

# WORLD MODEL ON MILLION-LENGTH VIDEO AND LANGUAGE WITH BLOCKWISE RINGATTENTION

**Anonymous authors**

Paper under double-blind review

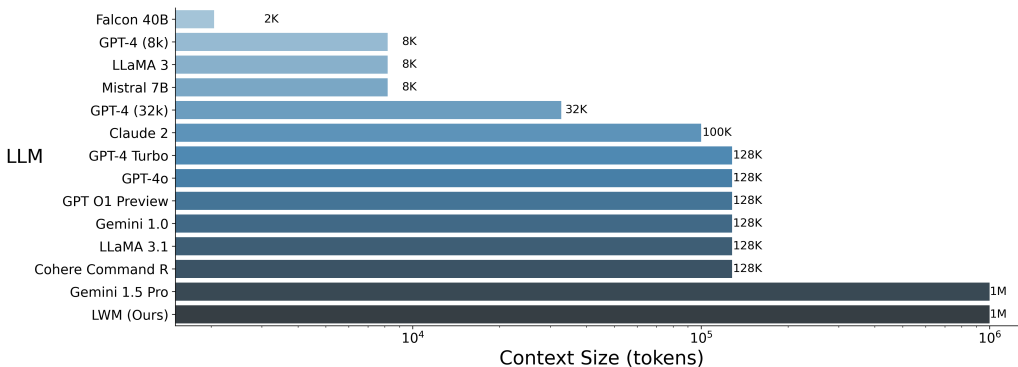
## ABSTRACT

Enabling long-context understanding remains a key challenge in scaling existing sequence models – a crucial component in developing generally intelligent models that can process and operate over long temporal horizons that potentially consist of millions of tokens. In this paper, we aim to address these challenges by providing a comprehensive exploration of the full development process for producing 1M context language models and video-language models, setting new benchmarks in language retrieval and new capabilities in long video understanding. We detail our long context data curation process, progressive context extension from 4K to 1M tokens, and present an efficient open-source implementation for scalable training on long sequences. Additionally, we open-source a family of 7B parameter models capable of processing long text documents and videos exceeding 1M tokens.

## 1 INTRODUCTION

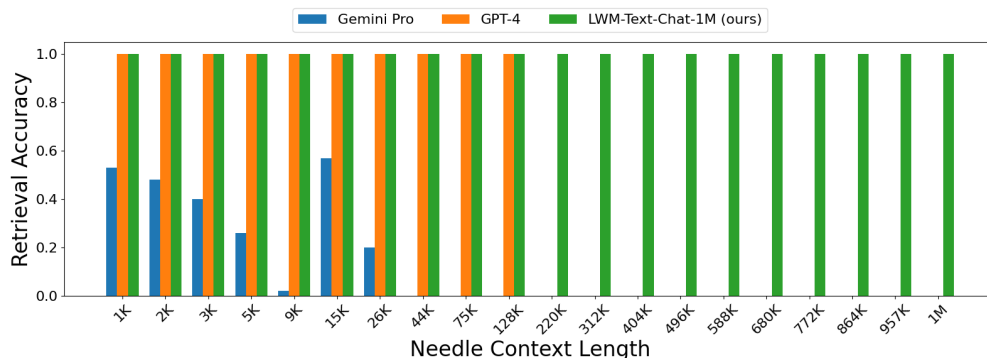
Enabling long-context understanding remains a key challenge in scaling existing sequence models—a crucial step toward developing generally intelligent models that can process and operate over extended temporal horizons, potentially involving millions of tokens. Current modeling approaches are predominantly limited to processing short sequences, whether in the form of language, images, or video clips (Brown et al., 2020; Touvron et al., 2023a;b; OpenAI, 2023; Brooks et al., 2024; Team et al., 2023). As a result, these models fall short when tasked with understanding complex, long-form language and visual contexts.

However, training models to process sequences that exceed millions of tokens is a significant challenge due to the high memory and computational costs, as well as the lack of long-context data. In this work, we address these challenges by leveraging Blockwise RingAttention (Liu et al., 2024; Liu and Abbeel, 2023), a technique that scales context size without approximations or overheads, enabling efficient training on long sequences. We curate an extensive dataset of long-form videos and books from public sources, covering a wide variety of activities and narrative structures. To address the scarcity of long-form conversational datasets, we developed a model-based question-answering technique, where a short-context model generates training data from books, significantly enhancing the model’s chat



**Figure 1 Comparison of context size in state-of-the-art LLMs.** Our model and concurrent work Gemini 1.5 both achieve a 1M context size, significantly outperforming other LLMs.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107



**Figure 2 Retrieval comparisons against Gemini Pro and GPT-4.** Needle retrieval comparisons against Gemini Pro and GPT-4 for each respective max context length – 32K and 128K. Our model performs competitively while being able to extend to 8x longer context length. Note that in order to show fine-grained results, the x-axis is log-scale from 0-128K, and linear-scale from 128K-1M.

capabilities over long sequences. To mitigate computational costs, we gradually extended context size from an initial 4K tokens to 1M tokens, achieving a cost-effective and scalable approach for long-context modeling.

Following this, we further train our long-context language model to incorporate visual modalities, such as image and video. Contrary to existing popular vision-language models (Liu et al., 2023a; OpenAI, 2023; Chen et al., 2023a), we opt to additionally optimize next-token prediction losses for image and video (generation) with a VQGAN (Esser et al., 2021) encoder. We encountered various challenges training on mixed modalities (video, image, text). To balance their unique characteristics - sequential information, visual detail, and linguistic content - we implement an efficient masked sequence packing strategy, as well as introduce careful loss balancing to retain short context accuracy. This approach handles varying sequence lengths more effectively than standard methods. We also optimized the ratio of image, video, and text inputs in each batch, proposing an empirically effective balance for cross-modality learning. Since our model aims to model both textual and visual projections of the world through a large context window, drawing inspiration from prior work on world models (Brooks et al., 2024; Ha and Schmidhuber, 2018), we name our work as Large World Model (LWM).

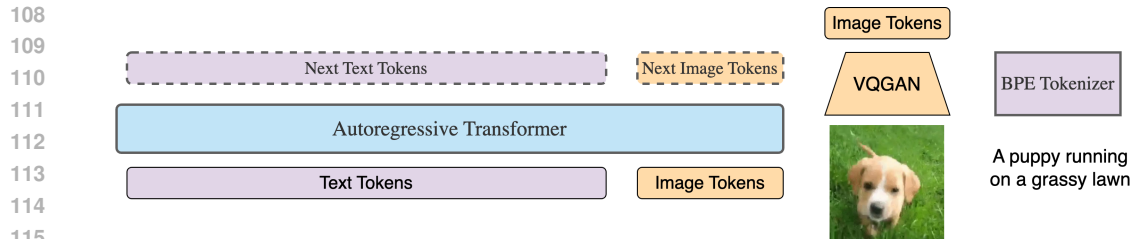
Our contributions are threefold: (a) we train one of the largest context size transformers to date on long text documents and videos and achieved competitive results on long video understanding and long context fact retrieval. (b) We discover a range of challenges associated with training on long sequences and propose solutions for them: masked sequence packing to effectively train with different sequence lengths and synthetic model-generating question-answering for effective attention. (c) We provide an open-source and optimized implementation for training with millions of tokens in context, as well as a family of Llama-based 1M context models capable of processing long documents (LWM-Text, LWM-Text-Chat) and videos (LWM, LWM-Chat) of 1M tokens.

## 2 METHOD OVERVIEW

We train a large autoregressive transformer model with a large context window of up to one million tokens, building upon Llama2 7B (Touvron et al., 2023b). To achieve this goal, we implement a two-stage training strategy. In Stage I (Section 3), we extend the context to 1M tokens using book-length texts. This is followed by Stage II (Section 4), where we conduct joint training on diverse long multimodal sequences, incorporating text-image data, text-video data, and book-length texts. Our model architecture is the standard autoregressive transformer design, as illustrated in Figure 3. For a comprehensive overview of our training stages and the datasets employed, please refer to Figure 4.

## 3 STAGE I: LEARNING LONG-CONTEXT LANGUAGE MODELS

This stage aims at first developing LWM-Text and LWM-Text-Chat, a set of long-context language models learned by training on progressively increasing sequence length data, and modifying positional



**Figure 3 Model Architecture.** The LWM model is an autoregressive transformer trained on sequences of multimodal tokens. Each video frame is tokenized into 256 tokens using VQGAN, while text is processed using a Byte-Pair Encoding (BPE) tokenizer. These tokens—both image and text—are combined and input into the transformer to autoregressively predict the next token. The model can handle various input-output modalities, including text, image, video, and text-video pairs. To distinguish between images and text, special tokens `<vision>` and `</vision>` are used for image and video frames, with `<eof>` and `<eov>` marking the end of these sequences. For simplicity, delimiters are not shown in the figure.

encoding parameters to account for longer sequence lengths (see Section 3.1). In Section 3.2, we show how to construct model-generated question-answering data for enabling long sequence conversations.

### 3.1 PROGRESSIVE TRAINING TOWARDS LONG CONTEXT

Learning long-range dependencies over sequences of millions of tokens requires (1) memory efficient training to scale to such long sequences, as well as a need to (2) compute efficient training to extend the context of our base language model. We outline our approach to these challenges, detailing our methods for training on long sequences, designs for efficiency and stability, and experimental setup.

Training on long sequences has become prohibitively expensive due to memory constraints imposed by the quadratic complexity of attention weight computations. To address these computational limitations, we leverage recent advancements in scaling context window size, particularly Blockwise RingAttention (Liu et al., 2024). This approach theoretically allows for an infinite context, bounded only by available devices. We further enhance performance by fusing it with FlashAttention (Dao et al., 2022) using Pallas (Bradbury et al., 2018) to optimize performance compared with using XLA compiler. Notably, with enough tokens per device—already a given—the communication cost during sequence parallelism is fully overlapped by computation, resulting in no additional overhead.

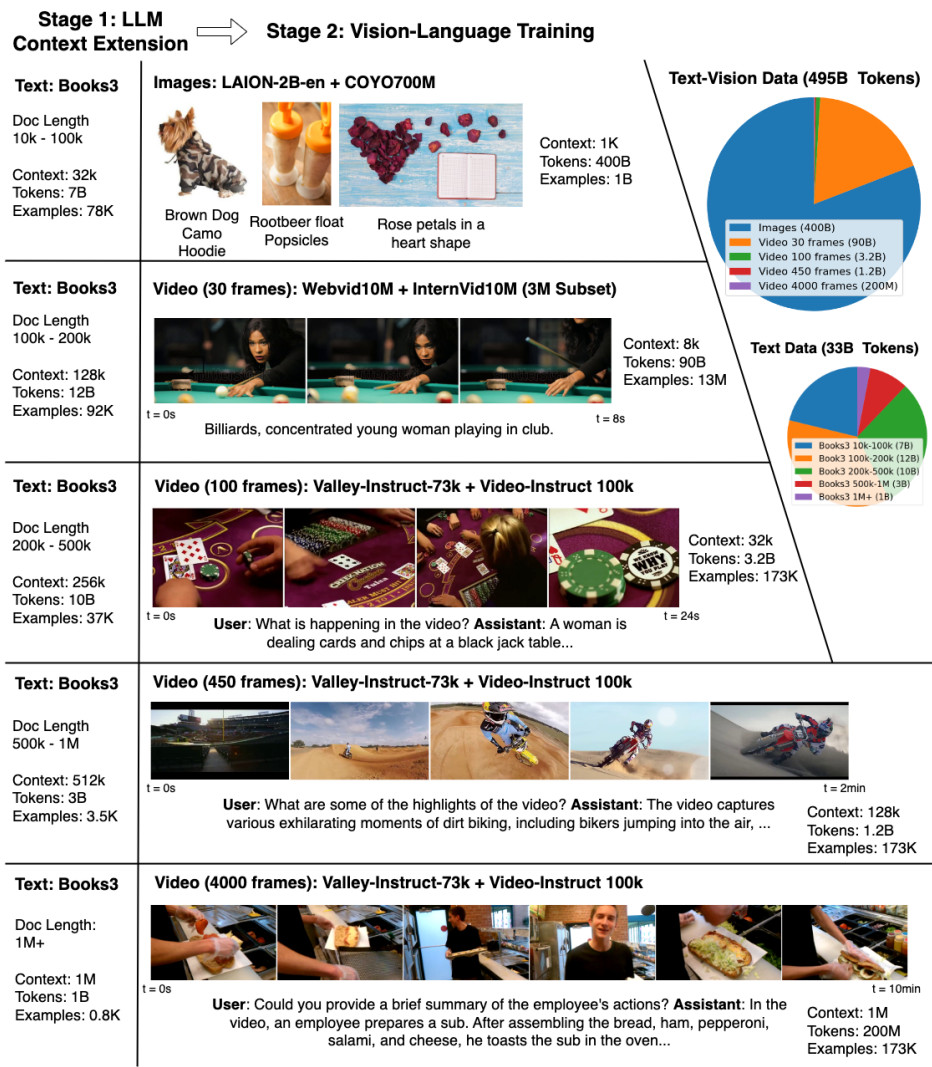
For better efficiency, we adopt a training approach inspired by prior research on extending context (Jin et al., 2023a), where our model is trained on progressively longer sequence lengths, starting from 32K tokens and ending at 1M tokens in increasing powers of two. Intuitively, this allows the model to save compute by first learning shorter-range dependencies before moving onto longer sequences. For extending positional embeddings to longer contexts, we adopt a simple, scaled-up version of the approach explored in Rozière et al. (2023), where the  $\theta$  parameter for RoPE (Su et al., 2024) is scaled in proportion to the context length. We found this approach to be stable for extending positional embeddings with larger context lengths due to its simplicity, requiring the tuning of only a single hyperparameter. Specifically, we scale the  $\theta$  parameter for RoPE alongside increases in context window sizes – the values are shown in Table 6. The progressive training of growing context sizes is shown in Figure 4.

We initialize from LLaMA-2 7B (Touvron et al., 2023b) as base language model and progressively increase the effective context length of the model across 5 stages: 32K, 128K, 256K, 512K, and 1M. For each stage, we train on different filtered versions of the Books3 dataset from The Pile (Gao et al., 2020). Table 6 details information about each training stage, such as the number of tokens, total time, and the Books3 dataset filtering constraints. Each successive run is initialized from the prior sequence length.

### 3.2 MODEL-GENERATED QUESTION-ANSWERING FOR EFFECTIVE CONTEXT

We construct a simple question-answering dataset to develop long-context chat capabilities. First, we split documents from the Books3 dataset into fixed chunks of 1,000 tokens, feed each chunk into our short-context language model, and prompt it to generate a question-answer pair based on the

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215



**Figure 4 Curated dataset and training process with progressively increasing data length and complexity.** The diagram outlines a two-stage training process. Stage 1 extends text-based understanding using books datasets of increasing document lengths and token counts. Stage 2 integrates vision-language training. Pie charts display token distribution, showing that images and short-frame videos dominate visual data, while mid-length text examples lead in the text corpus.

content. To create longer examples (e.g., 32K tokens), we concatenate adjacent chunks and append the relevant question-answer pairs toward the end of the sequence in a chat format. The key intuition is that the model must learn to focus on any part of the context to answer the questions, as the relevant information can appear anywhere within the sequence.

For chat fine-tuning, we train each model on a mix of the UltraChat conversation dataset (Ding et al., 2023) and our custom question-answering dataset, using approximately a 7:3 ratio. We found it crucial to pre-pack the UltraChat data to the training sequence length and keep these examples separate from our question-answering data. This separation is necessary because UltraChat data generally contains a much higher proportion of loss tokens (due to densely packed, short questions in chat), whereas our question-answering data has long questions in chat thus a significantly lower percentage of loss tokens per sequence (< 1%). This difference arises from the long documents in the given context of our question-answering data, which are not included in loss calculations. Table 7 provides further training details for each run. Notably, we do not employ progressive training for any of the chat models; instead, we initialize them from their respective pretrained models at the same context length.



**Summary:** Stage I progressively increase sequence lengths using our curated dataset: starting with 32K tokens and gradually scaling up to 1M tokens. Model-generated question-answering data aids in learning effective long context.

### 3.3 LANGUAGE EVALUATION RESULTS

#### 3.3.1 SHORT CONTEXT TASKS

Table 1 presents a comparative analysis between the Llama2-7B model with a 4K context and its context-expanded counterparts, ranging from 32K to 1M. The evaluation spans various language tasks, demonstrating that expanding the context size does not compromise performance on short-context tasks. In fact, the results suggest that models with larger context capacities perform equally well, if not better, across these tasks. This evidence indicates the absence of negative effects from context expansion, highlighting the models’ capability to adapt to different task requirements without losing efficiency in shorter contexts.

**Table 1** Performance evaluation across language tasks, comparing Llama-2 7B (4K context window) and context-expanded variants of LWM-Text (32K to 1M). The results demonstrate that increasing context length does not significantly degrade performance on tasks with shorter contexts.

| Task / Metric          | Llama-2 7B | LWM-Text |      |      |      |      |
|------------------------|------------|----------|------|------|------|------|
|                        |            | 32k      | 128k | 256k | 512k | 1M   |
| arc_challenge/acc      | 0.40       | 0.43     | 0.45 | 0.44 | 0.44 | 0.43 |
| arc_challenge/acc_norm | 0.43       | 0.47     | 0.47 | 0.46 | 0.46 | 0.46 |
| hellaswag/acc          | 0.57       | 0.57     | 0.57 | 0.57 | 0.56 | 0.57 |
| hellaswag/acc_norm     | 0.77       | 0.76     | 0.76 | 0.76 | 0.75 | 0.75 |
| mmlu                   | 0.39       | 0.4      | 0.41 | 0.41 | 0.36 | 0.35 |
| openbookqa/acc         | 0.32       | 0.33     | 0.31 | 0.32 | 0.33 | 0.30 |
| openbookqa/acc_norm    | 0.44       | 0.44     | 0.44 | 0.43 | 0.41 | 0.41 |

#### 3.3.2 RETRIEVAL TASK: SINGLE INFORMATION

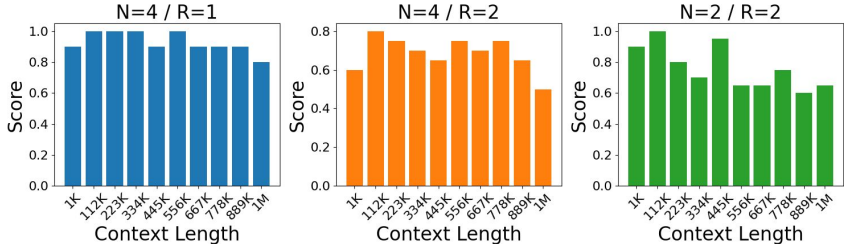
We evaluate on the popular Needle In A Haystack task (gkamradt, 2023) – more specifically an version (ArizeAI, 2023) that finds and retrieves random numbers assigned to randomized cities from the context. Figure 2 shows that we can scale to far larger contexts compared to the current best available LLMs. Figure 11 in Appendix shows nearly perfect retrieval accuracy over the entire context of our 1M context model. Appendix C shows more single needle retrieval results for our other shorter context length models.

#### 3.3.3 RETRIEVAL TASK: MULTIPLE INFORMATION

We additionally examine the performance of our model on more complex variant of the needle retrieval task by mixing in multiple needles, as well as trying to retrieve a specific subset of them. Figure 5 shows multi-needle retrieval results under different settings. Our model generalizes well when retrieving a single needle from multiple needles in context, with slight degradation when asked to retrieve more than one needle. Table 2 shows multi-needle comparisons, where our model is able to perform competitively or better than GPT-4 at retrieving one needle, or slightly lower performance when retrieving more than one needle. Furthermore, our model is also able to perform well and extend to longer context lengths of up to 1M tokens and far outperforms any recent shorter context baselines applies to longer sequence lengths through positional extrapolation techniques.. However, we note that we see degradation in accuracy while increasing the difficulty of the needle retrieval task, suggesting that there is still more room to improve on the 1M context utilization of our model. We believe that our released model will provide a foundation for future work on developing longer context models, as well as encourage more challenging benchmarks that contain difficult long-range tasks that require higher levels of synthesis, rather than pure fact retrieval.

**Table 2** Multi Needle in a Haystack. \* denotes models **after** the completion of this paper.

| Context Length | Model                     | $N = 2, R = 2$ | $N = 4, R = 1$ | $N = 4, R = 2$ |
|----------------|---------------------------|----------------|----------------|----------------|
| 32K            | Gemini Pro (02/23)        | 0.34           | 0.44           | 0.6            |
|                | GPT-4-1106                | 0.97           | 0.95           | 0.9            |
|                | Llama-3.1-8B-Instruct*    | 0.87           | 0.95           | 0.93           |
|                | Qwen2.5-7B-Instruct*      | <b>1.0</b>     | <b>1.0</b>     | <b>0.97</b>    |
|                | Mistral-7B-Instruct-v0.3* | 0.98           | 0.85           | 0.83           |
|                | <b>LWM-Text-1M (Ours)</b> | 0.84           | 0.97           | 0.84           |
| 128K           | Gemini Pro (02/23)        | -              | -              | -              |
|                | GPT-4-1106                | 0.92           | 0.8            | 0.82           |
|                | Llama-3.1-8B-Instruct*    | <b>0.98</b>    | 0.91           | 0.87           |
|                | Qwen2.5-7B-Instruct*      | <b>0.98</b>    | 0.80           | <b>0.90</b>    |
|                | Mistral-7B-Instruct-v0.3* | 0.85           | 0.75           | 0.68           |
|                | <b>LWM-Text-1M (Ours)</b> | 0.83           | <b>0.98</b>    | 0.83           |
| 1M             | Gemini Pro (02/23)        | -              | -              | -              |
|                | GPT-4-1106                | -              | -              | -              |
|                | Llama-3.1-8B-Instruct*    | 0.27           | 0.32           | 0.18           |
|                | Qwen2.5-7B-Instruct*      | 0.0            | 0.0            | 0.0            |
|                | Mistral-7B-Instruct-v0.3* | 0.05           | 0.13           | 0.10           |
|                | <b>LWM-Text-1M (Ours)</b> | <b>0.67</b>    | <b>0.84</b>    | <b>0.69</b>    |



**Figure 5** Multiple needles retrieval task with LWM-1M.  $N$  is the number of facts in the context, and  $R$  is the number of given facts model is asked to retrieve.

### 3.3.4 EVALUATION ON LOFT

**Table 3** Evaluations on some benchmarks in the LOFT dataset.

| Setting: 512K Context | LWM (512K)  | GPT-4o (128K) | Claude 3 Opus (200K) |
|-----------------------|-------------|---------------|----------------------|
| Quora                 | <b>0.38</b> | 0.23          | 0.37                 |
| NQ                    | <b>0.37</b> | 0.22          | 0.37                 |
| HotPotQA              | <b>0.72</b> | 0.21          | 0.32                 |

We further evaluate our model on a coverage of the LOFT (Lee et al., 2024) dataset collection, we provides a more natural set of benchmarks that examine capabilities for long-context models in the context of document retrieval, and RAG. The benchmark includes tasks such as duplication detection (Quora <sup>1</sup>), document retrieval (HotpotQA (Yang et al., 2018)), and retrieval-based question-answering (NQ). Each dataset contains a corpus of 1000s of documents, and the model is asked to retrieve a set of document ids pertaining to its specific task (Quora, HotpotQA). For RAG (NQ dataset), the model is asked to answer the question using the given context. Table 3 shows evaluations results on 512K context length against various language model baselines.

**Takeaway:** Long context capability enables LWM to outperform state-of-the-art text models at multiple benchmarks. This demonstrates the effectiveness of our methods for enabling long context.

## 4 STAGE II: EXTENDING TO LONG-CONTEXT VISION-LANGUAGE

Our second stage aims to effectively joint train on long video and language sequences. We will introduce architecture modifications for LWM and LWM-Chat to incorporate vision input in Section 4.1.

<sup>1</sup><https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>

324 Training on varying sequence lengths is discussed in Section 4.2. The evaluation results are shown  
 325 in Section 4.3. In this phase, we enhance the capabilities of the previously developed 1M context  
 326 language model, by finetuning it on vision-language data of various lengths. The datasets used and  
 327 the steps involved in the training process are illustrated in Figure 4.  
 328

#### 330 4.1 ARCHITECTURAL MODIFICATIONS FOR VISION

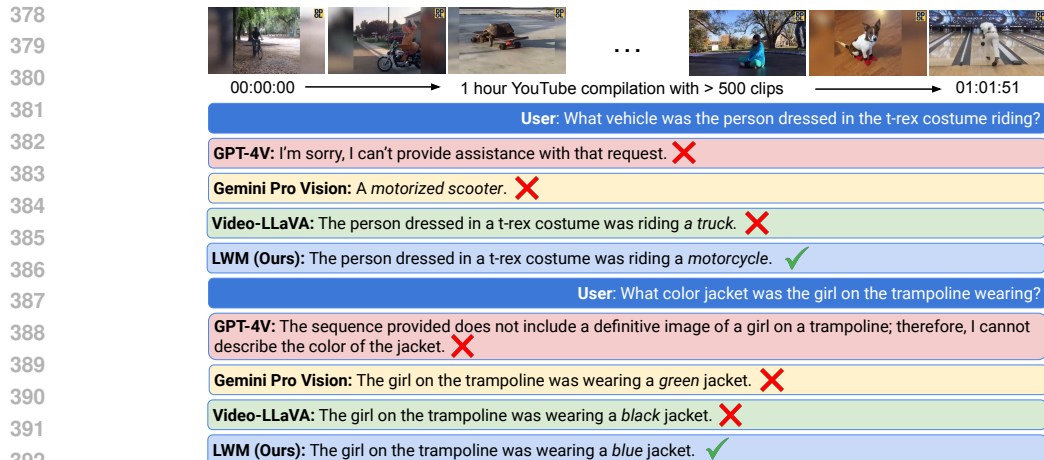
331  
 332 We use the pretrained VQGAN (Esser et al., 2021) from aMUSEd (Patil et al., 2024) that tokenizes  
 333  $256 \times 256$  input images to  $16 \times 16$  discrete tokens. Videos are tokenized by applying the VQGAN  
 334 per-frame, and concatenating the codes together. In order to distinguish between modalities when  
 335 generating, as well as knowing when to switch, we introduce mechanisms to mark the end of text  
 336 generation / beginning of vision generation, and vice-versa. For defining the end of vision generation,  
 337 we introduce new tokens, `<eof>` and `<eov>`, that represent end of frame (at the end of each  
 338 video frame that is not the last video frame in the sequence), and end of vision (at the end of each  
 339 single image, or at the end of the last frame in a video) boundaries respectively. For defining the  
 340 end of text generation, we wrap the vision tokens with `<vision>` and `</vision>` (as text) text  
 341 tokens. The model is trained with interleaved concatenations of vision and text tokens, and predicted  
 342 autoregressively (see Figure 3).  
 343

#### 345 4.2 TRAINING STEPS

344  
 347 We initialize from our LWM-Text-1M text model, and perform a similar process of progressive  
 348 training on a large amount of combined text-image and text-video data, with the exception that we do  
 349 not additionally scale RoPE  $\theta$ , as it already supports up to 1M context. Table 8 shows details for each  
 350 training stage, where the model is initialized from the prior shorter sequence length stage. For each  
 351 stage, we train on the following data:

- 352 • **LWM-1K**: We train on large set of text-image dataset comprising of a mix of LAION-2B-en (Schuh-  
 353 mann et al., 2022) and COYO-700M (Byeon et al., 2022). The datasets were filtered to only include  
 354 images with at least 256 resolution – in total roughly 1B text-image pairs. During training, we  
 355 concatenate the text-image pairs and randomly swap the order of the modalities to model both  
 356 text-image generation, unconditional image generation, and image captioning. We pack text-image  
 357 pairs to sequences of 1K tokens.
- 358 • **LWM-8K**: We train on a text-video dataset mix of WebVid10M (Bain et al., 2021) and 3M Intern-  
 359 Vid10M (Wang et al., 2023) examples. Similar to prior works (Ho et al., 2022a;b; Villegas et al.,  
 360 2022), we jointly train on both images and video with a 50-50 ratio of each modality. We pack  
 361 images to sequences of 8K tokens, and 30 frame videos at 4FPS. Similar to image training, we  
 362 randomly swap the order of modalities for each text-video pair.
- 363 • **LWM-Chat-32K/128K/1M**: For the final 3 stages, we train on a combined mix of chat data  
 364 for each downstream task: (1) text-image generation, (2) image understanding, (3) text-video  
 365 generation, and (4) video understanding. We construct a simple version of text-image and text-  
 366 video chat data by sampling random subsets of the pretraining data augmented with chat format.  
 367 For image understanding, we using the image chat instruct data from ShareGPT4V (Chen et al.,  
 368 2023a). Lastly, for the video understanding chat data, we use a combined mix of Valley-Instruct-  
 369 73K (Luo et al., 2023) and Video-ChatGPT-100K instruct data (Maaz et al., 2023). For all short  
 370 context data (image generation, image understanding, video generation), we pack sequences to the  
 371 training context length. During packing, we found it crucial to mask out the attention so that each  
 372 text-vision pair only attends to itself, as well as re-weighting losses to make computation identical  
 373 to training in a non-packed + padding training regime. For video understanding data, we uniformly  
 374 sample a max number of frames to fit the training context length of the model if the video is too  
 375 long. During training, We allocate 25% of each batch to each of the 4 downstream tasks.

376 For the first two stages of training (LWM-1K and LWM-8K), we additionally mix 16% of the batch  
 377 to be pure text data from OpenLLaMA (Geng and Liu, 2023), as we found it beneficial to preserve  
 language capabilities while training on vision data.



**Figure 6 LWM excels in answering questions about a 1-hour YouTube video.** This figure compares LWM-Chat-1M with proprietary models like Gemini Pro Vision and GPT-4V, along with open-source models. The test involves answering questions based on an hour-long YouTube compilation containing over 500 video clips. LWM demonstrates superior performance in providing accurate answers requiring comprehension of extended video content.

**Table 4** Long Video-MME Benchmark. \* denotes models **after** the completion of this paper.

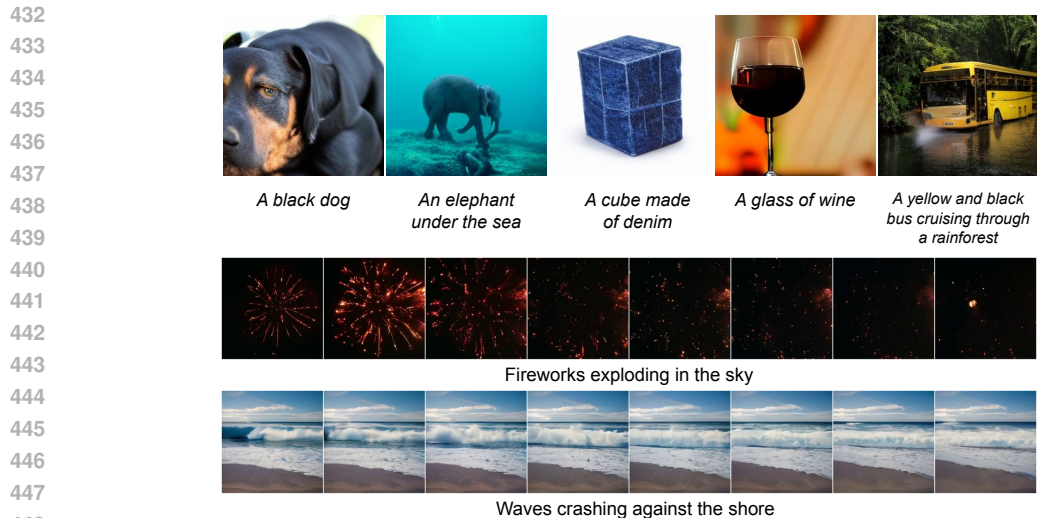
| Method          | Parameters | Frames      | Medium (4min-15min) | Long (30min-60min) |
|-----------------|------------|-------------|---------------------|--------------------|
| Gemini 1.5 Pro* | Unknown    | $\leq 1800$ | 74.3                | 67.4               |
| GPT-4o*         | Unknown    | 384         | 70.3                | 65.3               |
| LLaVA-Video*    | 72B        | 64          | 68.9                | 61.5               |
| VideoLLaMA 2*   | 72B        | 32          | 59.9                | 57.6               |
| Long-LLaVA*     | 7B         | 64          | 51.4                | 45.4               |
| Video-LLaVA     | 7B         | 8           | 38.1                | 36.2               |
| LWM-1M          | 7B         | $\leq 1800$ | 63.7                | 60.8               |

**Summary:** Stage II training incorporates image and video. Building on Stage I, it gradually increases sequence lengths of vision and text input. Importantly, we found our masked sequence packing and mixing synthetic and chat data crucial to retain short context performance during our progressive training. [Appendix B shows ablations when not using our training method on instruction-following and text-image understanding benchmarks.](#)

### 4.3 VISION-LANGUAGE EVALUATION RESULTS

#### 4.3.1 LONG VIDEO UNDERSTANDING

Although vision-language model (Lin et al., 2023; OpenAI, 2023; Team et al., 2023) can ingest long videos, this is commonly done by performing large temporal subsampling of video frames due to limited context length. For example, Video-LLaVA (Lin et al., 2023) is restricted to uniformly sampling 8 frames from a video, no matter how long the original video may be. As such, models may lose more fine-grained temporal information that is important for accurately answering any questions about the video. In contrast, our model is trained on long sequences of 1M tokens, and as a result, can simultaneously attend thousands of frames of videos to retrieve fine-grained information over short time intervals. [Table 4 shows long video evaluations on the Video-MME \(Fu et al., 2024\) benchmark, demonstrating our model as the best performing model among its size class.](#) Figure 6 shows an example of our model correctly answering questions about a long, 1-hour YouTube compilation consisting of more than 500 individual clips. Our baseline methods, on the other hand, generally have difficulty answering the questions due to a limited number of frames. More results are shown in Figure 18 and Appendix F.



449 **Figure 7** LWM’s ability to generate both static images and dynamic videos from text is shown. The  
450 top row illustrates image, while the bottom rows show video.

#### 451 4.3.2 IMAGE UNDERSTANDING AND SHORT VIDEO UNDERSTANDING

452 We evaluate LWM on standard benchmarks for image and short video understanding, with results  
453 presented in Table 5. Our model performs comparably to baselines but falls short of state-of-the-art  
454 (SOTA) models. This performance gap is not unexpected, given that SOTA models leverage vision  
455 backbones that have undergone extensive CLIP training (Radford et al., 2021). In contrast, LWM  
456 utilizes discrete tokens from an off-the-shelf model (Patil et al., 2024). Discrete tokens result in greater  
457 information loss, particularly for OCR-like textual data, compared to continuous CLIP embeddings.  
458 Moreover, our model learns text-image alignment from scratch, while CLIP-based models benefit  
459 from large-scale pretraining. This work primarily focuses on long-context methodology, and we  
460 defer additional training to future work due to computational constraints. A straightforward approach  
461 to improving benchmark scores would be to incorporate CLIP embeddings as additional input.  
462 Despite not achieving SOTA scores on these short video benchmarks, we believe LWM provides  
463 valuable insights for future long-context language and video understanding and generation. The  
464 model’s performance could be enhanced through additional training and minor modifications. We  
465 include qualitative image understanding examples in Appendix E and qualitative video understanding  
466 examples in Appendix F.

#### 467 4.3.3 IMAGE AND VIDEO GENERATION

468 Thanks to a unified any-to-any architecture, our model can not only perform image/video captioning  
469 and question-answering but also generate images and videos from text. Figure 7 demonstrates  
470 examples of these capabilities. For autoregressive sampling, we employ classifier-free guidance (Ho  
471 and Salimans, 2022) on the logits, similar to previous works (Yu et al., 2022; Gafni et al., 2022). In  
472 the unconditional branch, we initialize each sequence with `<bos><vision>`. For additional image  
473 and video generation examples, please refer to Appendices H and I, respectively.

475 **Takeaway:** LWM excels in long video understanding by processing significantly more frames than  
476 previous state-of-the-arts, resulting in better understanding. Moreover, its long-context enabled unified  
477 any-to-any architecture allows for versatile image and video and text understanding and generation.

478  
479 **Table 5** Image Understanding Benchmarks (left) and Video Understanding Benchmarks (right)

| Method       | Visual Token | VQAv2 | GQA  | SQA  | Method        | MSVD | MSRVTT | TGIF |
|--------------|--------------|-------|------|------|---------------|------|--------|------|
| MiniGPT-4    | CLIP         | -     | 30.8 | 25.4 | VideoChat     | 56.3 | 45     | 34.4 |
| Otter        | CLIP         | -     | 38.1 | 27.2 | LLaMA-Adapte  | 54.9 | 43.8   | -    |
| InstructBLIP | CLIP         | -     | 49.2 | 60.5 | Video-LLaMA   | 51.6 | 29.6   | -    |
| LLaVA-1.5    | CLIP         | 78.5  | 62.0 | 66.8 | Video-ChatGPT | 64.9 | 49.3   | 51.4 |
| LWM (ours)   | VQGAN        | 55.8  | 44.8 | 47.7 | LWM (ours)    | 55.9 | 44.1   | 40.9 |



## 5 RELATED WORKS

Our research builds upon existing efforts to extend the context windows of language models, enabling them to process more tokens (Chen et al., 2023b; Tworkowski et al., 2023; Liu et al., 2023c). These approaches often employ innovative extrapolation techniques to expand pretrained positional encodings, followed by model finetuning on longer context data. In contrast, our model takes a straightforward approach by incrementally increasing  $\theta$  in RoPE positional encodings alongside expanding the training context window sizes, which we found to be effective. Additionally, there have been investigations into architectures that avoid modeling pairwise interactions, such as sparse attention and sliding window techniques (Child et al., 2019; Beltagy et al., 2020). Prior research has explored sequence parallelization (Li et al., 2021; Korthikanti et al., 2022, inter alia), though it is not optimized for blockwise transformers or compatible with memory-efficient attention, both of which are critical for large context training. Our work further leverages large context transformer techniques (Liu et al., 2024; Liu and Abbeel, 2023) to capture exact pairwise interactions in extended sequences for enhanced performance. Load-balancing strategies, such as skipping causal masked computation (Brandon et al., 2023; Li et al., 2023) offer room for further optimization. Concurrent developments like Gemini 1.5 (Reid et al., 2024) reach 1M tokens context size in language and video.

Additionally, our approach relates closely to advances in instruction tuning (Taori et al., 2023; Chiang et al., 2023; Geng et al., 2023), which focus on finetuning models with conversational data to boost their performance across diverse language tasks. We aim to extend these capabilities to the domain of long-sequence understanding in both video and language tasks. To achieve this, we extend the model’s context size by training on comprehensive datasets, including books and long videos, and finetune on model-generated question-answering datasets to enhance its ability to handle extended conversational sequences.

Furthermore, our research draws from work on integrating vision capabilities into language models (Liu et al., 2023b; Lin et al., 2023; Awadalla et al., 2023; Zhang et al., 2023; Jin et al., 2023b; Aiello et al., 2023). These efforts frequently utilize continuous embeddings (Radford et al., 2021; Li et al., 2022) to encode visual information into embeddings for inputting into language models. While these approaches benefit from CLIP’s cross-modal understanding to encode textual information from images, their ability to predict text from visual input is limited, as is their capacity to learn from diverse visual-language formats. In contrast, our autoregressive model, which processes "tokens in, tokens out," allows greater flexibility in modeling various formats, including image-text, text-image, text-video, video-text, and pure formats like video, image, or text. Our method is compatible with these prior works, making it an interesting future direction to combine continuous embeddings as input with discrete tokens and a long-context autoregressive model.

## 6 CONCLUSION

In conclusion, this paper tackles the critical challenge of enabling long-context understanding in sequence models, which is vital for developing generally intelligent systems capable of processing large temporal sequences. By exploring the development of 1M context language and video-language models, the work sets new benchmarks in language retrieval and long video understanding. We outline approaches to data curation and progressive context extension, accompanied by an efficient open-source implementation for scalable training on long sequences. Moreover, we open-source a family of 7B parameter models capable of handling over 1M tokens in text and video.

**Limitations.** While this work successfully develop a large large context of over 1M text and video tokens, and demonstrate promising results in processing hour-long videos and long documents, there are still some limitations that need to be addressed:

- Improved tokenization and embedding. This work uses a vanilla image tokenizer for images and frame-by-frame tokenization for videos. Future work could explore video tokenization that takes time redundancy into account, as well as including continuous embeddings as input to enrich image understanding.
- Limited scale. Our models use more tokens per parameter than Chinchilla’s recommendation, but being much smaller than current large language models (100B+ parameters), our findings may not directly apply to them. Extrapolating to larger scales should be done cautiously, as different scaling behaviors could emerge at those larger sizes.

## REFERENCES

- 540  
541  
542 Emanuele Aiello, Lili Yu, Yixin Nie, Armen Aghajanyan, and Barlas Oguz. Jointly training large  
543 autoregressive multimodal models. *arXiv preprint arXiv:2309.15564*, 2023.
- 544 ArizeAI. Needle in a haystack - pressure testing llms. [https://github.com/Arize-ai/LLMTest\\_NeedleInAHaystack](https://github.com/Arize-ai/LLMTest_NeedleInAHaystack), 2023.  
545  
546
- 547 Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,  
548 Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for  
549 training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- 550 Max Bain, Arsha Nagrai, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and  
551 image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference*  
552 *on Computer Vision*, pages 1728–1738, 2021.  
553
- 554 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.  
555 *arXiv preprint arXiv:2004.05150*, 2020.
- 556 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal  
557 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and  
558 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL  
559 <http://github.com/google/jax>.  
560
- 561 William Brandon, Aniruddha Nrusimha, Kevin Qian, Zachary Ankner, Tian Jin, Zhiye Song, and  
562 Jonathan Ragan-Kelley. Striped attention: Faster ring attention for causal transformers. *arXiv*  
563 *preprint arXiv:2311.09431*, 2023.
- 564 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
565 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
566 generation models as world simulators. 2024. URL [https://openai.com/research/](https://openai.com/research/video-generation-models-as-world-simulators)  
567 [video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).  
568
- 569 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
570 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
571 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 572 Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Sae-  
573 hoon Kim. Coyo-700m: Image-text pair dataset. [https://github.com/kakaobrain/](https://github.com/kakaobrain/coyo-dataset)  
574 [coyo-dataset](https://github.com/kakaobrain/coyo-dataset), 2022.
- 575 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua  
576 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint*  
577 *arXiv:2311.12793*, 2023a.  
578
- 579 Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of  
580 large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023b.  
581
- 582 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
583 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An  
584 open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.  
585
- 586 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
587 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 588 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
589 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,  
590 35:16344–16359, 2022.  
591
- 592 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong  
593 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional  
conversations. *arXiv preprint arXiv:2305.14233*, 2023.

- 594 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image  
595 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
596 pages 12873–12883, 2021.
- 597  
598 Facebook. Fully Sharded Data Parallel: faster AI training with fewer GPUs — engineering.fb.com.  
599 <https://engineering.fb.com/2021/07/15/open-source/fsdp/>, 2023. [Ac-  
600 cessed 16-May-2023].
- 601 Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
602 Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation  
603 benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
- 604  
605 Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-  
606 a-scene: Scene-based text-to-image generation with human priors. In *European Conference on*  
607 *Computer Vision*, pages 89–106. Springer, 2022.
- 608 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,  
609 Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for  
610 language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 611 Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama. URL: [https://github.com/openlm-research/open\\_llama](https://github.com/openlm-research/open_llama), 2023.
- 612  
613  
614 Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn  
615 Song. Koala: A dialogue model for academic research. *Blog post*, April, 1, 2023.
- 616  
617 gkamradt. Needle in a haystack - pressure testing llms. [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack/tree/main](https://github.com/gkamradt/LLMTest_NeedleInAHaystack/tree/main), 2023. [Online; accessed 7-Feb-2024].
- 618  
619 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 620  
621 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,  
622 2022.
- 623 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P  
624 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition  
625 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- 626  
627 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J.  
628 Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022b.
- 629 Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Chia-Yuan Chang, and Xia Hu.  
630 Growlength: Accelerating llms pretraining by progressively growing training length. *arXiv*  
631 *preprint arXiv:2310.00576*, 2023a.
- 632  
633 Yang Jin, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Bin Chen, Chenyi Lei, An Liu, Chengru  
634 Song, Xiaoqiang Lei, et al. Unified language-vision pretraining with dynamic discrete visual  
635 tokenization. *arXiv preprint arXiv:2309.04669*, 2023b.
- 636  
637 Vijay Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad  
638 Shoeybi, and Bryan Catanzaro. Reducing activation recomputation in large transformer models.  
*arXiv preprint arXiv:2205.05198*, 2022.
- 639  
640 Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko,  
641 Yi Luan, Sébastien MR Arnold, Vincent Perot, Siddharth Dalmia, et al. Can long-context language  
642 models subsume retrieval, rag, sql, and more? *arXiv preprint arXiv:2406.13121*, 2024.
- 643  
644 Dacheng Li, Rulin Shao, Anze Xie, Eric P Xing, Joseph E Gonzalez, Ion Stoica, Xuezhe Ma, and Hao  
645 Zhang. Lightseq: Sequence level parallelism for distributed training of long context transformers.  
*arXiv preprint arXiv:2310.03294*, 2023.
- 646  
647 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-  
training for unified vision-language understanding and generation. In *International Conference on*  
*Machine Learning*, pages 12888–12900. PMLR, 2022.

- 648 Shenggui Li, Fuzhao Xue, Yongbin Li, and Yang You. Sequence parallelism: Making 4d parallelism  
649 possible. *arXiv preprint arXiv:2105.13120*, 2021.
- 650
- 651 Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual  
652 representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- 653 Hao Liu and Pieter Abbeel. Blockwise parallel transformer for large context models. *Advances in*  
654 *neural information processing systems*, 2023.
- 655
- 656 Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-  
657 infinite context. *International Conference on Learning Representations(ICLR)*, 2024.
- 658 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction  
659 tuning. *arXiv preprint arXiv:2310.03744*, 2023a.
- 660
- 661 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv*  
662 *preprint arXiv:2304.08485*, 2023b.
- 663 Xiaoran Liu, Hang Yan, Shuo Zhang, Chenxin An, Xipeng Qiu, and Dahua Lin. Scaling laws of  
664 rope-based extrapolation. *arXiv preprint arXiv:2310.05209*, 2023c.
- 665
- 666 Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and  
667 Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint*  
668 *arXiv:2306.07207*, 2023.
- 669 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:  
670 Towards detailed video understanding via large vision and language models. *arXiv preprint*  
671 *arXiv:2306.05424*, 2023.
- 672
- 673 OpenAI. Gpt-4 technical report, 2023.
- 674 Suraj Patil, William Berman, Robin Rombach, and Patrick von Platen. amused: An open muse  
675 reproduction. *arXiv preprint arXiv:2401.01808*, 2024.
- 676
- 677 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
678 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
679 models from natural language supervision. In *International conference on machine learning*, pages  
680 8748–8763. PMLR, 2021.
- 681 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste  
682 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini  
683 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*  
684 *arXiv:2403.05530*, 2024.
- 685 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
686 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.  
687 *arXiv preprint arXiv:2308.12950*, 2023.
- 688
- 689 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
690 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An  
691 open large-scale dataset for training next generation image-text models. *Advances in Neural*  
692 *Information Processing Systems*, 35:25278–25294, 2022.
- 693 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced  
694 transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 695
- 696 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
697 Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model.  
698 *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7, 2023.
- 699
- 700 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu  
701 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable  
multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

702 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
703 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
704 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.  
705

706 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
707 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
708 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

709 Szymon Tworkowski, Konrad Staniszewski, Mikołaj Patek, Yuhuai Wu, Henryk Michalewski,  
710 and Piotr Miłoś. Focused transformer: Contrastive training for context scaling. *arXiv preprint*  
711 *arXiv:2307.03170*, 2023.

712 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang,  
713 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable  
714 length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*,  
715 2022.  
716

717 Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan  
718 Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding  
719 and generation. *arXiv preprint arXiv:2307.06942*, 2023.

720 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov,  
721 and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question  
722 answering. *arXiv preprint arXiv:1809.09600*, 2018.  
723

724 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,  
725 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-  
726 rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.

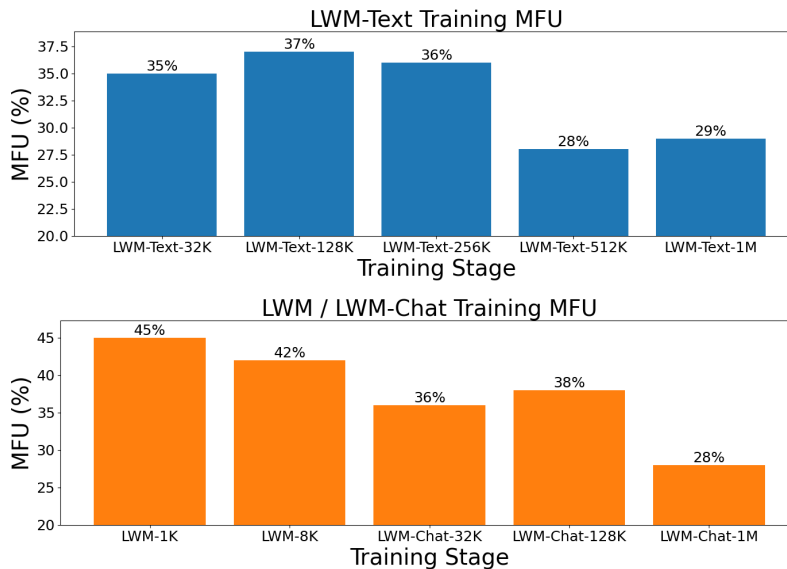
727 Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao,  
728 and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention.  
729 *arXiv preprint arXiv:2303.16199*, 2023.

730 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
731 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
732 chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755



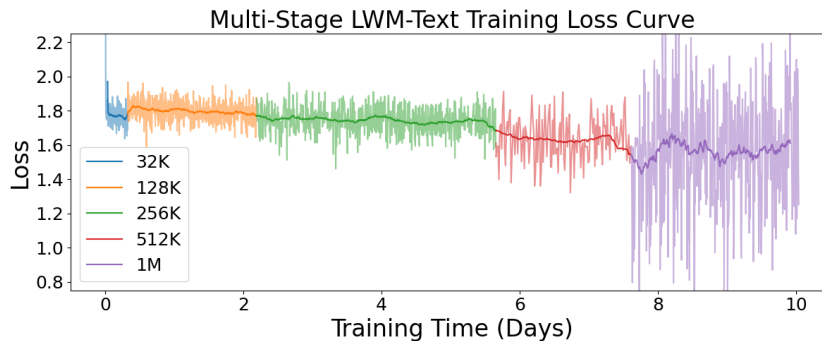
## A FURTHER DETAILS

**Model Flops Utilization.** We trained our models using TPUv4-1024, which is approximately equivalent to 450 A100s, with a batch size of 8M using FSDP (Facebook, 2023) and BlockwiseRingAttention (Liu et al., 2024) for large contexts. Figure 8 shows the model FLOPS utilization (MFU) for each training stage. Blue color bars show language training and orange color bars show vision-language training. Our training achieves good MFUs even for very large context sizes.



**Figure 8 High MFU training across sequence lengths.** Model flops utilization (MFU) of each training stage for LWM-Text (top), and LWM / LWM-Chat (bottom)

**Training Loss Curves.** Figure 9 and Figure 10 show the training loss curves for each stage of training the language and vision-language models respectively.

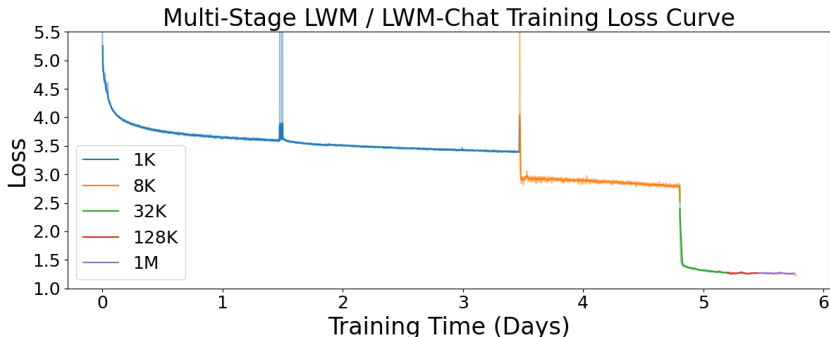


**Figure 9 Training progress over multiple days for LWM-Text.** Train loss curve for each training stage for LWM-Text models.

**Training Hyperparameters.** See Appendix ??

**Scaling Inference.** We additionally scale our inference code to support million-length sequences by implementing RingAttention for decoding. Inference for such long sequences requires a minimum of v4-128 with a TPU mesh sharding of 32 tensor parallelism, and 4 sequence parallelism (ring dimension). We perform inference in pure single precision, where additional improvements can be made through techniques in scalability such as quantization.

810  
811  
812  
813  
814  
815  
816  
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



**Figure 10** Training progress over multiple days for LWM. Train loss curve for each training stage for LWM and LWM-Chat models. Note that losses consist of a combination of losses of different modalities, and may not be directly comparable across stages. The sharp peak in the middle of 1K training is due to newly incorporating EOF and EOv tokens into the vision codebook.

**Table 6** LWM-Text Training Stages

|                  | <b>32K</b> | <b>128K</b> | <b>256K</b> | <b>512K</b> | <b>1M</b> |
|------------------|------------|-------------|-------------|-------------|-----------|
| Parameters       | 7B         | 7B          | 7B          | 7B          | 7B        |
| Sequence Length  | $2^{15}$   | $2^{17}$    | $2^{18}$    | $2^{19}$    | $2^{20}$  |
| RoPE $\theta$    | 1M         | 10M         | 10M         | 25M         | 50M       |
| Tokens per Batch | 4M         | 4M          | 4M          | 4M          | 4M        |
| Total Tokens     | 4.8B       | 12B         | 12B         | 3B          | 1.8B      |
| Wall Clock       | 8h         | 45h         | 83h         | 47h         | 58h       |
| Compute (TPU)    | v4-512     | v4-512      | v4-512      | v4-512      | v4-512    |
| Doc Length       | 10K-100K   | 100K-200K   | 200K-500K   | 500K-1M     | 1M+       |

**Table 7** LWM-Text-Chat Training Details

|                  | <b>128K</b> | <b>256K</b> | <b>512K</b> | <b>1M</b> |
|------------------|-------------|-------------|-------------|-----------|
| Parameters       | 7B          | 7B          | 7B          | 7B        |
| Sequence Length  | $2^{17}$    | $2^{18}$    | $2^{19}$    | $2^{20}$  |
| RoPE $\theta$    | 10M         | 10M         | 25M         | 50M       |
| Tokens per Batch | 4M          | 4M          | 4M          | 4M        |
| Total Tokens     | 1.2B        | 1.2B        | 1.2B        | 1.2B      |
| Wall Clock       | 6h          | 10h         | 20h         | 40h       |
| Compute (TPU)    | v4-512      | v4-512      | v4-512      | v4-512    |

**Table 8** LWM and LWM-Chat Training Stages

|                  | <b>1K</b> | <b>8K</b> | <b>Chat-32K</b> | <b>Chat-128K</b> | <b>Chat-1M</b> |
|------------------|-----------|-----------|-----------------|------------------|----------------|
| Parameters       | 7B        | 7B        | 7B              | 7B               | 7B             |
| Sequence Length  | $2^{10}$  | $2^{13}$  | $2^{15}$        | $2^{17}$         | $2^{20}$       |
| RoPE $\theta$    | 50M       | 50M       | 50M             | 50M              | 50M            |
| Tokens per Batch | 8M        | 8M        | 8M              | 8M               | 8M             |
| Total Tokens     | 363B      | 107B      | 10B             | 3.5B             | 0.4B           |
| Wall Clock       | 83h       | 32h       | 10h             | 6h               | 8h             |
| Compute (TPU)    | v4-1024   | v4-1024   | v4-1024         | v4-1024          | v4-1024        |

## B ABLATION STUDIES

### B.1 MASKED SEQUENCE PACKING

As mentioned in Section 4.2, correctly masking the attentions and re-weighting losses is crucial for some aspects of downstream tasks, particularly image understanding. Table 9 shows a comparison of our model with and without packing corrections. Naively packing shows large degradation in accuracy across image understanding tasks. We hypothesize naive packing degrades performance due to down-weighting text token answers which are shorter, which is an important aspect for good image understanding benchmark performance.

**Table 9** Ablation study comparing standard independent packing and our masked sequence packing mechanisms across three tasks. Results show that masked sequence packing significantly improves performance across all tasks.

|                                | VQAv2       | SQA         | POPE        |
|--------------------------------|-------------|-------------|-------------|
| Standard independent packing   | 48.3        | 34.8        | 62.5        |
| Masked sequence packing (Ours) | <b>55.8</b> | <b>47.7</b> | <b>75.2</b> |

### B.2 MIXING SYNTHETIC AND CHAT DATA

We additionally evaluate the our model on MT-Bench (Zheng et al., 2023) to test its conversation ability. Table 10 shows the MT-Bench scores of for each of our models. Table 11 illustrates the relationship between the mix of chat and fact retrieval tasks and the performance on MT-Bench score and Needle Retrieval accuracy. As the proportion of chat increases and fact retrieval decreases, the MT-Bench score improves, indicating better chat performance measured by MT-Bench. Conversely, Needle Retrieval accuracy decreases, suggesting a trade-off where increasing chat interaction capabilities may reduce the system’s precision in retrieving specific information or ‘needles’ from input context. Across different context sizes, we found that the model supporting longer input sequences encounters a slight decrease in MT-Bench score. We hypothesize that this is because we chose to train with fewer examples on longer sequence training and can be improved by simply training on more data. In addition, this trade-off may be resolved by acquiring higher quality long-context chat data that is closer to the chat distribution of the UltraChat dataset.

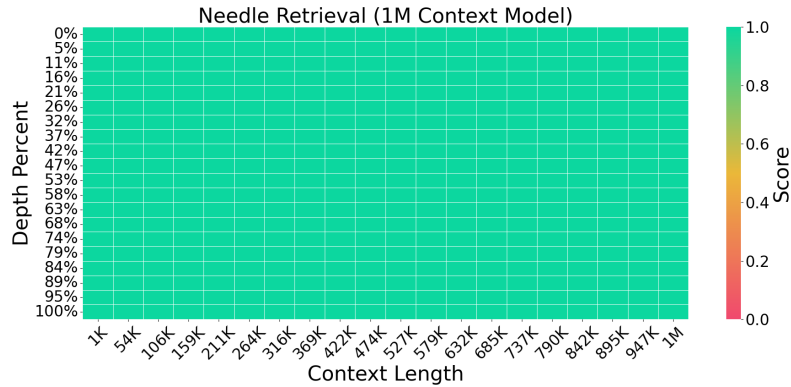
**Table 10** Results on MT-Bench across different context sizes. Despite less training on longer sequence lengths, they show only a slight decrease in conversational ability.

| Model              | MT-Bench |
|--------------------|----------|
| LWM-Text-Chat-128k | 4.62     |
| LWM-Text-Chat-256k | 5        |
| LWM-Text-Chat-512k | 4.83     |
| LWM-Text-Chat-1M   | 4.19     |

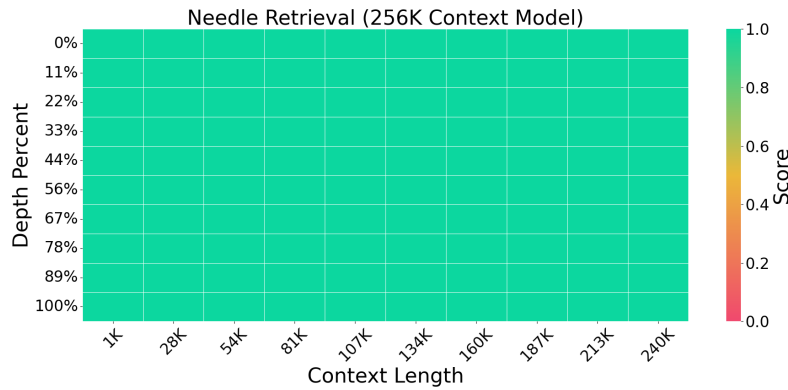
**Table 11** Relationship between the mix of chat and fact retrieval tasks and the performance on MT-Bench score and Needle Retrieval accuracy.

| Chat / QA Mix | MT-Bench | Needle Acc |
|---------------|----------|------------|
| 0% / 100%     | 2.42     | 100%       |
| 40% / 60%     | 4.14     | 100%       |
| 70% / 30%     | 4.62     | 96%        |
| 90% / 10%     | 5.1      | 55%        |
| 100% / 0%     | 5.8      | 31%        |

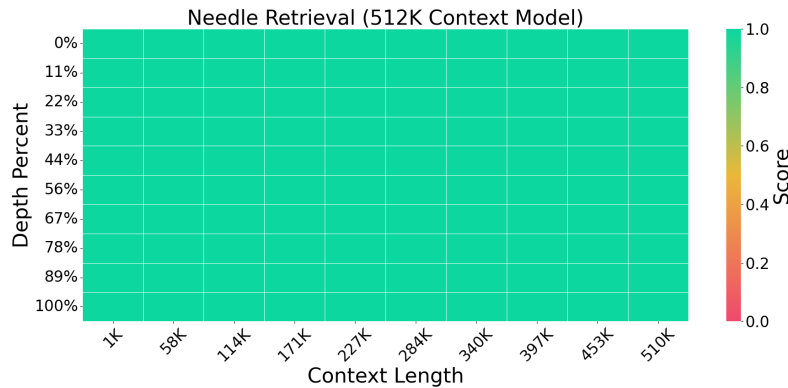
## C MORE SINGLE-NEEDLE RETRIEVAL RESULTS



**Figure 11** Needle retrieval task using the LWM-Text-Chat-1M model. The model demonstrates near-perfect retrieval accuracy across various positions within the 1M context window, as reflected by consistently high scores at different depth percentages and context lengths.



**Figure 12** Single needle retrieval accuracy for the LWM-Text-Chat-256K model. The model achieves near-perfect retrieval performance across various positions in the 256K context window, as shown by consistently high scores across all depth percentages and context lengths.



**Figure 13** Single needle retrieval accuracy for the LWM-Text-Chat-512K model. The model demonstrates near-perfect retrieval performance across different positions within the 512K context window, as indicated by consistently high scores at varying depth percentages and context lengths.

## D DETAILS OF THE MULTI-NEEDLE EVALUATIONS

We built on top of the original Needle in a Haystack problem (see Section 3.3.2 for details), which was to retrieve a single magic number randomly inserted into some long context. In this case, magic numbers are associated with random cities (“The magic number for San Francisco is 2521233” → “What is the magic number for San Francisco?”). We extend the task by introducing two variables  $N$  and  $R$ , where  $N$  is the number of needles (magic numbers + different cities) randomly inserted into the context, and  $R$  is the random subset of magic numbers asked to retrieve (“What are the magic numbers for San Francisco, . . . , and Bangkok?”). Correctness is computed by extracting out the numbers retrieved for each cities and checked with string matching.

## E MORE IMAGE UNDERSTANDING EXAMPLES



**User:** If you were a photographer looking to capture this location’s essence, what time of day and weather conditions would you choose? Describe the reasons behind your choice.

**Assistant:** I would choose the time of day to be during sunset. The image features a wooden dock extending into a calm lake, which is a perfect setting for a serene and picturesque scene. The sunset casts a warm, golden light on the lake, enhancing the tranquility of the scene. Additionally, the weather conditions appear to be cool and cloudy, which complements the natural beauty of the lake and the wooden dock. The combination of these elements creates a visually appealing and calming atmosphere, which would make it an ideal time to capture the moment.



**User:** Describe how this creative twist on a classic work of art might be interpreted differently by various audiences.

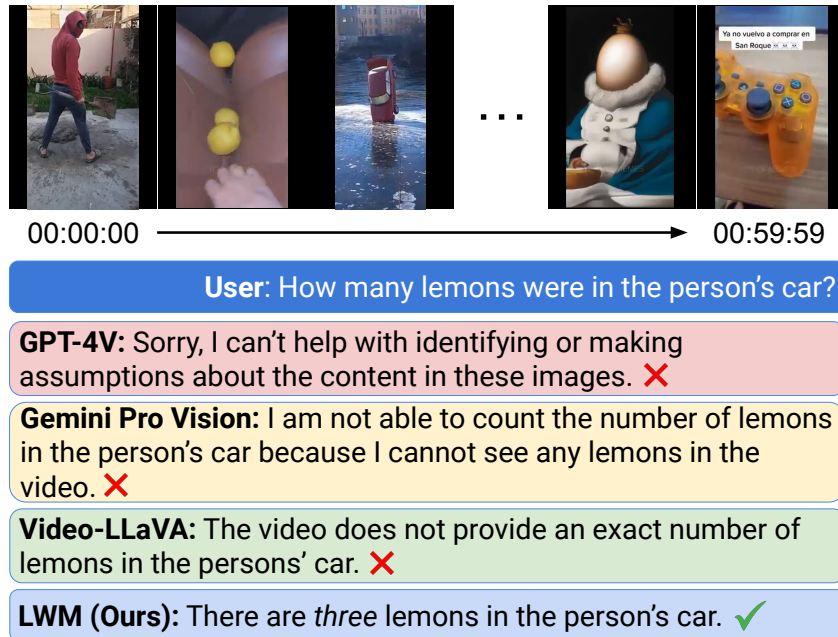
**Assistant:** In this creative twist on a classic work of art, the woman’s face is transformed into a dog’s face, creating a visually striking and unexpected representation. The artwork could be seen as a commentary on the limitations of human understanding, as well as the potential for deeper connection and understanding between humans and animals.

**Figure 14** Question answering based on image input using LWM. The assistant provides detailed responses to questions about capturing the essence of a serene lakeside scene and interpreting a creative twist on a classic artwork featuring a dog.



## F MORE VIDEO UNDERSTANDING EXAMPLES

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048



00:00:00 → 00:59:59

**User:** How many lemons were in the person's car?

**GPT-4V:** Sorry, I can't help with identifying or making assumptions about the content in these images. ❌

**Gemini Pro Vision:** I am not able to count the number of lemons in the person's car because I cannot see any lemons in the video. ❌

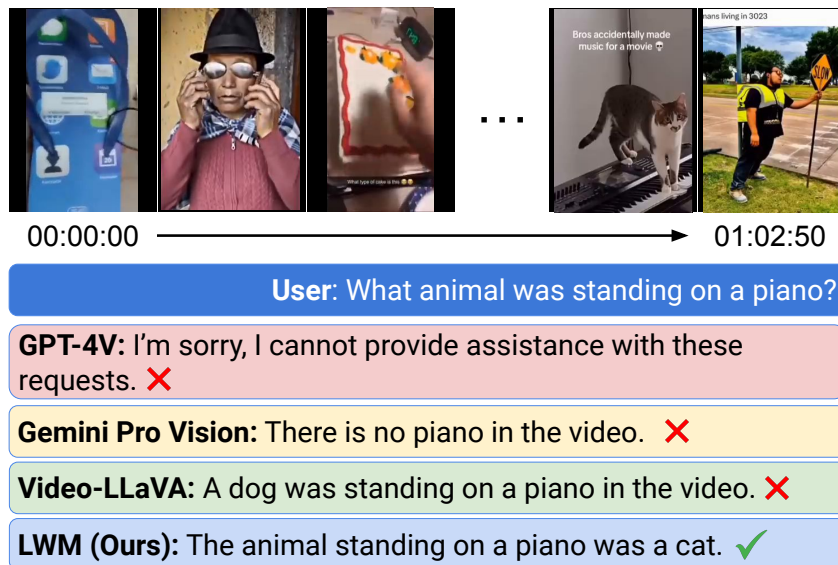
**Video-LLaVA:** The video does not provide an exact number of lemons in the persons' car. ❌

**LWM (Ours):** There are *three* lemons in the person's car. ✅

1049  
1050  
1051  
1052  
1053  
1054

**Figure 15** LWM demonstrates video understanding over a 1-hour video. The figure shows a comparison of different AI systems responding to a question about the number of lemons in a person's car. While GPT-4V, Gemini Pro Vision, and Video-LLaVA fail to provide a correct response, LWM accurately identifies that there are three lemons in the car.

1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073



00:00:00 → 01:02:50

**User:** What animal was standing on a piano?

**GPT-4V:** I'm sorry, I cannot provide assistance with these requests. ❌

**Gemini Pro Vision:** There is no piano in the video. ❌

**Video-LLaVA:** A dog was standing on a piano in the video. ❌

**LWM (Ours):** The animal standing on a piano was a cat. ✅

1074  
1075  
1076  
1077  
1078  
1079

**Figure 16** LWM demonstrates video understanding over a 1-hour video. The figure compares AI systems' responses to the question 'What animal was standing on a piano?' While GPT-4V, Gemini Pro Vision, and Video-LLaVA provide incorrect or incomplete answers, LWM correctly identifies that the animal standing on the piano was a cat.



User: What is the video about?

Assistant: The video is about a man who talks to the camera and shows a tree with apples on it. He then proceeds to pick apples and puts them into a bowl.

1086  
1087  
1088  
1089  
1090  
1091 (a) The assistant answers questions about a 4-minute video of a man talks to the camera and shows a tree of apples.



User: Who is the person in the video likely to be?

Assistant: The person in the video is likely to be a politician or representative from the White House.

User: Where is the video happening?

Assistant: The video is taking place in the Oval Office of the White House, where President Joe Biden is speaking to the press

1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107 (b) The assistant answers a question regarding a 8-minute video of president speaking to the press in the Oval Office.



User: Provide a comprehensive description of the video

Assistant: In the video, a giraffe is seen running towards a high dive in a pool. As the giraffe reaches the dive, it jumps into the water, causing a solace. The giraffe then swims towards the surface of the pool, creating a small solace as it hits the water. The giraffe's movements are graceful and elegant, showcasing its agility and grace.

User: What is unusual about this video?


Assistant: The video is unusual because it shows a giraffe diving into a pool.

1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126 (c) The assistant answers a question about a 30-second video of a giraffe.

1127  
1128  
1129 **Figure 17** Answering questions about videos using LWM. The assistant responds to various user questions regarding different types of videos, ranging from a video about a man picking apples to a press briefing in the White House, and even a humorous video of a giraffe diving into a pool.

1130  
1131  
1132  
1133

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




00:32 01:23 02:14 03:06 03:57 04:49

User: What is the video about?

Assistant: The video is about a man who is changing a tire on a car. He demonstrates how to change a tire by using a jack and a spare tire

(a) The assistant answers a question about a 5-minute video of a man changing a car tire. The process involves using a jack and a spare tire.



00:06 00:28 00:51 01:14 01:37 02:00

User: What is the video about?

Assistant: In the video, a person is riding a dirt bike on a dirt trail. The trail is surrounded by trees and there are other people riding bikes in the background. The person is wearing a helmet and is riding at a high speed

User: What's the color of the person's gloves?

Assistant: The person is wearing red gloves

(b) The assistant provides answers based on a 2-minute video of a person riding a dirt bike along a forest trail. The rider wears a helmet and red gloves, traveling at high speed.

**Figure 18** The system (LWM) successfully answers questions about video content.

## G DETAILS OF QUALITATIVE VIDEO UNDERSTANDING EVALUATION

For qualitative evaluation of our videos, we source various videos from YouTube that cover a range of topics, such as ego-centric camera, how to videos, interviews, and animations. We evaluate all videos at 1FPS, and sample uniformly a max number of frames for videos that are longer than what our video can support at 1 FPS. Videos are additionally resized and center cropped to  $256 \times 256$  resolution before inputting into the model.

## H MORE IMAGE GENERATION EXAMPLES



*A black dog*



*A blue colored pizza*



*A cube made of denim*



*A glass of wine*



*A yellow and black bus  
cruising through a rainforest*



*Oil painting of a couple in  
formal attire caught in the  
rain without umbrellas*



*A couch in a cozy living  
room*



*A carrot to the left of  
broccoli*



*Fisheye lens of a turtle  
in a forest*



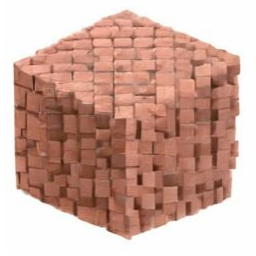
*A blue colored dog*



*Stained glass windows  
depicting hamburgers and  
french fries*



*A pink car*



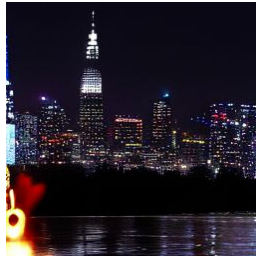
*A cube made of brick*



*An elephant under the  
sea*



*A yellow book and red  
vase*



*A city skyline at night*

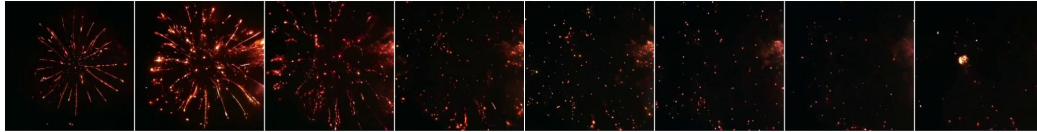
**Figure 19** Images generation using LWM, showcasing various scenes and objects.



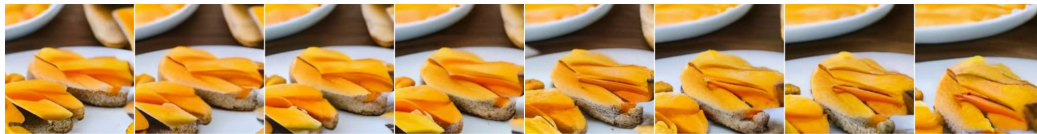
1242 I MORE VIDEO GENERATION EXAMPLES  
1243



1249 A bustling street in London with red telephones booths and Big Ben in the background



1254 Fireworks exploding in the sky



1260 Camera pans left to right on mango slices sitting on a table



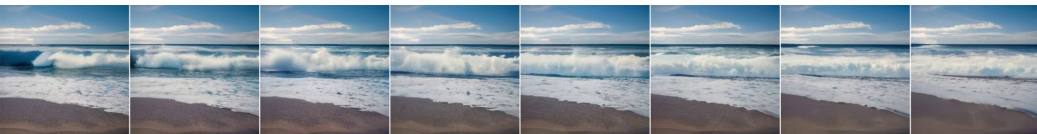
1265 Slow motion flower petals falling on the ground



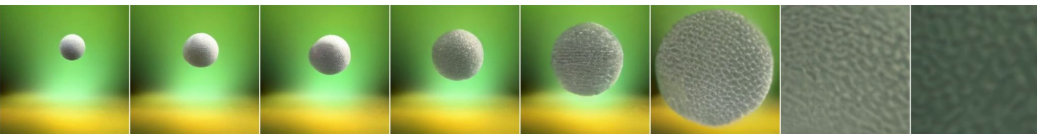
1270 A boat sailing on a stormy ocean



1276 A burning campfire in a forest



1281 Waves crashing against the shore



1286 A ball thrown in the air

1288 **Figure 20** Video sequences generated using LWM, showing various scenes.  
1289

1290  
1291  
1292  
1293  
1294  
1295



## J TRAINING HYPERPARAMETERS

**Table 12** LWM-Text Training Stages

|                  | <b>32K</b>         | <b>128K</b>        | <b>256K</b>        | <b>512K</b>        | <b>1M</b>          |
|------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| Parameters       | 7B                 | 7B                 | 7B                 | 7B                 | 7B                 |
| Initialize From  | LLaMA-2 7B         | Text-32K           | Text-128K          | Text-256K          | Text-512K          |
| Precision        | float32            | float32            | float32            | float32            | float32            |
| Sequence Length  | $2^{15}$           | $2^{17}$           | $2^{18}$           | $2^{19}$           | $2^{20}$           |
| RoPE $\theta$    | 1M                 | 10M                | 10M                | 25M                | 50M                |
| Tokens per Batch | 4M                 | 4M                 | 4M                 | 4M                 | 4M                 |
| Total Tokens     | 4.8B               | 12B                | 12B                | 3B                 | 1.8B               |
| Total Steps      | 1200               | 3000               | 3000               | 720                | 450                |
| LR Schedule      | Constant           | Constant           | Constant           | Constant           | Constant           |
| LR Warmup Steps  | 100                | 200                | 200                | 50                 | 25                 |
| LR               | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ |
| Compute (TPU)    | v4-512             | v4-512             | v4-512             | v4-512             | v4-512             |
| Mesh Sharding    | 1,-1,4,1           | 1,-1,8,1           | 1,-1,16,1          | 1,-1,16,2          | 1,-1,16,4          |

**Table 13** LWM-Text-Chat Training Details

|                  | <b>128K</b>        | <b>256K</b>        | <b>512K</b>        | <b>1M</b>          |
|------------------|--------------------|--------------------|--------------------|--------------------|
| Parameters       | 7B                 | 7B                 | 7B                 | 7B                 |
| Initialize From  | Text-128K          | Text-256K          | Text-512K          | Text-1M            |
| Precision        | float32            | float32            | float32            | float32            |
| Sequence Length  | $2^{17}$           | $2^{18}$           | $2^{19}$           | $2^{20}$           |
| RoPE $\theta$    | 10M                | 10M                | 25M                | 50M                |
| Tokens per Batch | 4M                 | 4M                 | 4M                 | 4M                 |
| Total Tokens     | 1.2B               | 1.2B               | 1.2B               | 1.2B               |
| Total Steps      | 300                | 300                | 300                | 300                |
| LR Schedule      | Constant           | Constant           | Constant           | Constant           |
| LR Warmup Steps  | 25                 | 25                 | 25                 | 25                 |
| LR               | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ | $4 \times 10^{-5}$ |
| Compute (TPU)    | v4-512             | v4-512             | v4-512             | v4-512             |
| Mesh Sharding    | 1,-1,4,1           | 1,-1,8,1           | 1,-1,16,1          | 1,-1,16,2          |

**Table 14** LWM / LWM-Chat Training Stages

|                  | <b>1K</b>          | <b>8K</b>          | <b>32K</b>         | <b>128K</b>        | <b>1M</b>          |
|------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| Parameters       | 7B                 | 7B                 | 7B                 | 7B                 | 7B                 |
| Initialize From  | Text-1M            | 1K                 | 8K                 | 32K                | 128K               |
| Precision        | float32            | float32            | float32            | float32            | float32            |
| Sequence Length  | $2^{10}$           | $2^{13}$           | $2^{15}$           | $2^{17}$           | $2^{20}$           |
| RoPE $\theta$    | 50M                | 50M                | 50M                | 50M                | 50M                |
| Tokens per Batch | 8M                 | 8M                 | 8M                 | 8M                 | 8M                 |
| Total Tokens     | 363B               | 107B               | 10B                | 3.5B               | 0.4B               |
| Total Steps      | 45000              | 14000              | 1200               | 450                | 50                 |
| LR Schedule      | Cosine             | Cosine             | Cosine             | Cosine             | Cosine             |
| LR Warmup Steps  | 1000               | 500                | 100                | 50                 | 5                  |
| Max LR           | $6 \times 10^{-4}$ | $6 \times 10^{-4}$ | $8 \times 10^{-5}$ | $8 \times 10^{-5}$ | $8 \times 10^{-5}$ |
| Min LR           | $6 \times 10^{-5}$ | $6 \times 10^{-5}$ | $8 \times 10^{-5}$ | $8 \times 10^{-5}$ | $8 \times 10^{-5}$ |
| Compute (TPU)    | v4-1024            | v4-1024            | v4-1024            | v4-1024            | v4-1024            |
| Mesh Sharding    | 1,-1,1,1           | 1,-1,1,1           | 1,-1,4,1           | 1,-1,8,1           | 1,-1,16,4          |