810 7 **APPENDIX**

811 812 813

RELATED WORKS А

814 Here we provide an additional discussion of related works that were omitted from the main paper 815 816

due to lack of space. The recently released Perception Test (Patraucean et al., 2024) consists of script-based recorded videos with manual annotations focusing on 4 broad skill areas - Memory, 817 Abstraction, Physics, Semantics, however videos are only 23s long (avg). Like Neptune, ActivityNet-818 RTL (Huang et al., 2024) was constructed in a semi-automatic fashion by querying GPT-4 to generate 819 comparative temporal localization questions from the captions in ActivityNet-Captions (Krishna 820 et al., 2017). CinePile (Rawal et al., 2024) was generated by prompting an LLM to generate multiple-821 choice questions. Because it is based on movie clips, it can leverage available human-generated 822 audio descriptions. Both ActivityNet-RTL and CinePile cover only limited domains and rely on existing annotations while Neptune covers a much broader spectrum of video types and its pipeline 823 is applicable to arbitrary videos. Our rater stage is lightweight, unlike other works that are entirely 824 manual (Zhou et al., 2024; Fang et al., 2024; Wang et al., 2024). In LVBench (Wang et al., 2024), 825 even the video selection is done manually, and for MoVOA (Zhang et al., 2023b), only the decoys 826 are generated automatically. Another recently released dataset (concurrent with our submission) is 827 the Video-MME dataset (Fu et al., 2024). The motivation of this dataset is similar to ours, namely 828 it covers videos of variable lengths, with 2,700 QADs covering a wide range of different question 829 types. The main difference between Video-MME and Neptune is that the former is entirely manually 830 annotated by the authors, while we propose a scalable pipeline which can be applied to new videos and 831 domains automatically, and can be tweaked to include different question types with reduced manual 832 effort. EgoSchema is the closest work to ours in motivation, but there are some key differences: (i) 833 it is limited to egocentric videos of exactly 3 minutes each, while Neptune covers many domains and follows a more natural length distribution for online videos (16s to 15min); (ii) it relies heavily 834 on manually obtained dense captions for egocentric videos, while our method generates captions 835 automatically too and hence can be easily applied to any video online; and more importantly (iii) 836 EgoSchema also has strong image and linguistic biases, while Neptune mitigates these. 837

838 Table 5: Comparison to Existing VideoQA datasets: Ann. Type: Annotation Type, QAD: Question, Answer 839 and Decoys, Rater V: Rater verified manually. † Movies are no longer available. ‡ Annotations are hidden 840 behind a test server, 500 are public. *average/max length. **short/medium/long. 841

842	Name	Ann	Rater V	Avg. len (s)	# Vids (total/test)	# Samples (total/test)	Available
843	MovieQA (Tapaswi et al., 2016)	QAD	1	200	6,771/1,288	6,462/1,258	X †
	MSRVTT-QA (Xu et al., 2017)	QA	×	15	10,000/2,990	243,680/72,821	1
844	ActivityNet-QA (Yu et al., 2019a)	QA	1	180	5,800/1,800	58,000/18,000	1
845	NExTQA (Xiao et al., 2021)	QAD	1	44	5,440/1,000	52,044/8,564	1
040	IntentQA (Li et al., 2023a)	QAD	1	44	4,303/430	16,297/2,134	1
846	EgoSchema (Mangalam et al., 2023)	QAD	1	180	5,063/5,063	5,063/5,063	√ ‡
0/7	Perception Test (Patraucean et al., 2024)	QAD	1	23	11,600	38,000	1
047	MVBench (Li et al., 2023c)	QAD	×	16	3,641	4,000	1
848	Video-Bench (Ning et al., 2023)	QAD	1	56	5,917	17,036	1
0.40	AutoEval-Video (Chen et al., 2023c)	QA	1	14.6	327	327	1
849	1H-VideoQA (Reid et al., 2024)	QAD	1	6,300 (max)	125	125	×
850	MLVU (Zhou et al., 2024)	QAD	1	720	2K	2593	1
000	Video-MME Fu et al. (2024)	QAD	1	82.5/562.7/2,385.5**	900	2,700	1
851	LongVideoBench Wu et al. (2024)	QAD	 Image: A second s	473	3,763	6,678	1
852	Neptune	QAD	1	150/901*	2,405	3,268	1
853	Neptune-MMH	QAD	1	159/901*	1,000	1,171	1

854

855

856 858

859

В THE NEPTUNE DATASET

B.1 ADDITIONAL INFORMATION ON QUESTION TYPES

860 Neptune covers a broad range of long video reasoning abilities, which are summarised below. These question types are obtained in the Question and Answer generation stage, for which the prompt is 861 provided in Sec. C.2.3. We provide further insights into the motivations of some of the question areas 862 provided in the prompt below.

Video Summarisation: Summarise and compare long parts of the video, as well as identify the most



917 (Fig. 5). The strongest difference was for counting questions, as LLM-proposed questions were often too easy, e.g. counting the number of times a certain word is mentioned.



B.2 DOMAINS IN NEPTUNE

A full graph of the domains in Neptune are provided in Fig. 6.

B.3 COMPARISON TO OTHER BENCHMARKS

We measure the complexity of Neptune compared to other benchmarks by analyzing the progression of model performance as we add more frames to the context. We use Gemini-1.5-Flash for this comparison since it is capable of handling very large contexts. Fig. 7 shows the results of this experiment, comparing Neptune with CinePile (Rawal et al., 2024), Perception Test (Patraucean et al., 2024) and Video-MME (Fu et al., 2024). We find that most benchmarks saturate at about 50 frames, including Video-MME, which has much longer videos than Neptune. While we included Perception Test here as it is new, it does not claim to be a long video benchmark and saturates at 16 frames.

- B.4 PER-TASK PERFORMANCE
- 971 We provide detailed per-task model performance in Tab. 6. See Fig. 4 (bottom left) for a graphical representation of a subset of these results. Overall, closed-source MLLMs perform best across all

Table 6: Per-task model performance. Tasks are abbreviated as follows: TO: Temporal Ordering, CE: Cause
And Effect, SC: State Changes, VN: Visual Narrative, CI: Creator Intent, CT: Counting, PR: Predictive, GR:
Goal Reasoning, CMP: Comparison, ID: Identification, SUM: Summarization, OTH: Other. The best accuracy
per task is printed in bold and the second best underlined.

076															
970	Method	Modalities	то	CE	SC	VN	CI	СТ	PR	GR	CMP	ID	SUM	ОТН	Task-avg
977	Image models BLIP2 (Li et al., 2023b)	RGB (center frame)	24.97	48.18	40.09	47.51	71.88	33.06	40.30	33.33	39.34	24.68	38.64	18.11	38,34
978	Short Context MLLMs														
979	Video-LLaVA (Lin et al., 2023) VideoLLaMA2 (Cheng et al., 2024a)	RGB (8 frames) RGB (16 frames)	22.95 35.71	36.06 57.27	28.30 48.36	30.79 57.31	46.88 78.13	19.35 33.06	35.82 53.73	53.33 60.00	31.15 50.54	20.78 42.86	20.06 47.35	23.40 29.81	30.74 49.51
090	VideoLLaMA2 (Cheng et al., 2024a)	RGB (16 frames) + ASR	34.08	60.00	53.77	59.06	90.63	38.71	56.72	73.33	61.96	54.55	59.00	35.09	56.41
981	Long Context MLLMs - open-source MA-LMM (He et al., 2024a) MiniGPT4-Video (Ataallah et al., 2024)	RGB (120 frames) RGB (45 frames)	19.34 20.43	22.12 34.24	18.87 24.06	22.58 30.79	25.00 34.38	15.32 21.77	16.42 31.34	20.00	17.49 21.31	20.78 32.47	19.32 23.16	16.60 21.89	19.49 27.43
982	MiniCPM-V 2.6 (Yao et al., 2024)	RGB (100 frames) RGB (50 frames)	41.32	65.15	67.3	70.38	84.38 75.0	37.9	82.09 67.16	86.67	60.66	66.23	66.22	46.42	62.53
983	Closed-source MLLMs JCEF (Min et al., 2024)	VLM captions (16 frames)	48.78	63.03	64.79	70.76	78.13	43.55	62.69	60.00	60.87	50.65	64.45	55.47	60.26
984	GPT-40 ⁴ (Achiam et al., 2023) Gemini-1.5-pro (Reid et al., 2024)	RGB (8 frames) + ASR RGB (all frames) + ASR	71.25 69.39	91.21 91.52	77.25 81.69	76.83 84.21	100.0 100.0	62.90 66.94	89.55 86.57	93.33 93.33	87.98 90.22	85.71 87.01	91.30 90.41	72.45 70.19	83.31 84.29
985	Gemini-1.5-flash (Reid et al., 2024)	RGB (all frames) + ASR	63.87	88.18	77.00	<u>80.99</u>	<u>96.88</u>	56.45	82.09	<u>86.67</u>	86.96	88.31	88.79	68.30	80.37

tasks, with Gemini-1.5-pro ranking best overall and GPT-40 ranking second. Even though their average scores are close, there are significant differences in per-task scores, showing the different capabilities of each model.

C IMPLEMENTATION DETAILS

993 C.1 VIDEO SELECTION

We choose the YT-Temporal-1Bn dataset (Zellers et al., 2022b) as the source for Neptune, because of its large and diverse corpus, and because of the high correlation between vision and audio transcripts.

997 Safety & Content Filters: We filter out videos with less than 100 views, that are uploaded within 998 90 days, and those tagged by YouTube content filters to contain racy, mature or locally controversial 999 content. We then identify and remove static videos (eg. those that consist of a single frame with a 1000 voiceover) by clustering similar frames in a video and ensure that there is more than 1 cluster. We 1001 also identify and remove videos comprising primarily of "talking heads". To achieve this, we apply a 1002 per-frame frontal-gazing face-detector at 1 fps and mark the frames where the bounding box height 1003 is greater than 20% as *talking head frames*. Then, we filter out videos where more than 30% of the 1004 frames are talking head frames. These thresholds are chosen based on an F1-score on a small dev set of 50 manually annotated videos. 1005

1006 **Diversity Sampling:** From the filtered set of videos, we sub-sample 100,000 videos to boost both 1007 semantic and demographic diversity. First, we cluster the videos based on video-level semantic 1008 embeddings and tag each video with a cluster id. Second, we tag each video with the perceived age 1009 and gender demographic information contained in the video. Third, we obtain a joint distribution 1010 of semantics (cluster id) and demographics (perceived age and gender) and apply a diversity boost 1011 function (Kim et al., 2022) on the joint distribution. Finally, we sample from videos from this distribution. Fig. 8, shows the down-sampling of over-represented cluster ids before and after 1012 applying the filter. We then uniformly sub-sample the videos further to reach the desired dataset size. 1013

1014

986 987

988

989

990 991

992

1015 C.2 PROMPTS FOR DATA GENERATION

1016 In this section we provide some of the prompts used for generating Neptune.

1018 C.2.1 PROMPT FOR FRAME CAPTIONING

¹⁰²⁰ We use the following prompt to obtain a caption for each video frame:

1021
1022 Answer the following questions about the given image. Then use the
1023 information from the answers only, and write a single sentence as caption.
1024

1025 Question(Mood): Describe the general mood in the image as succinctly as possible. Avoid specifying detailed objects, colors or text.



video using the "{head_topic}" in Scenes. A part of the video script



1134 **Partial Script:** doc_segment 1135 **Visual Support Caption** 1136 To extract better visual description of the segment that will be used for QA generation in the next 1137 phase, an extra step is performed to get visual support for each segment. That visual support is stored 1138 separately in conjunction with the dense caption for the segment. For this purpose, the dense caption 1139 from the previous step is used alongside the shot level visual captions. The following LLM prompt is 1140 used to extract the visual support: 1141 ******Task:****** I provide video scene information and your job is to summarize 1142 the exact elements from "Visual Captions" that directly support the "Scene 1143 Story" of the scene below. The visuals of the scene is broken down to 1144 shots and each shot is described in a line of text in the Visual Captions. 1145 **Scene Story:** dense caption for the segment 1146 **Visual Captions:** visual_captions_of_the_segment 1147 1148 **Output Format:** Plain text with at most 200 words summarizing the 1149 supporting visual elements. 1150 1151 1152 C.2.3 GENERATING QUESTIONS AND ANSWERS 1153 1154 I want you to act as a rigorous teacher in the "Long-term Video 1155 Understanding" class. Let's test your students' in-depth comprehension! 1156 Understanding: I'll provide you with the following: 1157 - Dense Captions: A detailed breakdown of the video, including key moments 1158 and timestamps. Analyze this carefully. 1159 1160 Your Task: Craft {target_number} Challenging Short-Answer Questions 1161 Requirement: 1162 - Challenge: Demonstrate your ability to create challenging, insightful 1163 short-answer questions about the video. These shouldn't test simple recall 1164 only. Aim to probe understanding of relationships, motives, subtle details, 1165 and the implications of events within the video. 1166 - Diversity: Design a variety of question types (more on this below). 1167 - Specificity: Each question must be self-contained and laser-focused on 1168 a single concept or event from the video. Avoid compound or overly broad 1169 questions. 1170 - Answers: Model the ideal answer format: Brief, accurate, and rooted 1171 directly in evidence from the video's content. 1172 - Video-Centric: Stay true to what's explicitly shown or stated in 1173 the video. Avoid relying on outside knowledge or speculation. Design 1174 questions so the correct answer cannot be easily determined without carefully analyzing the video. 1175 - Minimize Information Leakage: For question types like ranking or 1176 ordering, ensure that the order of candidates or options listed in the 1177 question doesn't inadvertently reveal the correct answer. Shuffle them to 1178 maintain neutrality. 1179 - Content-First: Timestamps and section titles within the captions are 1180 there for guidance. Do not explicitly refer to those markers in your 1181 questions or answers. Focus on the events and elements themselves. 1182 - Unambiguous: Ensure each question has a single, clearly defined correct 1183 answer. Avoid questions that are open to multiple interpretations (e.g., 1184 counting elements where viewers might disagree). 1185 - Visual Elements: Questions focused on visual reasoning or visual narratives should emphasize the interpretation of the visuals. Keep the 1186 question minimal, letting the answer describe the specific visual elements 1187 in detail.

1188 You want to test students' capabilities of understanding the video, 1189 including but not limited to the following aspects: 1190 Ability: Summarize and compare long parts of the video. 1191 Ability: Compress information from the video rather than just listing the 1192 actions that happened in the video. 1193 Ability: Identify the most important segments of the video. 1194 Ability: Recognize and understand the visual elements in different parts 1195 of the video. 1196 Ability: Understand the timeline of events and the plot in the video. 1197 Ability: Count objects, actions and events. Focus on higher-level 1198 counting where the same instance does not occur in all/every frame and actions are sufficiently dissimilar. 1199 Ability: Understand and reason about cause and effect in the video. 1200 Ability: Understand the unspoken message that the audience may perceive 1201 after watching the video, which may require common sense knowledge to 1202 infer. 1203 Ability: Understand the visual reasoning of why and how important visual 1204 content is shown in the video. 1205 Ability: Understand the visual narrative of the video and the mood of the 1206 video and which visual elements do contribute to that. 1207 Ability: Understand object states change over time, such as door opening 1208 and food being eaten. 1209 Presentation 1210 1211 - QUESTION: Introduce each question as "QUESTION 1, 2, 3: (capability) full question". - ANSWER: Follow the format "CORRECT ANSWER: correct answer". 1212 1213 Good example questions: - Question (counting): How many ingredients are 1214 added to the bowl in total throughout the video? Correct Answer: 3. 1215 - Question (goal reasoning): What is the purpose of the man standing in 1216 front of the whiteboard with a diagram on it? Correct Answer: To explain 1217 the features and capabilities of the vehicle. 1218 - Question (cause and effect): How does the document help people to be 1219 happier? Correct Answer: It helps people to identify and focus on the 1220 things that make them happy, and to develop healthy habits. 1221 1222 - Question (timeline event): In what order are the following topics 1223 discussed in the video: history of pantomime, importance of pantomime, 1224 mime as a tool for communication, benefits of pantomime? Correct Answer: Mime as a tool for communication, history of pantomime, importance of 1225 pantomime, benefits of pantomime. 1226 1227 - Question (predictive): What happens after the man jumps up and down on 1228 the diving board? Correct Answer: He jumps into the pool. 1229 - Question (summarization): What is the overall opinion of the reviewers 1230 about Hawaiian Shaka Burger? Correct Answer: The food is good, but the 1231 patties are frozen. 1232 - Question (creator intent): What message does the video creators try 1233 to send to the viewers? Correct Answer: Nature is essential for human 1234 well-being. 1235 1236 - Question (visual-temporal): What color is the scarf that Jessica wears 1237 before she enters the restaurant? Correct Answer: Red. 1238 - Question (visual narrative): How does John's overall facial expression 1239 contribute to the explanation of the financial situation that is described 1240 in the video? Correct Answer: He shows sad feelings and expression when 1241

1242 he described the financial collapse of the company which adds to the sense 1243 of empathy that video describes. 1244 - Question (visual reasoning): What was shown to support the effects of a 1245 high cholesterol diet in the video? Correct Answer: Video demonstrates 1246 how cholesterol gradually clogs blood vessels, using an animation to 1247 illustrate the cross-section of vessels and the buildup of plaque. 1248 Bad example questions because it can be answered by common sense. 1249 Question (counting): How many players are there in a soccer team? Correct 1250 Answer: 11. 1251 1252 Bad example questions because it asks for trivial details. - Question 1253 (counting): How many times the word 'hurricane' is said in the video? 1254 Correct Answer: 7. 1255 Bad example questions because the summary of topics are subjective and 1256 ambiguous. - Question (timeline event): List the sequence of topics 1257 Grace discusses in the video, starting with the earliest. Correct Answer: 1258 Getting ready for a photoshoot, attending a baseball game, showing off her 1259 new outfit, playing a Wayne's World board game, and discussing her upcoming week. 1260 1261 Dense Caption with Timestamps: {video_inputs_str} 1262 1263 C.2.4 GENERATING DECOYS FROM QUESTIONS AND ANSWERS 1264 1265 Role: You are a rigorous teacher in a "Long-term Video Understanding" class. You will assist students in developing strong critical thinking 1267 skills. This requires creating sophisticated test questions to accompany 1268 video content. 1269 Understanding: I will provide: 1270 1271 - Dense Captions: A breakdown of the video, including structure, key 1272 events, and timestamps. - Target Questions & Answers: A set of 1273 {target_number} questions about the video, along with their correct answers. 1274 Task: Generate High-Quality Multiple-Choice Questions 1275 1276 1. Analyze: Carefully study the dense captions, questions, and correct answers. Familiarize yourself with the nuanced details of the video 1277 content. 1278 1279 2. Decoy Design: For each target question, generate {decoy_number} 1280 incorrect answers (distractors). These distractors must be: 1281 - Challenging: Plausible to the point where students need deep content understanding and critical thinking to choose the correct answer. 1282 - Stylistic Match: Mimic the style, tone, and complexity of the correct 1283 answer. 1284 - Similar Length: Keep length close to that of the correct answer, 1285 preventing students from eliminating choices based on length differences. 1286 - Factually Relevant: Related to the video content, even if slightly 1287 incorrect due to a detail change, misinterpretation, or logical fallacy. 1288 - Reasonable: Each decoy should be something that could be true, making 1289 simple elimination impossible. 1290 Specific Techniques for Distractor Creation 1291 1292 - Subtle Tweaks: Alter a minor detail from the correct answer (e.g., change 1293 a time, location, or name). - Confusing Similarity: Use a concept from elsewhere in the video that 1294 seems related but applies to a different context. 1295

- Misdirection: Introduce a true statement related to the video's theme but

1296 not directly answering the question. 1297 - Order Shuffling: If the question involves the order of events, subtly 1298 rearrange the order within the distractors. 1299 Presentation: 1300 1301 - QUESTION: Repeat the provided question faithfully (e.g., "QUESTION 1 1302 (Capability): ...") - CORRECT ANSWER: Repeat the correct answer (e.g., "CORRECT ANSWER: ...") 1303 - WRONG ANSWERS: List each wrong answer on a separate line without using 1304 letters to label choices (e.g., "WRONG ANSWER 1: ...", "WRONG ANSWER 2: 1305 ...") 1306 1307 1308 *GOOD* Example: Question: What are the three main challenges that the college is taking on? Correct Answer: Food scarcity, pollution, and 1309 disease. Wrong Answer 1: Global warming, deforestation, and poverty. 1310 Wrong Answer 2: Hunger, homelessness, and crime. Wrong Answer 3: Obesity, 1311 malnutrition, and food insecurity. Wrong Answer 4: Food waste, water 1312 shortages, and air pollution. 1313 *BAD* examples where the decoys format is different from correct answer: 1314 Question: What color is the shirt that the woman is wearing? Correct 1315 Answer: Black. Wrong Answer 1: The woman is wearing a white shirt. 1316 Wrong Answer 2: The woman is wearing a blue shirt. Wrong Answer 3: The 1317 woman is wearing a green shirt. Wrong Answer 4: The woman is wearing a 1318 red shirt. 1319 *BAD* examples because only the correct answer is in positive sentiment. 1320 Question: What is the overall sentiment of the man in the video? Correct 1321 Answer: He is overjoyed with his new gift. Wrong Answer 1: He is upset 1322 his gift is not big enough. Wrong Answer 2: He is sad about life in 1323 general. Wrong Answer 3: He is upset the gift is not great. Wrong Answer 1324 4: He seems down and unhappy. 1325 1326 Dense Caption with Timestamps: {video_inputs_str} 1327 Question and Correct Answer: {question_and_answer_str} 1328 1329 C.2.5 QAD FILTERING 1330 The following prompt is used to filter out questions that can solve from QADs alone. 1331 1332 Instructions: 1333 Carefully analyze the following question and options. Rank the options 1334 provided below, from the most likely correct answer to the least likely 1335 correct answer. Please respond with "ANSWER" and "EXPLANATION". 1336 1337 Your response should be in the following format: 1338 * ANSWER: [Letter of the ranking, split by greater than symbol. (e.g., "ANSWER: A > B > C > D > E")]. 1339 * EXPLANATION: [Provide a brief explanation of your choice. Do not repeat 1340 the option.] 1341 1342 QUESTION: {question_str} 1343 Options: {options_str} 1344 1345 Please provide your response below. 1347 C.3 HUMAN RATING AND CORRECTION OF QADS 1348

We provide a screenshot of the UI used by raters to annotate automatically generated QADs in Fig.10. Note that if any of the four options under the 'Is the question valuable' field are not selected,

then the question is discarded from the dataset. We made sure to train raters using training raters (with detailed decks and feedback rounds), as well as applying rater replication (we used 3 raters per question independently), and rater pipelining (having an experienced rater verify the answer from a previous rater) in order to correct hallucinations and other mistakes, and discard QADs that were inappropriate. Overall, of the total 11,030 QADs that we obtained automatically, 7,762 (70%) were discarded by raters.



Figure 10: Screenshot of rater UI.

1382 C.4 FILTERING SUBSETS

1380 1381

Here we provide details for how we select the thresholds used to create the NEPTUNE-MMH and NEPTUNE-MMA subsets. For both subsets, we filtered NEPTUNE-FULL with the QAD filter described in Sec. 4.4. For NEPTUNE-MMA, we additionally filtered out QADs that human raters marked as requiring only the audio modality and answer (see Sec. 4.5). We refer to this as the "rater test". For NEPTUNE-MMH, we instead applied the ASR filter (Sec. 4.4). Both QAD and ASR filters were run by prompting an LLM (Gemini 1.0 Pro) three times, each time with a different random seed and then removing QADs that the LLM answered correctly at least X out of three times, where X is the threshold for the test.

1391 Fig. 11 shows how choosing different thresholds affects dataset size and accuracy scores. The top 1392 row shows the choices for the NEPTUNE-MMH subset. Raters marked almost half of the questions 1393 as answerable from audio only, so the rater filter already cuts the dataset size in half. Successively 1394 applying the QAD filter with increasing thresholds reduces data size up until less than 25%. We 1395 benchmark three models on the different subsets that have access to ASR only, vision only, or both vision and ASR, respectively. As expected, all three models show declining performance, with the 1396 ASR-only model showing the biggest losses. This suggests that all models were inferring the correct answer from the QAD only, which the filter successfully mitigates. The vision-only model gains 1398 slightly from removing QADs that fail the rater rest, which is expected as the test removes QADs 1399 that rely on audio, which the model does not have access to. However, like for the other models, its 1400 accuracy declines when adding the QAD test. 1401

1402 The bottom row of Fig. 11 shows the choices for the NEPTUNE-MMA subset where we use the ASR 1403 filter and the QAD filter with identical thresholds. This filter set has a stronger effect on the dataset size, reducing it to less than 15% of its original size at the highest threshold. Because the ASR-only



1428Figure 11: Effect of filtering thresholds for the NEPTUNE-MMH (top row) and NEPTUNE-MMA (bottom row)1429subsets.

model was used for the ASR filter, we exclude it from the accuracy comparison. The vision-only
and vision+ASR models both show declining accuracy with increasing thresholds. As expected, the
accuracy of the vision+ASR model declines faster. The effect of this filter set on the accuracy is
much stronger than that of the above filter set, suggesting that it increases the difficulty of the dataset
more strongly. Even the vision-only model declines faster than above, suggesting that this filter set
generally removes easier questions, even those that rely on vision only.

For both filtered sets, we opted to set the threshold to two, which in both cases significantly increases the dataset difficulty while still preserving enough QADs for statistically meaningful evaluation metrics. We noticed that when setting the threshold to three, there were less than five QADs left for some question types, preventing robust accuracy estimation for these tasks.

- 1442
- 1443 C.5 IMPLEMENTATION DETAILS FOR BENCHMARKS
- 1444 1445 C.5.1 BLIND BASELINES

For the Gemini-1.5-pro baseline with text only the prompt used was: "Carefully analyze the question and all available options then pick the most probable answer for this question"

- 1448
- 1449 C.5.2 VIDEO-LLAVA

For Video-LLaVA the following prompt was used - "Pick a correct option to answer the question.Question: question Options: options ASSISTANT:".

- 1453
- 1454 C.5.3 VIDEOLLAMA2 1455
- During inference, we uniformly sampled 8 frames from each video. Each frame undergoes padding
 and resizing to a standardized dimension. The pre-processed frames are then fed into the image
 encoder. These steps are set as default in the inference script provided by videoLlama2.

1458 **QAD Prompt:** _PROMPT_TEMPLATE = """Pick a correct option number to answer the question. 1459 Question: {question} Options: {options}:""" 1460 **OE Prompt:** Question: {question} 1461 1462 Output post processing: We eliminated extra characters and spaces using regex to get the final ID of 1463 the predicted option. 1464 1465 C.5.4 MINIGPT4-VIDEO 1466 We set the 300 maximum number of output tokens to be 300 for the open-ended task and 10 for the 1467 multiple choice eval. The prompts are as follows: 1468 PROMPT TEMPLATE MCQ = """Question: select the correct option for this task: question 1469 Options: options. Output format: [OPTION]: [Reason]""" 1470 1471 _PROMPT_TEMPLATE_OPEN_ENDED = """Question: question Answer:""" 1472 1473 C.5.5 MA-LMM 1474 1475 We set the 300 maximum number of output tokens to be 300 for the open-ended task and 300 for the multiple choice eval. The prompts are as follows: 1476 1477 _PROMPT_TEMPLATE_MCQ = """Question: select the best choice for this task: question Options: 1478 options Answer:""" 1479 PROMPT TEMPLATE OPEN ENDED = """Question: question Answer:""" 1480 1481 C.5.6 GPT-40 PROMPTS 1482 1483 **Open-ended evaluation with transcript** 1484 You are an expert in video understanding and question answering. You can 1485 analyze a video given its image sequence and and transcript and answer 1486 questions based on them. 1487 {video_frames} 1488 1489 Video Transcript: {transcript} 1490 Answer the question using the image sequence. Do not describe the frames 1491 just answer the question. Question: {question} 1492 1493 **Open-ended evaluation without transcript** 1494 You are an expert in video understanding and question answering. You can 1495 analyze a video given its image sequence and answer questions based on 1496 them. 1497 {video frames} 1498 1499 Answer the question using the image sequence. Do not describe the frames 1500 just answer the question. Question: {question} 1501 Multiple-choice evaluation with transcript 1502 You are an expert in video understanding and question answering. You can 1503 analyze a video given its image sequence and and transcript and answer 1504 questions based on them. 1505 1506 {video_frames} 1507 Video Transcript: {transcript} 1508 Answer the question using the image sequence. Do not describe the frames 1509 just answer the question by identifying the choice. Question: {question} 1510 Choices: {choices} Please identify the correct CHOICE and explain your 1511 reasoning concisely. Output Format: [CHOICE]: [REASON]

1512 Multiple-choice evaluation without transcript

You are an expert in video understanding and question answering. You can analyze a video as an image sequence and answer questions based on that.

1516 {video_frames}

Answer the question using the image sequence. Do not describe the frames just answer the question by identifying the choice. Question: {question} Choices: {choices} Please identify the correct CHOICE and explain your reasoning concisely. Output Format: [CHOICE]: [REASON]

1521 1522

1523

C.6 COMPUTE RESOURCES

The compute heavy part of the project was image frame captioning (as this involves reading high dimensional pixel data). The rest of the pipeline involves largely text-only LLMs and hence was less compute heavy. We estimate that the entire project in total took roughly 256 TPU v5e running over a period of 50 days.

1528

1530

¹⁵²⁹ D ADDITIONAL DETAILS FOR GEM

1531 D.1 CREATION OF GEM EQUIVALENCE DEV SET 1532

To create a development set that allows us to estimate the accuracy of different open-ended question answering metrics on Neptune, we sampled 97 question-answer pairs from the dataset and generated 3 candidate answers per question by prompting VideoLLAVA (Lin et al., 2023), Gemini-1.5-pro (Reid et al., 2024) and MA-LMM (He et al., 2024b) to write a free-form answer for each question without looking into the decoys or ground truth. We then manually annotated these responses between 0 and 1 by comparing it to the ground truth answer. We made sure that the annotators are blind to the model to avoid any bias. The resulting set has 292 equivalence pairs with an average score of 0.32, with 85 examples having score greater 0.5 and 206 examples with score less than 0.5

1540

1541 D.2 BENCHMARKING ON THE DEV SET

1543 In Table. 1, we evaluate several open-ended metrics on our dev set. The task of the metric is to classify 1544 whether the open-ended response and ground-truth answer are equivalent or not. We report F1-scores to balance false-positives and false-negatives. We evaluate both traditional rule-based metrics such 1545 as CIDEr and ROUGE-L, as well as established model-based metrics such as BEM(Bulian et al., 1546 2022). We also try using Gemini-1.5-pro (Reid et al., 2024) as an LLM based equivalence metric 1547 (by prompting it to estimate equivalence). First, we note that as expected, Gemini-1.5-pro correlates 1548 well with the human ground-truth annotation of the set, achieving a high F1-score of 72.5. However, 1549 given that Gemini is not open-source and proprietary, any change in the model can affect all the prior 1550 results in an external leader-board making it challenging as a metric. Traditional rule-based metrics 1551 perform much worse than Gemini-1.5-pro on this dev set as they are n-gram based and struggle to 1552 handle the diversity of domains and styles in the open-ended responses. The BERT model based 1553 BEM metric (Bulian et al., 2022) performs similarly, achieving an F1-score of 61.5.

Next, we evaluate lightweight open-source language models Gemma-2B (Team et al., 2024a), Gemma-7B (Team et al., 2024a) and Gemma-9B (Team et al., 2024b) in a zero-shot setting and find that performance improves with model size, with Gemma-9B bridging the gap well between traditional metrics and the Gemini-1.5-pro based metric. Finally, we fine-tune Gemma-9B on the open-source BEM answer equivalence dataset (Bulian et al., 2022), and find that Gemma-9B finetuned on the BEM dataset performs the best on our dev-set. We name this metric *GEM*.

1560 1561

1561 D.3 IMPLEMENTATION DETAILS

We use instruction-tuned variants of the Gemma models (gemma-it-2b, gemma-it-7b and gemma-it-9b) for our experiments. To develop a prompt, we experiment with several variations in a zero-shot setting and measure the performance on the dev-set. Our final prompt is shown below. To ensure responses occur in a standard format, we simply measure the softmax-probability over "TRUE"

response indicating the statements are equivalent and "FALSE" response indicating the statements are 1567 not equivalent. For each model, the threshold over probability is chosen to maximize the F-1 score 1568 on dev set. To finetune Gemma models on BEM dataset, we tokenize the same prompt as used in 1569 the zero-shot setting and train it using prefix-LM tuning for 10000 iterations using a learning rate of 1570 1e-6. For evaluation, we truncate the open-ended responses to 100 words, use a decode cache size of 1024 and threshold the softmax probability of the LM using the chosen threshold from dev-set. 1572 <start_of_turn>user Answer Equivalence Instructions: 1574 1575 Carefully consider the following question and answers. 1576 You will be shown a "gold-standard" answer from a human annotator, 1577 referred to as the "Reference Answer" and a "Candidate Answer". 1578 Your task is to determine whether the two answers are semantically 1579 equivalent. 1580 1581 In general, a candidate answer is a good answer in place of the "gold" reference if both the following are satisfied: 1582 1. The candidate contains at least the same (or more) relevant information 1583 as the reference, taking into account the question; in particular it 1584 does not omit any relevant information present in the reference. 1585 2. The candidate contains neither misleading or excessive superfluous 1586 information not present in the reference, taking into account the 1587 question. 1589 Your response should be one word, "TRUE" or "FALSE", in the following format: 1591 ANSWERS_ARE_EQUIVALENT: [TRUE or FALSE] 1592 Question: 1593 " { } " 1594 1595 Candidate Answer: 1596 " { } " 1597 1598 Reference Answer: 1599 " { } " Please provide your response below. <end_of_turn> 1603 <start_of_turn>model ANSWERS ARE EQUIVALENT: 1604 1605 1606 D.4 QUALITATIVE EXAMPLES FOR METRIC Below, we provide some qualitative examples of our metric applied to the predictions provided by 1608 the baseline models Video LLaMA2(Cheng et al., 2024a), Gemini-1.5-flash (Reid et al., 2024) and 1609 Gemini.1.5-pro (Reid et al., 2024). We find that Gemini-1.5-flash responses are qualitatively slightly 1610 better than Gemini-1.5-pro responses as they contain less hallucinations. Furthermore, we find that 1611 with a single frame as input, Gemini-1.5-pro usually responds with a variation of "This information 1612 is not present is the provided frame" instead of making an educated guess. This is unlike the case of 1613 multiple choice questions where the model commits to one of the choices. These aspects are reflected 1614 in the results presented in Table.4 and Table.3. We also find that the metric fails often in the questions 1615 about listing items / events in order. In such cases, the model is biased towards accepting the response 1616 (false positives) even if the order is incorrect provided the items / events are present in the candidate 1617 response. 1618

1619 Question: Compare the Pampered Chef Veggie Strip Maker to a food processor in terms of their effectiveness in making cauliflower

```
1620
          rice.
1621
     Answer: The Veggie Strip Maker is more effective than a food
1622
          processor in making cauliflower rice because it prevents the
1623
          cauliflower from becoming mushy.
1624
     Video LLaMA2: Both the Pampered Chef Veggie Strip Maker and a food
1625
          processor can be effective in making cauliflower rice, but they
          differ in their methods of operation and the resulting texture of
1626
                       The Pampered Chef Veggie Strip Maker is a kitchen
          the rice.
1627
          appliance designed specifically for making vegetable strips,
1628
          including cauliflower rice. It uses a grating mechanism to create
1629
          rice-like pieces by shredding the cauliflower. This method can
1630
          produce a relatively uniform and rice-like texture.
                                                                  On the other
1631
          hand, a food processor is a versatile kitchen appliance that can
1632
          also be used to make ...
1633
      GEM Equivalence: False
1634
      Gemini-1.5-flash: The Pampered Chef Veggie Strip Maker makes
1635
          cauliflower rice that is less mushy than a food processor.
1636
      GEM Equivalence: True
      Gemini-1.5-pro: The Pampered Chef Veggie Strip Maker is more
1637
          effective than a food processor because it doesn't make the
1638
          cauliflower mushy.
1639
      GEM Equivalence: True
1640
1641
1642
      Question: How does the use of behavioral training help in the
1643
          process of taking foot x-rays on elephants?
1644
     Answer: Behavioral training encourages elephants to willingly
1645
          participate in their healthcare, making the process of taking foot
1646
          x-rays less stressful for both the elephants and the caretakers.
      Video LLaMA2: The use of behavioral training plays a crucial role in
1647
          the process of taking foot x-rays on elephants. By training the
1648
          elephants to follow commands and stand on a platform, it allows the
1649
          caretakers to safely and effectively capture the x-rays of the
1650
          elephants' feet. This training helps to ensure that the elephants
1651
          remain still and cooperative during the x-ray process, minimizing
1652
          any potential harm or discomfort. Additionally, the behavioral
1653
          training helps to build a positive relationship between the
1654
          elephants and their caretakers, fostering trust and cooperation.
1655
          Overall, the use of behavioral training enhances the accuracy and
1656
          efficiency of the
1657
      GEM Equivalence: True
1658
      Gemini-1.5-flash: Behavioral training allows the elephants to
          willingly participate in the training sessions and participate in
1659
          their own health care.
1660
      GEM Equivalence: True
1661
      Gemini-1.5-pro: Behavioral training allows the elephants to
1662
          willingly participate in the x-ray process by presenting their feet
1663
          for imaging.
1664
      GEM Equivalence: True
1665
1666
      Question: What are the key features of the Adam Purvis Progeny
1667
          knife's blade?
1668
     Answer: The Adam Purvis Progeny knife's blade is made of 3.4-inch
          M390 steel. The blade shape is useful, simple, and highly effective,
1669
          the edge is very sharp, featuring a stonewash finish on the primary
1670
          bevels and a horizontal satin on the flats. It is exceptionally
1671
          sharp and has a 150 thousandths stock thickness, making it suitable
1672
          for everyday carry tasks.
1673
     Video LLaMA2: The Adam Purvis Progeny knife's blade is made of S35VN
```

```
1674
          steel, which is known for its high strength and excellent edge
1675
          retention. The blade has a flat ground and a hollow grind, which
1676
          allows for a thin and sharp edge. Additionally, the blade has a
1677
          black and silver color scheme, giving it a sleek and modern
1678
          appearance.
      GEM Equivalence: False
1679
      Gemini-1.5-flash: The Adam Purvis Progeny knife's blade is made of
1680
          M390 steel, has a modified sheepsfoot style, and is ground very
1681
          thinly for a sharp edge.
1682
      GEM Equivalence: False
1683
      Gemini-1.5-pro: The Adam Purvis Progeny knife features a 3.4-inch
1684
          M390 steel blade with a modified sheepsfoot shape, thin grind, sharp
1685
          edge, and a stonewashed finish on the bevels with a satin finish on
1686
          the flats.
1687
      GEM Equivalence: True
1688
1689
1690
      Question: How does interacting with civilians in nearby villages
          enhance cadets' leadership skills?
1691
     Answer: Engaging with civilians in nearby villages allows cadets to
1692
          develop respectful communication and relationship-building skills,
1693
          which are essential for effective leadership in various situations.
1694
     Video LLaMA2: Interacting with civilians in nearby villages enhances
1695
          cadets' leadership skills by providing them with real-world
1696
          experiences and opportunities to apply their knowledge and training
1697
          in a practical setting. By engaging with the local community, cadets
1698
          learn to communicate effectively, build trust, and foster positive
1699
          relationships with people from diverse backgrounds. This hands-on
1700
          experience helps them develop their interpersonal skills, cultural
          awareness, and problem-solving abilities, all of which are crucial
1701
          for successful leadership in various situations. Additionally,
1702
          interacting with civilians allows cadets to gain a deeper
1703
          understanding of the challenges and needs of the community, which
1704
          can inform their decision-making and ...
1705
      GEM Equivalence: True
1706
      Gemini-1.5-flash: Interacting with civilians in nearby villages
1707
          teaches cadets how to speak to people respectfully and build trust.
1708
      GEM Equivalence: True
1709
      Gemini-1.5-pro: It teaches them how to speak to people.
1710
      GEM Equivalence: False
1711
1712
      Question: List the sequence of topics discussed in the video,
1713
          starting with the earliest.
1714
     Answer: 1. Introduction of the GAN356 AIR cube and its unique
1715
          features. 2. Discussion of the GAN356 AIR's weight and GES (Gan
1716
          Elasticity System) springs. 3. Lubrication, tensioning, and personal
1717
          experience with the GAN356 AIR's performance. 4. Assessment of the
1718
          GAN356 AIR's performance, including its smoothness, speed, corner-
1719
          cutting, and tendency to lock up. 5. Overall evaluation of the
1720
          GAN356 AIR, including its strengths and weaknesses, and
1721
          recommendations for potential buyers.
1722
     Video LLaMA2: 1. Introduction to the video and the Rubik's Cube. 2.
          Explanation of the different types of Genuine Handcubes. 3.
1723
          Demonstration of how to solve the Rubik's Cube, including the
1724
          various algorithms and techniques used. 4. Showcasing of the
1725
          different colors and types of Rubik's Cubes available. 5.
1726
          Conclusion and closing credits.
1727
      GEM Equivalence: False
```

1728	Gemini-1.5-flash: The video starts with an introduction of the GAN
1729	356 Air cube, then discusses its features, including the GES system
1730	and its light weight. The reviewer then shares their opinion on the
1731	different GES springs and their effect on the cube's feel. The video
1732	continues with a discussion on the cube's performance, including its
1733	smoothness, speed, and tendency to lock up. Finally, the reviewer
1734	concludes with a recommendation for viewers who might be interested
1735	in purchasing the cube.
1736	GEM Equivalence: True
1737	Gemini-1.5-pro: The video discusses the main features of the Gan 356
1738	Air, the different types of springs, the weight and price of the
1739	cube, the reviewer's experience with the cube, and the reviewer's
1740	overall opinion of the cube.
1741	GEM Equivalence: False
1742	
1743	
1744	Question: In what order are the following events mentioned in the
1745	video: introduction of the wallet, features of the wallet, and care
1746	instructions?
1747	Answer: Introduction of the wallet, features of the wallet, care
1748	instructions.
17/10	Video LLaMA2: The video begins with the introduction of the wallet,
1750	followed by a description of its features, and concludes with
1751	instructions on how to care for the wallet.
1751	GEM Equivalence: True
1752	Gemini-1.5-flash: The speaker introduces the wallet, then describes
1754	its features, and finally mentions the care instructions.
1734	GEM Equivalence: Irue
1700	features of the wallet
1756	GEM Equivalence: False
1/5/	GER Equivalence. raise
1758	
1759	
1760	
1761	

E MODEL-BASED TEMPORAL CERTIFICATE

1764 The idea of a temporal certificate was introduced by EgoSchema (Mangalam et al., 2023) as a way 1765 of capturing the intrinsic temporal understanding length for long video QA datasets. It is defined 1766 as 'the length of the video a human verifier needs to observe to be convinced of the veracity of the marked annotation'. While the authors used it to uncover flaws in existing long video QA datasets, as 1767 well as to provide a difficulty measure independent of video length, we find that is has the following 1768 drawbacks: (i) it does not take into account the *length of time* or the *effort* taken by the annotator 1769 themselves, to find the correct time span in videos; (ii) it requires manual annotation from expert 1770 annotators to measure; and finally (iii) is subjective. 1771

As an attempt to mitigate these issues, we introduce a slightly modified version of the temporal certificate, which is *Model-Based*. We calculate this certificate using 129 samples from Neptune and EgoSchema, respectively. For this experiment we used Gemini 1.5 Pro, with one "driver" model run to answer the question and two other model runs with different random seeds to verify if the answer was not correct by random chance. Along with the question and options, we provided video clips of various lengths from the center of the video, and at various fps, as shown in Fig. 12.

Since this experiment queried a set of frames over various clip lengths, we defined it as the "needle in haystack" problem. Here, the needle is defined as a frame or set of frames needed to answer the question correctly, matching a human's ground truth response, while the haystack is a set of frames which need to be watched to find the needle frames. Iteratively, we increase the video length and fps for the query until the model achieves the correct response.



Figure 12: **Model-based Temporal Certificate:** Illustration of video clip querying for the model-based temporal certificate experiment. The red clip is the clip length that resulted in an incorrect response. As we increased the clip length wider, and the model correctly answered the question, we logged the frame count for incorrect response and correct response, and stopped querying. Besides clip length, we vary the fps of the query clip.



Figure 13: Frame level temporal certificate: We compared our dataset sample with EgoSchema to evaluate the number of frames needed by model to answer questions correctly. The figures above show the distribution of the minimum number of frames required to achieve the correct response.

1791

1792

1793

1795

As shown in Fig. 13, we find that the model needs more frames to answer the question correctly for the Neptune dataset as compared to EgoSchema. This resulted in a mean of 5.39 as certificate frames for Neptune which is 3.37 times the mean certificate frame number of 1.6 for EgoSchema. On the clip length level this translated to a mean of 21.22s of clip needed to respond correctly on the Neptune dataset, whereas for EgoSchema the mean was 9.07s. The model-based certificate lengths turn out to be much smaller than the certificate lengths reported by EgoSchema, where humans needed close to 100s to answer the questions for EgoSchema.

In addition, we define the *effort score* as the fraction of the maximum number of frames needed to be watched before answering the question correctly, as defined in Equation 1. An effort score closer to 0 suggests that the needle isn't very small compared to the haystack, i.e. most of the frames contain the answer to the question; while a high effort score means a high percentage of haystack frames needs to be included before we cover all frames required to answer correctly.

1827 1828

1829
1830EFFORT SCORE = $\frac{MAX NUMBER OF FRAMES RESULTING IN AN INCORRECT RESPONSE}{MIN NUMBER OF FRAMES RESULTING IN A CORRECT RESPONSE}$ (1)1831
1832For Neptune, the mean effort score was 0.47, whereas for EgoSchema, it was 0.19. This suggests that
Neptune requires 2.47 times the effort compared to EgoSchema according to the definition above,

1834 Neptune requires 2.47 times the effort compared to EgoSchema according to the definition above,
 1835 which closely corroborates the above results for the mean clip lengths needed to solve the questions from the respective datasets.

