperature τ . The temperature controls the sharpness of the softmax allocation (Eq. 6).

- $\tau \rightarrow 0$ (**Greedy**): The budget concentrates entirely on the single chunk with the highest utility. This risks overfitting to one region and failing on multi-hop tasks.
- $\tau \rightarrow \infty$ (**Uniform**): The policy degenerates into uniform sampling (equivalent to qTTT), losing the efficiency benefits of selection.
- $\tau \approx 1.0$ (**Balanced**): Empirically, values in $[0.5, 1.5]$ perform robustly. We adopt $\tau = 1.0$ as a tuning-free default.

Minimum Coverage k_{\min} . We set $k_{\min} = 1$ to ensure no chunk is entirely starved of gradients. Increasing k_{\min} reduces the budget available for utility-based redistribution, pushing the behavior closer to uniform sampling. $k_{\min} = 1$ represents the minimal constraint necessary to prevent mode collapse while maximizing the freedom of the allocation policy.

H Testable Predictions from Theorem 1

Theorem 1 yields three empirically verifiable predictions:

1. **Gradient Norm Variance:** The sequence $\{\|g_t\|\}$ under **GDWM** should exhibit lower variance than under global uniform sampling.
2. **Within-Chunk Cosine Similarity:** The cosine similarity $\cos(g_i, g_j)$ for positions i, j within the same chunk should be higher than for positions in different chunks.
3. **Loss Curve Monotonicity:** The training loss curve should converge more monotonically with fewer oscillations under chunk-restricted sampling.

These predictions are consistent with our empirical observations and provide a testable framework for validating the theoretical analysis.

I Largest Remainder Allocation Method

After the initial floor allocation in Eq. (7), there may be residual steps due to rounding. We distribute these using the *Largest Remainder Method*:

1. Compute fractional remainders: $r_c = K_{\text{rem}} \cdot w_c - \lfloor K_{\text{rem}} \cdot w_c \rfloor$
2. Sort chunks by r_c in descending order
3. Assign one additional step to each of the top- R chunks, where $R = K_{\text{total}} - \sum_c \lfloor K_{\text{rem}} \cdot w_c \rfloor - M \cdot k_{\min}$

This guarantees that the total allocation exactly equals K_{total} while respecting the utility-based priority ordering.

J Evidence Span Analysis

We formalize the relationship between chunk size and task performance through the concept of *evidence span*.

Definition 2 (Evidence Span). For a task \mathcal{T} with query q and context X , the evidence span $E_{\mathcal{T}}$ is the minimum contiguous token range required to contain all information necessary for correct response generation:

$$E_{\mathcal{T}} = \min_{[i,j]} \{j - i : X_{[i,j]} \text{ suffices for } \mathcal{T}\} \quad (16)$$

Evidence Fragmentation Problem. When chunk size $S < E_{\mathcal{T}}$, the evidence is partitioned across multiple chunks. Let the evidence span $[a, b]$ be split into $k = \lceil (b - a) / S \rceil$ chunks. For each fragment C_i , the CPPI estimate becomes:

$$\hat{U}(C_i) = \text{CPMI}(C_i) < \text{CPMI}([a, b]) = U_{\text{true}} \quad (17)$$

This underestimation occurs because each fragment, viewed in isolation, appears to have low contextual dependency—the long-range signal is diluted across fragments.

Optimal Chunk Size Selection. The optimal chunk size satisfies:

$$S^* = \min\{S : S \geq E_{\mathcal{T}}, \forall \mathcal{T} \in \text{target tasks}\} \quad (18)$$

For a diverse benchmark like ZeroSCROLLS containing both local and multi-hop tasks, $S = 1024$ emerges as the robust choice: it satisfies $S \geq E_{\mathcal{T}}$ for most multi-hop instances while maintaining reasonable granularity for local tasks.

Empirical Validation. Table 4 in the main text validates this analysis:

- At $S = 512$: MuSiQue collapses to 10.0 (evidence fragmentation)
- At $S = 1024$: MuSiQue recovers to 25.8 ($S \geq E_{\mathcal{T}}$)
- At $S = 2048$: Marginal improvement (25.3) with reduced local precision

This provides a principled explanation for why $S = 1024$ achieves the best cross-task robustness: it is the minimum chunk size that avoids catastrophic evidence fragmentation across the task distribution. For tasks with intermediate evidence spans like Qasper, performance naturally peaks at the corresponding granularity ($S = 512$), confirming that $S \approx E_{\mathcal{T}}$ is the theoretical optimum.

1690 J.1 Information-Theoretic Lower Bound on 1691 Fragmented Utility

1692 The empirical observation that fragmented evi-
1693 dence leads to utility underestimation admits a formal information-theoretic explanation. We show
1694 that chunking an evidence span introduces a *systematic negative bias* in utility estimation due to
1695 the loss of *synergistic information*.
1696
1697

1698 **Definition 3** (Contextual Utility as Mutual Informa-
1699 tion). *Let G denote the global context (long-range
1700 prefix) and L denote the local context (sliding win-
1701 dow). For a token x_t , the Contextual Utility can
1702 be interpreted as the absolute mutual information
1703 gain:*

$$1704 \begin{aligned} U(x_t) &= |I(x_t; G) - I(x_t; L)| \\ 1705 &\approx |I(x_t; G \setminus L \mid L)| \end{aligned} \quad (19)$$

1706 where $G \setminus L$ represents the distant context beyond
1707 the local window. For a chunk C , the aggregate
1708 utility is $U(C) = \sum_{x_t \in C} U(x_t)$.

1709 **Proposition 1** (Utility Underestimation Under
1710 Fragmentation). *Let $E = C_1 \cup C_2$ be an evidence
1711 span partitioned into two adjacent chunks. Then
1712 the sum of individual chunk utilities is bounded
1713 above by the utility of the unified span:*

$$1714 U(C_1) + U(C_2) \leq U(E) + \epsilon \quad (20)$$

1715 where $\epsilon \leq 0$ when the chunks exhibit informa-
1716 tion synergy (i.e., positive interaction information).
1717 Equality holds if and only if C_1 and C_2 are infor-
1718 mationally independent given the global context.

1719 *Proof.* We leverage the chain rule of mutual infor-
1720 mation. Let G denote the global context. The joint
1721 utility of the unified evidence span $E = C_1 \cup C_2$
1722 satisfies:

$$1723 I(E; G) = I(C_1; G) + I(C_2; G \mid C_1) \quad (21)$$

1724 The fragmented utility estimation treats C_1 and
1725 C_2 independently:

$$1726 \hat{U}(E) = I(C_1; G) + I(C_2; G) \quad (22)$$

1727 The *fragmentation gap* is therefore:

$$1728 \begin{aligned} \Delta_{\text{frag}} &= I(E; G) - \hat{U}(E) \\ 1729 &= I(C_2; G \mid C_1) - I(C_2; G) \\ 1730 &= -I(C_1; C_2; G) \end{aligned} \quad (23)$$

1731 where $I(C_1; C_2; G)$ is the *interaction information*
1732 (also known as co-information or multivariate mu-
1733 tual information), defined as:

$$1734 I(C_1; C_2; G) = I(C_2; G) - I(C_2; G \mid C_1) \quad (24)$$

1735 **Interpretation.** The interaction information
1736 $I(C_1; C_2; G)$ measures the degree to which C_1 and
1737 C_2 *synergistically* inform G :

- 1738 • **Positive interaction** ($I > 0$): C_1 and C_2 are
1739 *redundant*—knowing one reduces the infor-
1740 mation gain from the other. Fragmentation
1741 causes *overestimation* (rare).
- 1742 • **Negative interaction** ($I < 0$): C_1 and C_2
1743 are *synergistic*—together they provide more
1744 than the sum of parts. This is characteristic
1745 of **reasoning chains**. Fragmentation causes
1746 *underestimation*.

1747 For multi-hop reasoning tasks like MuSiQue,
1748 where the answer requires synthesizing facts from
1749 multiple locations, the evidence exhibits strong neg-
1750 ative interaction information. Thus:

$$1751 \begin{aligned} I(C_1; C_2; G) < 0 &\implies \Delta_{\text{frag}} > 0 \\ &\implies U(E) > \hat{U}(E) \end{aligned} \quad (25)$$

1753 This proves that fragmented utility estimation is
1754 a *systematic lower bound* on the true utility when
1755 evidence is synergistic. \square

1756 **Consequence for Chunk Size Selection.** Propo-
1757 sition 1 rigorously justifies the catastrophic perfor-
1758 mance collapse at small chunk sizes on multi-hop
1759 tasks. The fragmentation gap Δ_{frag} is not merely
1760 noise but a *structured bias* scaling with inter-chunk
1761 synergy. The optimal S^* must satisfy $S \geq E_{\mathcal{T}}$ to
1762 ensure synergistic evidence is not partitioned.

1763 **Connection to Cognitive Science.** This parallels
1764 the *binding problem* in cognitive neuroscience: a
1765 reasoning chain stored in working memory must
1766 be represented as a unified chunk to preserve logi-
1767 cal coherence. Fragmenting it destroys emergent
1768 meaning, analogous to how fragmenting a sentence
1769 into words loses compositional semantics.

1770 K Theoretical Justification for CPMI vs. 1771 Surprisal

1772 A natural alternative to our CPMI-based utility is a
1773 simpler, non-uniform gating policy based solely on
1774 *Surprisal* (or Perplexity), i.e., prioritizing chunks
1775 with high $\mathcal{L}_{\text{local}}(x) = -\log P(x \mid x_{\text{local}})$. While
1776 intuitively appealing (allocating compute to “hard”
1777 regions), this approach is mathematically subopti-
1778 mal for long-context adaptation.

Proposition. Surprisal measures *difficulty*, whereas CPMI measures *dependency*.

Let the information content of a token x_t be decomposed as:

$$I(x_t | x_{\text{global}}) = I(x_t | x_{\text{local}}) - \underbrace{I(x_t; x_{\text{distant}} | x_{\text{local}})}_{\text{CPMI}} \quad (26)$$

High surprisal ($I(x_t | x_{\text{local}})$ is large) can arise from two sources:

1. **Aleatoric Uncertainty:** The token is inherently unpredictable (e.g., a random name or number), regardless of context.
2. **Epistemic Uncertainty (Contextual):** The token is predictable given long-range context but unpredictable locally.

Gradient adaptation on Type 1 tokens is wasteful—it forces the model to memorize noise. Adaptation on Type 2 tokens is high-value—it retrieves recoverable information. CPMI ($|\log P_{\text{full}} - \log P_{\text{local}}|$) specifically isolates Type 2 uncertainty by cancelling out the intrinsic difficulty. A surprisal-based baseline would confuse noise with signal, allocating budget to intrinsically hard but context-irrelevant tokens. Thus, CPMI is the theoretically correct objective for *context-dependent* memory consolidation.

L Detailed Related Work

Context Engineering Context engineering optimizes information payloads for large language models, spanning techniques from foundational prompt design to advanced management strategies (Mei et al., 2025; Huang, 2025; Hua et al., 2025; Sahoo et al., 2024; Haider et al., 2024). Core methodologies include in-context learning and chain-of-thought prompting for adaptive reasoning (Weng, 2024; Allingham et al., 2023; Brown et al., 2020a; Chen et al., 2025b; Wei et al., 2022; Ohalete et al., 2025). Recent approaches focus on retrieval-augmented generation (RAG) and automated prompt optimization (Lewis et al., 2020; Nogueira and Cho, 2019; Shin et al., 2020; Li et al., 2025a; Karpukhin et al., 2020), while efficient context management techniques address quadratic scaling limitations through window extension, compression, and hierarchical processing (Li et al., 2023; Vaswani et al., 2017; Mao et al., 2024; Fu et al., 2024; Duan et al., 2025; Zhu et al., 2025;

Tan et al., 2024; Wang et al., 2024b; Song et al., 2024; Zhou et al., 2024; Hou et al., 2024a; Ratner et al., 2022; Zhang et al., 2024f; Sun et al., 2025; Wu et al., 2025b; Wang et al., 2024a). Parallel efforts tackle robustness and knowledge conflicts via context-aware decoding, representation engineering, and activation alignment (He et al., 2024b,a; Govindan et al., 2025; Park et al., 2025a; Wang et al., 2025a; Zhao et al., 2024b; Longpre et al., 2021; Khandelwal et al., 2025; Shi et al., 2023; Zhou et al., 2023; Zhao et al., 2024a; Shen et al., 2025; Porretta et al., 2025; Katrix et al., 2025; Jukić et al., 2025; Houlisby et al., 2019; Park et al., 2025b). Domain-specific applications demonstrate these methods in tabular analysis, translation, clinical NLP, and traffic prediction (Hollmann et al., 2023; Wu and Hu, 2023; Sivarajkumar et al., 2024; Zheng et al., 2024; Yang et al., 2023), often utilizing filtering mechanisms to enhance retrieval precision (Shi et al., 2025; Cheng et al., 2024; Chakraborty et al., 2025).

Memory Management Memory management in large language models has evolved to address the critical bottlenecks imposed by model weights and, increasingly, the ephemeral activation memory required for inference. While parametric memory stores implicit knowledge in model weights (Li et al., 2025b; Zhang et al., 2025b; Hsieh et al., 2023), the key-value (KV) cache has emerged as the dominant constraint, often consuming significantly more memory than parameters and scaling linearly with sequence length (Liu et al., 2023b; Zhang et al., 2024d; Kwon et al., 2023; Dai et al., 2024; Cai et al., 2025; Roy et al., 2025; Shutova et al., 2025; Barua, 2024; Modarressi et al., 2023; Li et al., 2025c; Cheng et al., 2025; Brown et al., 2020b; Gao et al., 2018; Dao et al., 2022). To mitigate these limitations, research has introduced virtual memory and paging techniques akin to operating systems (Kwon et al., 2023; Xue et al., 2024; Koilia and Kachris, 2024; Zhang et al., 2025a; Chitty-Venkata et al., 2025), as well as advanced compression and quantization methods to reduce footprints without quality loss (Xie et al., 2025; Srinivas and Runkana, 2025; Zhang et al., 2024c). Dynamic eviction strategies, including heavy-hitter identification and attention-based pruning, selectively discard less critical tokens to maintain fixed budgets (Yuan et al., 2025; Zhang et al., 2023; Zeng et al., 2025; Shinwari and Usama, 2025; Yi et al., 2024). Parallelization and offloading frameworks

further optimize resource utilization across devices (Yang et al., 2024; Ren et al., 2021; Sheng et al., 2023; Wu and Tu, 2024; Banasik, 2025), while hierarchical and biologically inspired architectures offer structured approaches to long-term memory and personalization (Fang et al., 2024; Rezazadeh et al., 2024; Kang et al., 2025a; Hou et al., 2024b; Li et al., 2024; Metinov et al., 2025; Zhong et al., 2023a; Huang et al., 2024).

Agent Memory Research on agent memory in large language models addresses the constraints of limited context windows and static training data by enabling persistent information storage and retrieval. Existing approaches categorize memory into distinct types, differentiating between short-term working memory and long-term storage (Han et al., 2024; Wang and Chen, 2025; Sridhar et al., 2023), as well as distinguishing parametric memory in model weights from explicit contextual memory (Du et al., 2025a; Shan et al., 2025). More granular taxonomies identify specialized forms such as episodic, consensus, semantic, and procedural memory (Han et al., 2024; Terranova et al., 2025). Structurally, methods range from knowledge-organization and retrieval-oriented mechanisms to architecture-driven hierarchies (Kang et al., 2025b; Zhong et al., 2023b). While early work relied heavily on Retrieval Augmented Generation (RAG) and static vector databases (Lewis et al., 2020; Xu et al., 2025b; Vishwakarma et al., 2025), recent advancements explore graph-based systems (Anokhin et al., 2024; Vishwakarma et al., 2025) and diverse formats like natural language or embeddings (Wang et al., 2023; Shinn et al., 2023). Contemporary frameworks implement core operations—consolidation, updating, indexing, and forgetting—to manage memory dynamics, often incorporating biologically inspired mechanisms like forgetting curves (Xiong et al., 2025; Park et al., 2023; Zhao et al., 2023; Cao et al., 2025). Despite these advances, challenges remain in autonomous memory management, with many systems relying on manual predefinitions (Zhang et al., 2025d; Wang et al., 2025b; Zhang et al., 2024e) or suffering from structural rigidity and catastrophic forgetting (Xu et al., 2025b; Zhang et al., 2025c; Guo et al., 2023). Applications span gaming, dialogue, and procedural tasks (Hu et al., 2024b; Zeng et al., 2024; Hu et al., 2024a; Fang et al., 2025; Mohammed, 2025). Future research directions emphasize unified theo-

retical frameworks (Wu et al., 2025a), episodic memory for single-shot learning (Pink et al., 2025), and scalable, autonomous systems capable of self-evolution (Salama et al., 2025; Zhang et al., 2024e; Hu et al., 2025; Huang et al., 2025).

Test-Time Adaptation. A growing line of work treats deployment as an optimization phase, improving robustness and adaptation by updating model states at inference time under shift or novelty, often via self-supervised objectives and safeguards against forgetting or noisy updates (Niu et al., 2022; Tang et al., 2023; Park et al., 2024; Zhang et al., 2024a). Within LLMs, inference-time optimization spans gradient-based test-time training on task instances or auxiliary data, including adaptation driven by retrieved neighbors, contextual streams, active/verification-guided sample selection, and selective test-time learning for evaluation models (Hardt and Sun, 2024; Muhtar et al., 2024; Hübötter et al., 2024; Hübötter et al., 2025; Moradi et al., 2025; Jwa et al., 2025; Akyürek et al., 2025). In parallel, test-time compute can be optimized at the policy level—allocating computation across candidate solutions or update actions—through principled test-time scaling, meta-learned compute control, or reinforcement learning at inference (Snell et al., 2024; Qu et al., 2025; Zuo et al., 2025). Complementary inference-time strategies reduce effective context cost by compressing or distilling long inputs and intermediate representations, ranging from selective augmentation/compression in retrieval pipelines to learned and training-free prompt compression, dynamic allocation of soft tokens, activation-based beacons, and near-lossless KV compression (Xu et al., 2023; Chevalier et al., 2023; Jiang et al., 2024; Pan et al., 2024; Fei et al., 2025; Chen et al., 2025a; Zhang et al., 2024b; Chari et al., 2025).

M Datasets

We evaluate GDWM on two complementary long-context benchmarks: **ZeroSCROLLS**—a *zero-shot* suite adapted from SCROLLS with reliable, task-specific automatic metrics—and **LongBench v2**—a *realistic* long-context benchmark using multiple-choice questions for robust evaluation.

M.1 ZeroSCROLLS

Benchmark overview. ZeroSCROLLS is a zero-shot benchmark for long-text understanding that contains *no training split* and only small validation

Task	Type	Metric	Avg #Words
GovReport	Summarization	ROUGE	7,273
QMSum	QB-Summ	ROUGE	10,839
Qasper	QA	F1	3,531
NarrativeQA	QA	F1	49,384
QuALITY (Quality)	MC-QA	Accuracy	4,248
MuSiQue	QA (Multi-hop)	F1	1,749

Table 5: ZeroSCROLLS tasks used in this work and their official metrics/statistics (Shaham et al., 2023).

sets, with each task capped at 500 examples to keep evaluation affordable. It extends SCROLLS by adapting six long-document tasks and adding four new tasks, covering summarization, question answering, and information aggregation.

Tasks used in this work. Following prior long-context evaluation practice, we select six representative ZeroSCROLLS tasks spanning both *sparse-evidence* and *dense-coverage* regimes: (i) **GovReport** and **QMSum** for (query-based) summarization, (ii) **Qasper** and **NarrativeQA** for long-context QA, (iii) **QuALITY** (denoted as **Quality** in our paper) for multiple-choice comprehension, and (iv) **MuSiQue** for multi-hop QA.

Evaluation metrics. ZeroSCROLLS uses task-aligned automatic metrics: ROUGE for summarization, F1 for extractive/abstractive QA-style tasks, and Accuracy for multiple-choice QA. We report the official metrics with the benchmark’s evaluation scripts.

M.2 LongBench v2

Benchmark overview. LongBench v2 is a realistic long-context benchmark designed to test *deep understanding and reasoning* over long inputs. It contains **503** challenging **multiple-choice** questions with contexts ranging from **8k to 2M words**, organized into **6** major categories and **20** subtasks. All examples are in English, and each instance includes a long context, a question, four options, a ground-truth answer, and annotated evidence for verification. The benchmark emphasizes evaluation *reliability* by using accuracy-based scoring (rather than free-form generation metrics).

Task categories and statistics. Table 6 summarizes the six categories and their median context lengths. Notably, LongBench v2 includes long-dialogue history understanding and code-repository understanding, which directly stress *memory* and *retrieval under long context*.

Category	#Questions	Median #Words
Single-Document QA	175	51k
Multi-Document QA	125	34k
Long In-context Learning	81	71k
Long-dialogue History Understanding	39	25k
Code Repository Understanding	50	167k
Long Structured Data Understanding	33	49k

Table 6: LongBench v2 category-level statistics reported in Bai et al. (2025).

Evaluation protocol. We follow the official LongBench v2 evaluation protocol and report **accuracy** (overall and category-wise). For consistency across settings, all methods are evaluated under the same context-length budget of the underlying base model; when an instance exceeds the model limit, we follow the benchmark-recommended pre-processing/truncation behavior (Bai et al., 2025).

2014
2015
2016
2017
2018
2019
2020
2021