

# ASSESSING EPISODIC MEMORY IN LLMs WITH SEQUENCE ORDER RECALL TASKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Current LLM benchmarks focus on evaluating models’ memory of facts and semantic relations, primarily assessing semantic aspects of long-term memory. However, in humans, long-term memory also includes episodic memory, which links memories to their contexts, such as the time and place they occurred. The ability to contextualize memories is crucial for many cognitive tasks and everyday functions. Existing benchmarks have poor coverage of episodic memory. To address the gap in evaluating memory in LLMs, we define episodic memory for LLMs and introduce Sequence Order Recall Tasks (SORT), which we adapt from tasks used in cognitive psychology. SORT requires *causal* LLMs to recall the correct order of text segments, and provides a general framework that is both easily extendable and does not require any additional annotations. We present an initial evaluation dataset, Book-SORT, comprising 36k pairs of segments extracted from 9 books recently added to the public domain. Based on a human experiment with 155 participants, we show that humans can recall sequence order based on long-term memory of a book. We find that models can perform the task with high accuracy when relevant text is given in-context during the SORT evaluation. However, when presented with the book text only during training, LLMs’ performance on SORT falls short. By evaluating a new aspect of memory, we believe that SORT will aid in the emerging development of memory-augmented models.

## 1 INTRODUCTION

Large language models (LLMs) have impressive performance on many benchmarks that test factual or semantic knowledge learned during training or in-context (Hendrycks et al., 2020; Ryo et al., 2023; Logan IV et al., 2019; Petroni et al., 2019; Yu et al., 2023; Sun et al., 2023). While these advances are noteworthy, the type of long-term knowledge that these datasets test is only one of several types that naturally intelligent systems store, retrieve, and update continuously over time (Norris, 2017; Izquierdo et al., 1999; McClelland et al., 1995). Current evaluation tasks do not assess episodic memory, which is a form of long-term knowledge thought to be important for cognitive function in humans and animals. Below we propose a definition of episodic memory in LLMs, which is based on the human literature (see Appendix A for a discussion).

### Definition 1.1: Episodic Memory in LLMs

Episodic memory refers to knowledge in a language model that:

- (1) is specific and unique to a particular sequence;
- (2) is acquired through a single exposure to that sequence (single-shot learning);
- (3) contains information about relations between parts (e.g. encountered events/items, incl. more abstract items) within that sequence;
- (4) can still be retrieved when arbitrarily many tokens are processed in between encoding and retrieval;
- (5) has functional implications, meaning the knowledge can be used by the model to answer explicit queries.

054 In contrast to semantic memory, episodic memory links memories to their contexts, such as the  
055 time and place they occurred. Research on human memory also originally focused on semantic,  
056 rather than episodic memory – however, researchers realized that one could distinguish the ‘what’  
057 (semantic) content from the ‘where’ (spatial context) and ‘when’ (temporal context) Tulving (2002).  
058 This ability to organize memory based on spatial and temporal details enables us to reconstruct events  
059 that occurred in the possibly distant past, predict the future, and relate information across multiple  
060 events that are separated by time windows spanning a lifetime, capabilities crucial for many cognitive  
061 tasks and everyday functions. We propose SORT as a first benchmark to assess an important aspect  
062 of episodic memory.

063 The ability to link contextual details to stored information—particularly, details about temporal context—  
064 may be key to improving LLM performance on several tasks. More human-like episodic memory  
065 may improve models’ continual learning and adaptation to shifting data distributions, performance  
066 on tasks requiring long contexts (e.g., long chat exchanges with a user), and source attribution via  
067 knowledge of where and when a memory was acquired, which could help to reduce or identify  
068 hallucinations.

069 To address the gap in evaluating crucial attributes of memory in causal LLMs, we propose the  
070 Sequence Order Recall Task (SORT), which we adapt from tasks in cognitive psychology that are  
071 used to assess long-term episodic memory in humans and animals (Eichenbaum, 2013; Davachi &  
072 DuBrow, 2015). Specifically, SORT requires a model to recall the correct order of sequential data,  
073 such as segments of text. We hope that SORT will be the first of many benchmarks that assess various  
074 aspects of episodic memory in LLMs.

075 We provide a specific instantiation of SORT that requires causal language models to recall the correct  
076 order of two segments sampled from text, along with a corresponding evaluation dataset—Book-SORT.  
077 Book-SORT contains over 36k pairs of text segments from 9 books, with variations in segment length  
078 (20 and 50 words) and distance between segments (up to 16k words). We chose books that were very  
079 recently released from U.S. copyright to minimize the possibility that LLMs were pre-trained on these  
080 texts. This allowed us to test three common methods of giving a causal language model access to a  
081 specific text: (1) during inference in-context, (2) during inference via retrieval augmented generation  
082 (RAG), and (3) during training via fine-tuning with a language modeling objective. Furthermore,  
083 we provide a human evaluation from 155 participants who had finished reading a whole book and  
084 were tested with no additional access to the book, showing that humans can recall segment order with  
085 up to 70% accuracy based on their long-term memory. While the ceiling performance on SORT is  
086 100% (assuming that texts do not contain duplicate segments), our human data provides an important  
087 reference point to compare and contrast long-term memory across models and humans.

088 When given access to excerpts from the books in-context, we find that models achieve up to 95%  
089 accuracy with relevant 250-word excerpts but degrade quickly as longer excerpts are presented. Using  
090 Retrieval Augmented Generation, models can recall sequence order with limited performance. Finally,  
091 models fine-tuned with a language modeling objective on the book texts do not significantly improve  
092 their SORT performance, showing that parametric memory in current transformer models supports  
093 semantic but not episodic long-term memory.

094 Our main contributions can be summarized as follows:

- 095 • definition of episodic memory in the context of LLMs
- 096 • proposal of the self-supervised task SORT, which requires LLMs to recall the correct order  
097 of segments from a sequence and can be used to assess capabilities in causal LLMs that  
098 would be supported by episodic memory in humans
- 099 • a new dataset Book-SORT comprised of 36k samples from 9 public domain books and an  
100 evaluation framework that is easily extendable to new datasets
- 101 • first-of-its-kind human evaluation ( $N = 155$ ) showing that humans are capable of recalling  
102 the order of text from an entire book based on long-term memory
- 103 • a comprehensive evaluation of open-source and closed language models on Book-SORT,  
104 showing that current models: i) have good in-context memory performance, when all  
105 necessary information is presented in the prompt and the prompt is short; ii) quickly lose  
106 the ability to recall sequence order as the excerpt provided in-context gets longer, though  
107 still far below their advertised context-lengths; iii) fail to recall segment order based on

parametric memory formed via fine-tuning with a language modeling objective; (iv) perform worse on SORT with retrieval augmented memory than with in-context memory.

## 2 RELATED WORK

**Evaluation of parametric semantic memory in LLMs.** Benchmarks such as MMLU (Hendrycks et al., 2020), T-REx (Elsahar et al., 2018), LAMA (Petroni et al., 2019), WICE (Ryo et al., 2023), KoLA (Yu et al., 2023), and others (Sun et al., 2023) test models’ retrieval and reasoning ability on different domains, such as recalling a chemistry fact.

Other benchmarks that partially evaluate LLM semantic memory are those that require reasoning using temporal (Ning et al., 2020; Zhou et al., 2021; Feng et al., 2023) (e.g. lunch happens before dinner), causal (Srivastava et al., 2023) (e.g. she is eating, therefore she is hungry), or other commonsense knowledge (e.g. food is edible) (Ismayilzada et al., 2023) acquired during pretraining. In contrast to these benchmarks, our work proposes a task that involves judgments regarding temporal context information about text segments that either (a) are available through in-context memory or (b) were otherwise previously presented to the model, e.g. via fine-tuning or Retrieval Augmented Generation.

**Evaluation of in-context memory in LLMs.** We evaluate in-context memory, in which the model has in-context access to all relevant text for the task. This relates to works that evaluate a model’s ability to retrieve information from its context input, such as Needle In A Haystack (Kamradt, 2023) and FLenQA (Levy et al., 2024). These requirement 3 in our definition only minimally by testing the ability to retrieve an atomic piece of information regardless of its context.

Previous datasets and benchmarks that evaluate performance over long context lengths, such as Long Range Arena (Tay et al., 2021), SCROLLS (Shaham et al., 2022), and MULD (Hudson & Al Moubayed, 2022), are also relevant. The evaluation of in-context memory with SORT differs from these works by focusing on order information, which is key to episodic memory in humans. In NarrativeXL (Moskvichev & Mai, 2023) and NarrativeQA (Kočíský et al., 2017), models have to perform reading comprehension and free recall tasks when given entire books. SORT is different in that it places the focus on memory of segments where the complete context (all other segments) matters for the evaluation. This is not always the case in reading comprehension tasks where questions can be about atomic parts of the context that need to be retrieved.

**Tasks related to SORT.** Previously proposed tasks that most closely relate to SORT are BART’s denoising training objective (Lewis et al., 2020), which permutes the order of sentences in a document and learns to reconstruct the correct order, and BERT’s next sentence prediction objective (Devlin et al., 2019), which learns to predict whether two sentences follow each other in a text. SORT differs from these tasks, as it is not intended as a training objective, and it can include text segments with an arbitrary distance between each other in a document, possibly exceeding the context input length of the model. In ChapterBreak (Sun et al., 2022), long segments ending at a chapter boundary taken from a book are presented to an LLM along with multiple segments of chapter beginnings from the same book. The task for the LLM is then to tell which one is the directly following chapter and which are not. This suffix-identification task aims to evaluate narrative-understanding based reasoning about books, while we propose SORT as an evaluation for episodic memory in LLMs, involving both a model and a memory-insertion method. By evaluating a SORT baseline in which the models do not have access to relevant source texts, we show that memory is needed for SORT and general narrative-reasoning ability is not enough.

## 3 SEQUENCE ORDER RECALL TASK

We introduce a novel evaluation task: recalling the order of parts of a sequence, which we term the Sequence Order Recall Task (SORT). SORT is adapted from recency judgment tasks used in cognitive psychology to evaluate episodic memory in humans and animals (Eichenbaum, 2013; Davachi & DuBrow, 2015). In this task, a sequence is presented to a participant. Then, after some delay, the participant is asked to judge the order in which two segments of the sequence appeared. We adapt this task to test memory in models. The general task can be applied to any sequential domain, including video and audio. Here we focus on the text domain to evaluate LLMs (Fig. 1).

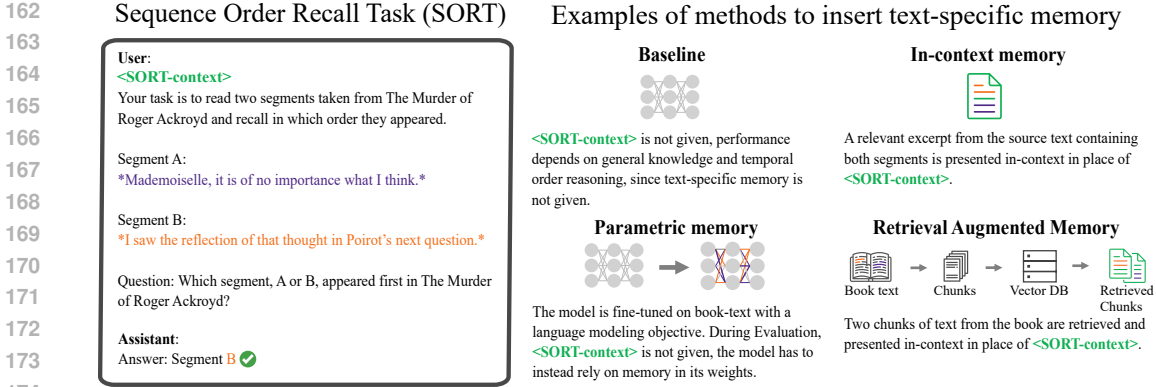


Figure 1: Overview of the Sequence Order Recall Task (SORT) to evaluate how models can access memory of temporal order. Left: Example task prompt for SORT. A prefix to the prompt can be given to assess in-context forms of memory. Right: Examples of methods to insert memory of specific texts into a model.

**Formal description of SORT.** The general form of the task can be described as follows. Let  $\mathbf{X} \in \mathbb{R}^{T \times F}$  be sequential data, where  $T$  is the number of time-steps (e.g. token in a text) and  $F$  is the number of features (e.g. vocabulary size). We define start indices  $t_j$  and  $t_k$  for pairs of segments of length  $L \in \mathbb{N}^+$  in  $\mathbf{X}$ , such that both  $t_j < t_k$  and  $t_j + L \leq t_k$ . Using these, we extract non-overlapping segments from the original sequence  $\mathbf{X}$  as  $\tilde{\mathbf{X}}_i = \mathbf{X}[t_i : t_i + L - 1, :]$ . The order of segments  $\tilde{\mathbf{X}}_j$  and  $\tilde{\mathbf{X}}_k$  is randomized, yielding  $[\tilde{\mathbf{X}}_A \tilde{\mathbf{X}}_B]$ , which is then given as part of a model’s input. The task for a model  $\mathcal{M}_\theta$  is to infer whether  $t_A < t_B$ , i.e. in SORT, the task of a model is to predict which of two non-overlapping subsequences  $\tilde{\mathbf{X}}_A$  and  $\tilde{\mathbf{X}}_B$  has the lower starting index in  $\mathbf{X}$ . The task can be used to evaluate a variety of methods to include document-specific memory in models. To assess in-context memory, i.e. memory based on text presented in-context, the segments are preceded by  $\mathbf{X}$  in the model’s input. When assessing retrieval-augmented generation methods, instead of prepending  $\mathbf{X}$ , segments of  $\mathbf{X}$  are retrieved and prepended. For the assessment of parametric long-term memory,  $\tilde{\mathbf{X}}$  is not part of a model’s input, instead the model’s parameters (or a subset thereof)  $\theta$  are a function of  $\mathbf{X}$  via pre-training or fine-tuning:  $\theta = f(\mathbf{X})$ .

The general form of SORT is the following input, which can be preceded by additional context to insert a memory:

$$I_{SORT} = [P_{context} P_{task} P_{label_A} \tilde{\mathbf{X}}_A P_{label_B} \tilde{\mathbf{X}}_B P_{question} P_{answer}], \tag{1}$$

where  $P_{context}$  can either be relevant context, such as (parts of) the source sequence  $\mathbf{X}$  to assess in-context memory (stored in activation slots), or an empty string when parametric memory (stored in weights) is assessed;  $P_{task}$  instructs the model for the sequence order recall task to read two segments and describes the objective: answering which of the two labeled segments appears first in  $\tilde{\mathbf{X}}$ ;  $P_{label_A}$  and  $P_{label_B}$  are the labels (e.g. the characters “A” and “B”) for the first and second segment presented in the task  $\tilde{\mathbf{X}}_A$  and  $\tilde{\mathbf{X}}_B$ ;  $P_{question}$  repeats the SORT objective as a question; finally,  $P_{answer}$  provides the beginning of the answer string as “Answer: Segment”.

### 3.1 EVALUATING LARGE LANGUAGE MODELS ON SORT

We greedily sample an answer token  $\mathbf{a} = \text{argmax}(\mathcal{M}_\theta(\mathbf{I}))$  from the model  $\mathcal{M}_\theta$ , which is parameterized by  $\theta$ , and decode the sampled answer token  $\mathbf{a}$  as either “A” or “B”.

The answer is evaluated as correct if it corresponds to the segment that truly appears first in  $\mathbf{X}$ . For proprietary (OpenAI) models that do not allow completing assistant responses with prepended text, we omit  $P_{answer}$ . In this case we resort to generating a sequence of 25 tokens, and parse the generated text for A or B responses.

**Prompt selection.** Using a single prompt formulation across all models may bias the results. To prevent this, we compiled a set of 12 prompts that vary formulations in  $P_{context}$  and  $P_{task}$ . For each

model, we evaluate each prompt on a held-out dataset of 400 samples and used the best performing prompt for each model. The full prompts and further details on prompt selection are given in Appendix C.2-C.3.

**Baseline without book-specific memory.** We want to ensure that performance on SORT is due to text-specific memory and not due to temporal order reasoning supported by more semantic forms of memory such as commonsense knowledge (e.g. lunch happens before dinner). We isolate the effects on SORT that are due to text-specific memory by contrasting performance between a baseline model that does not have access to the specific text and a model that has access to the sequences in one of various ways in which memory can be inserted.

### 3.2 INSERTING TEXT-SPECIFIC MEMORY INTO MODELS

We evaluate three examples of methods to insert text-specific memory into models: (1) via in-context presentation, (2) via fine-tuning with a language modeling objective, and (3) via retrieval augmented generation of short chunks of text in a book.

**In-context presentation.** When assessing in-context memory,  $\mathbf{P}_{\text{context}}$  in Eq. 1 contains relevant excerpts from the source text along with the book title. The prompt includes the instruction to carefully read the text from the book (a list of used prompts is shown in Appendix 6). To test in-context memory, We make sure that excerpts contain both segments and vary the length of excerpts in our experiments.

**Finetuning with a language modeling objective.** Instead of presenting text from the books in the same prompt in which the SORT task is given, we are interested in parametric memory of the texts. In this condition,  $\mathbf{P}_{\text{context}}$  in Eq. 1 is an empty string. To insert parametric memory of the source texts into a model, we fine-tune the model with a next-token prediction objective on the books, split into chunks of 5000 words and contextualized by the books’ titles. Since we need to preserve the models’ ability to understand and follow the task instructions (Allen-Zhu & Li, 2024), we fine-tune on a dataset which additionally includes 3,500 random instruction-following examples that are unrelated to SORT. This helps to prevent catastrophic forgetting during continued finetuning (Luo et al., 2024). We finetune on 8 A100 GPUs with an initial learning rate of  $5e-6$  and a batch size of 192. Full details of the fine-tuning setup are given in Appendix F and our code will be available. SORT should be informative about episodic memory more broadly, we did not train models on SORT (see Appendix F.4).

**Retrieval Augmented Generation.** To include memory of text via retrieval augmented generation (RAG), we built a typical naive RAG pipeline that relies on two separately pretrained models for the retriever and the reader (Gao et al., 2024). The retriever returns text passages from a database to serve as task context for the LLM (i.e. as  $\mathbf{P}_{\text{context}}$ , Eq. 1).

The retrieval database contained text embeddings of all passages from Book-SORT (Sec. 4). We used the LangChain recursive text splitter to chunk Book-SORT text into  $\sim 1024$  character, non-overlapping passages (average 183 words). Each passage was then encoded into a 1024-d vector using a high-performing, open-source text retrieval model (BGE-v1.5, (Xiao et al., 2024)). To retrieve the passages, we conduct an exact nearest neighbor search. The search returns the  $k = 2$  nearest neighbors. We maintained this similarity order when inserting the retrieved passages into the prompt, i.e. the most similar passage appears first in  $\mathbf{P}_{\text{context}}$ .

## 4 BOOK-SORT DATASET AND EVALUATION

We created an English language dataset to evaluate episodic memory in humans and LLMs. The selected sequence data considered several factors: (1) we chose long texts (mean length = 72,700 words) that exceed the context windows of most transformer LLMs; (2) we used books to enhance memorability for human readers and facilitate our human evaluation experiment; (3) we selected books from *Project Gutenberg* that recently entered the U.S. public domain to avoid ethical and copyright issues, and minimize pre-training contamination in LLMs. Within these constraints, we aimed to maximize content diversity, including narrative fiction novels, a physics text, and an extended essay. Further details on the 9 books in the Book-SORT dataset are available in Appendix B.1.

#### 270 4.1 BOOK-SORT CREATION

271  
272 We constructed a dataset that varies across factors that can affect human or model performance  
273 on SORT. Based on prior reports on LLMs (Liu et al., 2024), we first varied (1)  $L_E$ , the length of  
274 the text excerpt presented in context. Since the typical standard context length of the LLMs in our  
275 study was 4096 tokens, we set  $L_E = \{250, 1000, 2500\}$  words. For models with extended context  
276 windows, we also created datasets where  $L_E = \{10000, 20000\}$  words, which excluded one book  
277 that was too short. Our pilot experiments on humans suggested two other factors that would affect  
278 task performance: (2)  $L_S$ , the length of the segments from the text, and (3)  $D_S$ , the distance between  
279 the segments in the original text. To mirror the human experiments, we set  $L_S = \{20, 50\}$  words.  
280 We then created 4 different distance bins  $D_S = \{d_0, d_1, d_2, d_3\}$ , whose values were bounded by the  
281 excerpt length  $L_E$  (Appendix Table 4).

282 Within each unique combination of the first two factors  $L_E$  and  $L_S$ , we randomly sampled 110  
283 excerpts from each of the 9 books (i.e. 100 samples for SORT evaluation, and 10 samples for  
284 prompt selection per book). All excerpts and segments began at a sentence boundary. Within each  
285 combination of  $L_E, L_S$ , we randomly sampled 4 different segment pairs, one from each distance bin  
286  $D_S$ . This minimized the possibility that observing an effect of distance on SORT performance would  
287 be due to differences in the semantic content of the text segments. Finally, for all 110 trials within  
288 each of these 3 factors, we counterbalanced the correct answer. This yielded a well-controlled and  
289 easily extendable dataset of about 36K text segment pairs for SORT evaluation.

#### 290 4.2 HUMAN LONG-TERM MEMORY EVALUATION

291  
292 As a reference point (but not a performance ceiling), we further provide a human evaluation from  
293 155 participants who had recently finished reading one of the 9 books in the Book-SORT dataset,  
294 *The Murder of Roger Ackroyd* (Christie, 1927). This evaluation assessed long-term memory, as the  
295 average time between reading and testing was 7.5 days, far surpassing short-term memory duration  
296 (Hasson et al., 2015). There is no previously reported data on long-term memory for entire books from  
297 large samples, so we designed an experiment to collect this data. Given the difficulty of recruiting  
298 participants to read lengthy books specifically for an experiment, we used a creative recruiting  
299 strategy: inviting members of the online reading community *Goodreads* who had recently finished  
300 *The Murder of Roger Ackroyd*. Participants completed an online survey within 30 days of finishing  
301 the book. The expected compensation for participation was \$12 and the study was approved by  
302 the IRB at Anonymized University. We provide 1570 segment pair samples from 155 participants.  
303 Further details about this one-of-a-kind study are provided in Appendix B.3.

#### 304 4.3 MODELS

305  
306 We evaluate a selection of open models covering a broad range of scores on popular benchmarks  
307 such as MMLU (see Table 5) ranging from 7b to 8x22b parameter transformer models. Initial  
308 experiments with non-instruction-tuned models resulted in chance performance on Book-SORT (see  
309 Appendix E), which we attribute to the lack of instruction tuning<sup>1</sup>, and thus focus on evaluating  
310 instruction-tuned models in this work. We have selected models from different model families  
311 including Llama3 (AI@Meta, 2024), Llama2 (Touvron et al., 2023), Mistral (Jiang et al., 2023),  
312 Mixtral (Jiang et al., 2024), Gemma (Team et al., 2024) and OpenAI GPTs (Achiam et al., 2023). For  
313 our experiments on finetuning as a method for inserting memory into models, we focus on two models  
314 Mistral-v0.2-7b-Instruct and Llama3-8b-Instruct because they allow full-parameter fine-tuning with  
315 8 A100 GPUs.

### 316 5 RESULTS

317  
318 We present empirical findings for a baseline without text-specific memory of the books in Book-SORT,  
319 as well as three methods to include memory, using 9 open-source models and 2 closed language  
320 models.  
321  
322

323 <sup>1</sup>(Zhang et al., 2024b) provides an overview of instruction tuning approaches

Table 1: Baseline: SORT performance before models are exposed to the books in Book-SORT.

|                       | Segment length 20 | Segment length 50 |
|-----------------------|-------------------|-------------------|
| Llama3-70b-inst       | $0.52 \pm 0.007$  | $0.54 \pm 0.007$  |
| Llama3-8b-inst        | $0.51 \pm 0.008$  | $0.52 \pm 0.007$  |
| Mixtral-8x22b-inst    | $0.52 \pm 0.007$  | $0.55 \pm 0.007$  |
| Mixtral-8x7b-DPO-inst | $0.52 \pm 0.008$  | $0.54 \pm 0.008$  |
| Llama2-70b-inst       | $0.51 \pm 0.007$  | $0.51 \pm 0.008$  |
| Gemma-1.1-7b-inst     | $0.51 \pm 0.008$  | $0.51 \pm 0.007$  |
| Mistral-v0.2-7b-inst  | $0.51 \pm 0.007$  | $0.51 \pm 0.008$  |
| Mistral-v0.1-7b-inst  | $0.50 \pm 0.008$  | $0.50 \pm 0.008$  |
| Llama2-7b-inst        | $0.50 \pm 0.008$  | $0.49 \pm 0.008$  |
| GPT-3.5-turbo         | $0.52 \pm 0.009$  | $0.52 \pm 0.012$  |
| GPT-4                 | $0.53 \pm 0.008$  | $0.57 \pm 0.007$  |

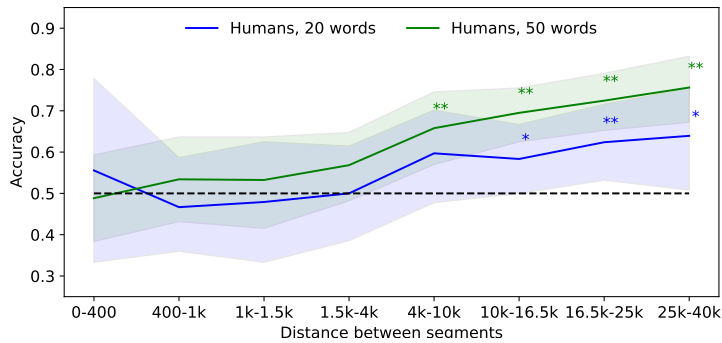


Figure 2: Human long-term memory performance on SORT for different segment lengths and distances between segments. Shaded areas depict bootstrapped 95% confidence intervals. Significant difference from chance is marked with asterisks (\* p-value<0.05, \*\* p-value<0.01).

## 5.1 BASELINE

**SORT requires memory specific to books in Book-SORT.** To validate that it is not possible to achieve high performance on Book-SORT without memory of the specific books that are included in the dataset, we evaluate models before they have access to the books. This shows that SORT requires memory of particular books and cannot be easily solved via temporal order reasoning (Hendrycks et al., 2020). We find that segment pairs with a very short and with a very long distance in the book allow a higher performance than chance (see Appendix D.1), indicating that some of these segment pairs can be ordered based not on memory but on reasoning or common-sense. However, none of the models have a high performance for any segment distance bin, as would be expected if SORT requires book-specific memory.

## 5.2 HUMAN EXPERIMENT

**Humans can perform in SORT based on long-term memory.** The results from human long-term memory (LTM) experiments, depicted in Figure 2, demonstrate that humans can perform in SORT based on long-term memory. The average accuracy is 0.64 for segments of 50 words and 0.56 for segments of 20 words). Human performance is higher for pairs of segments that have a greater distance in the book, with a peak accuracy of 0.76 for distances greater than 25,000 words and 50-word segments. Binomial tests show that beyond a distance of 4000 words, humans perform statistically significantly better than chance. Note that we present these results as evidence that one possible information processing system—a human—can perform SORT based on long-term memory. Importantly, these results do not present the ceiling performance on the memory task that we propose. The expected ceiling performance on SORT is 100%, given that the books do not contain duplicated segments of text, which is less probable for longer segment lengths.

Table 2: Mean of in-context memory performance with 95% bootstrapped confidence interval. SORT-extend shows performance with excerpts of lengths 10000 and 20000 words, which exceeds most models’ context lengths.

| Model name            | Parameters | Max context | SORT             | SORT-extend      |
|-----------------------|------------|-------------|------------------|------------------|
| Llama3-70b-inst       | 70b        | 8k          | $0.92 \pm 0.020$ | /                |
| Llama3-8b-inst        | 8b         | 8k          | $0.93 \pm 0.007$ | /                |
| Mixtral-8x22b-inst    | 8x22b      | 64k         | $0.95 \pm 0.020$ | $0.79 \pm 0.038$ |
| Mixtral-8x7b-DPO-inst | 8x7b       | 32k         | $0.89 \pm 0.030$ | $0.56 \pm 0.058$ |
| Llama2-70b-inst       | 70b        | 8k          | $0.77 \pm 0.040$ | /                |
| Gemma-1.1-7b-inst     | 7b         | 8k          | $0.85 \pm 0.010$ | /                |
| Mistral-v0.2-7b-inst  | 7b         | 32k         | $0.85 \pm 0.032$ | $0.65 \pm 0.045$ |
| Mistral-v0.1-7b-inst  | 7b         | 8k          | $0.77 \pm 0.013$ | /                |
| Llama2-7b-inst        | 7b         | 4k          | $0.56 \pm 0.014$ | /                |
| GPT-3.5-turbo         | unknown    | 16k         | $0.86 \pm 0.010$ | /                |

### 5.3 IN-CONTEXT MEMORY

**Models generally perform well on SORT based on in-context memory.** Nearly all models achieve above 77% accuracy when given in-context access to relevant excerpts from the books, reaching up to 95% (Table 2). This indicates that very large models are not necessary to perform this task effectively, as demonstrated by the Llama3-8b model outperforming larger models such as Llama3-70b and Mixtral-8x7b-DPO.

**In-context memory performance increases with greater distance between segments.** We further evaluate the effect of another factor which may influence the model performance—the distance between the text segments in the excerpt. Figure 3b shows an increasing trend in accuracy as the distance between segments increases. This improvement in accuracy is consistent across excerpt lengths and is observed across all models (see Appendix D.2).

**In-context memory performance decreases with increasing excerpt length.** Average performance on longer excerpts (Table 2, SORT-extend) is substantially lower than in the standard context lengths, despite the presence of longer segment distances. For increasing excerpt lengths, we see a consistently monotonic decrease in average accuracy (Figures 3a and 3).

**Additional analyses.** Further analyses are presented in Appendix D.2. Models handle longer segments (50 words) slightly more effectively than shorter segments (20 words), with an improvement of up to 4%. We found no significant differences across books from different domains (Table 11-12).

### 5.4 PARAMETRIC MEMORY VIA FINETUNING

**Full parameter fine-tuning on books with a language modeling objective did not improve SORT performance.** For Llama3-8b-Instruct and Mistral-7b-v0.2-Instruct, we do not observe any difference in performance on SORT after memory is inserted via fine-tuning on large chunks of book-text. A pairwise statistical analysis across epochs of fine-tuning, relative to two baselines that either exclude the books from the fine-tuning dataset or instead include only summaries of the books, shows no substantial improvement (see Appendix F).

### 5.5 RETRIEVAL AUGMENTED MEMORY

**RAG based memory leads to worse performance than in-context memory.** Due to the fact that the order of multiple passages from the same document is not preserved in a standard RAG setting, the performance is lower than in in-context memory and does not reach 70% accuracy for any distance between segments (Figure 4a). Curiously, we find that bigger models (i.e. Mixtral-8x22b-Instruct and Llama3-70b-Instruct) do not substantially outperform smaller models with RAG. Even when passages containing both segments are retrieved and presented in the correct order, we find that Llama3-8b-Instruct outperforms two much larger models on SORT (Figure 4b).



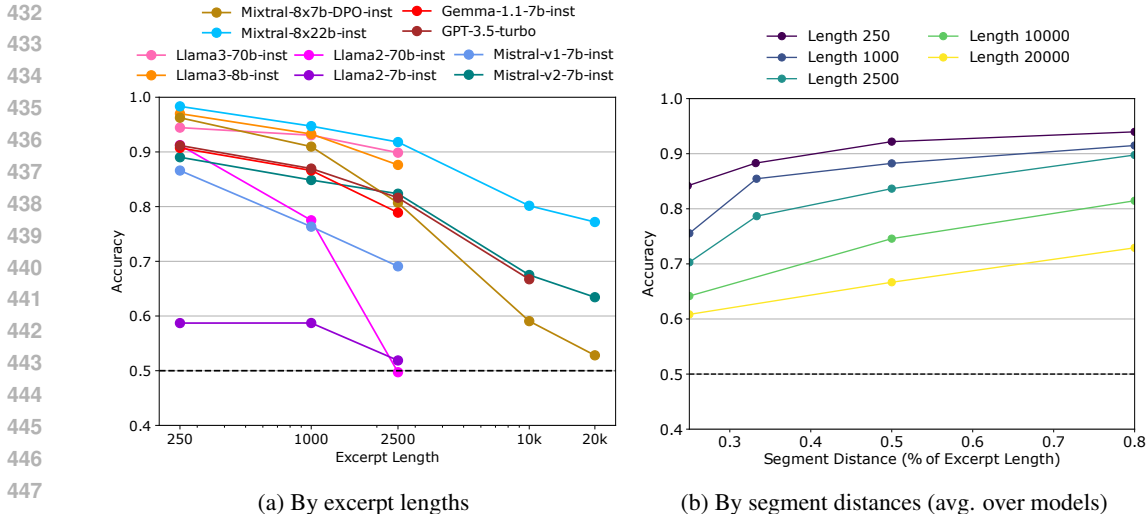


Figure 3: Factors affecting SORT performance based on in-context memory. (a) SORT accuracy by excerpt length. (b) Average over SORT performance of different models across segment distances for different excerpt lengths.

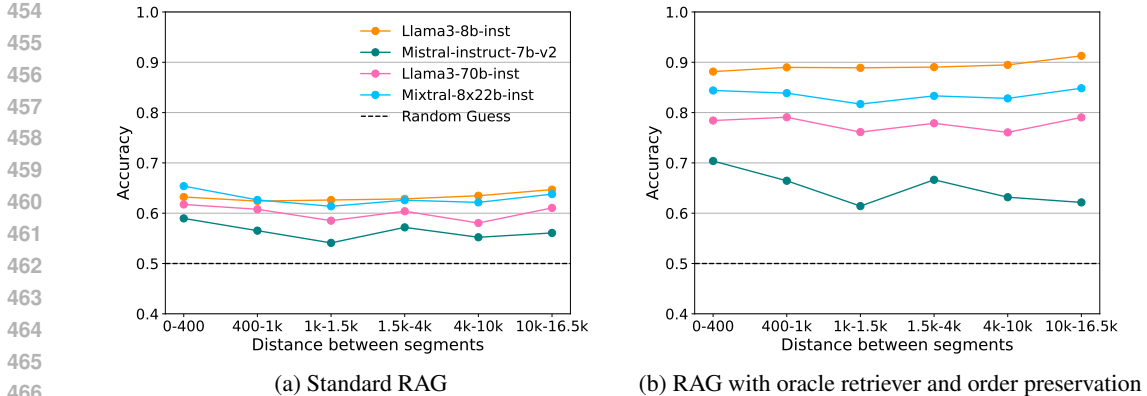


Figure 4: SORT performance based on RAG memory. (a) Accuracy with standard RAG memory. (b) Accuracy with RAG memory given that the correct passages of text are retrieved and presented in the order in which they appeared in the books.

## 6 DISCUSSION

We provide a new evaluation task, SORT, for assessing episodic memory in causal large language models, that can be used with any text data and without the need for annotation. We created Book-SORT, a dataset for SORT based on books that were recently added to the public domain and we validated that book-specific memory is indeed needed to achieve high performance on Book-SORT. We evaluated three different ways to include memory of specific texts in a model to assess whether they support a key function of episodic memory. Below, we discuss our results for these methods in relation to episodic memory in humans.

**Is in-context memory a form of episodic memory?** Several links have been drawn between in-context memory in transformers and multiple models of episodic memory in humans (Ji-An et al., 2024; Whittington et al., 2022; 2024; Ellwood, 2024) and our results suggest that it does support sequence order recall. However, our in-context memory results suggest that performance degrades for sequence order recall, as it does for other tasks (Liu et al., 2024; Levy et al., 2024). We believe that the result of decreasing performance with more context supports a view of in-context memory as

486 an extended form of working memory. The two key problems with in-context memory that make it  
487 unlike long-term episodic memory are that it does not generalize well to arbitrarily long sequences  
488 and its cost increases as the context gets longer. In terms of our definition of episodic memory  
489 in LLMs (1), in-context memory disqualifies as supporting episodic memory because the fourth  
490 requirement is not fulfilled for current models.

491 **Is parametric memory in transformers a form of episodic memory?** High performance on  
492 benchmarks including MMLU suggests that parametric memory in LLMs learned via a language  
493 modeling objective can support semantic forms of memory (e.g. when recalling knowledge to answer  
494 factual questions). Our evaluation on SORT showing close to chance performance suggests that  
495 current forms of parametric memory insertion might not support functions similar to those of episodic  
496 memory.

497 **Is retrieval augmented memory a form of episodic memory?** Since it avoids the problems of  
498 context-length generalization, Retrieval Augmented Generation presents a potentially strong way to  
499 include memory of episodes via a retrieval process and subsequent in-context presentation. However,  
500 our results suggest that there is a lot of room for improvement over the performance of vanilla RAG.  
501 However, in vanilla RAG, retrieved segments are presented without surrounding context information  
502 (since chunks of the same document are independent). Order-preserving (OP) RAG (Yu et al.,  
503 2024) presents one way to retain relative positional information between retrieved passages as one  
504 kind of temporal context and can thereby increase performance on SORT. The episodic memory  
505 system in animals does not only bind temporal order information to memories but its context-binding  
506 generalizes to more abstract types of context (Eichenbaum, 2015a; Qiu et al., 2024) that would not be  
507 given in OP-RAG since memories are encoded independently of each other (i.e. the third criterion in  
508 our definition (1) is not properly fulfilled).

509 **Limitations.** Current high performing causal LLMs do not disclose their training data, which means  
510 that care needs to be taken in selecting suitable data to include in a SORT dataset. To minimize  
511 the probability that models have been trained on books used for our SORT evaluation, we curated  
512 Book-SORT based on books that were not publicly available when models were trained. However we  
513 cannot rule out that no copyrighted material was used in training of a model, which would require us  
514 to interpret results as indicating the effectiveness of additional rather than initial memory-insertion.  
515 Furthermore the reliance on instruction-following can limit the applicability to both non-instruction-  
516 tuned models and models that have poor instruction-following ability. While we provide a few  
517 examples of memory-insertion methods, we leave more extensive studies on how to induce episodic  
518 memories without relying on complete in-context presentation to future work.

519 **Future work.** Improving long-term memory in LLMs is an emerging area of research (Liu et al.,  
520 2023; Borgeaud et al., 2022; Fournier et al., 2023; Phang et al., 2023; Wang et al., 2024; Zhong  
521 et al., 2022; 2024), and SORT can be used to assess improvement in an crucial aspect of an important  
522 form of memory in new models. Specifically, improving episodic memory in models may improve  
523 models' continual learning, performance on tasks at long contexts such as extended chat exchanges  
524 with a user, and source attribution via knowledge of where and when a memory was acquired. Recent  
525 efforts have highlighted the potential of augmenting causal LLMs with additional episodic memory  
526 mechanisms (Fountas et al., 2024; Das et al., 2024), and we expect that SORT can be used to evaluate  
527 these classes of models, once such a model with a sufficiently strong instruction-following ability is  
528 released. Another possibility is to identify new and better methods to insert episodic memory of texts  
529 into existing models. Additionally, SORT can be extended to other types of inputs, such as audio and  
530 video, which can be used to evaluate episodic memory in multimodal models in the future.

531 **Conclusion.** The ability of LLMs to retain and retrieve long-term knowledge is crucial for their  
532 continued integration in many applications. Therefore, a more comprehensive and systematic  
533 evaluation of these abilities is needed. We believe that the new evaluation framework SORT offers a  
534 promising path for future research aimed at better understanding and improving these capabilities in  
535 foundation models.

536 **Ethics Statement.** To avoid ethical issues concerning copyright, we based Book-SORT on books  
537 that were recently added to the public domain. Our human experiment with 155 participants was  
538 approved by the IRB at Anonymized University and participants were compensated.

540 **Reproducibility Statement.** We will publicly release the Book-SORT dataset as well as all code  
 541 to generate new SORT datasets and evaluate models on SORT. For open models, evaluation on  
 542 Book-SORT is deterministic due to greedy sampling and the use of an answer prefix.  
 543

## 544 REFERENCES

- 546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
 547 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
 548 *arXiv preprint arXiv:2303.08774*, 2023.
- 549 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/  
 550 blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 552 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and  
 553 extraction, 2024. URL <https://arxiv.org/abs/2309.14316>.
- 554 Per Andersen, Richard GM Morris, David G Amaral, Tim Bliss, and John O’Keefe. *The Hippocampus  
 555 Book*. Oxford University Press, New York, 2006. ISBN 9780195100273.
- 557 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican,  
 558 George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al.  
 559 Improving language models by retrieving from trillions of tokens. In *International conference on  
 560 machine learning*, pp. 2206–2240. PMLR, 2022.
- 561 Agatha Christie. *The Murder of Roger Ackroyd*. Cosimo Classics, 1927.
- 563 Laura Lee Colgin, Edvard I Moser, and May-Britt Moser. Understanding memory through hippocam-  
 564 pal remapping. *Trends in neurosciences*, 31(9):469–477, 2008.
- 565 Martin A Conway. Sensory–perceptual episodic memory and its context: Autobiographical memory.  
 566 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 356  
 567 (1413):1375–1384, 2001.
- 569 Payel Das, Subhajt Chaudhury, Elliot Nelson, Igor Melnyk, Sarath Swaminathan, Sihui Dai, Aurélie  
 570 Lozano, Georgios Kollias, Vijil Chenthamarakshan, Jiří, Navrátil, Soham Dan, and Pin-Yu Chen.  
 571 Larimar: Large language models with episodic memory control, 2024. URL [https://arxiv.  
 572 org/abs/2403.11901](https://arxiv.org/abs/2403.11901).
- 573 Lila Davachi and Sarah DuBrow. How the hippocampus preserves order: the role of prediction and  
 574 context. *Trends in cognitive sciences*, 19(2):92–99, 2015.
- 575 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
 576 bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–  
 577 4186, 2019.
- 579 Howard Eichenbaum. Memory on time. *Trends in cognitive sciences*, 17(2):81–88, 2013.
- 580 Howard Eichenbaum. The hippocampus as a cognitive map ... of social space. *Neuron*, 87(1):  
 581 9–11, 2015a. ISSN 0896-6273. doi: <https://doi.org/10.1016/j.neuron.2015.06.013>. URL [https://  
 582 //www.sciencedirect.com/science/article/pii/S0896627315005267](https://www.sciencedirect.com/science/article/pii/S0896627315005267).
- 584 Howard Eichenbaum. The hippocampus as a cognitive map ... of social space. *Neuron*, 87  
 585 (1):9–11, July 2015b. ISSN 0896-6273. doi: 10.1016/j.neuron.2015.06.013. URL [http:  
 586 //dx.doi.org/10.1016/j.neuron.2015.06.013](http://dx.doi.org/10.1016/j.neuron.2015.06.013).
- 587 Howard Eichenbaum and Neal J. Cohen. Can we reconcile the declarative memory and spatial naviga-  
 588 tion views on hippocampal function? *Neuron*, 83(4):764–770, August 2014. ISSN 0896-6273. doi:  
 589 10.1016/j.neuron.2014.07.032. URL [http://dx.doi.org/10.1016/j.neuron.2014.  
 590 07.032](http://dx.doi.org/10.1016/j.neuron.2014.07.032).
- 591 Ian T. Ellwood. Short-term hebbian learning can implement transformer-like attention. *PLOS  
 592 Computational Biology*, 20(1):1–18, 01 2024. doi: 10.1371/journal.pcbi.1011843. URL [https:  
 593 //doi.org/10.1371/journal.pcbi.1011843](https://doi.org/10.1371/journal.pcbi.1011843).

- 594 Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon Hare, Frederique  
595 Laforest, and Elena Simperl. T-rex: A large scale alignment of natural language with knowledge  
596 base triples. In *Proceedings of the Eleventh International Conference on Language Resources and*  
597 *Evaluation (LREC 2018)*, 2018.
- 598  
599 Yu Feng, Ben Zhou, Haoyu Wang, Helen Jin, and Dan Roth. Generic temporal reasoning with  
600 differential analysis and explanation. In *Proceedings of the 61st Annual Meeting of the Association*  
601 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 12013–12029, 2023.
- 602 Zafeirios Fountas, Martin A Benfeghoul, Adnan Oomerjee, Fenia Christopoulou, Gerasimos Lam-  
603 pouras, Haitham Bou-Ammar, and Jun Wang. Human-like episodic memory for infinite context  
604 llms, 2024. URL <https://arxiv.org/abs/2407.09450>.
- 605  
606 Quentin Fournier, Gaétan Marceau Caron, and Daniel Aloise. A practical survey on faster and lighter  
607 transformers. *ACM Computing Surveys*, 55(14s):1–40, 2023.
- 608 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng  
609 Wang, and Haofen Wang. Retrieval-Augmented Generation for Large Language Models: A Survey,  
610 2024. URL <http://arxiv.org/abs/2312.10997>.
- 611  
612 Robert R Hampton and Bennett L Schwartz. Episodic memory in nonhumans: what, and where, is  
613 when? *Current opinion in neurobiology*, 14(2):192–197, 2004.
- 614 Uri Hasson, Janice Chen, and Christopher J Honey. Hierarchical process memory: memory as an  
615 integral component of information processing. *Trends in cognitive sciences*, 19(6):304–313, 2015.
- 616  
617 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
618 Steinhardt. Measuring massive multitask language understanding. In *International Conference on*  
619 *Learning Representations*, 2020.
- 620 George Hudson and Noura Al Moubayed. Muld: The multitask long document benchmark. In  
621 *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 3675–3685,  
622 2022.
- 623  
624 Mete Ismayilzada, Debjit Paul, Syrielle Montariol, Mor Geva, and Antoine Bosselut. Crow: Bench-  
625 marking commonsense reasoning in real-world tasks. In *Proceedings of the 2023 Conference on*  
626 *Empirical Methods in Natural Language Processing*, pp. 9785–9821, 2023.
- 627 Iván Izquierdo, Jorge H Medina, Mônica RM Vianna, Luciana A Izquierdo, and Daniela M Barros.  
628 Separate mechanisms for short-and long-term memory. *Behavioural brain research*, 103(1):1–11,  
629 1999.
- 630  
631 Li Ji-An, Corey Y. Zhou, Marcus K. Benna, and Marcelo G. Mattar. Linking in-context learning  
632 in transformers to human episodic memory, 2024. URL <https://arxiv.org/abs/2405.14992>.
- 633  
634 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
635 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
636 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 637  
638 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
639 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
640 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 641  
642 Greg Kamradt. Llmtest\_needleinahaystack, 2023. URL [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack). Accessed: 2024-06-03.
- 643  
644 Oliver Kliegl and Karl-Heinz T Bäuml. The mechanisms underlying interference and inhibition: A  
645 review of current behavioral and neuroimaging research. *Brain Sciences*, 11(9):1246, 2021.
- 646  
647 Tomáš Kočíský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,  
and Edward Grefenstette. The narrativeqa reading comprehension challenge, 2017. URL <https://arxiv.org/abs/1712.07040>.

- 648 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
649 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
650 serving with pagedattention, 2023.  
651
- 652 Janet Levin. Functionalism. In Edward N. Zalta and Uri Nodelman (eds.), *The Stanford Encyclopedia*  
653 *of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2023 edition, 2023.  
654
- 655 Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length on  
656 the reasoning performance of large language models. *arXiv preprint arXiv:2402.14848*, 2024.  
657
- 658 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,  
659 Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for  
660 natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual*  
661 *Meeting of the Association for Computational Linguistics*, pp. 7871–7880, 2020.  
662
- 663 Zhenrui Liao and Attila Losonczy. Learning, fast and slow: Single- and many-shot learning in  
664 the hippocampus. *Annual Review of Neuroscience*, 47(1):187–209, August 2024. ISSN 1545-  
665 4126. doi: 10.1146/annurev-neuro-102423-100258. URL [http://dx.doi.org/10.1146/  
666 annurev-neuro-102423-100258](http://dx.doi.org/10.1146/annurev-neuro-102423-100258).  
667
- 668 Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang.  
669 Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint*  
*arXiv:2311.08719*, 2023.  
670
- 671 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and  
672 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*  
673 *Association for Computational Linguistics*, 12:157–173, 2024.  
674
- 675 Robert L Logan IV, Nelson F Liu, Matthew E Peters, Matt Gardner, and Sameer Singh. Barack’s  
676 wife hillary: Using knowledge-graphs for fact-aware language modeling. *arXiv preprint*  
*arXiv:1906.07241*, 2019.  
677
- 678 Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study  
679 of catastrophic forgetting in large language models during continual fine-tuning, 2024. URL  
680 <https://arxiv.org/abs/2308.08747>.  
681
- 682 Andrew R Mayes and Neil Roberts. Theories of episodic memory. *Philosophical Transactions of the*  
683 *Royal Society of London. Series B: Biological Sciences*, 356(1413):1395–1408, 2001.  
684
- 685 James L McClelland, Bruce L McNaughton, and Randall C O’Reilly. Why there are complementary  
686 learning systems in the hippocampus and neocortex: insights from the successes and failures of  
687 connectionist models of learning and memory. *Psychological review*, 102(3):419, 1995.  
688
- 689 Arseny Moskvichev and Ky-Vinh Mai. Narrativexl: A large-scale dataset for long-term memory  
690 models, 2023. URL <https://arxiv.org/abs/2305.13877>.  
691
- 692 Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. Torque: A reading  
693 comprehension dataset of temporal ordering questions. In *Conference on Empirical Methods*  
694 *in Natural Language Processing*, 2020. URL [https://api.semanticscholar.org/  
695 CorpusID:218470560](https://api.semanticscholar.org/CorpusID:218470560).  
696
- 697 Dennis Norris. Short-term memory and long-term memory are still different. *Psychological bulletin*,  
698 143(9):992, 2017.  
699
- 700 John O’Keefe and Lynn Nadel. *The Hippocampus as a Cognitive Map*. Oxford: Clarendon Press,  
701 Oxford, UK, 1978. ISBN 9780198572060.
- 702 Oded Ovadia, Menachem Brief, Moshik Mishaelli, and Oren Elisha. Fine-tuning or retrieval? compar-  
ing knowledge injection in llms, 2024. URL <https://arxiv.org/abs/2312.05934>.

- 702 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi  
703 Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Gv,  
704 Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyra,  
705 Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song,  
706 Xiangru Tang, Johan Wind, Stanisław Woźniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and  
707 Rui-Jie Zhu. RWKV: Reinventing RNNs for the transformer era. In Houda Bouamor, Juan  
708 Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP*  
709 *2023*, pp. 14048–14077, Singapore, December 2023. Association for Computational Linguistics.  
710 doi: 10.18653/v1/2023.findings-emnlp.936. URL [https://aclanthology.org/2023.  
711 findings-emnlp.936](https://aclanthology.org/2023.findings-emnlp.936).
- 712 Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller,  
713 and Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*,  
714 2019.
- 715 Jason Phang, Yao Zhao, and Peter J Liu. Investigating efficiently extending transformers for long  
716 input summarization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural*  
717 *Language Processing*, pp. 3946–3961, 2023.
- 718 Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson.  
719 Fine-tuning aligned language models compromises safety, even when users do not intend to!, 2023.  
720
- 721 Yidan Qiu, Huakang Li, Jiajun Liao, Kemeng Chen, Xiaoyan Wu, Bingyi Liu, and Ruiwang  
722 Huang. Forming cognitive maps for abstract spaces: the roles of the human hippocampus and  
723 orbitofrontal cortex. *Communications Biology*, 7(1), May 2024. ISSN 2399-3642. doi: 10.1038/  
724 s42003-024-06214-5. URL <http://dx.doi.org/10.1038/s42003-024-06214-5>.
- 725 Kamoji Ryo, Goyal Tanya, and Rodriguez Juan Diego. Wice: Real-world entailment for claims in  
726 wikipedia. *arXiv preprint arXiv: 2303.01432 v1*, 2023.  
727
- 728 Bennett L Schwartz and Siân Evans. Episodic memory in primates. *American Journal of Primatology:*  
729 *Official Journal of the American Society of Primatologists*, 55(2):71–85, 2001.
- 730 Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong,  
731 Mor Geva, Jonathan Berant, et al. Scrolls: Standardized comparison over long language sequences.  
732 In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,  
733 pp. 12007–12021, 2022.
- 734 Larry R Squire and Stuart M Zola. Structure and function of declarative and nondeclarative memory  
735 systems. *Proceedings of the National Academy of Sciences*, 93(24):13515–13522, 1996.  
736
- 737 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
738 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the  
739 imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions*  
740 *on Machine Learning Research*, 2023.
- 741 Jørgen Sugar and May-Britt Moser. Episodic memory: Neuronal codes for what, where, and when.  
742 *Hippocampus*, 29(12):1190–1205, 2019.  
743
- 744 Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. Head-to-tail: How knowledge-  
745 able are large language models (llm)? aka will llms replace knowledge graphs? *arXiv preprint*  
746 *arXiv:2308.10168*, 2023.
- 747 Simeng Sun, Katherine Thai, and Mohit Iyyer. ChapterBreak: A challenge dataset for long-range  
748 language models. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz  
749 (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for*  
750 *Computational Linguistics: Human Language Technologies*, pp. 3704–3714, Seattle, United States,  
751 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.271.  
752 URL <https://aclanthology.org/2022.naacl-main.271>.
- 753 Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,  
754 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient  
755 transformers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qVyeW-grC2k>.

- 756 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,  
757 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models  
758 based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.  
759
- 760 Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023.  
761 URL <https://huggingface.co/datasets/teknium/OpenHermes-2.5>.
- 762 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
763 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
764 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.  
765
- 766 Endel Tulving. Episodic memory: From mind to brain. *Annual Review of Psychology*, 53(1):  
767 1–25, February 2002. ISSN 1545-2085. doi: 10.1146/annurev.psych.53.100901.135114. URL  
768 <http://dx.doi.org/10.1146/annurev.psych.53.100901.135114>.
- 769 Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei.  
770 Augmenting language models with long-term memory. *Advances in Neural Information Processing*  
771 *Systems*, 36, 2024.  
772
- 773 James C. R. Whittington, Joseph Warren, and Tim E.J. Behrens. Relating transformers to models and  
774 neural representations of the hippocampal formation. In *International Conference on Learning*  
775 *Representations*, 2022. URL <https://openreview.net/forum?id=B8DV09B1YE0>.
- 776 James C.R. Whittington, William Dorrell, Timothy E.J. Behrens, Surya Ganguli, and Mohamady  
777 El-Gaby. On prefrontal working memory and hippocampal episodic memory: Unifying memories  
778 stored in weights and activity slots. *bioRxiv*, 2024. doi: 10.1101/2023.11.05.565662. URL <https://www.biorxiv.org/content/early/2024/03/04/2023.11.05.565662>.  
779
- 780 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
781 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von  
782 Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama  
783 Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art  
784 natural language processing, 2020.  
785
- 786 Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. C-  
787 Pack: Packaged Resources To Advance General Chinese Embedding. In *Proceedings of the 47th*  
788 *International ACM SIGIR Conference on Research and Development in Information Retrieval*.  
789 arXiv, 2024. URL <http://arxiv.org/abs/2309.07597>.  
790
- 791 Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun  
792 Yao, Xiaohan Zhang, Hanming Li, et al. Kola: Carefully benchmarking world knowledge of large  
793 language models. *arXiv preprint arXiv:2306.09296*, 2023.  
794
- 795 Tan Yu, Anbang Xu, and Rama Akkiraju. In defense of rag in the era of long-context language  
796 models, 2024. URL <https://arxiv.org/abs/2409.01666>.
- 797 Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang.  
798 Removing rlhf protections in gpt-4 via fine-tuning, 2024.  
799
- 800 Peitian Zhang, Ninglu Shao, Zheng Liu, Shitao Xiao, Hongjin Qian, Qiwei Ye, and Zhicheng Dou.  
801 Extending llama-3’s context ten-fold overnight, 2024a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2404.19553)  
802 [2404.19553](https://arxiv.org/abs/2404.19553).
- 803 Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi  
804 Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. Instruction tuning for large language models: A  
805 survey, 2024b.  
806
- 807 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing  
808 large language models with long-term memory. *Proceedings of the AAAI Conference on Artificial*  
809 *Intelligence*, 38(17):19724–19731, 2024. doi: 10.1609/aaai.v38i17.29946. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29946>.

810 Zexuan Zhong, Tao Lei, and Danqi Chen. Training language models with memory augmentation. In  
811 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.  
812 5657–5673, 2022.

813 Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. Temporal  
814 reasoning on implicit events from distant supervision. In *Proceedings of the 2021 Conference of*  
815 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
816 *Technologies*, pp. 1361–1371, 2021.

817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



## A DEFINING EPISODIC MEMORY FOR LLMs

Episodic memory is a new concept for LLMs that has not previously been explored in-depth. As such, it might be helpful to provide more explanations about the criteria that we adopted to define this form of memory in LLMs based on extensive research and mature theories about episodic memory in humans Tulving (2002); Andersen et al. (2006). Unlike human definitions of episodic memory, we do not require any notion of conscious personal experience, or a specific neural implementation (e.g. something like hippocampal dependence) Tulving (2002). Therefore this section will go through each criterion to provide additional background for our reasoning behind each item.

**(1) Episodic memory is specific to a single sequence.** Episodic memory is a form of memory that is always specific to a single sequence and its unique temporal context. For humans, two experiences are easily distinguishable due to their high-dimensional nature, which makes it possible to find a number of features in which they differ. In the neuroscience literature, this is often described as memory representations that do not interfere with one another Sugar & Moser (2019); Colgin et al. (2008). This differentiates it from semantic memory, which can not necessarily be linked to a particular, unique sequence.

**(2) Episodic memory is learned through a single exposure to that sequence.** Unlike semantic memory, episodic memories in humans and other animals can be based on single experiences, i.e. acquired through single-shot learning Liao & Losonczy (2024) or single-trial learning Schwartz & Evans (2001).

**(3) Episodic memory binds context to memory content.** This is based on decades of research establishing that episodic memory binds the 'what', 'where', and 'when' of specific memories Sugar & Moser (2019). That is, episodic memories encode the spatial and temporal context associated with the encoded information, establishing a "cognitive map" in episodic memory O'Keefe & Nadel (1978). A contemporary theory of episodic memory posits that it is a more general "relational processing mechanism" Eichenbaum & Cohen (2014). It can include more abstract relations, e.g. functioning as a map of social space of memory (Eichenbaum, 2015b). This criterion does not just take the conservative view that episodic memory spans cognitive maps for space and time of events and items, but it can include more abstract relations between parts of a sequence, allowing for more abstractive relational mapping between and within parts of a sequence. The most conservative, well-studied, and easily testable aspects however remain temporal order and spatial memory, and we suggest that it is sensible to evaluate these types of contextual relations in LLMs before including more abstract relations, which we did not do in this work.

**(4) Episodic memory can potentially persist for an arbitrarily long time.** In humans and other biological systems, episodic memory is specifically a form of long-term memory Squire & Zola (1996) that stores knowledge which can persist up to the span of a lifetime Mayes & Roberts (2001); Conway (2001). In LLMs this means episodic memory needs to allow for (up to) arbitrarily many tokens in between the memory-sequence and a query about that sequence.

**(5) Episodic memory is generally accessible for explicit reasoning and communication.** In classic views of human memory systems, episodic memory is a type of declarative, or explicit memory Squire & Zola (1996). This criterion is partly based on this classic view. While explicit memory is often characterized as consciously accessible memory, for LLMs, we map this to whether an LLM is able to answer explicit questions, suggesting that such information could explicitly be requested and then used in internal reasoning, which does not invoke any notion of consciousness. This is a *functionalist* criterion, as defined in the philosophy of mind Levin (2023). That is, our definition of episodic memory in deep learning models requires them to actually use the memory for general task solving. This criterion is also shared by neuroscientists studying episodic memory in non-human animals Hampton & Schwartz (2004).

## B ADDITIONAL DETAILS ON BOOK-SORT DATA SET

**Preprocessing book text.** We wrote custom Python code to only retain the book text that formed a continuous narrative. We stripped the front and back matter of the book, and extracted chapter titles

if they existed. 8 of the 9 books contained individual section or chapter breaks. For these 8 books, we parsed the text corresponding to each chapter. Chapter titles or section headings (e.g. ‘VI’ to indicate section six) were removed, and all remaining text was concatenated. This string was split into words (assuming simple whitespace separators with Python `string.split()`) to produce a final text array for each book. This text array was sampled for the Book-SORT dataset.

## B.1 BOOK SELECTION

We provide details about the 9 books in Book-SORT in Table 3.

Table 3: Project Gutenberg metadata on Book-SORT books.

| ID    | Title                              | Author                             | Word count | Release   | Pub  | LoCC* | Subjects  |
|-------|------------------------------------|------------------------------------|------------|-----------|------|-------|---|
| 69087 | The Murder of Roger Ackroyd        | Christie, Agatha                   | 69,720     | 10/2/2022 | 1926 | PR    | Detective and mystery stories; Fiction: Private investigators - England, Murder - Investigation, Belgians - England |
| 72578 | Tom Swift and His Talking Pictures | Appleton, Victor                   | 43,853     | 1/1/2024  | 1928 | PZ    | Adventure stories; Motion pictures  |
| 72600 | The Trumpeter of Krakow            | Kelly, Eric Philbrook              | 59,081     | 1/2/2024  | 1928 | PZ    | Juvenile fiction: Middle Ages, Poland - History - Casimir IV, 1447-1492   |
| 72869 | Meet the Tiger                     | Charteris, Leslie                  | 79,946     | 2/4/2024  | 1928 | PR    | Fiction: Private investigators - England; Detective and mystery stories   |
| 72958 | Hunting for Hidden Gold            | Dixon, Franklin W.                 | 42,354     | 2/14/2024 | 1928 | PZ    | Juvenile fiction: Brothers, Gold mines and mining, Montana, Robbers and outlaws; Mystery and detective stories      |
| 72963 | The Nature of the Physical World   | Eddington, Arthur Stanley, Sir     | 104,530    | 2/15/2024 | 1928 | Q     | Physics - Philosophy; Science - Philosophy  |
| 72972 | Money for Nothing                  | Wodehouse, P.G. (Pelham Grenville) | 82,331     | 2/16/2024 | 1928 | PR    | Humorous stories; Fiction: Swindlers and swindling, Greed   |
| 73017 | Pomona; or, the Future of English  | De Selincourt, Basil               | 9,273      | 2/22/2024 | 1928 | PE    | English language  |
| 73042 | The Well of Loneliness             | Hall, Radclyffe                    | 163,217    | 2/26/2024 | 1928 | PR    | Fiction: Lesbians - England - Social conditions   |

\*LoCC = Library of Congress classification.

## B.2 BETWEEN-SEGMENT DISTANCES

The segment distance  $L_S$  for Book-SORT is sampled from one of four distance bins. The right edge of each bin is given in Table 4. Distance is computed between the beginning of the first segment and the beginning of the second segment. The minimum distance  $L_S$  therefore produces adjacent, non-overlapping segments.

Table 4: Right edge of each distance bin used to create samples for Book-SORT.

|                   | Minimum | Bin0    | Bin1    | Bin2    | Bin3      |
|-------------------|---------|---------|---------|---------|-----------|
| $L_E \leq 2,500$  | $L_S$   | $L_E/4$ | $L_E/3$ | $L_E/2$ | $L_E/0.8$ |
| $L_E \geq 10,000$ | $L_S$   | 1000    | $L_E/4$ | $L_E/2$ | $L_E/0.8$ |

## B.3 HUMAN STUDY DETAILS

**Participant compensation.** Participants were compensated via a lottery system with a chance to win a gift card to a popular book store. The expected value of the compensation came out to \$12 per hour.

**Study design.** Each participant completed an online survey. First, the participant consented to the study, read a brief set of instructions, and completed a brief survey, including a question regarding when the participant finished reading the book. The complete set of survey questions is listed below. Each participant was then asked to answer "Which segment occurred first in the book?" for 10 randomly chosen text segment pairs from a total set of 540 unique segment pairs sampled from the whole book. We chose to present a sample number of trials to each participant to minimize interference effects from repeated memory retrieval (Kliegl & Bäuml, 2021). The presentation order of the text segments was randomized across participants. In the end, each participant was asked 4

972 simple questions about the book plot to verify that the participant had indeed read the book. Each  
 973 participant was only allowed to participate in the study once.  
 974

975 **Demographics questions.** The human participants were asked the following set of demographics  
 976 questions before beginning the experiment:  
 977

- 978 1. I have finished the book The Murder of Roger Ackroyd [Options: True/False]
- 979 2. On what date did you finish the book? [Calendar question type]
- 980 3. Did you read or listen to the book? [Options: Read/Listen]
- 981 4. Was this your first time reading / listening to the book? [Options: Yes / No]
- 982 5. What is your age? [Options: 18-25, 25-35, 35-45, 45-55, 55-65, 65+]
- 983 6. What gender do you identify with? [Options: Female/Male/Other]
- 984 7. What is your experience with the English language? [Options: Native / Fluent / Advanced /  
 985 Intermediate / Beginner]
- 986 8. How many books did you read or listen to in the past year? [Options: 1-2 / 3-5 / 6-10 / 10+]

989 We use the responses above to determine the number of days that have passed since finishing the  
 990 book, and make this information available in the human dataset together with the responses.  
 991

992 **Inclusion criteria.** We include data from participants who answered at least 3 of 4 plot questions  
 993 correctly, and finished reading the book within 30 days of participating in the study. These inclusion  
 994 criteria result in 155 participants.  
 995

## 996 C MODEL AND PROMPTING DETAILS

### 997 C.1 MODEL DETAILS

1000 We listed all models we used in this paper and their download links from HuggingFace in Table  
 1001 5. For the OpenAI models, we used the gpt-3.5-turbo-0125 version of GPT-3.5, and gpt-4-turbo-  
 1002 2024-04-09 for GPT-4. Models were selected to cover a broad range of performance on more  
 1003 semantic/knowledge-based tasks such as those included in MMLU.  
 1004

1005 Table 5: Model Details

| 1006                           | 1007                  | 1008       | 1009 |
|--------------------------------|-----------------------|------------|------|
| Name in HuggingFace            | Name in Paper         | MMLU score |      |
| Llama-3-70B-Instruct           | Llama3-70b-inst       | 80.06      |      |
| Llama-3-8B-Instruct            | Llama3-8b-inst        | 66.60      |      |
| Mixtral-8x22B-Instruct-v0.1    | Mixtral-8x22b-inst    | 77.77      |      |
| Nous-Hermes-2-Mixtral-8x7B-DPO | Mixtral-8x7b-DPO-inst | 72.28      |      |
| Mistral-7B-Instruct-v0.1       | Mistral-v1-7b-inst    | 60.10      |      |
| Mistral-7B-Instruct-v0.2       | Mistral-v2-7b-inst    | 60.07      |      |
| Llama-2-70b-chat               | Llama2-70b-inst       | 68.90      |      |
| Llama-2-7b-chat                | Llama2-7b-inst        | 45.30      |      |
| gemma-1.1-7b-inst              | Gemma-1.1-7b-inst     | 64.30      |      |

### 1020 C.2 PROMPTING

1021 For our experiments with Book-SORT, we created a total of 12 prompts that are composed of two  
 1022 parts. The prompts differ in how they phrase the tasks. The first part contains instructions to read  
 1023 the text excerpt from the book as well as a placeholder for the actual excerpt. The second part of  
 1024 the prompt contains the description of SORT, including a mention of the book or document title as  
 1025 well as two segments from that document. We found that current open LLMs fail at the task even

1026 with in-context access to the text, if they are asked to tell which segment appeared second or last.  
 1027 For this reason, we ran all experiments with the placeholder <position> set to "first". All of these  
 1028 prompts were preceded by the same generic system prompt: "You are a helpful, respectful and honest  
 1029 assistant."

1030 Table 6: Selection of 12 prompts used for prompt validation

| No. | Reading instruction  | SORT instruction   |
|-----|--|--|
| 1   | "Please take some time to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>" | "You will be shown pairs of text fragments from <booktitle>. Please select which of two fragments appeared <position> in the book. You will be shown 10 such pairs. <segments> Which fragment appeared <position> in the book, <label_0> or <label_1>?"  |
| 2   | "I need you to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>"            | "In this exercise, your objective is to identify the text segment, either <label_0> or <label_1>, that appeared <position> in <booktitle>. Please read the segments carefully to determine their order of appearance in <booktitle> and respond with either <label_0> or <label_1>: <segments> Which of these, <label_0> or <label_1>, was <position> in <booktitle>?" |
| 3   | "I need you to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>"            | "Your task is to recall which text segment, either <label_0> or <label_1>, appeared <position> in the book <booktitle>. Please read the segments carefully to remember in which order they appeared in <booktitle> and respond with either <label_0> or <label_1>: <segments> Which of these, <label_0> or <label_1>, was <position> in the book <booktitle>?"         |
| 4   | "I need you to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>"            | "You will be shown two text segments, labeled as <label_0> and <label_1>. Please recall in which order they appeared in the book <booktitle> and tell me which one came <position>. Please read the segments carefully: <segments> Which of these two parts of the book, <label_0> or <label_1>, came <position> in the book <booktitle>?"                             |
| 5   | "I need you to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>"            | "I will show you two short parts from a book, labeled as <label_0> or <label_1>. Your task is to tell me which of them appeared <position> in the book <booktitle>. Please read both segments carefully and try to remember where in the book they come from: <segments> Which of these, <label_0> or <label_1>, appeared <position> in the book <booktitle>?"         |

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

Table 6: Selection of 13 prompts used for prompt validation

| No. | Reading instruction   | SORT instruction   |
|-----|---|--|
| 6   | "I need you to thoroughly read and comprehend this extract from the book <booktitle>. The passage is as follows: <excerpt>" | "This is your task: Given two segments from a book, labeled as <label_0> and <label_1>, please tell me which of them appeared <position> in <booktitle>. Read both segments carefully and try to remember where in <booktitle> they appeared: <segments> Which of these, <label_0> or <label_1>, comes <position> in the book <booktitle>?"                            |
| 7   | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>"                     | "You will be shown pairs of text fragments from <booktitle>. Please select which of two fragments appeared <position> in the book. You will be shown 10 such pairs. <segments> Which fragment appeared <position> in the book, <label_0> or <label_1>?"  |
| 8   | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>"                     | "In this exercise, your objective is to identify the text segment, either <label_0> or <label_1>, that appeared <position> in <booktitle>. Please read the segments carefully to determine their order of appearance in <booktitle> and respond with either <label_0> or <label_1>: <segments> Which of these, <label_0> or <label_1>, was <position> in <booktitle>?" |
| 9   | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>"                     | "Your task is to recall which text segment, either <label_0> or <label_1>, appeared <position> in the book <booktitle>. Please read the segments carefully to remember in which order they appeared in <booktitle> and respond with either <label_0> or <label_1>: <segments> Which of these, <label_0> or <label_1>, was <position> in the book <booktitle>?"         |
| 10  | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>"                     | "You will be shown two text segments, labeled as <label_0> and <label_1>. Please recall in which order they appeared in the book <booktitle> and tell me which one came <position>. Please read the segments carefully: <segments> Which of these two parts of the book, <label_0> or <label_1>, came <position> in the book <booktitle>?"                             |
| 11  | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>"                     | "I will show you two short parts from a book, labeled as <label_0> and <label_1>. Your task is to tell me which of them appeared <position> in the book <booktitle>. Please read both segments carefully and try to remember where in the book they come from: <segments> Which of these, <label_0> or <label_1>, appeared <position> in the book <booktitle>?"        |

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Table 6: Selection of 13 prompts used for prompt validation

| No. | Reading instruction   | SORT instruction  |
|-----|---|---|
| 12  | "Please carefully read this excerpt from the book <booktitle>. This is the relevant passage: <excerpt>" | "This is your task: Given two segments from a book, labeled as <label_0> and <label_1>, please tell me which of them appeared <position> in <booktitle>. Read both segments carefully and try to remember where in <booktitle> they appeared: <segments> Which of these, <label_0> or <label_1>, comes <position> in the book <booktitle>?" |

### C.3 PER-MODEL RESULTS ON PROMPT SELECTION SWEEP

To identify the prompts that work best for each model, we take 400 segment-pair samples that we excluded from the main evaluation and evaluate models’ in-context memory with all prompts shown in Table 6. To select the best prompt we considered both the proportion of A and B responses, which should be around 0.5, and the accuracy. We report the best selected prompts in Table 10 with numbers referring to the prompts presented in Table 6.

Table 7: Selected prompts for each model.

| Model Name           | Best Prompt |
|----------------------|-------------|
| Llama3-70b-inst      | 4           |
| Llama3-8b-inst       | 3           |
| Mixtral-8x22b-inst   | 4           |
| Llama2-70b-inst      | 7           |
| Gemma-1.1-7b-inst    | 8           |
| Mistral-v0.2-7b-inst | 3           |
| Mistral-v0.1-7b-inst | 2           |
| Llama2-7b-inst       | 10          |
| GPT-3.5-turbo        | 7           |
| GPT-4                | 7           |

### C.4 RAG PROMPT SELECTION

There were two different prompts to select for the retrieval-augmented generation experiments: the retrieval prompt (i.e. the search query), and the LLM prompt.

#### C.4.1 RETRIEVAL PROMPT (SEARCH QUERY)

The goal of retrieval in our RAG experiments is to find the text passages that will provide the most information about the segments for the sequence ordering task. After we created the vector database of all the text passages from Book-SORT, we formulated several different search queries (Table 8). We then ran retrieval using a validation subset of Book-SORT (50-word segments, 250-word excerpts from all books). The retrieval used the same database and text embedding model as described in the RAG portion of Section 3.2. The best search query was simple and only consisted of the segment text (query 8, Table 8). This search query is used for all RAG experiments.

#### C.4.2 RAG LLM PROMPTS

We followed a procedure similar to the one outlined in Section C.2. We created a total of 10 modifications to the reading instructions from Table 6.

Table 8: The search queries for the RAG experiment and their average retrieval recall@10 on a validation subset of Book-SORT (250 word excerpts, 50 word segments).

| No. | Search Query Text   | Recall@10 |
|-----|---|-----------|
| 0   | "Please determine the order in which the following text segments appeared in <booktitle>: <segments>" | 0.728     |
| 1   | "We need to put text segments from <booktitle> in order. These are the segments: <segments>"          | 0.817     |
| 2   | "Please find these text segments from <booktitle>: <segments>"  | 0.869     |
| 3   | "Please find these text segments from <booktitle> to provide context for the next task: <segments>"   | 0.875     |
| 4   | "Which text chunks from <booktitle> contain the following segments? <segments>"                       | 0.802     |
| 5   | "Which text excerpts from <booktitle> contain the following segments? <segments>"                     | 0.799     |
| 6   | "Which text chunks from <booktitle> overlap with these text segments: <segments>"                     | 0.782     |
| 7   | "<booktitle> contains this text: <segments>"  | 0.865     |
| 8   | "<segments>"  | 0.906     |
| 9   | "<booktitle> <segments>"  | 0.858     |

Table 9: RAG prompt modifications.

| No. | RAG Reading Instruction  |
|-----|--|
| 0   | "Here are some relevant excerpts from the book <booktitle>: <context>"                                     |
| 1   | "The following excerpts from the book <booktitle> may be helpful context for the task. Context: <context>" |
| 2   | "Context: <context>"   |
| 3   | "Searching a book database found these relevant text snippets: <context>"                                  |
| 4   | "The following search results may be useful context: <context>"  |
| 5   | "I will show you some relevant text found by searching a database of books: <context>"                     |
| 6   | "Please read some text deemed relevant for the task before performing the task. Relevant text: <context>"  |
| 7   | "Please read these search results carefully to help you perform the task. Search results: <context>"       |
| 8   | "Your objective may become easier with the use of these search results: <context>"                         |
| 9   | "This context may be helpful: <context>"   |

#### C.4.3 PER-MODEL RESULTS ON RAG PROMPT SELECTION

For a given LLM, we modified the reading instruction of the best prompt from Table 10 with each of the 10 options in Table 9. We then ran a sweep over the same 400 segment-pair samples detailed in Section C.3 and found the instruction that resulted in the highest performance on this held-out dataset.

Table 10: Best RAG instruction prompts for each model.

| Model Name           | Best RAG Instruction No. |
|----------------------|--------------------------|
| Llama3-70b-inst      | 7                        |
| Llama3-8b-inst       | 7                        |
| Mixtral-8x22b-inst   | 3                        |
| Mistral-v0.2-7b-inst | 6                        |

## D ADDITIONAL DETAILS ON BOOK-SORT RESULTS

### D.1 MEMORY-LESS BASELINE RESULTS

Figure 5 shows performance on Book-SORT without any memory-insertion of the books used in Book-SORT. We find that performance is higher in segment pairs that are very proximal or very distant in the book, indicating that it might be easier to sort these pairs based on temporal order reasoning. Performance without additional memory-insertion is generally low, showing that memory is needed for SORT.

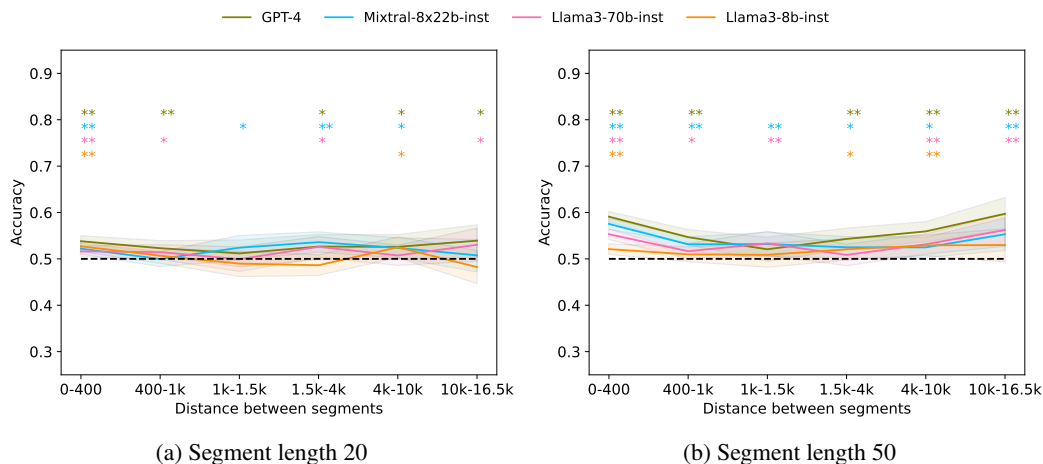


Figure 5: Baseline SORT performance without memory of books in Book-SORT. Significant difference from chance is marked with asterisks (\*p-value<0.05, \*\*p-value<0.01).

### D.2 IN-CONTEXT MEMORY FULL RESULTS

In this section, we provide a comprehensive overview of the in-context memory results across various models in Table 11 and Table 12. The table below illustrates the accuracy of different models on multiple books at segment lengths of 20 and 50 words. We observe that, while models generally perform slightly better with longer segments (50 words) compared to shorter ones (20 words), the improvement is modest, averaging up to 4%.



Table 11: Accuracy and Difference of Various Models on Multiple Books at Excerpt Lengths of 20 and 50, with in-context memory (Part 1)

| Model name      | Book  | SORT S20  | SORT S50  | SORT-Extend S20 | SORT-Extend S50 |
|-----------------|-------|-----------|-----------|-----------------|-----------------|
| Llama3-8b-inst  | 69087 | 0.89±0.03 | 0.92±0.03 | /               | /               |
| Llama3-8b-inst  | 72578 | 0.91±0.02 | 0.93±0.02 | /               | /               |
| Llama3-8b-inst  | 72600 | 0.92±0.03 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 72869 | 0.92±0.03 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 72958 | 0.92±0.02 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 72963 | 0.92±0.03 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 72972 | 0.92±0.03 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 73017 | 0.91±0.03 | 0.94±0.02 | /               | /               |
| Llama3-8b-inst  | 73042 | 0.92±0.03 | 0.94±0.02 | /               | /               |
| Llama2-70b-inst | 69087 | 0.74±0.12 | 0.90±0.08 | /               | /               |
| Llama2-70b-inst | 72578 | 0.75±0.12 | 0.90±0.09 | /               | /               |
| Llama2-70b-inst | 72600 | 0.71±0.13 | 0.91±0.09 | /               | /               |
| Llama2-70b-inst | 72869 | 0.71±0.13 | 0.91±0.09 | /               | /               |
| Llama2-70b-inst | 72958 | 0.71±0.13 | 0.90±0.09 | /               | /               |
| Llama2-70b-inst | 72963 | 0.72±0.13 | 0.89±0.10 | /               | /               |
| Llama2-70b-inst | 72972 | 0.70±0.13 | 0.88±0.10 | /               | /               |
| Llama2-70b-inst | 73017 | 0.70±0.13 | 0.87±0.10 | /               | /               |
| Llama2-70b-inst | 73042 | 0.71±0.13 | 0.88±0.10 | /               | /               |
| Llama2-7b-inst  | 69087 | 0.56±0.05 | 0.56±0.05 | /               | /               |
| Llama2-7b-inst  | 72578 | 0.57±0.05 | 0.55±0.05 | /               | /               |
| Llama2-7b-inst  | 72600 | 0.57±0.05 | 0.56±0.04 | /               | /               |
| Llama2-7b-inst  | 72869 | 0.57±0.05 | 0.56±0.04 | /               | /               |
| Llama2-7b-inst  | 72958 | 0.57±0.05 | 0.56±0.04 | /               | /               |
| Llama2-7b-inst  | 72963 | 0.57±0.05 | 0.57±0.05 | /               | /               |
| Llama2-7b-inst  | 72972 | 0.57±0.05 | 0.56±0.05 | /               | /               |
| Llama2-7b-inst  | 73017 | 0.57±0.05 | 0.56±0.05 | /               | /               |
| Llama2-7b-inst  | 73042 | 0.57±0.05 | 0.56±0.05 | /               | /               |
| Llama3-70b-inst | 69087 | 0.90±0.08 | 0.92±0.09 | /               | /               |
| Llama3-70b-inst | 72578 | 0.92±0.08 | 0.92±0.09 | /               | /               |
| Llama3-70b-inst | 72600 | 0.92±0.08 | 0.93±0.09 | /               | /               |
| Llama3-70b-inst | 72869 | 0.93±0.07 | 0.93±0.08 | /               | /               |
| Llama3-70b-inst | 72958 | 0.93±0.07 | 0.94±0.08 | /               | /               |
| Llama3-70b-inst | 72963 | 0.92±0.08 | 0.93±0.09 | /               | /               |
| Llama3-70b-inst | 72972 | 0.91±0.08 | 0.93±0.09 | /               | /               |
| Llama3-70b-inst | 73017 | 0.92±0.08 | 0.94±0.09 | /               | /               |
| Llama3-70b-inst | 73042 | 0.91±0.09 | 0.94±0.08 | /               | /               |

### D.3 RESULTS PER BOOK

In Fig. 6, we provide the baseline results without text-specific memory separately for each of the 9 books in Book-SORT.

In Fig. 7, we provide the in-context memory results separately for each of the 9 books in Book-SORT.

### D.4 RELATIONSHIP BETWEEN IN-CONTEXT MEMORY RESULTS AND DISTANCE BETWEEN SEGMENTS ACROSS EXCERPT LENGTHS

In Fig. 8 and Fig 9, we show the average accuracy by the distance between segments for all the excerpt lengths and segment lengths.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

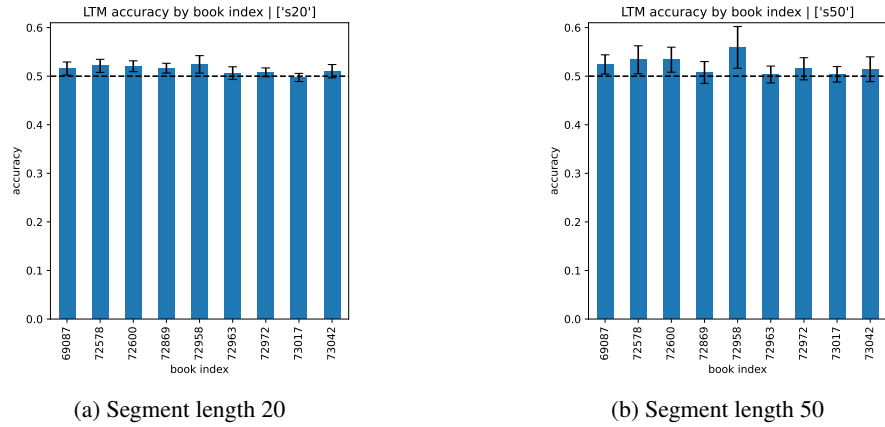


Figure 6: Models' baseline performance by book (error bars indicate standard deviation)

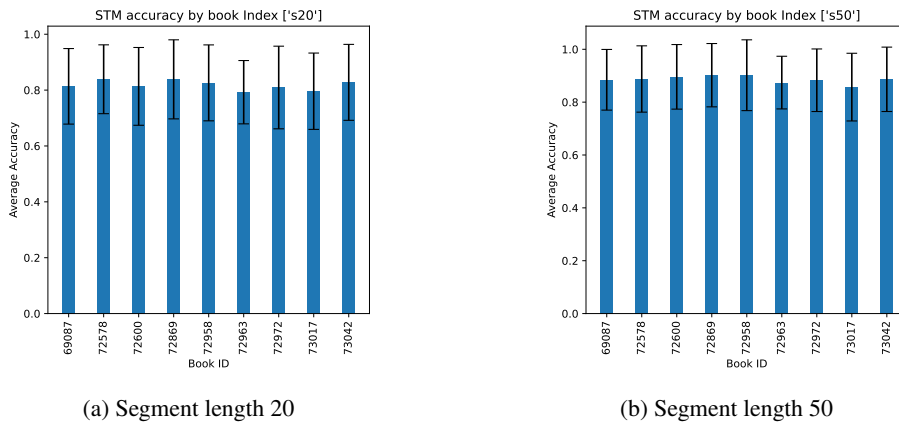


Figure 7: Models' in-context memory performance by book (error bars indicate standard deviation)

## D.5 IN-CONTEXT MEMORY: FULL BOOK EVALUATION

Recent LLMs have longer context windows for which they can perform well. This opens the possibility to test in-context memory for these books when presented with a complete book in their context window. It is straightforward to create an adequate SORT-dataset to test how well an LLM can recall sequence order when presented with a full book instead of just an excerpt. Initial findings from evaluation of Llama3.1-8b-Instruct when presented with the complete book *The Murder of Roger Ackroyd* suggest that the model is performing worse than humans (see Figure 10). While the human performance shown as a reference in Figure 10 is based on reading the book up to 30 days prior to testing, the model only has the book-text in its context window without the presence of task-irrelevant additional episodes that occurred for the humans. To match the difficulty to human episodic memory testing, future versions of SORT datasets could introduce the addition of irrelevant documents, similar to a needle-in-a-haystack task (Kamradt, 2023).

The SORT-dataset for this evaluation is based on 300 pairs of 50-word segments for each of the segment-distance bins that we report. Because we test on a complete book and not an excerpt from a book, we modified the reading instruction of prompt 2 in Table 6 by replacing “excerpt from the book” with “the complete book”.

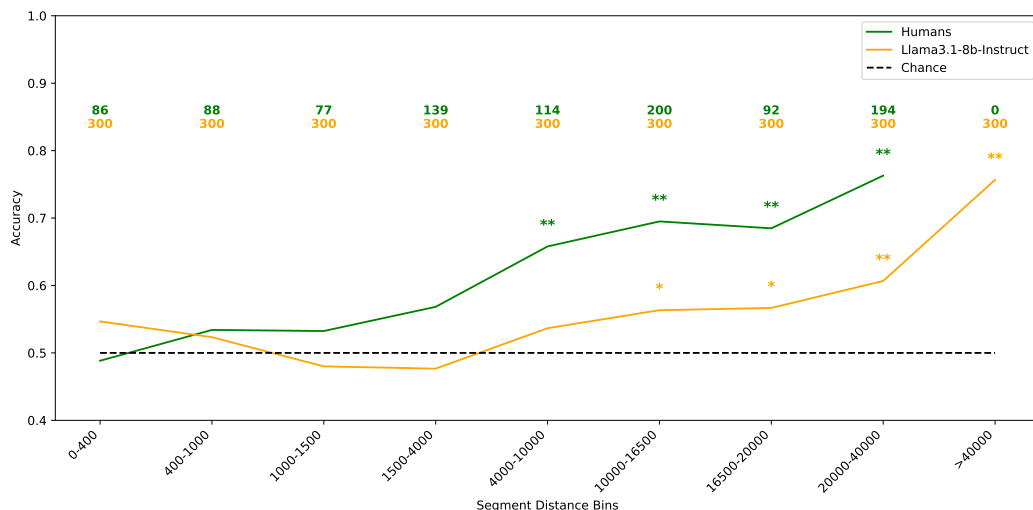


Figure 10: SORT evaluation of Llama3.1-8b-Instruct (128k context window) on a SORT dataset based on *The Murder of Roger Ackroyd*. Numbers indicate the sample-number within the bins. Asterisks indicate statistical significance: \* $p < 0.05$  and \*\* $p < 0.01$

## D.6 IN-CONTEXT MEMORY: LOST-IN-THE-MIDDLE EFFECT

The lost-in-the-middle effect is present in the in-context memory condition of SORT (Fig. 11). Both segments are considered to be in the middle section of the excerpt when the first word of each segment is in the middle one-third of the excerpt. Using logistic regression, we examined the interaction between whether both segments are in the middle section of the excerpt and the excerpt length, after controlling for the distance between the two segments. For nine out of ten models, we found a significant lost-in-the-middle effect, where accuracy is lower when both segments are in the middle section. For seven out of ten models, there is also a significant interaction between whether both segments are in the middle section of the excerpt and the excerpt length, suggesting that the degree of the lost-in-the-middle effect varies across excerpt lengths.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

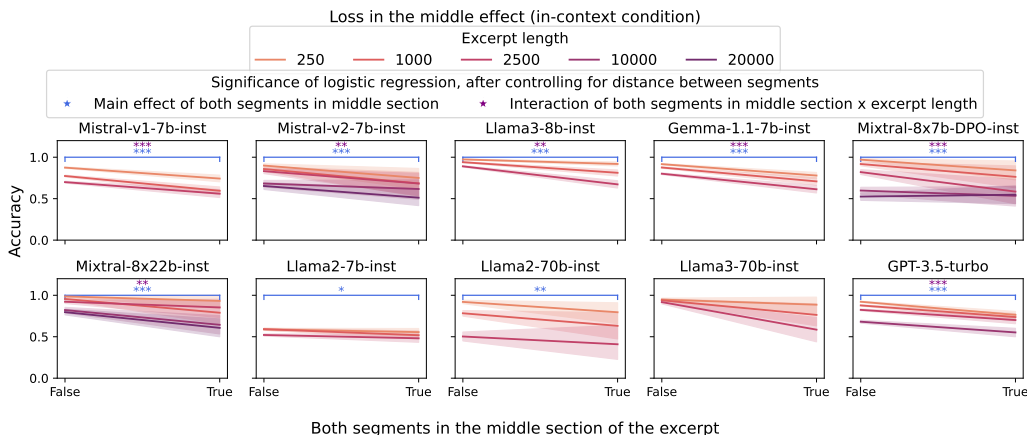


Figure 11: Lost-in-the-middle effect for the in-context memory condition

D.7 IN-CONTEXT MEMORY: CONTEXT-EXTENDED LLAMA3-8B-INSTRUCT MODEL

We evaluate a single context-extended Llama3-8b-Instruct model (Zhang et al., 2024a) on SORT that reports high performance on other benchmarks after context-extension to see how performance changes. Our findings (Figure 12) show a strong degradation in performance with low performance on 20k word excerpts, despite a claimed support for 80k context windows and high performance on needle-in-a-haystack tests for long sequences. This highlights the need for high quality data during context-extension, and shows an additional use for SORT: to benchmark context-extension methods.

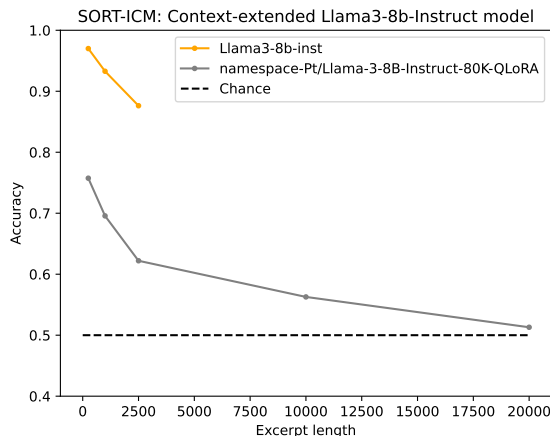


Figure 12: Evaluation of a context-extended model with a supposed 80k context support shows degraded performance.

D.8 BASELINE PERFORMANCE

In Fig. 13, we provide the SORT results based on parametric memory for all models across various segment distances. Due to the recent addition of the texts in Book-SORT to the public domain, we expect that models were not trained on these texts, i.e. they should not have text-specific memory. Performance is higher for segment pairs that have a short distance and a high distance in the books, indicating that these are more likely to be sort-able without episodic memory, based on temporal order reasoning.

## 1512 E BOOK-SORT RESULTS FROM ADDITIONAL MODELS

1513

1514

1515

### E.1 BASE MODELS

1516

1517

1518

1519

1520

We chose 2 base models to evaluate, Llama3-8b and Mistral-7b, whose fine-tuned versions (Llama3-8b-inst and Mistral-v2-7b-inst) performed well on SORT based on in-context memory. Figure 14 shows that both the base models got around chance performance across all the excerpt lengths and segment lengths.

1521

1522

### E.2 STATE-SPACE MODELS

1523

1524

1525

1526

1527

1528

1529

1530

We tested an instruction-tuned version of the state space model RWKV (Peng et al., 2023), available in Huggingface as RWKV/rwkv-raven-7b. The results of the prompt sweep on SORT with in-context memory yielded a performance of 51% – very close to chance levels. A possibility for this is a larger sensitivity to prompting, e.g. this model might require instructions to be given in a different order. We assume that this is due to insufficient instruction tuning. While it could be interesting to see the performance of a state-space model with memory other than in-context, we leave this question to future work.

1531

1532

## F FINETUNING OF LLAMA3-8B-INSTRUCT

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

**Fine-tuning details.** We fine-tuned Llama3-8b-Instruct and Mistral-7b-v0.2-Instruct on a single node with 8 A100 GPUs. The books (without pre-processing beyond removing Project Gutenberg related text, i.e. including chapter signifiers) are split into chunks of 5000 words and contextualized in the same way in which excerpts are presented in-context in our experiments, i.e. together with the book-title in a user prompt along with a preceding system prompt. For the instruction data, we exclude the following task types: "experience", "stylized\_response", "joke", "trivia", "roleplay", "riddle" and "greeting". Samples containing both book-chunks and instruction-following examples are padded to the maximum length in a batch. The effective batch size in our experiments is 192. We choose a moderately low initial learning rate of  $5e-6$  with cosine decay and a small amount of weight decay set to  $1e-4$ . The chunks of books comprise a total of 116 independent samples. Together with 3500 instruction samples from the OpenHermes dataset (Teknum, 2023), this means 19 steps of gradient descent are taken in one epoch. We fine-tuned both models for a total of 5 epochs, however only the first epoch would qualify for episodic memory testing, since one of the characteristics of episodic memory is that it is single-shot learned (see definition 1 and Appendix A for a discussion). This precludes complete memorization of texts, for which multiple repetitions are needed (Ovadia et al., 2024).

1550

1551

1552

1553

1554

1555

1556

1557

1558

**Inclusion of instruction data to avoid catastrophic forgetting.** Fine-tuning an instruction-tuned model on specific data can lead to catastrophic forgetting (Luo et al., 2024), such that only a few steps of gradient descent can be enough to undo previous behavioral alignment (Qi et al., 2023; Zhan et al., 2024). To retain the general ability to follow instructions, and to allow for control condition fine-tuned models in which the book text is not part of the training data, we include 3,500 instruction samples from the OpenHermes2.5 dataset on Huggingface (Teknum, 2023). Therefore the baseline without text-specific memory to compare with is not only the respective initial model before fine-tuning, but the same model when fine-tuned on the same 3,500 instruction samples but excluding the 116 samples of book chunks.

1559

1560

1561

1562

1563

1564

1565

**Overfitting does not lead to better performance on SORT.** To test whether overfitting on the texts would lead to better performance, we finetuned a Llama3-8b-Instruct model for 120 and 300 repetitions of the same chunks of book text (mixed with 300k and 30k instruction samples respectively). We found that training with many repetitions of the book text, resulting in a sharper decrease in perplexity, also did not support the ability to recall the order of segments. This still holds when a SORT-task readout layer is trained on 8 out of 9 books and then evaluated on the held-out book (see Figure 15). Even though the perplexity for *The Murder of Roger Ackroyd* dropped from 9.9 to 2.4, the performance for SORT did not increase.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

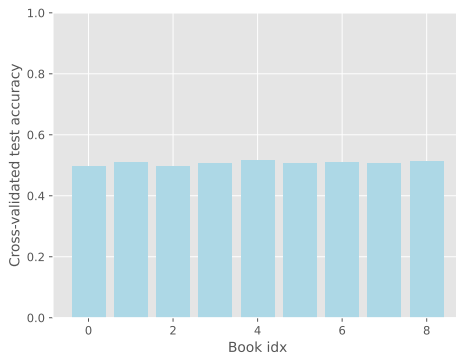


Figure 15: Crossvalidated task-readout performance for Llama3-8b-Instruct after finetuning on 120 repetitions of the book chunks, along with 300k instruction samples.

F.1 PERPLEXITY ANALYSIS OF FINE-TUNED MODELS

To confirm that fine-tuning on the books makes a model learn about the segments, we compare the perplexities of the two segments shown in SORT without source text presented in-context. We find that when the models are finetuned on data that includes the chunks of the books, they have a substantially lower perplexity for both segments, compared with the models fine-tuned only on the instruction data (see figure 16). Note that the scale of these perplexity values highlights that our task is likely out of distribution, presumably with little to no similar instruction data seen during pre-training and fine-tuning.

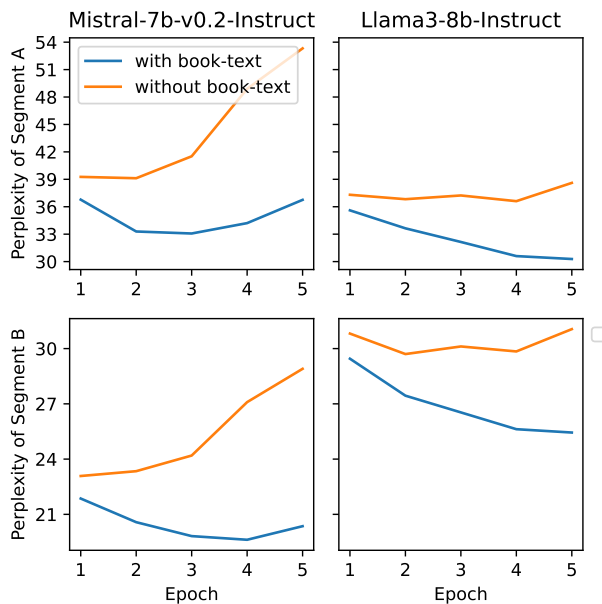


Figure 16: Perplexity of the two segments after fine-tuning of Mistral-7b-v0.2-Instruct and Llama3-8b-Instruct, when presented in the absence of in-context access to source excerpts.

F.2 COMPARISON OF SORT PERFORMANCE AFTER FINE-TUNING USING McNEMAR’S TEST

We find that even though the book-text finetuned Llama3-8b model has a form of memory of the books’ texts, the epoch-matched performance between the models fine-tuned without the book-chunks does not differ statistically for any epoch (Figure 17). For this analysis we use McNemar’s test

since we have an exact match of presented samples for both the memory-finetuned model and the baseline that does not form any memory of the text (Figure 16). We find high p-values, indicating no difference in performance between models fine-tuned with and without the book text (Figure 18), neither for Llama3-8b-Instruct, nor for Mistral-7b-v0.2-Instruct. Note that only epoch 1 qualifies to be tested for episodic memory since the books are only seen once (requirement 2 in Definition 1).

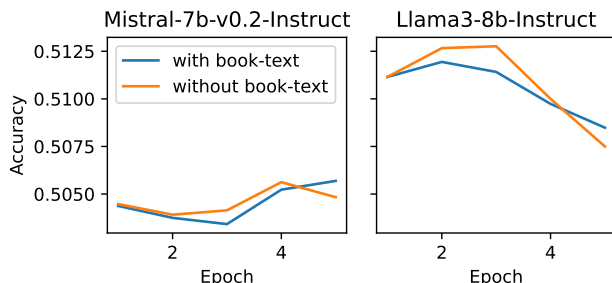


Figure 17: Accuracy of Llama3-8b-Instruct and Mistral-7b-v0.2-Instruct across epochs of finetuning on data including and excluding relevant book-text. Figure 18 shows that differences between accuracies shown here are not statistically significant ( $p>0.05$ ).

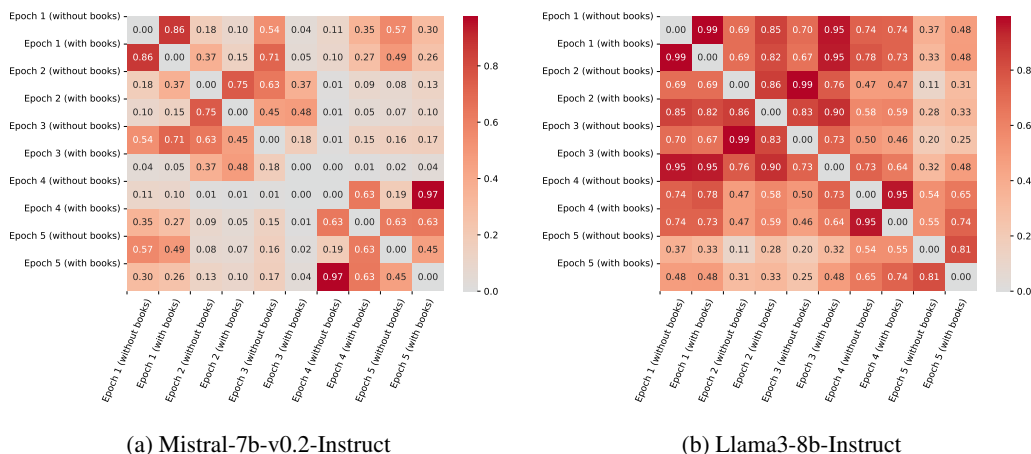


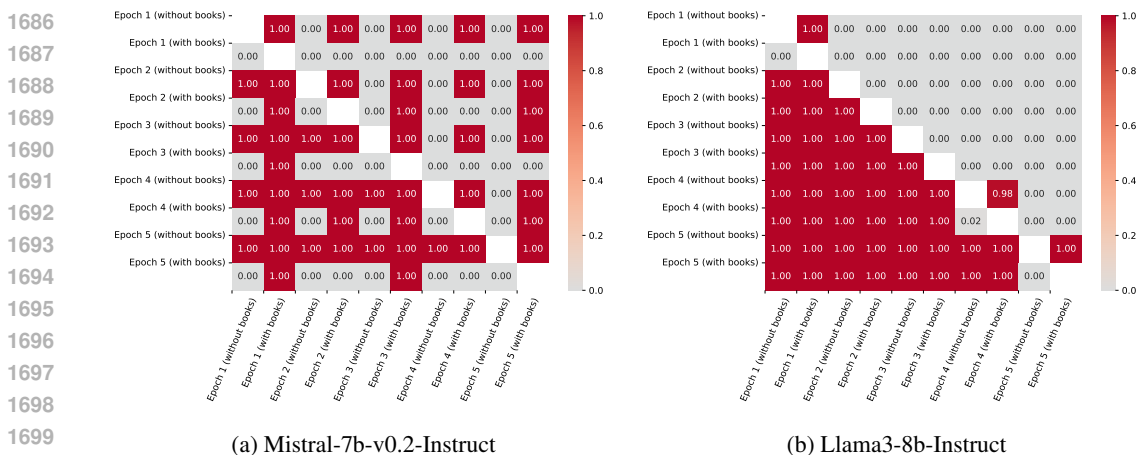
Figure 18: McNemar’s Test matrix of fine-tuned models performance. Shown are p-values indicating whether a model checkpoint (row) is different in its accuracy compared to another checkpoint (columns) with statistical significance. We fine-tuned with and without the books used in Book-SORT. There is no statistically significant difference between the models finetuned without and with book text. The effect of fine-tuning seems insignificant even without correcting these p-values for multiple comparisons.

### F.3 COMPARISON OF SORT PERFORMANCE AFTER FINE-TUNING USING A PAIRWISE T-TEST

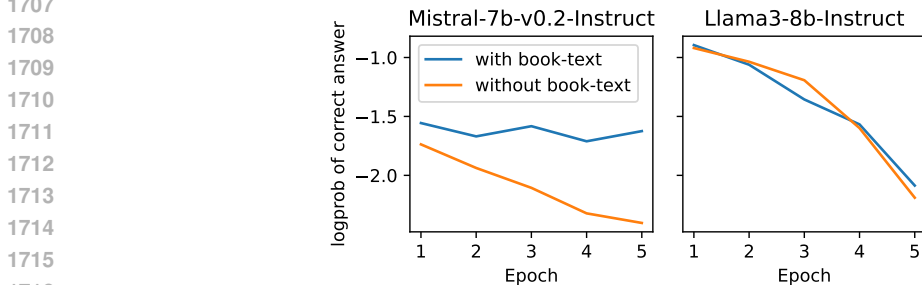
Testing the binary correctness evaluated based on a greedily sampled token does not allow us to draw conclusions about sub-threshold effects of fine-tuning on task performance. To test whether the models fine-tuned on the books is better than the models that are fine-tuned without chunks from the books, we performed a pairwise t-test on a continuous measure of accuracy based on the token log-probabilities. We compute the likelihood of the correct answer by taking the log ratio of the correct answer among all answers that can be mapped to either A or B, i.e. we are interested in  $\log\left(\frac{p(a=y)}{p(a=A)+p(a=B)}\right)$ , where  $y$  is the correct answer.

The results shown in figure 19 suggest that fine-tuned models do improve over the base model, with the book text condition performing better than the others after one epoch of training with statistical

1674 significance ( $p < 0.01$ ). Even though there is an effect, the magnitude is very small, as can be seen in  
 1675 Figure 20, and this positive effect could also be attributed to interleaving the instruction data with  
 1676 samples including longer texts (5,000 words) compared to just the instruction samples. Decreasing  
 1677 log-probability of the correct answer reflects catastrophic forgetting associated with training for  
 1678 multiple epochs on a small dataset, which violates requirement (2) of episodic memory (Definition  
 1679 1). Figure 19 (second row for both matrices) shows that log-probabilities after the first epoch is  
 1680 statistically significantly higher for the fine-tuned models that included the book-text compared to all  
 1681 other epochs with and without including the book-text in fine-tuning. However, Figure 20 shows that  
 1682 the difference in log-probabilities of the correct answer is small after the first epoch, indicating that  
 1683 inserted memory does not affect explicit question answering ability about temporal order memory  
 1684 (see Figure 17), which is one of the requirements for episodic memory (req. 5 in definition 1).  
 1685



1699 Figure 19: Pairwise t-test matrix of fine-tuned models. Shown are p-values indicating whether a  
 1701 model (row) has higher log probabilities of the correct answer compared to another model (columns)  
 1702 with statistical significance. Row 2 is significantly better than all other epochs with or without  
 1703 book-text.  
 1704



1718 Figure 20: Log-probability of the correct answer for fine-tuned models across epochs. Figure F.3  
 1719 shows statistical significance between conditions and epochs for this data.  
 1720

1721 F.4 TRAINING ON SORT?

1723 Our findings that fine-tuning with a language-modeling objective does not lead to temporal order  
 1724 memory might be due to the out-of-distribution property of the task - recalling the order of text is  
 1725 something that is not well covered in the models' training data distribution. Rather than implying that  
 1726 SORT is not suitable to test this type of memory, this highlights a lacking ability in current specific  
 1727 models to use parametric memory for this aspect of episodic memory: if they are asked about the  
 order of a document that they have been trained on, they can in fact not recall it. Directly training a



1728 model on parametric sequence order recall tasks might improve its performance, however we would  
1729 not expect this to generalize to other aspects of episodic memory (Appendix A). In this work, we  
1730 only evaluated existing models' capabilities (i.e. we did not modify models to become better at the  
1731 task), and we suggest that for SORT to be maximally informative about general episodic memory  
1732 capabilities of a model, it should only be used as a test-set task (not to be included in training and not  
1733 to be optimized for specifically in terms of architecture and hyperparameters). Training on SORT  
1734 risks that a less generalizable mechanism is used to specifically infer temporal order information,  
1735 without generalizing to other types of relations between parts of a sequence (see Appendix A). This  
1736 mirrors how we can explicitly include order information in RAG, which could also easily "solve"  
1737 SORT, but would not be informative about RAGs suitability to function as an episodic memory  
1738 insertion method for a model more broadly (see our discussion on RAG in section 6).

1739

#### 1740 F.5 IN-CONTEXT MEMORY PERFORMANCE OF FINE-TUNED MODELS

1741 Despite the inclusion of instruction data in fine-tuning, the accuracy with source excerpts presented  
1742 in-context of SORT decreased from 0.93 to 0.90 after a single epoch and to 0.88 after three epochs of  
1743 fine-tuning for Llama3-8b-Instruct. For the instruction-data only baseline of Llama3-8b-Instruct, the  
1744 performance degraded slightly less with an accuracy of 0.91 after the first epoch of fine-tuning.

1745

## 1746 G CODE AND DATA

1747

1748 Upon publication we will provide the code to create SORT datasets and evaluate models on SORT in  
1749 a public GitHub repository, along with the Book-SORT dataset used in this work. Our evaluation  
1750 code currently supports the OpenAI API, Huggingface Transformers (Wolf et al., 2020) and vLLM  
1751 (Kwon et al., 2023) for distributed inference.

1752

1753 Experiment data from our Book-SORT evaluation is located in a Google Drive folder, along with the  
1754 human experiment data. These will be made accessible openly through Huggingface datasets.

1755

1756 **License.** We make our code and data openly available under a permissive BSD-3 license for code.  
1757 Data including Book-SORT is available under a CC0 license.

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

1782 Table 12: Accuracy and Difference of Various Models on Multiple Books at Excerpt Lengths of 20  
 1783 and 50, with in-context memory (Part 2)  
 1784

| 1785 | <b>Model name</b>     | <b>Book</b> | <b>SORT S20</b> | <b>SORT</b> | <b>SORT-Extend S20</b> | <b>SORT-Extend S50</b> |
|------|-----------------------|-------------|-----------------|-------------|------------------------|------------------------|
| 1786 | Mixtral-8x7b-DPO-inst | 69087       | 0.86±0.10       | 0.87±0.13   | 0.63±0.18              | 0.49±0.14              |
| 1787 | Mixtral-8x7b-DPO-inst | 72578       | 0.88±0.10       | 0.90±0.10   | 0.63±0.18              | 0.57±0.14              |
| 1788 | Mixtral-8x7b-DPO-inst | 72600       | 0.89±0.10       | 0.91±0.10   | 0.63±0.18              | 0.58±0.15              |
| 1789 | Mixtral-8x7b-DPO-inst | 72869       | 0.90±0.09       | 0.92±0.10   | 0.61±0.17              | 0.55±0.15              |
| 1790 | Mixtral-8x7b-DPO-inst | 72958       | 0.90±0.09       | 0.93±0.09   | 0.57±0.16              | 0.57±0.15              |
| 1791 | Mixtral-8x7b-DPO-inst | 72963       | 0.89±0.10       | 0.92±0.10   | 0.56±0.16              | 0.55±0.15              |
| 1792 | Mixtral-8x7b-DPO-inst | 72972       | 0.89±0.10       | 0.91±0.09   | 0.55±0.16              | 0.54±0.15              |
| 1793 | Mixtral-8x7b-DPO-inst | 73017       | 0.87±0.10       | 0.91±0.10   | 0.55±0.14              | 0.54±0.14              |
| 1794 | Mixtral-8x7b-DPO-inst | 73042       | 0.87±0.10       | 0.91±0.09   | 0.57±0.14              | 0.55±0.14              |
| 1795 | Mixtral-8x22b-inst    | 69087       | 0.92±0.08       | 0.93±0.09   | 0.73±0.13              | 0.73±0.11              |
| 1796 | Mixtral-8x22b-inst    | 72578       | 0.92±0.08       | 0.95±0.07   | 0.76±0.12              | 0.76±0.12              |
| 1797 | Mixtral-8x22b-inst    | 72600       | 0.93±0.08       | 0.96±0.07   | 0.77±0.12              | 0.78±0.11              |
| 1798 | Mixtral-8x22b-inst    | 72869       | 0.93±0.08       | 0.97±0.07   | 0.78±0.12              | 0.80±0.11              |
| 1799 | Mixtral-8x22b-inst    | 72958       | 0.93±0.08       | 0.97±0.06   | 0.79±0.12              | 0.80±0.11              |
| 1800 | Mixtral-8x22b-inst    | 72963       | 0.92±0.09       | 0.97±0.06   | 0.78±0.12              | 0.78±0.12              |
| 1801 | Mixtral-8x22b-inst    | 72972       | 0.92±0.09       | 0.97±0.07   | 0.78±0.12              | 0.79±0.12              |
| 1802 | Mixtral-8x22b-inst    | 73017       | 0.93±0.09       | 0.97±0.07   | 0.78±0.12              | 0.79±0.12              |
| 1803 | Mixtral-8x22b-inst    | 73042       | 0.93±0.09       | 0.97±0.07   | 0.78±0.12              | 0.79±0.12              |
| 1804 | Mistral-v2-7b-inst    | 69087       | 0.85±0.10       | 0.87±0.11   | 0.64±0.15              | 0.66±0.13              |
| 1805 | Mistral-v2-7b-inst    | 72578       | 0.85±0.11       | 0.87±0.10   | 0.63±0.15              | 0.65±0.14              |
| 1806 | Mistral-v2-7b-inst    | 72600       | 0.86±0.11       | 0.87±0.10   | 0.64±0.14              | 0.67±0.14              |
| 1807 | Mistral-v2-7b-inst    | 72869       | 0.85±0.11       | 0.87±0.11   | 0.64±0.15              | 0.68±0.13              |
| 1808 | Mistral-v2-7b-inst    | 72958       | 0.86±0.10       | 0.88±0.11   | 0.65±0.15              | 0.68±0.14              |
| 1809 | Mistral-v2-7b-inst    | 72963       | 0.83±0.11       | 0.88±0.11   | 0.64±0.14              | 0.68±0.14              |
| 1810 | Mistral-v2-7b-inst    | 72972       | 0.84±0.11       | 0.88±0.10   | 0.63±0.14              | 0.68±0.14              |
| 1811 | Mistral-v2-7b-inst    | 73017       | 0.83±0.11       | 0.88±0.10   | 0.63±0.14              | 0.68±0.14              |
| 1812 | Mistral-v2-7b-inst    | 73042       | 0.83±0.11       | 0.88±0.10   | 0.63±0.14              | 0.68±0.14              |
| 1813 | Mistral-v1-7b-inst    | 69087       | 0.74±0.04       | 0.82±0.03   | /                      | /                      |
| 1814 | Mistral-v1-7b-inst    | 72578       | 0.75±0.04       | 0.81±0.03   | /                      | /                      |
| 1815 | Mistral-v1-7b-inst    | 72600       | 0.74±0.04       | 0.80±0.03   | /                      | /                      |
| 1816 | Mistral-v1-7b-inst    | 72869       | 0.74±0.04       | 0.81±0.03   | /                      | /                      |
| 1817 | Mistral-v1-7b-inst    | 72958       | 0.74±0.04       | 0.81±0.03   | /                      | /                      |
| 1818 | Mistral-v1-7b-inst    | 72963       | 0.74±0.04       | 0.80±0.03   | /                      | /                      |
| 1819 | Mistral-v1-7b-inst    | 72972       | 0.75±0.04       | 0.80±0.03   | /                      | /                      |
| 1820 | Mistral-v1-7b-inst    | 73017       | 0.74±0.04       | 0.80±0.03   | /                      | /                      |
| 1821 | Mistral-v1-7b-inst    | 73042       | 0.75±0.04       | 0.80±0.03   | /                      | /                      |
| 1822 | Gemma-1.1-7b-inst     | 69087       | 0.82±0.03       | 0.88±0.03   | /                      | /                      |
| 1823 | Gemma-1.1-7b-inst     | 72578       | 0.83±0.04       | 0.89±0.03   | /                      | /                      |
| 1824 | Gemma-1.1-7b-inst     | 72600       | 0.83±0.04       | 0.88±0.03   | /                      | /                      |
| 1825 | Gemma-1.1-7b-inst     | 72869       | 0.84±0.04       | 0.89±0.03   | /                      | /                      |
| 1826 | Gemma-1.1-7b-inst     | 72958       | 0.84±0.04       | 0.89±0.03   | /                      | /                      |
| 1827 | Gemma-1.1-7b-inst     | 72963       | 0.84±0.04       | 0.88±0.03   | /                      | /                      |
| 1828 | Gemma-1.1-7b-inst     | 72972       | 0.84±0.04       | 0.87±0.03   | /                      | /                      |
| 1829 | Gemma-1.1-7b-inst     | 73017       | 0.83±0.04       | 0.87±0.03   | /                      | /                      |
| 1830 | Gemma-1.1-7b-inst     | 73042       | 0.84±0.04       | 0.87±0.03   | /                      | /                      |
| 1831 | GPT-3.5-turbo         | 69087       | 0.86±0.03       | 0.88±0.03   | /                      | 0.69±0.04              |
| 1832 | GPT-3.5-turbo         | 72578       | 0.87±0.03       | 0.89±0.03   | /                      | 0.69±0.04              |
| 1833 | GPT-3.5-turbo         | 72600       | 0.87±0.03       | 0.89±0.03   | /                      | 0.67±0.04              |
| 1834 | GPT-3.5-turbo         | 72869       | 0.87±0.03       | 0.90±0.03   | /                      | 0.67±0.04              |
| 1835 | GPT-3.5-turbo         | 72958       | 0.87±0.03       | 0.90±0.03   | /                      | 0.67±0.04              |
|      | GPT-3.5-turbo         | 72963       | 0.86±0.03       | 0.89±0.03   | /                      | 0.67±0.04              |
|      | GPT-3.5-turbo         | 72972       | 0.86±0.03       | 0.88±0.03   | /                      | 0.67±0.04              |
|      | GPT-3.5-turbo         | 73017       | 0.85±0.03       | 0.88±0.03   | /                      | 0.67±0.04              |
|      | GPT-3.5-turbo         | 73042       | 0.85±0.03       | 0.88±0.03   | /                      | 0.67±0.04              |

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

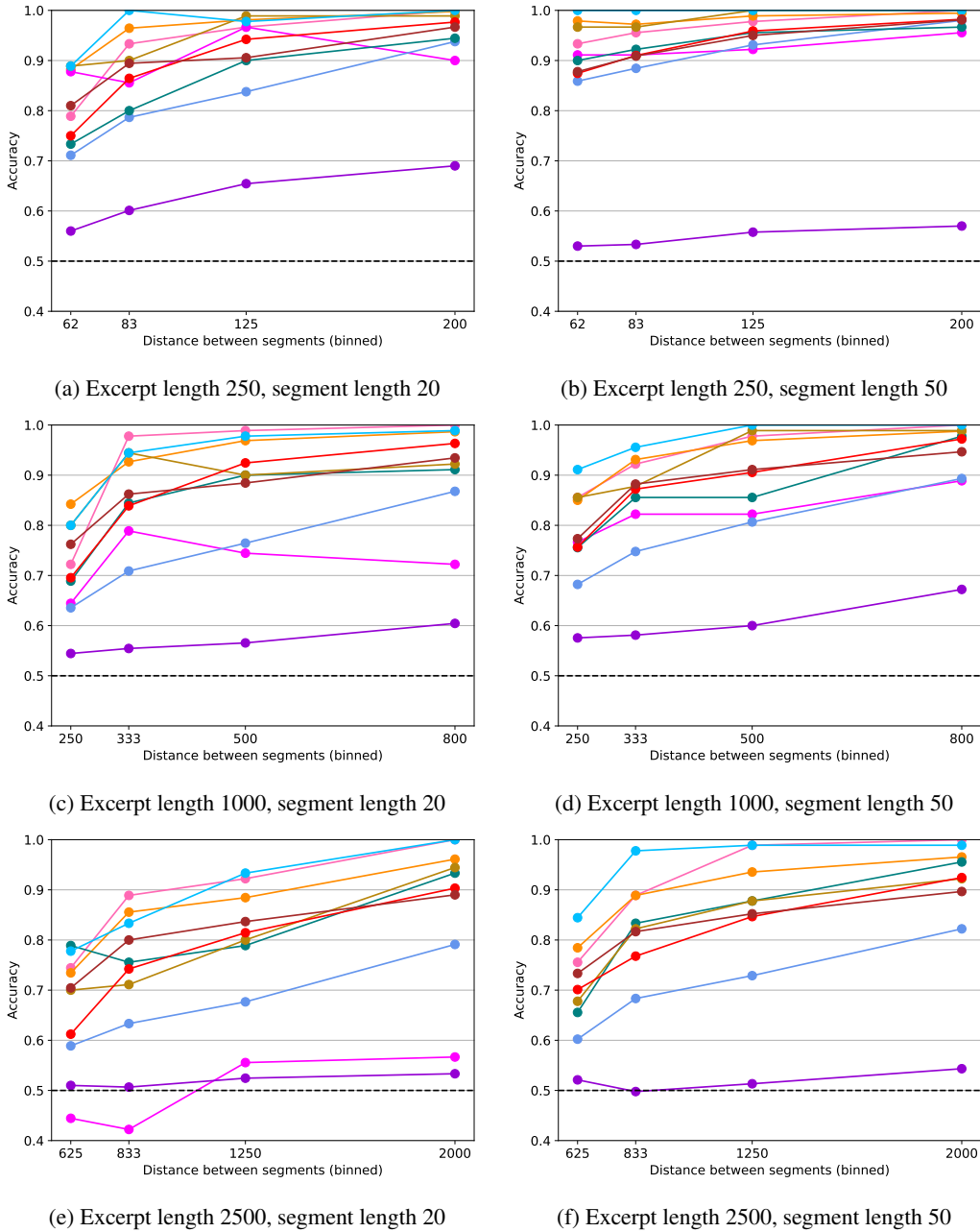
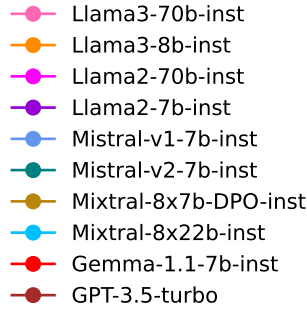


Figure 8: Average accuracy by distance between segments (All excerpt length), part A.

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

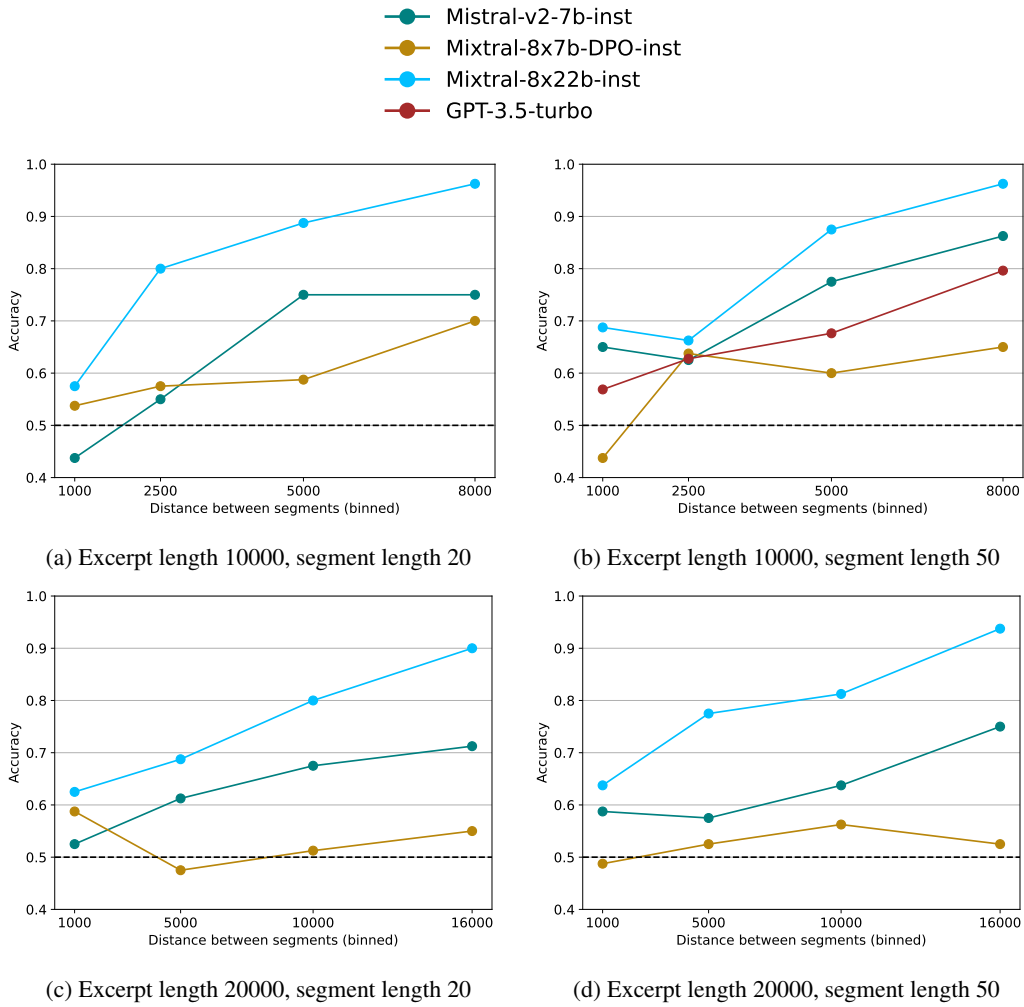


Figure 9: Average accuracy by distance between segments (All excerpt length), part B.

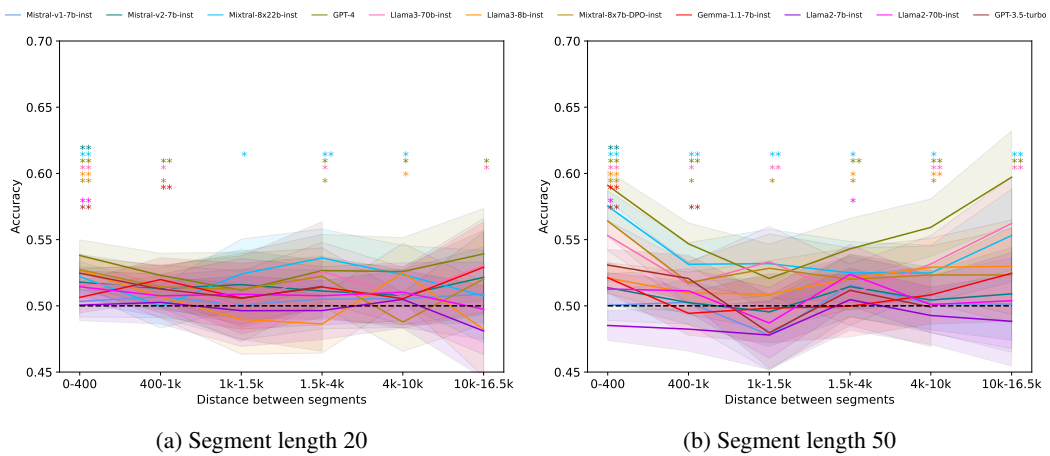
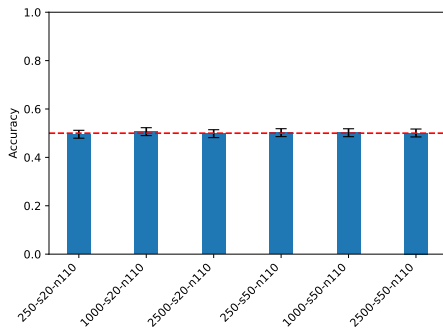
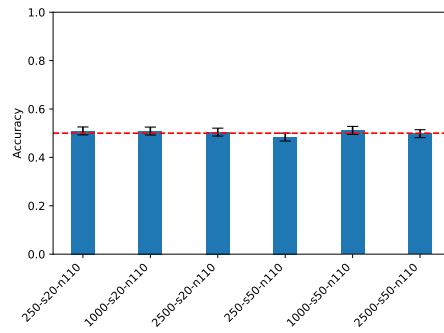


Figure 13: Baseline model performance on SORT without text-specific memory by segment distance (95% bootstrapped confidence interval). Significant difference from chance is marked with asterisks (\*p-value<0.05,\*\*p-value<0.01).

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997



(a) Llama3-8b



(b) Mistral-7b

Figure 14: Base model performance for SORT (in-context memory).