

# ARIA : An Open Multimodal Native Mixture-of-Experts Model

Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Guoyin Wang, Bei Chen, Junnan Li<sup>✉</sup>

Rhymes AI

## Abstract

Information comes in diverse modalities. Multimodal native AI models are essential to integrate real-world information and deliver comprehensive understanding. While proprietary multimodal native models exist, their lack of openness imposes obstacles for adoptions, let alone adaptations. To fill this gap, we introduce ARIA, an open multimodal native model with best-in-class performance across a wide range of multimodal, language, and coding tasks. ARIA is a mixture-of-expert model with 3.9B and 3.5B activated parameters per visual token and text token, respectively. It outperforms Pixtral-12B and Llama3.2-11B, and is competitive against the best proprietary models on various multimodal tasks. We pre-train ARIA from scratch following a 4-stage pipeline, which progressively equips the model with strong capabilities in language understanding, multimodal understanding, long context window, and instruction following. We open-source the model weights along with a codebase that facilitates easy adoptions and adaptations of ARIA in real-world applications.

Code: <https://github.com/rhymes-ai/Aria>

Website: <https://rhymes.ai/>

## 1 Introduction

In this report, we present ARIA, the first open mixture-of-experts (MoE) model that is **multimodal native**. The term multimodal native has been used in prior literature to refer to different model capabilities, without a clear consensus. Here, we provide a quantifiable definition: *A multimodal native model refers to a single model with strong understanding capabilities across multiple input modalities (e.g. text, code, image, video), that matches or exceeds the modality-specialized models of similar capacities.* Our definition aligns with the user experience of proprietary multimodal models such as GPT-4o or Gemini-1.5, where a user does not need to differentiate inputs from different modalities. Instead, the model is expected to seamlessly handle and integrate multiple modalities' input with a single model.

While proprietary multimodal native models are not uncommon, their training recipes remain largely undisclosed. As a result, most open-source models are modal-specialized or show subpar performance across modalities. In this research, we fill the gap and introduce training recipes for developing multimodal native models from scratch, which includes key aspects below:

- **Model Architecture.** The core of our model is a fine-grained mixture-of-experts decoder, which enables faster training and inference speed over dense decoders, due to more efficient parameter utilization through expert specialization. ARIA MoE activates 3.5B parameters per text token and has a total of 24.9B parameters. Visual input of variable length, size, and aspect is encoded as visual tokens using a lightweight visual encoder of 438M parameters. ARIA has a long multimodal context window of 64k tokens.

---

<sup>✉</sup> corresponding author: [junnanli@rhymes.ai](mailto:junnanli@rhymes.ai)

Category	Benchmark	ARIA	Pixtral-12B	Llama3.2-11B	GPT-4V	GPT-4o mini	GPT-4o	Gemini-1.5 Flash	Gemini-1.5 Pro
<b>Knowledge/Math (Multimodal)</b>	MMMU (val)	54.9	52.5	50.7	56.4	59.4	69.1	56.1	62.2
	MathVista (testmini)	66.1	58.0	51.5	-	54.7	63.8	58.4	63.9
<b>Document/Chart/Scene Text Understanding</b>	DocVQA (test)	92.6	90.7	88.4	88.4	-	92.8	89.9	93.1
	ChartQA (test)	86.4	81.8	83.4	78.4	-	85.7	85.4	87.2
	TextVQA (val)	81.1	-	-	78.0	-	-	78.7	78.7
<b>General Visual QA</b>	MMBench-1.1	80.3	-	-	79.8	76.0	82.2	-	73.9
<b>Long Video Understanding</b>	EgoSchema (test)	66.8	-	-	-	-	72.2	65.7	72.2
	LongVideoBench (test)	65.3	47.4	45.7	60.7	58.8	66.7	62.4	64.4
	VideoMME (w subs)	72.1	47.5	50.2	63.3	68.9	77.2	75.0	81.3
<b>Knowledge/Math/Reasoning (Language)</b>	MMLU (5-shot)	73.3	69.2	69.4	86.4	-	89.1	78.9	85.9
	MATH (CoT)	50.8	48.1	51.9	-	70.2	76.6	-	-
	ARC Challenge	91.0	-	83.4	-	96.4	96.7	-	-
<b>Coding</b>	HumanEval	73.2	72.0	72.6	67.0	87.2	90.2	74.3	84.1

Table 1: Performance comparison across various multimodal and language benchmarks. Results of competing models are collected from verified official sources or reruned with official settings.

- **Data.** ARIA is pre-trained on 6.4T language tokens and 400B multimodal tokens. We develop a rigorous procedure to curate high-quality data from a diverse set of sources. The multimodal pre-train data includes four major categories: interleaved image-text sequence from common crawl, synthetic image captions, documents transcriptions and question-answering pairs, synthetic video captions and question-answering pairs.
- **Training Pipeline.** We design a 4-stage training pipeline, including language pre-training, multimodal pre-training, multimodal long-context pre-training, and multimodal post-training. Each stage is designed to progressively enhance certain model capabilities while maintaining those acquired in early stages. Our pipeline efficiently and effectively exploits the data and compute resources to maximize model performance.

Following this recipe, ARIA demonstrates state-of-the-art performance as an open multimodal native model. Compared to Pixtral-12B [Mixtral, 2024] and Llama3.2-11B [Dubey et al., 2024], Aria demonstrates superior performance across a wide range of multimodal, language, and coding tasks, while enjoying lower inference cost due to the fewer number of activated parameters. In addition, ARIA also performs on par with proprietary models such as GPT-4o and Gemini-1.5 on various multimodal tasks. The detailed benchmark results are present in Table 1.

We release ARIA under the Apache 2.0 license, free for both academic and commercial use. To facilitate easier adoption, we open-source a training framework that enables finetuning ARIA on a wide variety of data sources and formats, using as few as one GPU.

## 2 Model

### 2.1 Fine-Grained Mixture-of-Experts

MoE has emerged as a preferred architecture over dense models for building compute-efficient large language models [Fedus et al., 2022, Jiang et al., 2024, Dai et al., 2024, Ludziejewski et al., 2024]. The core idea of MoE is to replace each feed-forward layer (FFN) in a Transformer with a set of experts, where each expert is structurally identical to

an FFN. Each input token is routed to only a subset of experts in each layer. The sparsity of expert activation ensures computational efficiency of an MoE layer.

Due to the vast diversity of multimodal data, we hypothesize that *expert specialization* is important for an multimodal MoE to understand input from different data distributions. To this end, we use a large number of fine-grained experts with smaller FFN hidden dimension than standard FFNs, similar to [Dai et al., 2024]. In particular, ARIA has 66 experts in each MoE layer, 2 of the 66 experts are shared among all inputs to capture common knowledge, whereas 6 more experts are activated for each token by a router module. Table 2 shows the detailed architectural configuration.

ARIA is significantly different from previous multimodal MoEs which either design modality-specific expert architectures or rely on upcycling from dense models [Lin et al., 2024b, Shen et al., 2023, Lin et al., 2024a]. Our multimodal native MoE is pre-trained from scratch with modality-generic experts. In Section 4.2, we show that multimodal expert specialization naturally arises after pre-training.

#total parameters	#activated parameters	#experts	#activated experts	expert FFN dim	hidden dim	#layers
24.9B	3.5B	$2^\Delta+64$	$2^\Delta+6$	1664	2560	28

Table 2: Architectural configuration of our MoE decoder.  $\Delta$  denotes shared experts.

## 2.2 Visual Encoder

We design a lightweight visual encoder to convert visual inputs (i.e. images or video frames) into continuous visual tokens with the same feature dimension as word embeddings, which enables the MoE to seamlessly integrate visual and language inputs.

Drawing inspiration from previous work [Li et al., 2023, Bai et al., 2023, Laurençon et al., 2024], our visual encoder consists of a Vision Transformer (ViT) and a projection module. The ViT accepts images in their native aspect ratio as variable-length sequences of patches [Lee et al., 2023, Dehghani et al., 2023], which preserves the inherent information structure in images. We categorize image size into three ranges: (1) medium-resolution images, where the longer edge is resized to 490 pixels; (2) high-resolution images, where the longer edge is resized to 980 pixels; (3) ultra-high-resolution images, where an image is dynamically decomposed into multiple high-res images, following a strategy similar to Liu et al. [2024]. We initialize the weights of our ViT using the SigLIP-SO400M model [Zhai et al., 2023] and continue pre-train the ViT on our multimodal data.

Our projection module transforms the sequence of image embeddings from the ViT into a sequence of visual tokens. It comprises a single cross-attention layer and a FFN layer. The cross-attention layer employs a set of trainable vectors as queries and the image embeddings as keys. Medium-resolution images are processed by 128 queries, whereas high-resolution images are processed by an additional 128 queries (256 queries in total). The outputs from the cross-attention layer are then fed to an FFN, which then outputs visual tokens for the MoE decoder to further process.

## 2.3 Infrastructure

ARIA is trained on an extensively modified Megatron framework [Shoeybi et al., 2019]. We eschew pipeline parallelism and instead implement a combination of expert parallelism [Lepikhin et al., 2020] and ZeRO-1 data parallelism [Rajbhandari et al., 2020] to optimize performance. Due to the carefully designed parallelism method and the small model size, ARIA can be effectively trained without using tensor parallelism, which significantly reduces communication overhead and enhances training efficiency.

We implement a load balancing loss to prevent routing collapse and encourage balanced expert activation. We find that the expert-level load balancing loss in previous work [Fedus et al., 2022, Dai et al., 2024] is overly restrictive for our MoE due to the large number of experts. Therefore, we relax the load balancing to groups of experts, where each group contains 8 fine-grained experts. We also employ z-loss [Zoph et al., 2022] to stabilize training.

## 3 Training

In this section, we delineate our 4-stage training pipeline. In each stage, the model aims to learn new capabilities while maintaining those acquired previously. We perform evaluation during each stage to ensure that such goal is achieved in a data-efficient and compute-efficient way.

### 3.1 Language Pre-training

The first stage pre-trains the MoE decoder with a large amount of curated language data converted into discrete text tokens, using a next-token prediction loss, which enables the MoE to learn general knowledge about the world. The context window length is 8K tokens.

**Language Data.** Our language pre-training data contains 6.4T tokens in total, curated from a variety of data sources containing knowledge until May 2024. We de-duplicate the data at different granularities and perform rigorous quality filtering, using a combination of rule-based approach and model-based quality classifiers. To enhance model’s in-context learning capability, we employ data clustering and pack similar data in the same sequence during training, akin to the approach in [Shi et al. \[2023\]](#). However, their original method is less scalable and likely to generate numerous long-tail structures when processing trillions of tokens. Instead, we utilize a minimum spanning tree algorithm for language data clustering, which resulted in a noticeable performance gain.

### 3.2 Multimodal Pre-training

The second stage pre-trains the MoE decoder and the visual encoder with a mixture of language and multimodal data, using the same next-token prediction loss. This stage aims to enable the model with broad multimodal understanding abilities, while maintaining or even improving its language understanding. To this end, the language data contains a high-quality subset of 1T tokens, covering topics including code, reasoning, and knowledge. The multimodal data contains 400B tokens from a diverse set of sources, which can be categorized into four major categories below.

**Interleaved image-text web data.** We extract and filter web pages from Common Crawl. The filtering process first removes web pages with low image or text quality. Then, it de-duplicate images, and removes web pages where the images and the text have low overall CLIP score [[Radford et al., 2021](#)]. Additionally, we adjust the position of the images in the sequence, by moving an image to the front of a sentence if the sentence has higher CLIP score and is in front of the image. In total, we curate 190B interleaved image-text tokens.

**Synthetic image captions.** Alt texts directly extracted for web images are generally short, less descriptive, and noisy. It has been shown in previous work that synthetic data at scale can improve multimodal pre-training [[Li et al., 2022](#)]. We thus synthesize image captions using a small model which has learned to generate longer and more descriptive image captions by re-writing the alt texts. We create synthetic captions for 300M images in the LAION-400M dataset [[Schuhmann et al. \[2021\]](#)], resulting in a total of 70B multimodal tokens.

**Document transcriptions and QA.** To improve the model’s capability of understanding text-heavy images, we transcribe document images into texts using public OCR methods. We also render images using plain text, chart json or table/equation latex code. In order to enhance the model’s ability to not only transcribe text but also understand its meaning, we use a language model to create synthetic question-answering pairs. In total, our multimodal document data contains 102B tokens.

**Video captions and QA.** We collect 4.4M videos of varying lengths from a diverse range of sources. We train a model to generate frame-level dense descriptions for the videos. Then, we use a language model to generate question-answering pairs and video summarizations based on the dense video descriptions. In total, our video data contains 35B tokens. We select samples within 8K length for multimodal pre-training.

Model	#Params activated (total)	LongVideoBench		VideoMME		MMLongBench-Doc	
		test	val	w subs	w/o subs	acc	f1
<i>Open-source</i>							
ARIA	3.9B (25.3B)	65.3	63.0	72.1	67.6	28.3	24.6
Qwen2-VL-7B	7B	56.8	55.6	69.0	63.3	21.3	22.7
Idefics2	8B	49.4	49.7	-	-	7.0	6.8
MiniCPM-V-2.6	8B	55.7	54.9	63.7	60.9	11.5	11.6
Llama3.2-11B	11B	45.7	45.5	49.5	46.0	13.8	11.3
Pixtral-12B	12B	47.4	44.9	47.5	40.7	6.4	6.0
InternVL-Chat-V1.5	26B	51.7	51.2	52.4	50.7	14.6	13.0
InternVL2-40B	40B	60.6	59.3	62.4	61.2	18.2	17.9
LLaVA-OneVision-72B	72B	63.2	61.3	69.6	66.3	-	-
Qwen2-VL-72B	72B	61.7	60.4	77.8	71.2	33.3	35.7
<i>Proprietary</i>							
Gemini-1.5-Flash	-	62.6	61.4	75.0	60.3	27.0	21.3
Gemini-1.5-Pro	-	64.4	64.0	81.3	75.0	28.2	20.6
GPT-4o mini	-	58.8	56.5	68.9	64.8	29.0	28.6
GPT-4o	-	66.7	66.7	77.2	71.9	42.9	44.9

Table 3: Evaluation of long-context multimodal understanding on videos and documents. Results of competing models are collected from verified official leaderboards or reruned with official settings.

### 3.3 Multimodal Long-Context Pre-training

In this stage, we pre-train on long sequences to extend the model’s context window to 64K tokens. Language long-sequence data is selected from the pre-train data source. Multimodal long-sequence data contains long videos, long documents and synthetic long sequences constructed from short multimodal data. In particular, we concatenate a sequence of independent images as input, and concatenate their image descriptions as target. This stage consumes 12B language tokens and 21B multimodal tokens, where 69% of the 33B tokens are long sequences. We increase the RoPE base frequency hyperparameter from 100K to 5M.

After this stage, the model perfectly solves the needle-in-a-haystack task [Kamradt, 2023] for up to 64K context window. It also demonstrates substantial performance improvement on long video understanding and long multimodal document understanding tasks.

### 3.4 Multimodal Post-training

The final post-training stage anneals the learning rate to converge the model. The learning focuses on improving the model’s question-answering and instruction-following capabilities, using a mixture of high-quality open-source datasets and human-annotated datasets, covering domains including multimodal, code, math, and reasoning. This stage digests 20B tokens in total.

## 4 Evaluation and Analysis

### 4.1 Benchmark Results

In Table 1, we compare ARIA against leading open models of similar scale and proprietary models across a variety of established benchmarks. In Table 3 and Table 4, we examine the long-context multimodal understanding and instruction following capability, respectively. Based on the evaluation result, we highlight the following key observations.

**ARIA is the best-in-class open multimodal native model**, showing clear advantages over Pixtral-12B and Llama3.2-11B across a wide range of multimodal, language, and coding tasks.

**ARIA is competitive against proprietary models on various multimodal tasks**, including document understanding, chart reading, scene text recognition, and video understanding.

**ARIA excels in long-context multimodal understanding**. Real-world multimodal data is complex by nature and often involves long sequences of interleaved vision-language input,

	ARIA	Phi-3 Vision	Qwen2-VL-7B	Idefics2	Pixtral-12B	InternVL-Chat-v1.5	LLaVA-NeXT-34B	MiniCPM-V-2.5	Gemini-1.0-Pro	Reka-Core	Claude-3-Sonnet	GPT-4o
MIA-Bench (Multimodal)	8.76	7.60	8.07	5.14	8.43	7.54	7.56	7.63	7.06	7.70	7.94	8.86
MT-Bench (Language)	8.53	6.27	6.41	-	7.68	-	-	-	-	-	-	-

Table 4: Evaluation of instruction following capabilities. Results of competing models are copied from Qian et al. [2024] for MIA-Bench and Mixtral [2024] for MT-Bench.

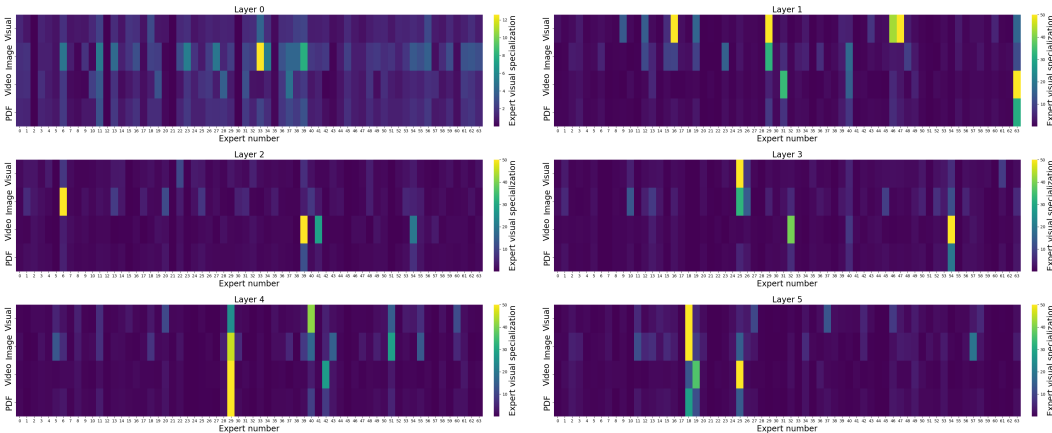
such as videos with subtitles or multi-page documents. ARIA excels in understanding such data, significantly outperforming open models such as Qwen2-VL-7B [Bai et al., 2023] and LLaVA-OneVision-72B [Li et al., 2024]. Compared to proprietary models, ARIA outperforms GPT-4o mini in long video understanding [Wu et al., 2024], and is superior to Gemini-1.5-Flash in long document understanding, making ARIA a preferable choice for processing long multimodal data in a compute-efficient and time-efficient manner.

**ARIA has strong instruction following capabilities**, outperforming other open models on both multimodal and language-only benchmarks. See Section 4.3 for qualitative examples.

## 4.2 Expert Modality Specialization

We analyze the expert activation behavior across all layers in ARIA MoE after the multimodal pre-training stage. We use multimodal data from three domains for analysis: natural image, video, and PDF-format image. For each expert, we first compute its activation rate for both visual tokens and text tokens, denoted as  $R_v$  and  $R_t$ .  $R_v$  refers to the number of visual tokens that activates the expert divided by the total number of visual tokens processed by all experts of that layer, and  $R_t$  refers to the same meaning for text tokens. Then we compute the ratio  $R_v/R_t$ , which represents the expert’s level of visual specialization. A higher visual specialization suggests that the expert is more frequently activated by visual tokens compared to text tokens.

Figure 1 shows a visualization of the expert visual specialization value (capped at 50). It is observed that a number of visual-specialized experts exist in most layers. Furthermore, multiple layers (layer 4, 5, 14, 15, 16, 17, 20) have a single visual expert specialized in all three visual domains. This analysis shows that despite its modality-generic architecture, our MoE learns effective expert utilization during pre-training.



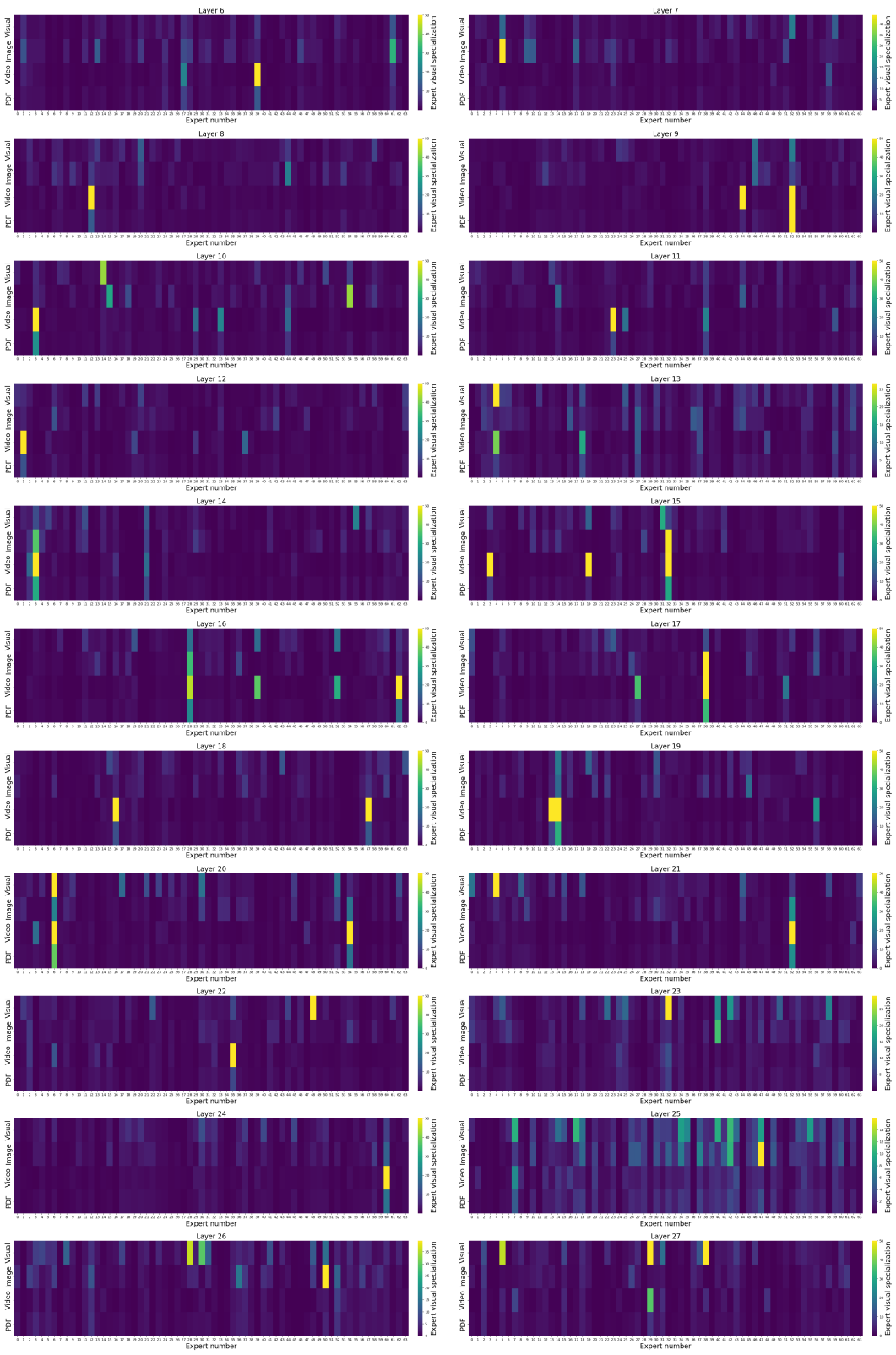


Figure 1: Visualization of expert specialization in each MoE layer.

## 4.3 Qualitative Analysis

### 4.3.1 Multimodal Native Reasoning with Vision, Language, Coding Capabilities

---

#### ARIA Example - Weather Forecast Information Extraction

---

**User** Organize the highest temperature, lowest temperature, and humidity for all even dates into a table.  
Write a python code to draw a line chart for the temperatures, and a bar chart for the humidity. Use red and green colors for the temperatures, and use blue color for the humidity.

2 Week Extended Forecast in Singapore, Singapore

Day	Conditions		Comfort		
	Temperature	Weather	Feels Like	Wind	Humidity
Fri 20 Sep	33 / 26 °C	Isolated ts Storms. Overcast.	38 °C	16 km/h	↑ 62%
Sat 21 Sep	32 / 26 °C	A few ts Storms. Overcast.	37 °C	6 km/h	↗ 68%
Sun 22 Sep	33 / 26 °C	Isolated ts Storms. Overcast.	39 °C	12 km/h	↗ 63%
Mon 23 Sep	31 / 26 °C	Isolated ts Storms. Overcast.	34 °C	6 km/h	↗ 73%
Tue 24 Sep	32 / 26 °C	Thundershowers. Overcast.	34 °C	12 km/h	↖ 70%
Wed 25 Sep	31 / 25 °C	Thunderstorms. Overcast.	36 °C	3 km/h	↗ 70%
Thu 26 Sep	29 / 24 °C	Passing showers. Overcast.	32 °C	8 km/h	↘ 76%
Fri 27 Sep	27 / 24 °C	Showery. Overcast.	29 °C	5 km/h	↗ 89%
Sat 28 Sep	27 / 26 °C	Passing showers. Overcast.	29 °C	10 km/h	↗ 69%
Sun 29 Sep	27 / 26 °C	Thunderstorms. Overcast.	31 °C	5 km/h	↗ 90%
Mon 30 Sep	27 / 26 °C	Passing showers. Overcast.	29 °C	11 km/h	↑ 66%

---

**ARIA** Table of values for even dates:

Day	Highest Temperature (°C)	Lowest Temperature (°C)	Humidity (%)
20 Sep	33	26	62
22 Sep	33	26	63
24 Sep	32	26	70
26 Sep	29	24	76
28 Sep	27	26	69
30 Sep	27	26	66

```
import matplotlib.pyplot as plt

# Data for even dates
days = ['20 Sep', '22 Sep', '24 Sep', '26 Sep', '28 Sep', '30 Sep']
high_temps = [33, 33, 32, 29, 27, 27]
low_temps = [26, 26, 26, 24, 26, 26]
humidity = [62, 63, 70, 76, 69, 66]

# Line chart for temperatures
plt.figure(figsize=(14, 6))

plt.subplot(1, 2, 1)
plt.plot(days, high_temps, 'r-', label='Highest Temperature', marker='o')
plt.plot(days, low_temps, 'g-', label='Lowest Temperature', marker='o')
plt.xlabel('Days')
plt.ylabel('Temperature (°C)')

plt.title('Temperature Variation')
plt.legend()
plt.grid(True)
```



---

## ARIA Example - Weather Forecast Information Extraction

---

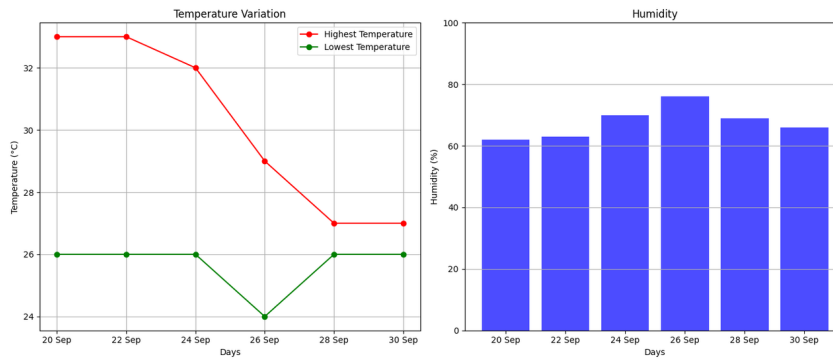
### ARIA (continued)

```
# Bar chart for humidity
plt.subplot(1, 2, 2)
plt.bar(days, humidity, color='blue', alpha=0.7)
plt.xlabel('Days')
plt.ylabel('Humidity (%)')
plt.title('Humidity')
plt.ylim(0, 100)
plt.grid(axis='y')

plt.tight_layout()
plt.show()
```

---

**Result** Figures generated by the code provided by ARIA is shown as below.



---

Weather forecast image source (14-day forecast, screenshot on Sep 20, 2024):  
<https://www.timeanddate.com/weather/singapore/singapore/hourly>

## ARIA Example - Financial Report Analysis

**User** How does the company's profit margin change from year 2020 to 2021? What factors cause such changes? Answer step by step.

Apple Inc.			
CONSOLIDATED STATEMENTS OF OPERATIONS			
(In millions, except number of shares which are reflected in thousands and per share amounts)			
	Years ended		
	September 25, 2021	September 26, 2020	September 28, 2019
<b>Net sales:</b>			
Products	\$ 297,392	\$ 220,747	\$ 213,883
Services	68,425	53,768	46,291
<b>Total net sales</b>	<b>365,817</b>	<b>274,515</b>	<b>260,174</b>
<b>Cost of sales:</b>			
Products	192,266	151,286	144,996
Services	20,715	18,273	16,786
<b>Total cost of sales</b>	<b>212,981</b>	<b>169,559</b>	<b>161,782</b>
<b>Gross margin</b>	<b>152,836</b>	<b>104,956</b>	<b>98,392</b>
<b>Operating expenses:</b>			
Research and development	21,914	18,752	16,217
Selling, general and administrative	21,973	19,916	18,245
<b>Total operating expenses</b>	<b>43,887</b>	<b>38,668</b>	<b>34,462</b>
<b>Operating income</b>	<b>108,949</b>	<b>66,288</b>	<b>63,930</b>
Other income/(expense), net	258	803	1,807
<b>Income before provision for income taxes</b>	<b>109,207</b>	<b>67,091</b>	<b>65,737</b>
Provision for income taxes	14,527	9,680	10,481
<b>Net income</b>	<b>\$ 94,680</b>	<b>\$ 57,411</b>	<b>\$ 55,256</b>
<b>Earnings per share:</b>			
Basic	\$ 5.67	\$ 3.31	\$ 2.99
Diluted	\$ 5.61	\$ 3.28	\$ 2.97
<b>Shares used in computing earnings per share:</b>			
Basic	16,701,272	17,352,119	18,471,336
Diluted	16,864,919	17,528,214	18,595,651

**ARIA** To determine how the company's profit margin changed from 2020 to 2021, we need to follow these steps:

1. Calculate the profit margins for both years:

• **2020:**

- Net income: \$57,411 million
- Total net sales: \$274,515 million
- Profit margin:

$$\text{Profit Margin} = \left( \frac{\text{Net Income}}{\text{Total Net Sales}} \right) \times 100 = \left( \frac{\$57,411}{\$274,515} \right) \times 100 \approx 20.91\%$$

• **2021:**

- Net income: \$94,680 million
- Total net sales: \$365,817 million
- Profit margin:

$$\text{Profit Margin} = \left( \frac{\text{Net Income}}{\text{Total Net Sales}} \right) \times 100 = \left( \frac{\$94,680}{\$365,817} \right) \times 100 \approx 25.88\%$$

---

## ARIA Example - Financial Report Analysis

---

ARIA (continued)

### 2. Analyze the factors causing the change:

- **Increase in Net Income:** The net income increased from \$57,411 million in 2020 to \$94,680 million in 2021, showing a significant rise of \$37,269 million.
- **Sales Growth:** Total net sales also increased from \$274,515 million in 2020 to \$365,817 million in 2021, indicating a strong growth in revenue.
- **Cost Management:** The cost of sales increased from \$151,286 million in 2020 to \$192,266 million in 2021, but the percentage increase in sales outpaced the increase in costs, leading to a higher profit margin.
- **Operating Expenses:** Operating expenses grew from \$38,668 million in 2020 to \$43,887 million in 2021, but the rise in sales was more substantial, positively impacting the profit margin.

Therefore, the company's profit margin improved from approximately 20.91% in 2020 to approximately 25.88% in 2021. This improvement is primarily due to a substantial increase in net income and effective cost management relative to the growth in total net sales.

---

---

Financial report image source:

<https://www.financestrategists.com/wealth-management/financial-statements/>

---

### ARIA Example - Handwritten Math Equation Explanation

---

**User** Explain the equation according to the notes, then give an example with numbers.

The image shows the handwritten equation  $P(A|C) = \frac{P(C|A)P(A)}{P(C)}$ . Annotations include: an arrow from  $P(C|A)$  to the text "Probability someone who definitely has an allergy would make the claim that they do."; an arrow from  $P(A)$  to the text "General probability that someone has a food allergy"; an arrow from  $P(A|C)$  to the text "Probability that someone actually has a food allergy give they say they do."; and an arrow from  $P(C)$  to the text "probability someone would claim to have a food allergy".

---

**ARIA** The equation in the image represents Bayes' Theorem, which is used to update the probability estimate of a hypothesis (in this case, having a food allergy) given new evidence (claiming to have a food allergy).

**Explanation:**

- $P(C|A)$ : Probability that someone who definitely has an allergy would claim that they do.
- $P(A)$ : General probability that someone has a food allergy.
- $P(C)$ : Probability that someone would claim to have a food allergy.
- $P(A|C)$ : Probability that someone actually has a food allergy given that they claim to have one.

The equation:

$$P(A|C) = \frac{P(C|A) \cdot P(A)}{P(C)}$$

**Example with numbers:**

Assume:

- $P(C|A) = 0.9$ , (90% of people with an allergy claim they do)
- $P(A) = 0.05$ , (5% of the population has an allergy)
- $P(C) = 0.1$ , (10% of the population claims to have an allergy)

Using the formula:

$$P(A|C) = \frac{0.9 \times 0.05}{0.1} = 0.45$$

So, the probability that someone who claims to have a food allergy actually does is 45%.

---

---

Image source:

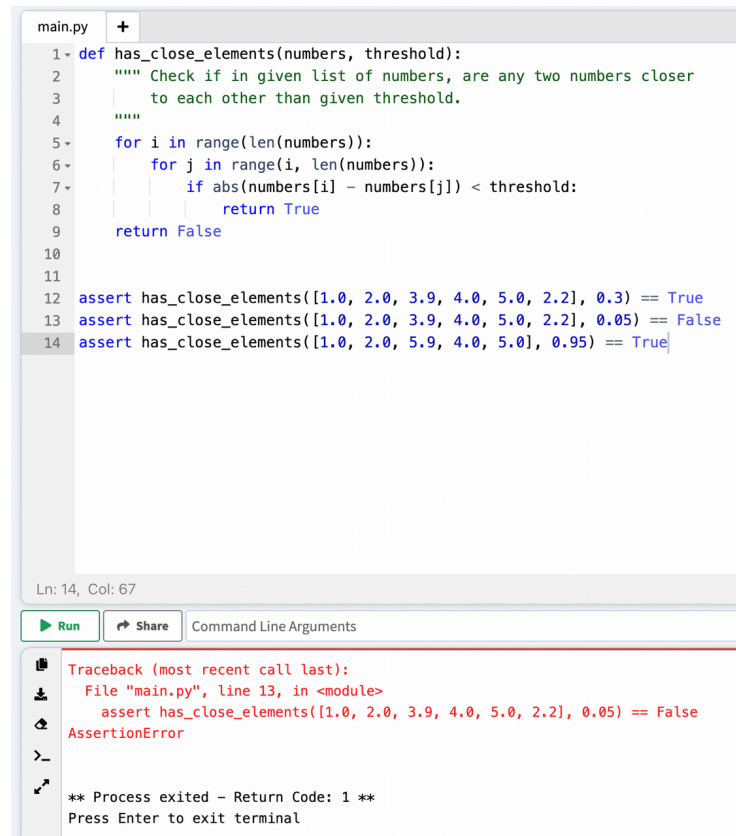
<https://ngathanblog.wordpress.com/2020/06/25/baysian-struggle/>.

---

## ARIA Example - Code Debugging from Screenshot

---

User Debug this code.



```
main.py +
1- def has_close_elements(numbers, threshold):
2-     """ Check if in given list of numbers, are any two numbers closer
3-         to each other than given threshold.
4-     """
5-     for i in range(len(numbers)):
6-         for j in range(i, len(numbers)):
7-             if abs(numbers[i] - numbers[j]) < threshold:
8-                 return True
9-     return False
10
11
12 assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.3) == True
13 assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.05) == False
14 assert has_close_elements([1.0, 2.0, 5.9, 4.0, 5.0], 0.95) == True

Ln: 14, Col: 67

Run Share Command Line Arguments

Traceback (most recent call last):
  File "main.py", line 13, in <module>
    assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.05) == False
AssertionError

** Process exited - Return Code: 1 **
Press Enter to exit terminal
```

**ARIA** The code has a logical error in the nested loop logic. The current implementation of the `has_close_elements` function does not correctly check if any two numbers in the list are closer to each other than the given threshold. Here's the issue:

1. The current condition `if abs(numbers[i] - numbers[j]) < threshold:` checks the absolute difference between each pair of numbers but should only check pairs where `i` and `j` are different indices.
2. The inner loop should start from `i + 1` instead of `i` to avoid comparing a number with itself.

Here's the corrected version:

```
def has_close_elements(numbers, threshold):
    """
    Check if in given list of numbers, are any two numbers closer to
    each other than the given threshold.
    """
    for i in range(len(numbers)):
        for j in range(i + 1, len(numbers)): # Start from i + 1 to avoid comparing with
            itself
                if abs(numbers[i] - numbers[j]) < threshold:
                    return True
    return False

# Tests
assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.3) == True
assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.05) == False
assert has_close_elements([1.0, 2.0, 5.9, 4.0, 5.0], 0.95) == True
```

This change ensures the function correctly checks all pairs of numbers and returns `True` if any pair is closer than the threshold, otherwise `False`.

---

---

### ARIA Example - Multi-image Food Chain Identification

---

**User** Based on the five species in the images, draw a food chain. Explain the role of each species in the food chain.



**ARIA**

1. **Mushrooms** (*Image 1*): Decomposers, breaking down dead organic matter and recycling nutrients.
2. **Grass** (*Image 4*): Primary producers, converting sunlight into energy via photosynthesis.
3. **Grasshopper** (*Image 2*): Primary consumers, herbivores feeding on plants.
4. **Mouse** (*Image 3*): Secondary consumers, herbivores that eat plants and possibly small insects.
5. **Hawk** (*Image 5*): Tertiary consumers, predators that eat mice and other small animals.

#### Food Chain

Grass → Grasshopper → Mouse → Hawk → Mushrooms

#### Explanation

1. Grass serves as the base, providing energy to grasshoppers.
  2. Grasshoppers eat the grass and are then consumed by mice.
  3. Mice are preyed upon by hawks.
  4. Hawks are apex predators, feeding on mice.
  5. Mushrooms decompose dead matter, recycling nutrients back into the ecosystem.
-

## 4.3.2 Long Multimodal Context Window

### ARIA Example - Paper Reading

User Summarize the key contribution of the paper, and explain how it uses various figures and tables to showcase such contribution.



**ARIA** The paper introduces LONGVIDEOBENCH, a benchmark for evaluating Large Multimodal Models (LMMs) in long-context video-language understanding. It addresses the challenge of processing longer video inputs by proposing a question-answering benchmark that tests the models' ability to understand and reason about video contexts interleaved with up to an hour-long video. This benchmark includes 3,763 videos with diverse themes and 6,678 human-annotated multiple-choice questions across 17 categories.

The key contribution is the comprehensive evaluation of LMMs' performance on long-form video understanding, highlighting the limitations of existing models in handling extended inputs. The paper demonstrates that while proprietary models like GPT-4 and Gemini-1.5-Pro show improvements with longer inputs, open-source models lag significantly behind.

---

### ARIA Example - Paper Reading

---

**ARIA** (continued) Figures and tables are used to illustrate these points:

- **Figure 1:** Demonstrates the performance of different models on the benchmark, showing how accuracy varies with the number of frames.
- **Figure 2:** Provides examples of the 17 categories of referring reasoning questions.
- **Table 1:** Lists the benchmark and previous benchmarks for comparison.
- **Table 2:** Details the 17 categories of referring reasoning questions.
- **Table 3:** Shows statistics of videos in LONGVIDEOBENCH by category groups and duration.
- **Table 5:** Validation set results categorized by duration groups.
- **Table 7:** Test set leaderboard of LMMs by duration groups and question categories.

These visual aids help in understanding the performance metrics, the diversity of the dataset, and the specific challenges posed by longer videos and complex questions, thus emphasizing the need for improved models in long-form video understanding.

---

**User** According to the paper, what are the two major difficulties in understanding long videos?

---

**ARIA** The two major difficulties in understanding long videos, as outlined in the paper, are:

1. **Retrieving details from long videos:** Existing Large Multimodal Models (LMMs) often struggle to extract specific details from long sequences. To accurately assess tasks in LONGVIDEOBENCH, there is a need for models to focus on granular details such as objects, events, or attributes, rather than providing a summary or topic overview.
2. **Reasoning contextual relations in long videos:** Questions in LONGVIDEOBENCH require models to analyze the interconnections among diverse contents. This involves understanding the relationships among objects, events, or attributes within the video, which is significantly challenging for extensive inputs. The tasks demand models to derive the correct answer by examining the context and relations across multiple moments in the video.

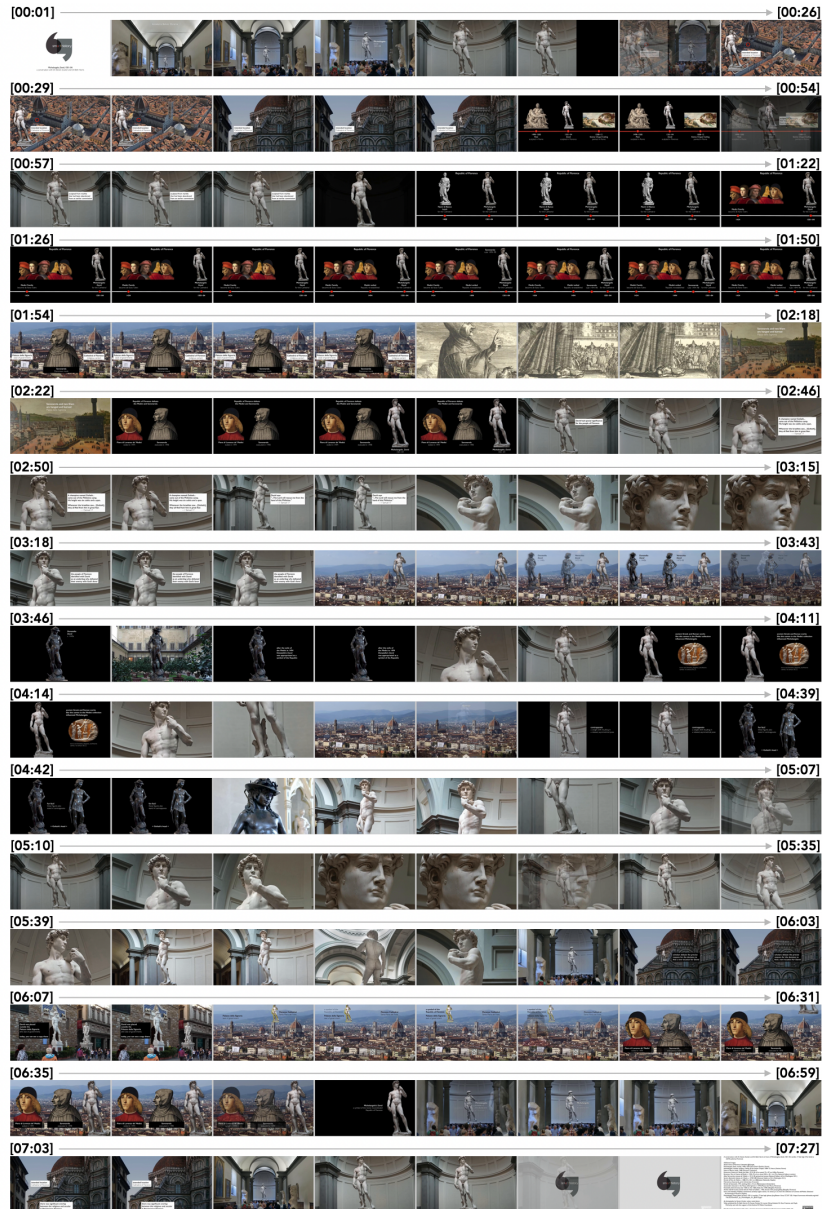


---

## ARIA Example - Long Video Understanding (Scene Splitting)

---

**User** Please split this video into scenes, providing start time, end time, a title and detailed descriptions for each scene.



---

**ARIA** Certainly! Here are the scenes from the video, along with their start and end times, and detailed descriptions:

1. Scene 1: Introduction

**Start Time:** 00:01

**End Time:** 00:05

**Title:** Introduction

**Description:** The video opens with a title screen featuring the *smarthistory* logo and the text "Michelangelo, David, 1501-04," indicating a conversation with Dr. Steven Zucker and Dr. Beth Harris.

---

---

## ARIA Example - Long Video Understanding (Scene Splitting)

---

ARIA (continued)

2. Scene 2: Accademia Gallery, Florence

**Start Time:** 00:05

**End Time:** 00:40

**Title:** Accademia Gallery, Florence

**Description:** The scene transitions to the interior of the Accademia Gallery in Florence, showcasing Michelangelo's David statue surrounded by visitors. The focus is on the statue and its historical context.

3. Scene 3: Historical Context

**Start Time:** 00:40

**End Time:** 01:08

**Title:** Historical Context

**Description:** The scene shifts to a timeline highlighting significant artworks and events from 1498 to 1512, including Michelangelo's David and the Sistine Chapel ceiling.

4. Scene 4: Comparison of Statues

**Start Time:** 01:08

**End Time:** 01:50

**Title:** Comparison of Statues

**Description:** The video presents a comparison between Michelangelo's David and Nanni di Banco's Isaiah for the Cathedral, emphasizing their historical significance.

5. Scene 5: Republic of Florence

**Start Time:** 01:50

**End Time:** 02:08

**Title:** Republic of Florence

**Description:** The scene highlights the Republic of Florence timeline, focusing on the Medici family and Savonarola, with a focus on their impact on the city.

6. Scene 6: Savonarola's Execution

**Start Time:** 02:08

**End Time:** 02:22

**Title:** Savonarola's Execution

**Description:** The scene depicts the execution of Savonarola and two friars in Piazza della Signoria, emphasizing the historical context.

7. Scene 7: Republic Defeats Medici and Savonarola

**Start Time:** 02:22

**End Time:** 02:36

**Title:** Republic Defeats Medici and Savonarola

**Description:** The scene shows the Republic of Florence defeating the Medici and Savonarola, with a focus on the historical figures involved.

---

---

## ARIA Example, Long Video Understanding (Scene Splitting)

---

ARIA (continued)

8. Scene 8: David's Significance  
**Start Time:** 02:36  
**End Time:** 03:29  
**Title:** David's Significance  
**Description:** The scene focuses on Michelangelo's David, highlighting its special significance for the people of Florence.
  9. Scene 9: David in Context  
**Start Time:** 03:29  
**End Time:** 03:42  
**Title:** David in Context  
**Description:** The scene provides a broader context of Florence, showing the cityscape and comparing different David statues.
  10. Scene 10: Donatello's David  
**Start Time:** 03:42  
**End Time:** 03:57  
**Title:** Donatello's David  
**Description:** The scene highlights Donatello's David, emphasizing its historical and artistic significance.
  11. Scene 11: Michelangelo's David  
**Start Time:** 03:57  
**End Time:** 05:49  
**Title:** Michelangelo's David  
**Description:** The scene returns to Michelangelo's David, focusing on its artistic details and historical context.
  12. Scene 12: David's Relocation  
**Start Time:** 05:49  
**End Time:** 06:28  
**Title:** David's Relocation  
**Description:** The scene discusses the relocation of David to its current position outside the Palazzo della Signoria.
  13. Scene 13: Historical Figures  
**Start Time:** 06:28  
**End Time:** 06:42  
**Title:** Historical Figures  
**Description:** The scene highlights historical figures like Piero de' Medici and Savonarola, emphasizing their impact on Florence.
  14. Scene 14: Michelangelo's David  
**Start Time:** 06:42  
**End Time:** 06:59  
**Title:** Michelangelo's David  
**Description:** The scene focuses on Michelangelo's David as a symbol of the newly reconstituted Republic of Florence.
-

---

## ARIA Example - Long Video Understanding (Scene Splitting)

---

ARIA (continued)

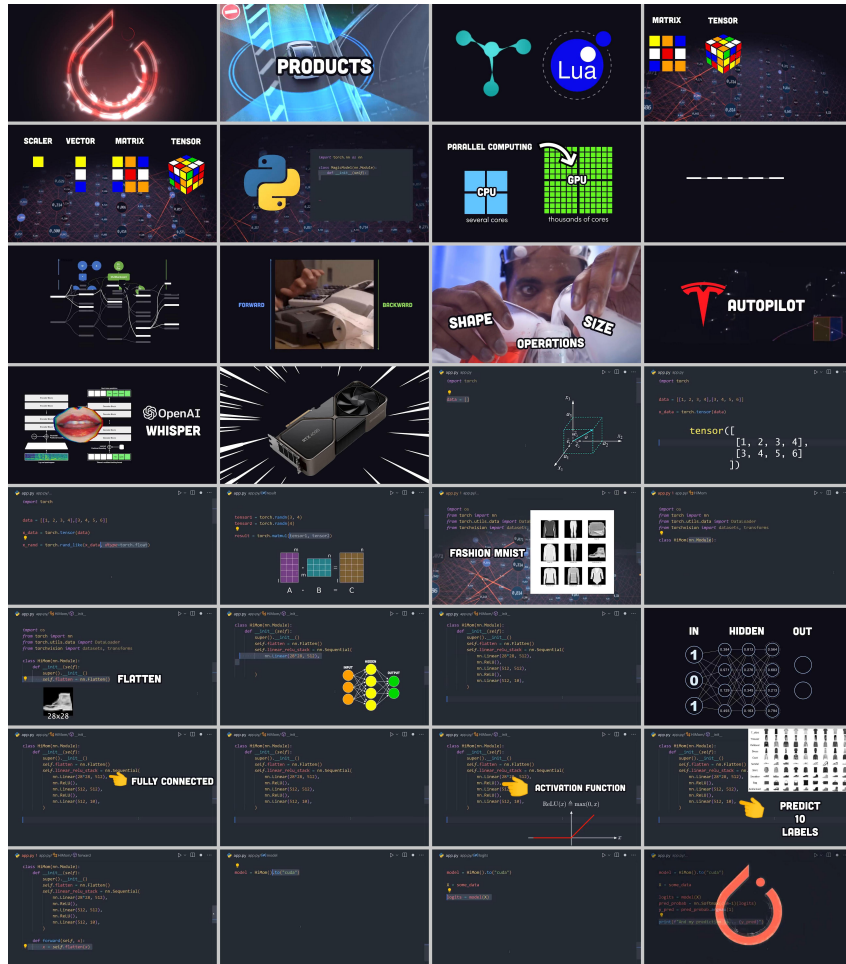
15. Scene 15: Accademia Gallery
  - Start Time:** 06:59
  - End Time:** 07:03
  - Title:** Accademia Gallery
  - Description:** The scene returns to the Accademia Gallery, showcasing the statue of David and its historical significance.
16. Scene 16: Cathedral of Florence
  - Start Time:** 07:03
  - End Time:** 07:10
  - Title:** Cathedral of Florence
  - Description:** The scene shifts to the Cathedral of Florence, highlighting its architectural significance.
17. Scene 17: David's Significance
  - Start Time:** 07:10
  - End Time:** 07:20
  - Title:** David's Significance
  - Description:** The scene emphasizes the significance of Michelangelo's David in the context of Renaissance Florence.
18. Scene 18: Conclusion
  - Start Time:** 07:20
  - End Time:** 07:24
  - Title:** Conclusion
  - Description:** The video concludes with the *smarthistory* logo, summarizing the key points discussed.
19. Scene 19: Credits
  - Start Time:** 07:24
  - End Time:** 07:27
  - Title:** Credits
  - Description:** The scene provides credits for the video, including contributors and sources.

The video has a total of 19 scenes, each with specific start and end times, titles, and detailed descriptions.

---

## ARIA Example - Code Understanding from Videos

**User** On the page with text 'Predict 10 Labels' on it, a Python class is defined. Please note down the code that defines this class. After that, explain each part of the code.



**ARIA** The Python class defined on the page with text Predict 10 Labels is as follows:

```
class HiMom(nn.Module):
    def __init__(self):
        super().__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
        )
    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

---

## ARIA Example - Code Understanding from Videos

---

ARIA (continued) Explanation of the code:

1. **Class Definition:** The class `HiMom` is defined, inheriting from `nn.Module`.
2. **Initialization (`__init__`):**
  - `super().__init__()`: Calls the initializer of the parent class.
  - `self.flatten = nn.Flatten()`: Creates a flatten layer to convert input tensors to 1D tensors.
  - `self.linear_relu_stack = nn.Sequential(...)`: Defines a sequence of layers:
    - `nn.Linear(28*28, 512)`: A fully connected layer with 28\*28 input features and 512 output features.
    - `nn.ReLU()`: ReLU activation function.
    - `nn.Linear(512, 512)`: Another fully connected layer with 512 input and output features.
    - `nn.ReLU()`: ReLU activation function.
    - `nn.Linear(512, 10)`: Output layer with 10 output features.
3. **Forward Method:**
  - `x = self.flatten(x)`: Flattens the input tensor.
  - `logits = self.linear_relu_stack(x)`: Passes the flattened tensor through the defined sequence of layers.
  - `return logits`: Returns the output logits.

---

Video source: <https://www.youtube.com/watch?v=ORMx45xqWkA>.

## 5 Acknowledgement

The authors appreciate Fan Zhou, Chengen Huang, Yanpeng Li, Chongyan Zhu, Xiaoyi Ren, Peng Liu, Chao Li, Tianhang Zhu, Chenxu Lv, Xiang He, Shenzhi Wang, Chuji Zheng, Wenhao Huang, Ziyang Luo, Shiming Yang, Fengji Zhang, Yuhao Yang, Linli Yao, Xiao Li, Lihuan Zhang, Liying Li, Wen Xie, Jun Tian, Yifan Ye and Hanshu Yan for their valuable input and suggestions.

## References

- J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-VL: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- D. Dai, C. Deng, C. Zhao, R. Xu, H. Gao, D. Chen, J. Li, W. Zeng, X. Yu, Y. Wu, et al. DeepSeekMoE: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- M. Dehghani, B. Mustafa, J. Djolonga, J. Heek, M. Minderer, M. Caron, A. Steiner, J. Puigcerver, R. Geirhos, I. M. Alabdulmohsin, A. Oliver, P. Padlewski, A. A. Gritsenko, M. Lucic, and N. Houlsby. Patch n’ pack: Navit, a vision transformer for any aspect ratio and resolution. In *NeurIPS*, 2023.
- A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Rozière, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. M. Kloumann, I. Misra, I. Evtimov, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, and et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23, 2022.
- A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. de Las Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud, L. Saulnier, M. Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao, T. Gervet, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mixtral of Experts: A sparse mixture of experts language model. *arXiv preprint arXiv:2401.04088*, 2024.
- G. Kamradt. Llmtest\_needleinahaystack, 2023. URL [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack/blob/main/README.md](https://github.com/gkamradt/LLMTest_NeedleInAHaystack/blob/main/README.md).
- H. Laurençon, L. Tronchon, M. Cord, and V. Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024.
- K. Lee, M. Joshi, I. R. Turc, H. Hu, F. Liu, J. M. Eisenschlos, U. Khandelwal, P. Shaw, M. Chang, and K. Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *ICML*, 2023.
- D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.

- B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, Y. Li, Z. Liu, and C. Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.
- J. Li, D. Li, C. Xiong, and S. C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022.
- J. Li, D. Li, S. Savarese, and S. C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023.
- B. Lin, Z. Tang, Y. Ye, J. Cui, B. Zhu, P. Jin, J. Zhang, M. Ning, and L. Yuan. MoE-LLaVA: Mixture of experts for large vision-language models. *arXiv preprint arXiv:2401.15947*, 2024a.
- X. V. Lin, A. Shrivastava, L. Luo, S. Iyer, M. Lewis, G. Ghosh, L. Zettlemoyer, and A. Aghajanyan. MoMa: Efficient early-fusion pre-training with mixture of modality-aware experts. *arXiv preprint arXiv:2407.21770*, 2024b.
- H. Liu, C. Li, Y. Li, and Y. J. Lee. Improved baselines with visual instruction tuning. In *CVPR*, 2024.
- J. Ludziejewski, J. Krajewski, K. Adamczewski, M. Pióro, M. Krutul, S. Antoniak, K. Ciebiera, K. Król, T. Odrzygózd, P. Sankowski, M. Cygan, and S. Jaszczur. Scaling laws for fine-grained mixture of experts. *arXiv preprint arXiv:2402.07871*, 2024.
- Mixtral. Pixtral 12b - the first-ever multimodal mistral model, 2024. URL <https://mistral.ai/news/pixtral-12b/>.
- Y. Qian, H. Ye, J. Fauconnier, P. Grasch, Y. Yang, and Z. Gan. Mia-bench: Towards better instruction following evaluation of multimodal llms. *arXiv preprint arXiv:2407.01509*, 2024.
- A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. In M. Meila and T. Zhang, editors, *ICML*, 2021.
- S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, 2020.
- C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, and A. Komatsuzaki. LAION-400M: open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- S. Shen, Z. Yao, C. Li, T. Darrell, K. Keutzer, and Y. He. Scaling vision-language models with sparse mixture of experts. In *EMNLP*, 2023.
- W. Shi, S. Min, M. Lomeli, C. Zhou, M. Li, G. Szilvassy, R. James, X. V. Lin, N. A. Smith, L. Zettlemoyer, S. Yih, and M. Lewis. In-context pretraining: Language modeling beyond document boundaries, 2023.
- M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- H. Wu, D. Li, B. Chen, and J. Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*, 2024.
- X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023.
- B. Zoph, I. Bello, S. Kumar, N. Du, Y. Huang, J. Dean, N. Shazeer, and W. Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.