

Deep Temporal Reasoning in Video Language Models: A Cross-Linguistic Evaluation of Action Duration and Completion through Perfect Times

Anonymous ACL submission

Abstract

Human perception of events is intrinsically tied to distinguishing between completed (perfect and telic) and ongoing (durative) actions, a process mediated by both linguistic structure and visual cues. In this work, we introduce the **Perfect Times** dataset, a novel, quadrilingual (English, Italian, Russian, and Japanese) multiple-choice question-answering benchmark designed to assess video-language models (VLMs) on temporal reasoning. By pairing everyday activity videos with event completion labels and perfectivity-tailored distractors, our dataset probes whether models truly comprehend temporal dynamics or merely latch onto superficial markers. Experimental results indicate that state-of-the-art models, despite their success on text-based tasks, struggle to mirror human-like temporal and causal reasoning grounded in video. This study underscores the necessity of integrating deep multimodal cues to capture the nuances of action duration and completion within temporal and causal video dynamics, setting a new standard for evaluating and advancing temporal reasoning in VLMs.

1 Introduction

Understanding how events unfold in time requires a detailed analysis of their both sequential and causal relationships. Sequential events are not simply arranged chronologically; rather, one event often triggers the next upon reaching its completion. Moens and Steedman (Moens and Steedman, 1988) highlighted that human memory organizes actions based on contingency, their cause-effect relationships, where a cause reaches its culmination before triggering an effect. This causal linkage is encoded in language through grammatical time and, critically, through aspect, specifically, perfectivity and its semantic correlate, telicity. Telicity refers to a property of verb phrases that denotes a definitive endpoint (e.g., “to put something somewhere”),

while atelic (or durative) expressions (e.g., “to hold something”) lack such clear termination. This property is deeply woven into the language (Tenny, 1994).

Temporal relations, therefore, are not conveyed solely through grammatical time; they are also robustly signaled by aspect. Yet, language often omits critical information that visual modality can provide. Observing how events unfold in time, like in videos, offer dynamic cues, such as key frame transitions and motion patterns, that decisively indicate whether an action has been completed or is still ongoing. This multimodal integration is indispensable for resolving ambiguities inherent in linguistic descriptions of events.

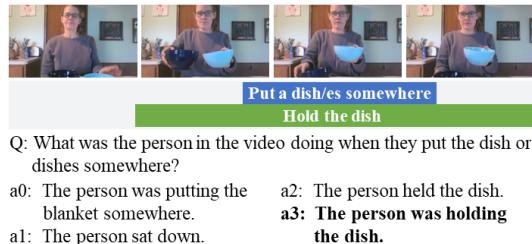


Figure 1: An example from the Perfect Times dataset illustrating the interaction between a completed action (*to put a dish or dishes somewhere*) and a durative action (*to hold the dish*). The blue and green stripes indicate temporal progression. The correct answer highlighted in bold.

To evaluate these phenomena, we introduce the **Perfect Times** dataset, a multiple-choice question-answering (MCQA) benchmark designed specifically for video-language models (VLMs). Our questions are constructed as complex sentences that juxtapose two actions, thereby probing the interplay between verb forms and temporal conjunctions in main and dependent clauses. Twelve carefully designed templates systematically cover all possible temporal relations between the two actions in accordance with the universal Allen’s in-

068 terval algebra (Allen, 1984). Moreover, our study
069 adopts a quadrilingual approach (English, Italian,
070 Russian, and Japanese) to capture the diverse ways
071 in which languages encode time and aspect. For
072 instance, Russian frequently encodes perfectivity
073 at the lexical level, making markers of completion
074 predefined in verbs, while Italian involves a sophis-
075 ticated interaction between time, mood, and aspect,
076 as it has more branched *concordanza dei tempi*
077 (the sequence of tenses) than English. In Japanese,
078 the commonplace understated and ambiguous lin-
079 guistic encoding necessitates stronger reliance on
080 visual context to disambiguate temporal relations.

081 By combining visual data with linguistically
082 complex questions, our work addresses the fol-
083 lowing questions:

- 084 • How do VLMs leverage linguistic markers
085 and visual cues to distinguish sequential from
086 simultaneous actions?
- 087 • Are there cross-linguistic differences in the
088 interpretation of temporal relations based on
089 grammatical encoding of perfectivity?
- 090 • How closely do model predictions align with
091 human understanding of event completion
092 expressed in grammatical time and aspect?

093 Figure 1 presents an example of an English
094 question-answer (QA) pair from the Perfect Times
095 dataset. In addition to determining the correct
096 action to answer the question in the video, the
097 model must catch the matching forms of seman-
098 tic or grammatical aspect in the question and the
099 answer option. For human native speakers, mak-
100 ing such comparisons is not difficult, but none of
101 the tested state-of-the-art models reach the 50%
102 accuracy threshold.

103 Our comprehensive, cross-linguistic approach
104 thus aims to set a new benchmark for evaluating
105 multimodal temporal reasoning in VLMs.

106 2 Related Works

107 2.1 Cognitive Approach

108 Several studies use the visual world paradigm to
109 examine how aspect influences real-time language
110 comprehension. Foppolo et al. (2021) shows that
111 visual and linguistic signals jointly shape the inter-
112 pretation of completed actions, guiding anticipatory
113 eye movements through verb aspect and visual
114 cues. Similarly, Foppolo et al. (2016) finds that
115 Italian adults rapidly focus on images of finished

116 events upon hearing the corresponding verb. In
117 an eye-tracking study, Bosch et al. (2021) reports
118 that while Italian children use verb semantics and
119 aspect to anticipate outcomes, their processing of
120 aspectual information lags behind basic lexical se-
121 mantics. van Hout (2008) challenges uniformity
122 in aspect acquisition by demonstrating that Italian
123 children acquire perfectivity later than their Polish
124 and Dutch peers. Minor et al. (2022) further show
125 that perfective and imperfective aspects influence
126 listeners’ expectations differently across Russian,
127 Spanish, and English, with Russian perfectives
128 strongly indicating completion and the English sim-
129 ple past being less reliable. Finally, Chang et al.
130 (2023) provide evidence from animation experi-
131 ments that visual cues of goal information affect
132 the choice of past versus progressive verb forms in
133 Japanese. These findings contributing to cognitive
134 science illustrate that VLMs may have greater poten-
135 tial to replicate human-like language processing
136 than text-only large language models (LLMs).

137 2.2 Time and Aspect in Transformers

138 Zhao et al. (2021) compared aspect interpretation
139 in humans and transformers using verb tenses and
140 resultative structures. Humans responded to all
141 cues, while transformers were sensitive to explicit
142 telicity but struggled with resultatives. Lombardi
143 and Lenci (2023) found that the transformer-based
144 (Vaswani et al., 2017) Italian model GilBERTo
145 performed similarly to humans on inherently telic
146 and atelic verbs, yet had difficulty with context-
147 dependent verbs, suggesting its limitations in han-
148 dling nuanced larger temporal contexts. Metheniti
149 et al. (2022) showed that transformer-based models
150 classify activities by duration and perfectivity in
151 English and French with over 80% accuracy.

152 However, all these studies are on text-only
153 LLMs, and the experiments lack a dedicated multi-
154 modal benchmark that fully captures all possible
155 temporal relations grounded in physical world.

156 2.3 Action Recognition and MCQA for Video

157 Numerous datasets offer short video clips for ac-
158 tion recognition (Kay et al., 2017; Liu et al., 2021;
159 Heilbron et al., 2015; Damen et al., 2020; Sigurdsson
160 et al., 2016). Kinetics (Kay et al., 2017) pro-
161 vides general action recognition (10-second clips);
162 FineAction (Liu et al., 2021) and ActivityNet (Heil-
163 bron et al., 2015) deliver fine-grained temporal an-
164 notations. EPIC-Kitchens (Damen et al., 2020) and
165 Charades (Sigurdsson et al., 2016) target specific

Template	MCA	DCA	Question	Answer
Precedence				
3	t	a	What had the person in the video done before holding the sandwich?	The person had opened the refrigerator.
4	t	t	What had the person in the video done before the other person walked through a doorway?	The person had closed the laptop.
6	a	a	What had the person in the video been doing before playing with a phone or camera?	The person had been watching the television.
7	a	t	What had the person in the video been doing before turning off the light?	The person had been playing with a phone or camera.
Succession				
1	t	a	What did the person in the video do after tidying up the table?	The person put the food somewhere.
2	t	t	What did the person in the video do after they had opened the box?	The person put the bag somewhere.
10	a	a	What was the person in the video doing after sitting in a chair?	The person was sitting on the floor.
11	a	t	What was the person in the video doing after taking the cup from somewhere?	The person was drinking from the cup.
9	a	t	What was the person in the video doing when the other person closed the door?	The person was watching the television.
Simultaneity				
5	t	a	What did the person in the video do while the other person was sitting on the floor?	The person opened the door.
8	a	a	What was the person in the video doing while holding the book?	The person was holding the bag.
12	t	t	What did the person in the video do when they grasped onto a doorknob?	The person opened the door.

Table 1: Examples of questions and answers in Perfect Times generated by temporal and aspectual templates with respect to the telicity markers (t: telic, a: atelic).

scenarios for classification and localization.

These datasets lay the foundation for the video MCQA datasets. Templated questions appear in MSVD-QA and MSRVTT-QA (Xu et al., 2017), STVQ (Jang et al., 2019), and STAR (Wu and Yu, 2021), while manually annotated formats are used in ActivityNet-QA (Yu et al., 2019), TVQA (Lei et al., 2018), NExT-QA (Xiao et al., 2021), Charades-SRL-QA (Sadhu et al., 2021), and Action Genome (Ji et al., 2019). All datasets cover content in English.

Our multilingual Perfect Times contains of the Charades videos, labelled with Action Genome action class markers, and uses STAR-like temporal templates to emphasize aspect, the dimension overlooked in previous work.

2.4 Video Language Models

Video Language Models integrate visual and textual processing through three key components: a pre-trained visual encoder, a pre-trained large language model (LLM), and a modality interface (Zhong et al., 2022). The visual encoder (Radford et al., 2021; Li et al., 2022, 2023) compresses raw video or audio into compact representations, while pre-trained LLMs (Chung et al., 2022; Chiang et al., 2023; Touvron et al., 2023) supply broad world knowledge for downstream tasks. Since LLMs cannot directly process encoder outputs, a learnable interface aligns the modalities, often via a Q-Former that integrates at the token (Lin et al., 2023) or feature (Alayrac et al., 2022) level.

While many models score high on popular benchmarks (not least because the benchmark data may appear in their training sets) they typically

struggle with novel scenarios. As we developed a completely new benchmark and set the baseline, for evaluation we focused on top models from the HuggingFace VLM leaderboard, including Pangea (Yue et al., 2024), Gemini (Team et al., 2024), PALO (Rasheed et al., 2025), Video-Llama2 (Cheng et al., 2024), Phi-3.5-vision (Abdin et al., 2024), Qwen2-VL (Wang et al., 2024), Llava-NEXT-Video (Zhang et al., 2024), InternVL2 (Chen et al., 2024), MiniCPM-V (Yao et al., 2024) and DeepSeek-VL (Lu et al., 2024). We filtered out the models that do not support the languages featured in Perfect Times.

3 Method

We anchor the main clause action, the queried action in the correct answer (*mca*), and position the dependent clause action which sets the temporal and causal context relative to it (*dca*). This sequential shift reflects their temporal progression and fully covers Allen’s interval algebra (Allen, 1984) that encodes all possible temporal relationships between two actions. Linguistically, these relations map to verb tenses and aspects and subordinate temporal conjunctions: for instance, *after* (Italian *dopo* (*che*), Russian *после того, как*, Japanese *後*) indicates that the *dca* precedes the *mca*; *while* (Italian *mentre*, Russian *пока*, Japanese *ながら*) denotes simultaneity; and *before* (Italian *prima* (*che*), Russian *перед тем, как*, Japanese *前に*) signals that the *dca* follows the *mca*.

For perfectivity, which captures an action’s continuous or complete (and causally linked) nature, we focus on the start and end points of both the main (*st_mca*, *et_mca*) and dependent (*st_dca*,

166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199

200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233

et_dca) clauses. We also label each action as telic or atelic based on its inherent properties. Accordingly, we group Allen’s temporal relations into three categories (precedence, succession, and simultaneity for the main clause action) and define the plausible aspect correlations within each group. For example, an incomplete action without a clear endpoint unlikely occurs *before* another action – a hardly plausible combination in sequential events even if grammatically acceptable in languages like Russian.

Next, inspired by STAR (Wu and Yu, 2021) and its sequence functional programs that exploit before/after relations, we developed templates in four languages to cover all temporal boundaries in conjunction with the semantics of action completion in both the main clause and dependent clause actions¹. The templates were developed with the help of native speakers who have linguistic or philological background and specialize in translations.

For instance, when *mca* is a completed action, we use as its base *What did the person in the video do...* (Italian Cosa ha fatto la persona nel video.../Cosa aveva fatto la persona nel video.../Cosa fece la persona nel video...; Russian Что сделал человек на видео...; Japanese ビデオに写っている人は...何をしましたか). The ongoing action as *mca* corresponds to *What was the person in the video doing...?* (Italian Cosa stava facendo la persona nel video... or Cosa faceva la persona nel video...; Russian Что делал человек на видео...; Japanese ビデオに写っている人は...何をしていましたか).

Below we outline temporal relations and templates by groups, while the complete list of templates in all languages is given in Appendix C. The examples of all the temporal questions in Perfect Times are given in Table 1².

3.1 Precedence

As shown in Figure 2, when mca happens before dca ($\text{et_mca} \leq \text{st_dca}$), all combinations of action completeness are possible. In the languages with tense agreement, English and Italian, the tense of

¹Our goal is not to exhaustively cover all surface-level expressions of temporal relations (e.g., alternative conjunctions, adverbial phrases, or nominal constructions). Instead, we focus on evaluating the model’s overall comprehension of temporal semantics with respect to perfectivity, regardless of paraphrasing.

²In the table, we provide examples in English for the general reference; some other examples in all the languages are in Appendix D.

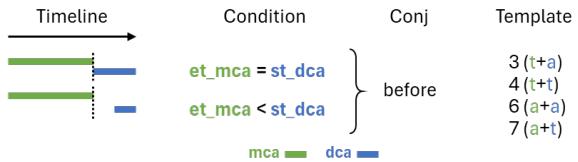


Figure 2: Precedence of mca in templates 3, 4, 6, and 7 with all combinations of mca and dca telicity markers.

mca is backshifted to a past form relative to dca. Russian and Japanese, in contrast, convey only the aspect of both mca and dca in place of shifting tenses.

3.2 Succession

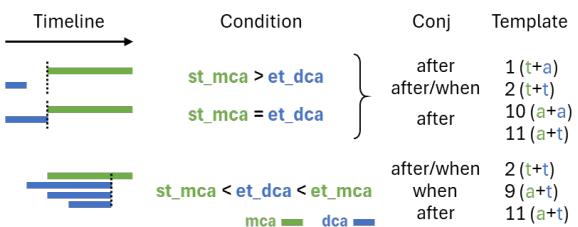


Figure 3: Succession of mca in templates 1, 2, 9, 10, 11. The strict sequence of actions allows for all combinations of telicity in mca and dca. Attention is focused on st_mca after or simultaneously with et_dca. If mca continues after dca, the sequence of events is focused on et_dca, which is perfective.

When mca follows dca, the trigger is the dca’s endpoint: at least some part of mca must take place after it. Symmetrically to precedence, English and Italian have agreement of tenses.

In case of partial following of the telic dca by the atelic mca the general temporal conjunction *when* (Italian quando; Russian когда; Japanese 時) is used the same way as the explicit sequential *after*, even though both actions may not strictly follow one another. This relationship is coded by Template 9 (Figure 4).



Figure 4: Example in Perfect Times made by template 9b: two actions with different agents, durative mca and telic dca.

3.3 Simultaneity

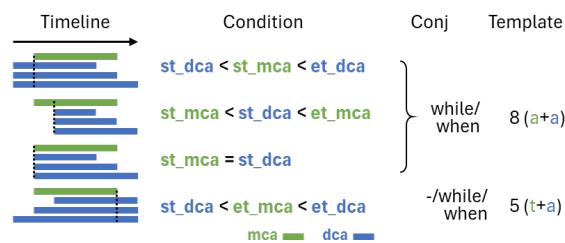


Figure 5: Overlapping mca and dca in templates 5 and 8.

The prototypical conjunctions for expressing overlapping actions, or simultaneity, are *while* (Italian mentre, Russian пока, Japanese ながら/間に) and the more general *when*. The continuous nature of dca (against which the mca of any perfectivity is questioned) underlies these relations.

When two telic actions finish at the same time, the ambiguous *when* emphasizes the simultaneity of the completion – as a synonym to *the moment when...* (Figure 6).

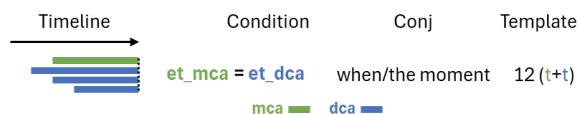


Figure 6: Template 12: simultaneity when both actions mca and dca end at the same moment in.

Each QA pair comes with three distractors that consider both the semantics and the grammatical form of the action.

3.4 Distractors

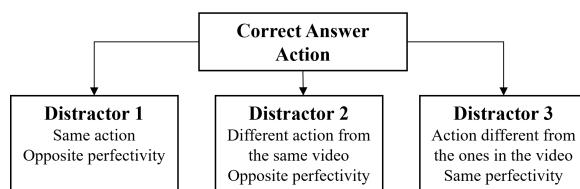


Figure 7: Distractors in relation to the Correct Answer.

Each distractor type is designed to probe different aspects of comprehension: linguistic nuance, context relevance, and discriminative ability within a shared scene. The diagram in Figure 7 illustrates the configuration of the distractors.

Distractor Type 1 is a variation of the correct answer rendered in a different form to convey the opposite completeness of the action (e.g., switching from past simple tense to a continuous/ing form

in English or taking the opposite aspect word in Russian). It is intended to be very close to the correct answer with the subtle difference in linguistic details. This tests whether a model (or human) can notice and correctly interpret the aspect mismatch between the question and the answer.

Distractor Type 2 comes from an alternative action that occurs in the same video, but uses the switched perfectivity/duration type the same way as in Type 1. Being contextually related, it forces the evaluator to differentiate between two plausible actions within the same scene. It tests the boundary between action comprehension and grammatical understanding.

Distractor Type 3 is taken from a completely different context. It comes from the list of actions that do not appear in the current video. Its purpose is to introduce an option that should be clearly out of context for the video at hand. A model with good action recognition should be able to dismiss such an option as completely implausible.

Among the four answer choices, both completed and durative actions are represented equally. Each answer option (one correct and three distractors) must be assigned to a random position for each question. To avoid option selection bias (Loginova et al., 2024), all distractor types and the correct answer should be distributed approximately equally across the answer options.

4 Experiments and Results

4.1 Dataset

Four annotators³ labeled 400 clips of the Charades (Sigurdsson et al., 2016) collection of approximately 30-second videos. All activities in these videos were categorised into 157 verb action classes derived from Action Genome (Rai et al., 2021), along with their respective time boundaries⁴.

³The annotators were recruited as volunteers and were all proficient in at least two target languages. Two of them have a background in linguistics, and their academic levels include two undergraduates and two postgraduates. Although none were directly involved in the project, each received detailed briefing and debriefing instructions from the professional linguist on the team. Their language proficiencies are as follows: Annotator 1 — English (fluent), Italian (intermediate), Russian (native), Japanese (fluent); Annotator 2 — English (native), Russian (native); Annotator 3 — English (advanced), Italian (native), Russian (native); Annotator 4 — Italian (native), Russian (native).

⁴The Charades dataset was initially annotated with activities and their time intervals for different research purposes. This previous annotation was inadequate for precisely determining the sequence of events or their concurrent occurrence due to its broad definition of action boundaries and a general

354
355 The statistics of videos and classes are given in Ap-
356 pendix A.

357 The action classes, which are essentially verb
358 phrases, were also annotated with telicity labels
359 irrespective of their context. Inter-annotator agree-
360 ment scores 0.67 by Fleiss' kappa, which is consid-
361 ered substantial (Landis and Koch, 1977), although
362 indicates that there is some inherent ambiguity or
363 subjectivity in assessing the verb classes for telicity
364 in isolation.

365 The algorithm begins by mixing actions within
366 the same video. It goes line by line in the shuf-
367 fled annotation file and assigns mca to the first
368 action and dca to the second one for each pair of
369 neighboring actions. Since the mca and dca an-
370 notations have telicity labels and time boundaries
371 used as conditions, the algorithm applies all cor-
372 responding templates and generates the questions
373 and correct answer options. Then the dataset is
374 populated with distractors. The correct answers
375 and distractor types are balanced across the answer
376 options. We end up with the dataset of 3,739 QA
377 pairs.

378 Subsequently, speakers of each language were
379 asked to take this MCQA test with instructions sim-
380 ilar to the prompts for VLMs: *Choose the correct*
381 *answer (a0, a1, a2, or a3) and respond only with*
382 *the option key (e.g., a0)*⁵. As a result, we obtained
383 an inter-annotator agreement of 0.8 according to
384 Fleiss' kappa (substantial agreement). This is how
385 the gold standard of 93.36% accuracy was devel-
386 oped. This high percentage is due to the fact that
387 distractors are designed in such a way that the cor-
388 rect answer is not difficult to intuitively predict.
389 Notably, most annotators made the most mistakes
390 in Distractor Type 2. Additional statistics on the
391 annotators' responses are presented in Appendix
E.

392 4.2 Models

393 We tested both the open-source and closed-source
394 multilingual VLMs: Qwen2-VL (Wang et al.,
395 2024), MiniCPM-V (Yao et al., 2024), InternVL2
396 (Chen et al., 2024), LLaVA-NeXT-Video (Zhang
397 et al., 2024), GPT-4o (Achiam et al., 2023), and

398 interpretation of actions without considering their causal and
399 temporal relationships. Additionally, the initial markup failed
400 to note whether different actions were performed by the same
401 character or different ones, a crucial detail for constructing
402 meaningful questions.

403 ⁵We deliberately did not provide additional instructions to
404 annotators, as we aimed to collect the data based on linguistic
405 intuition.

406 Gemini-2.0-Flash-Lite (Team et al., 2024). Ad-
407 ditional details on open-source models are in Ap-
408 pendix F.

409 We passed the video either directly (Qwen2-VL,
410 Gemini-2.0-Flash-Lite) or first extracted frames
411 every 3 seconds and passed the list of frames to
412 those models that can only process a sequence of
413 images (GPT-4o, InternVL2, LLaVA-NeXT-Video,
414 and MiniCPM-V).

415 4.3 Results

416 4.3.1 Evaluation

417 We measure **Accuracy**, **F1 Macro**⁶ and **Distractor**
418 **Rate by Type**, which is the percentage of choices
419 of each distractor type in case of error, calculated
420 as following:

421 Let us define D_i — the type of distractor chosen
422 by the model for question i , if it was predicted
423 incorrectly; T_k — a specific type of distractor (e.g.,
424 Type 1, Type 2, Type 3); n_k — the number of times
425 the model selected a distractor of type T_k .

426 The proportion of errors on distractors of type
427 T_k is calculated as:

$$P(T_k) = \frac{n_k}{\sum_j n_j} \times 100, \quad (1)$$

428 where $P(T_k)$ — the probability of selecting a
429 distractor of type T_k and $\sum_j n_j$ — the total num-
430 ber of cases where the model made an error⁷.

431 4.4 Quantitative and Qualitative Results

432 Table 2 shows the results of our tests on Perfect
433 Times dataset with the breakdown by the distractor
434 types.

435 Across all the languages, LLaVA-NeXT-Video
436 performs the worst. Despite its multilingual LLM
437 backbone (based on Qwen (Bai et al., 2023)), its
438 performance is nearly random for the languages
439 other than English. Moreover, answer analysis re-
440 vealed its severe selection bias⁸ towards the first
441 and second options in English and the first option
442 in the other languages. This indicates a system-
443 atic flaw in its logic and temporal reasoning, even
444 though it is able to handle basic contextual cues
445 (Distractor Type 3 is a minority). Its temporal

446 ⁶We also checked F1 Micro, but because the dataset is
447 rather balanced, this metric with TP/FP/FN does not differ
448 much from Accuracy.

449 ⁷A similar metric can be applied if we compute the error
450 frequency relative to all predictions, not just erroneous ones:
451 $P(T_k) = \frac{n_k}{N} \times 100$, where N is the total number of questions.

452 ⁸All other models revealed less prominent selection bias
453 in the first option, a0.

Table 2: VLM performance on Perfect Times across different languages (in percentage).

Model	Accuracy	F1 Macro	Distructor Type 1	Distructor Type 2	Distructor Type 3
English					
Gemini-2.0-flash-lite	43.41	42.19	22.15	30.42	4.01
GPT-4o	43.25	34.39	21.96	31.14	3.64
MiniCPM-V-2_6	36.19	35.5	27.68	31.08	5.05
Qwen2-VL-7B-Instruct	35.09	32.68	21.24	39.2	4.47
InternVL2-8B	34.86	33.36	27.34	29.09	8.7
LLaVA-NeXT-Video-7B	33.38	32.86	25.39	30.99	10.22
Italian					
Gemini-2.0-flash-lite	43.11	42.15	21.85	30.78	4.25
Qwen2-VL-7B-Instruct	41.59	39.63	21.95	32.09	4.35
GPT-4o	40.71	32.34	25.29	31.45	2.49
MiniCPM-V-2_6	37.42	35.73	29.55	28.96	4.07
InternVL2-8B	34.07	32.38	32.00	24.63	9.3
LLaVA-NeXT-Video-7B	25.92	17.65	25.40	29.74	18.88
Russian					
Gemini-2.0-flash-lite	46.99	46.33	19.04	29.85	4.12
GPT-4o	45.04	44.97	19.60	32.15	3.21
InternVL2-8B	36.87	34.37	28.82	25.54	8.77
MiniCPM-V-2_6	36.29	34.61	28.88	30.20	4.63
Qwen2-VL-7B-Instruct	34.13	32.75	20.40	39.5	5.96
LLaVA-NeXT-Video-7B	26.95	24.7	24.3	36.00	12.74
Japanese					
Gemini-2.0-flash-lite	43.06	41.77	22.19	31.43	3.32
MiniCPM-V-2_6	38.73	36.31	30.25	26.34	4.68
GPT-4o	38.49	30.65	23.78	35.23	2.49
Qwen2-VL-7B-Instruct	37.52	36.85	20.17	37.90	4.41
InternVL2-8B	35.05	31.65	31.22	23.92	9.81
LLaVA-NeXT-Video-7B	26.18	19.24	25.25	33.06	15.41

capabilities could be improved through enhanced multilingual action recognition of spatiotemporal features.

Gemini-2.0-flash-lite consistently performed better in all the languages with close accuracy in English, Italian and Japanese. The other models have a greater variance in accuracy, which is likely due to differing saturation of training data for each language. However, the accuracy of all the models in all the languages is significantly lower than the human gold standard. In combination with the consistent error trends across all the languages and all the models, it highlights that strongest current VLMs lack robust mechanisms for precise temporal fusion.

Gemini-2.0-flash-lite, GPT-4o and InternVL2-8B generally perform better on Russian data, where perfectivity is lexically encoded, and worse on Japanese, which demands stronger visual disam-

biguation. This reinforces our claim from Section 1 that language-specific encoding of time and aspect greatly affects temporal reasoning. It also echoes the findings in cognitive literature discussed in Section 2.1 that for native speakers perfectivity encoded in telicity is easier to interpret.

Another parallel between human cognitive processes of causality and the behavior of models comes from the qualitative analysis: all the tested VLMs, when mistaken, prefer perfect (or telic) answers to durative (or atelic) ones. This may stem from the specifics of the source data with strong causal connections in sequential actions.

All models tend to ignore the linguistic aspect, as the most errors are related to Distractors Type 1 and Type 2, where the predicted verb form does not match the verb form of the question. The majority of errors across all models except InternVL2 are of the Type 2 distractor. Apparently, the models do

477
478
479
480
481
not understand which moment in time is crucial for
the answer, even though they are able to distinguish
actions relevant to the video. Therefore, VLMs
should leverage multimodal fusion or specialized
architectures to capture temporality.

482
483
484
485
486
487
488
489
490
491
InternVL2, in contrast, seems to have a good
understanding of temporal context but struggles
with subtle linguistic differences in aspect, as its
majority of errors are of the Type 1 distractor in
every language except for English (possibly due to
greater data contamination). Its worst performance
in Italian, the language with the most versatile
range of verb forms for encoding temporal rela-
tions, proves that the model does not understand
the difference between grammar forms.

492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
Regarding specific temporal conjunctions, all
tested models struggle with disambiguating *when*
(and its equivalents in other languages). The high-
est number of correct answers is about simultaneity
marked with *while* (and its equivalents) by Tem-
plate 8a1. Conversely, in the similar pattern where
when is used instead of *while* to denote the same
temporal relationship with identical verb forms
(Template 8a2), there are significantly more errors⁹.
Similarly, the unambiguous conjunction *after* in
2a1 and 2b1 has a lower error rate than in 2a2
and 2b2 that use *when* to describe the same action.
Such behavior points to a weakness in aligning tem-
poral boundaries and semantics of action relations
with verb aspect. Following the cognitive find-
ings that underline human reliance on integrated
multimodal cues, it indicates that the separate lan-
guage stream that comes in multimodal models
along with the vision modality is insufficient for
temporal reasoning in videos.

5 Conclusion

512
513
514
515
516
517
518
519
520
521
522
523
524
We have demonstrated that a full understanding
of event dynamics requires an integrated approach
where both visual cues and grammatical structures
inform the interpretation of action completion. Our
cross-linguistic evaluation shows that the SoTA
VLMs rely on limited superficial cues, particu-
larly when disambiguating ambiguous temporal
conjunctions and distinguishing subtle aspectual
markers. Models favor telic responses, mirroring
human tendencies to find cause and effect depen-
dencies in events and exposing a reliance on lexical
over grammatical signals. These observations re-

525
526
527
528
529
inforce the need for enhanced multimodal fusion
and specialized temporal representations.

530
531
532
533
534
535
536
537
538
539
A key advantage of our **Perfect Times** dataset
is its semi-synthetic, template-based design. By
leveraging a unique methodology that uses univer-
sal temporal and aspectual templates, our dataset
can be easily augmented and adapted to any lan-
guage provided videos are annotated with action
labels and timestamps. This flexibility makes our
benchmark not only a novel tool for evaluating
temporal reasoning in VLMs but also a scalable
resource to drive further research. Our work thus
sets a new standard for probing the intricate inter-
play of time, causality, and aspect in multimodal
contexts.

6 Limitations

540
541
542
543
544
545
546
Our study is the first to link tense, aspect, and vi-
sual modality to test deep temporal reasoning in
models. However, it is not exhaustive regarding the
full range of synonymous temporal constructions.
Future work could augment the dataset with para-
phrases generated by LLMs to increase coverage.

547
548
549
550
551
552
553
Additionally, our semi-synthetic, template-
based approach may not fully reflect the variability
found in natural language. Although templates
ensure controlled testing of specific temporal re-
lations, they might not fully reflect spontaneous
speech. Expanding the dataset to include more nat-
uralistic examples could provide further insights.

554
555
556
557
558
559
560
Another limitation is the number of verb classes:
157 classes for 400 videos. While a broader range
of classes would capture more nuanced actions,
the current quantity is balanced by high-quality an-
notations and language-specific templates, which
were meticulously crafted by experts, a resource-
intensive process.

561
562
563
564
565
566
567
568
Finally, with more annotations per language, the
gold standard accuracy might decrease when av-
eraged across annotators, although it is unlikely
to reach the performance levels of current state-
of-the-art models. Future research should explore
scaling up the dataset, both in terms of linguis-
tic variety and video domains, to fully assess and
enhance multimodal temporal reasoning.

7 Ethics Statement

569
570
571
572
573
We introduce a new benchmark dataset derived
from the video datasets *Charades* and *Action
Genome*. We re-annotate their videos to support
the new evaluation of both open- and closed-source

⁹Some examples with quantitative data are provided in the Appendix G.

574 vision-language models (VLMs). In doing so, we
575 strictly adhere to the licensing terms of the original
576 datasets, ensuring that our derivative work com-
577 plies with all copyright and usage restrictions.

578 Our annotation process involved trained anno-
579 tators following the guidelines aimed at ensuring
580 consistency. Despite these efforts, we acknowledge
581 that any human annotation process may introduce
582 subjective interpretations. We therefore encourage
583 users of our dataset to consider potential annotation
584 biases when interpreting experimental results.

585 The evaluation of VLMs, particularly those that
586 generate open-ended text, carries inherent risks.
587 We have taken measures to mitigate such risks by
588 prompting and conducting evaluations. We pro-
589 mote transparency through the public release of
590 our dataset and code upon acceptance of this paper.
591 Our intention is to foster reproducible research and
592 to provide a resource that can contribute to improv-
593 ing the trustworthiness and robustness of VLMs.

594 References

- 595 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed
596 Awadallah, Ammar Ahmad Awan, Nguyen Bach,
597 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat
598 Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck,
599 Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav
600 Chaudhary, Dong Chen, Dongdong Chen, Weizhu
601 Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng,
602 Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen
603 Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao,
604 Min Gao, Amit Garg, Allie Del Giorno, Abhishek
605 Goswami, Suriya Gunasekar, Emman Haider, Jun-
606 heng Hao, Russell J. Hewett, Wenxiang Hu, Jamie
607 Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi,
608 Xin Jin, Nikos Karampatziakis, Piero Kauffmann,
609 Mahoud Khademi, Dongwoo Kim, Young Jin Kim,
610 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi
611 Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui
612 Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu,
613 Weishung Liu, Xiaodong Liu, Chong Luo, Piyush
614 Madan, Ali Mahmoudzadeh, David Majercak, Matt
615 Mazzola, Caio César Teodoro Mendes, Arindam Mi-
616 tra, Hardik Modi, Anh Nguyen, Brandon Norick,
617 Barun Patra, Daniel Perez-Becker, Thomas Portet,
618 Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang
619 Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy,
620 Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil
621 Salim, Michael Santacroce, Shital Shah, Ning Shang,
622 Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia
623 Song, Masahiro Tanaka, Andrea Tupini, Praneetha
624 Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan
625 Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel
626 Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia
627 Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu,
628 Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang,
629 Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu,
630 Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen
631 Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yu-
632 nan Zhang, and Xiren Zhou. 2024. Phi-3 technical
633 report: A highly capable language model locally on
634 your phone. *Preprint*, arXiv:2404.14219.
- OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, 635
Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale- 636
man, Diogo Almeida, Janko Altenschmidt, Sam Alt- 637
man, Shyamal Anadkat, Red Avila, Igor Babuschkin, 638
Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim- 639
ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, 640
Jake Berdine, Gabriel Bernadett-Shapiro, Christo- 641
pher Berner, Lenny Bogdonoff, Oleg Boiko, Made- 642
laine Boyd, Anna-Luisa Brakman, Greg Brockman, 643
Tim Brooks, Miles Brundage, Kevin Button, Trevor 644
Cai, Rosie Campbell, Andrew Cann, Brittany Carey, 645
Chelsea Carlson, Rory Carmichael, Brooke Chan, 646
Che Chang, Fotis Chantzis, Derek Chen, Sully 647
Chen, Ruby Chen, Jason Chen, Mark Chen, Ben- 648
jamin Chess, Chester Cho, Casey Chu, Hyung Won 649
Chung, Dave Cummings, Jeremiah Currier, Yunx- 650
ing Dai, Cory Decareaux, Thomas Degry, Noah 651
Deutsch, Damien Deville, Arka Dhar, David Do- 652
han, Steve Dowling, Sheila Dunning, Adrien Eco- 653
fet, Atty Eleti, Tyna Eloundou, David Farhi, Liam 654
Fedus, Niko Felix, Sim'on Posada Fishman, Just- 655
ton Forte, Isabella Fulford, Leo Gao, Elie Georges, 656
Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel 657
Goh, Raphael Gontijo-Lopes, Jonathan Gordon, Mor- 658
gan Grafstein, Scott Gray, Ryan Greene, Joshua 659
Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal- 660
lacy, Jesse Han, Jeff Harris, Yuchen He, Mike 661
Heaton, Johannes Heidecke, Chris Hesse, Alan 662
Hickey, Wade Hickey, Peter Hoeschele, Brandon 663
Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost 664
Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, 665
Angela Jiang, Roger Jiang, Haozhun Jin, Denny 666
Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer 667
Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kan- 668
itscheider, Nitish Shirish Keskar, Tabarak Khan, Lo- 669
gan Kilpatrick, Jong Wook Kim, Christina Kim, 670
Yongjik Kim, Hendrik Kirchner, Jamie Ryan Kiros, 671
Matthew Knight, Daniel Kokotajlo, Lukasz Kon- 672
draciuk, Andrew Kondrich, Aris Konstantinidis, 673
Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael 674
Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Le- 675
lung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly 676
Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, 677
Ryan Lowe, Patricia Lue, Anna Adeola Makanju, 678
Kim Malfacini, Sam Manning, Todor Markov, Yaniv 679
Markovski, Bianca Martin, Katie Mayer, Andrew 680
Mayne, Bob McGrew, Scott Mayer McKinney, 681
Christine McLeavey, Paul McMillan, Jake McNeil, 682
David Medina, Aalok Mehta, Jacob Menick, Luke 683
Metz, Andrey Mishchenko, Pamela Mishkin, Winnie 684
Monaco, Evan Morikawa, Daniel P. Mossing, Tong 685
Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin 686
Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Nee- 687
lakantan, Richard Ngo, Hyeonwoo Noh, Ouyang 688
Long, Cullen O'Keefe, Jakub W. Pachocki, Alex 689
Paino, Joe Palermo, Ashley Pantuliano, Giambattista 690
Parascandolo, Joel Parish, Emy Parparita, Alexandre 691
Passos, Mikhail Pavlov, Andrew Peng, Adam Perel- 692

693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723	<p>man, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Pondé de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Shepard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin D. Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas A. Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cer'on Uribe, Andrea Valalone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. Gpt-4 technical report.</p> <p>Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangoeei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. Flamingo: a visual language model for few-shot learning. <i>ArXiv</i>, abs/2204.14198.</p> <p>James F. Allen. 1984. Towards a general theory of action and time. <i>Artificial Intelligence</i>, 23(2):123–154.</p> <p>Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. <i>arXiv preprint arXiv:2309.16609</i>.</p> <p>Jasmijn E. Bosch, Mathilde Chailleux, and Francesca Foppolo. 2021. Incremental processing of telicity in italian children.</p> <p>Franklin Chang, Tomoko Tatsumi, Yuna Hiranuma, and Colin Bannard. 2023. Visual heuristics for verb production: Testing a deep-learning model with experiments in japanese. <i>Cognitive science</i>, 47 8:e13324.</p> <p>Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i>, pages 24185–24198.</p> <p>Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. 2024. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-lmms. <i>arXiv preprint arXiv:2406.07476</i>.</p> <p>Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.</p> <p>Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. <i>ArXiv</i>, abs/2210.11416.</p> <p>Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2020. The epic-kitchens dataset: Collection, challenges and baselines. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i>, 43:4125–4141.</p> <p>Francesca Foppolo, Jasmijn E. Bosch, Ciro Greco, Maria Nella Carminati, and Francesca Panzeri. 2021. Draw a star and make it perfect: Incremental processing of telicity. <i>Cognitive science</i>, 45 10:e13052.</p> <p>Francesca Foppolo, Francesca Panzeri, Cesare Greco, and Matteo Carminati. 2016. The incremental processing of accomplishment predicates.</p> <p>Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. <i>2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i>, pages 961–970.</p>	753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809
---	--	---

810	Y. Jang, Yale Song, Chris Dongjoo Kim, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2019. <i>Video question answering with spatio-temporal reasoning</i> . <i>International Journal of Computer Vision</i> , 127:1385 – 1412.	865
811		866
812		867
813		868
814		869
815	Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. 2019. <i>Action genome: Actions as compositions of spatio-temporal scene graphs</i> . <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 10233–10244.	870
816		
817		
818		
819		
820	Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Apostol Natsev, Mustafa Suleyman, and Andrew Zisserman. 2017. <i>The kinetics human action video dataset</i> . <i>ArXiv</i> , abs/1705.06950.	871
821		872
822		873
823		874
824		875
825		
826	J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. <i>Biometrics</i> , 33(1).	876
827		877
828		878
829	Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. 2018. <i>Tvqa: Localized, compositional video question answering</i> . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	879
830		880
831		881
832		882
833	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023. <i>Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models</i> . In <i>International Conference on Machine Learning</i> .	883
834		884
835		885
836		
837		
838	Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. 2022. <i>Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation</i> . In <i>International Conference on Machine Learning</i> .	886
839		887
840		888
841		889
842		890
843	Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. <i>Video-llava: Learning united visual representation by alignment before projection</i> . <i>ArXiv</i> , abs/2311.10122.	891
844		
845		
846		
847	Yi Liu, Limin Wang, Xiao Ma, Yali Wang, and Y. Qiao. 2021. <i>Fineaction: A fine-grained video dataset for temporal action localization</i> . <i>IEEE Transactions on Image Processing</i> , 31:6937–6950.	892
848		893
849		894
850		895
851	Olga Loginova, Oleksandr Bezrukov, and Alexey Kravets. 2024. <i>Addressing blind guessing: Calibration of selection bias in multiple-choice question answering by video language models</i> . <i>Preprint</i> , arXiv:2410.14248.	896
852		897
853		898
854		
855		
856	Agnese Lombardi and Alessandro Lenci. 2023. <i>Agen-tività e telicità in gilberto: implicazioni cognitive</i> . <i>Preprint</i> , arXiv:2307.02910.	902
857		903
858		904
859	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024. <i>Deepseek-vl: Towards real-world vision-language understanding</i> . <i>Preprint</i> , arXiv:2403.05525.	905
860		906
861		
862		
863		
864		
510	Eleni Metheniti, Tim Van De Cruys, and Nabil Hathout. 2022. <i>About time: Do transformers learn temporal verbal aspect?</i> In <i>Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics</i> , pages 88–101, Dublin, Ireland. Association for Computational Linguistics.	907
511		908
512		909
513		910
514		911
515		912
516		913
517		914
518		915
519		916
520		917
521		918
522		919
523		920
524		921

922		
923	Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun,	985
924	Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan,	986
925	Jeremiah Liu, Andras Orban, Fabian Güra, Hao	987
926	Zhou, Xinying Song, Aurelien Boffy, Harish Gana-	988
927	pathy, Steven Zheng, HyunJeong Choe, Ágoston	989
928	Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jar-	990
929	rod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah,	991
930	Emanuel Taropa, Majd Al Merey, Martin Baeuml,	992
931	Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Ol-	993
932	can Sercinoglu, George Tucker, Enrique Piqueras,	994
933	Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo	995
934	Danihelka, Becca Roelofs, Anaïs White, Anders	996
935	Andreassen, Tamara von Glehn, Lakshman Yagati,	997
936	Mehran Kazemi, Lucas Gonzalez, Misha Khalman,	998
937	Jakub Sygnowski, Alexandre Frechette, Charlotte	999
938	Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen,	1000
939	James Lottes, Nathan Schucher, Federico Lebron,	1001
940	Alban Rrustemi, Natalie Clay, Phil Crone, Tomas	1002
941	Kociský, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi	1003
942	Howard, Adam Bloniarz, Jack W. Rae, Han Lu,	1004
943	Laurent Sifre, Marcello Maggioni, Fred Alcober,	1005
944	Dan Garrette, Megan Barnes, Shantanu Thakoor, Ja-	1006
945	cob Austin, Gabriel Barth-Maron, William Wong,	1007
946	Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha,	1008
947	Arun Ahuja, Gaurav Singh Tomar, Evan Senter,	1009
948	Martin Chadwick, Ilya Kornakov, Nithya Attaluri,	1010
949	Ifñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan,	1011
950	Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang,	1012
951	Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia,	1013
952	Thanumalayan Sankaranarayana Pillai, Jacob Devlin,	1014
953	Michael Laskin, Diego de Las Casas, Dasha Valter,	1015
954	Connie Tao, Lorenzo Blanco, Adrià Puigdomènec	1016
955	Badia, David Reitter, Mianna Chen, Jenny Brennan,	1017
956	Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela	1018
957	Surita, Jane Labanowski, Abhi Rao, Stephanie Win-	1019
958	kler, Emilio Parisotto, Yiming Gu, Kate Olszewska,	1020
959	Ravi Addanki, Antoine Miech, Annie Louis, De-	1021
960	nis Teplyashin, Geoff Brown, Elliot Catt, Jan Bal-	1022
961	aguera, Jackie Xiang, Pidong Wang, Zoe Ashwood,	1023
962	Anton Briukhov, Albert Webson, Sanjay Ganapa-	1024
963	thy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang,	1025
964	Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur	1026
965	Bapna, Matthew Aitchison, Pedram Pejman, Henryk	1027
966	Michalewski, Tianhe Yu, Cindy Wang, Juliette Love,	1028
967	Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter	1029
968	Humphreys, Thibault Sellam, James Bradbury, Varun	1030
969	Godbole, Sina Samangooei, Bogdan Damoc, Alex	1031
970	Kaskasoli, Sébastien M. R. Arnold, Vijay Vasudevan,	1032
971	Shubham Agrawal, Jason Riesa, Dmitry Lep-	1033
972	ikhin, Richard Tanburn, Srivatsan Srinivasan, Hyeon-	1034
973	taek Lim, Sarah Hodkinson, Pranav Shyam, Johan	1035
974	Ferret, Steven Hand, Ankush Garg, Tom Le Paine,	1036
975	Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Za-	1037
976	heer Abbas, Sarah York, Machel Reid, Elizabeth	1038
977	Cole, Aakanksha Chowdhery, Dipanjan Das, Do-	1039
978	minika Rogozińska, Vitaliy Nikolaev, Pablo Sprech-	1040
979	mann, Zachary Nado, Lukas Zilka, Flavien Prost,	1041
980	Luheng He, Marianne Monteiro, Gaurav Mishra,	1042
981	Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Alla-	1043
982	manis, Clara Huiyi Hu, Raoul de Liedekerke, Justin	1044
983	Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou,	1045
984	Disha Shrivastava, Anirudh Baddepudi, Alex Goldin,	1046
	Adnan Ozturel, Albin Cassirer, Yunhan Xu, Daniel	1047
	Sohn, Devendra Sachan, Reinald Kim Amplayo,	1048
	Craig Swanson, Dessie Petrova, Shashi Narayan,	
	Arthur Guez, Siddhartha Brahma, Jessica Landon,	
	Miteyan Patel, Ruizhe Zhao, Kevin Villela, Luyu	
	Wang, Wenhao Jia, Matthew Rahtz, Mai Giménez,	
	Legg Yeung, James Keeling, Petko Georgiev, Diana	
	Mincu, Boxi Wu, Salem Haykal, Rachel Saputro,	
	Kiran Vodrahalli, James Qin, Zeynep Cankara,	
	Abhanshu Sharma, Nick Fernando, Will Hawkins,	
	Behnam Neyshabur, Solomon Kim, Adrian Hutter,	
	Priyanka Agrawal, Alex Castro-Ros, George van den	
	Driessche, Tao Wang, Fan Yang, Shuo yiin Chang,	
	Paul Komarek, Ross McIlroy, Mario Lučić, Guodong	
	Zhang, Wael Farhan, Michael Sharman, Paul Natsev,	
	Paul Michel, Yamini Bansal, Siyuan Qiao, Kris Cao,	
	Siamak Shakeri, Christina Butterfield, Justin Chung,	
	Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houldsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhwatty, Aditya Siddhant, Nenad Tomasev, Jin-wei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyana Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsilhas, Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas Fitzgerald, Keren Gu-Lemberg,	

1049	Mina Khan, Lisa Anne Hendricks, Marie Pellat,	1113
1050	Vladimir Feinberg, James Cobon-Kerr, Tara Sainath,	1114
1051	Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives,	1115
1052	Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd,	1116
1053	Le Hou, Qingze Wang, Thibault Sottiaux, Michela	1117
1054	Paganini, Jean-Baptiste Lespiau, Alexandre Mou-	1118
1055	farek, Samer Hassan, Kaushik Shivakumar, Joost van	1119
1056	Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh	1120
1057	Goyal, Matthew Tung, Andrew Brock, Hannah She-	1121
1058	ahan, Vedant Misra, Cheng Li, Nemanja Rakićević,	1122
1059	Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk	1123
1060	Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew	1124
1061	Lamm, Nicola De Cao, Charlie Chen, Sidharth Mud-	1125
1062	gal, Romina Stella, Kevin Brooks, Gautam Vasude-	1126
1063	van, Chenxi Liu, Mainak Chain, Nivedita Melink-	1127
1064	eri, Aaron Cohen, Venus Wang, Kristie Seymore,	1128
1065	Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krish-	1129
1066	nakumaran, Brian Albert, Nate Hurley, Motoki Sano,	1130
1067	Anhad Mohananey, Jonah Joughin, Egor Filonov,	1131
1068	Tomasz Kępa, Yomna Eldawy, Jiawern Lim, Rahul	1132
1069	Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang,	1133
1070	Sanil Jain, Sri Gayatri Sundara Padmanabhan, Subha	1134
1071	Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert	1135
1072	Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong	1136
1073	Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun	1137
1074	Wu, Sam Sobell, Andrea Siciliano, Alan Papir,	1138
1075	Robby Neale, Jonas Bragagnolo, Tej Toor, Tina	1139
1076	Chen, Valentin Anklin, Feiran Wang, Richie Feng,	1140
1077	Milad Gholami, Kevin Ling, Lijuan Liu, Jules Walter,	1141
1078	Hamid Moghaddam, Arun Kishore, Jakub Adamek,	1142
1079	Tyler Mercado, Jonathan Mallinson, Siddhinita Wan-	1143
1080	dekar, Stephen Cagle, Eran Ofek, Guillermo Gar-	1144
1081	rido, Clemens Lombriser, Maksim Mukha, Botu Sun,	1145
1082	Hafeezul Rahman Mohammad, Josip Matak, Yadi	1146
1083	Qian, Vikas Peswani, Paweł Janus, Quan Yuan, Leif	1147
1084	Schelin, Oana David, Ankur Garg, Yifan He, Olek-	1148
1085	sii Duzhyi, Anton Ålgmyr, Timothée Lottaz, Qi Li,	1149
1086	Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shiv-	1150
1087	anna, Aleksandr Chuklin, Josie Li, Carrie Spadine,	1151
1088	Travis Wolfe, Kareem Mohamed, Subhabrata Das,	1152
1089	Zihang Dai, Kyle He, Daniel von Dincklage, Shyam	1153
1090	Upadhyay, Akanksha Maurya, Luyan Chi, Sebas-	1154
1091	tian Krause, Khalid Salama, Pam G Rabinovitch,	1155
1092	Pavan Kumar Reddy M, Aarush Selvan, Mikhail	1156
1093	Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu	1157
1094	Gupta, Boyi Liu, Deepak Sharma, Idan Heimlich	1158
1095	Shtacher, Shachi Paul, Oscar Akerlund, François-	1159
1096	Xavier Aubet, Terry Huang, Chen Zhu, Eric Zhu,	1160
1097	Elico Teixeira, Matthew Fritze, Francesco Bertolini,	1161
1098	Liana-Eleonora Marinescu, Martin Bölle, Dominik	1162
1099	Paulus, Khyatti Gupta, Tejas Latkar, Max Chang,	1163
1100	Jason Sanders, Roopa Wilson, Xuewei Wu, Yi-	1164
1101	Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid	1165
1102	Lall, Swaroop Mishra, Wanming Chen, Thang Lu-	1166
1103	ong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk,	1167
1104	Dominik Rabiej, Vipul Ranjan, Krzysztof Styrc,	1168
1105	Pengcheng Yin, Jon Simon, Malcolm Rose Har-	1169
1106	riott, Mudit Bansal, Alexei Robsky, Geoff Bacon,	1170
1107	David Greene, Daniil Mirylenka, Chen Zhou, Obaid	1171
1108	Sarvana, Abhimanyu Goyal, Samuel Andermatt,	1172
1109	Patrick Siegler, Ben Horn, Assaf Israel, Francesco	1173
1110	Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici,	1174
1111	Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin	1175
1112	Guu, Roey Yogev, Xiaochen Cai, Alessandro Agos-	1176

1177	Evan Palmer, Paul Suganthan, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marcus Wu, Ricardo Aguilera, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnapalli, Caroline Kaplan, Jiri Simska, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylewicz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhhal, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Liting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Urias, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebecca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Põder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivière, Alanna Walton, Clément Crepy, Alicia Parrish, Zongwei Zhou, Clement Farabet, Carey Radabaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fidjeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Ramamohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Evgenii Eltyshev, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripraraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, Xiang-Hai Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurumurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikulik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu, Xuanyi Dong, Amruta Muthal, Senaka Butthpitiya, Sarthak Jauhari, Nan Hua, Urvashi Khandelwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Shahar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferring, Chenel Elkind, Aviel Atias, Paulina Lee, Vít Listík, Mathias Carlen, Jan	1241 1242 1243 1244 1245 1246 1247 1248 1249 1250 1251 1252 1253 1254 1255 1256 1257 1258 1259 1260 1261 1262 1263 1264 1265 1266 1267 1268 1269 1270 1271 1272 1273 1274 1275 1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292 1293 1294 1295 1296 1297 1298 1299 1300 1301 1302 1303 1304
------	---	--

1305	van de Kerkhof, Marcin Pukus, Krunoslav Zaher, Paul Müller, Sasha Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirnschall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai, Anca Stefanou, Bo Feng, Keshav Dhandhania, Manish Katyal, Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhargava Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal Majmundar, Michael Alverson, Michael Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandru, Xiangkai Zeng, Ben Bariach, Laura Weidinger, Tu Vu, Alek Andreev, Antoine He, Kevin Hui, Sheleem Kashem, Amar Subramanya, Sissie Hsiao, Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohman, Yonghui Wu, Slav Petrov, Jeffrey Dean, and Oriol Vinyals. 2024. Gemini: A family of highly capable multimodal models. <i>Preprint</i> , arXiv:2312.11805.	1336
1330	Carol Tenny. 1994. Aspectual roles and the syntax-semantics interface.	1331
1332	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poult, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>ArXiv</i> , abs/2307.09288.	1333
1335	Angeliac van Hout. 2008. Acquiring perfectivity and telicity in dutch, italian and polish. <i>Lingua</i> , 118:11740–1765.	1336
1337	Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Neural Information Processing Systems</i> .	1338
1339	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	1340
1341	Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	1342
1343	Bo Wu and Shoubin Yu. 2021. Star: A benchmark for situated reasoning in real-world videos. In <i>NeurIPS Datasets and Benchmarks</i> .	1344
1345	Junbin Xiao, Xindi Shang, Angela Yao, and Tat seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 9772–9781.	1346
1347	D. Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. <i>Proceedings of the 25th ACM international conference on Multimedia</i> .	1348
1349	Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. Minicpm-v: A gpt-4v level mllm on your phone. <i>arXiv preprint arXiv:2408.01800</i> .	1350
1351	Zhou Yu, D. Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. <i>ArXiv</i> , abs/1906.02467.	1352
1353	Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyaranayanan Ramamoorthy, and Graham Neubig. 2024. Pangea: A fully open multilingual multimodal lilm for 39 languages. <i>arXiv preprint arXiv:2410.16153</i> .	1354
1355	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. <i>Preprint</i> , arXiv:2303.15343.	1356
1357	Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024. Video instruction tuning with synthetic data. <i>Preprint</i> , arXiv:2410.02713.	1358
1359	Yiyun Zhao, Jian Gang Ngui, Lucy Hall Hartley, and Steven Bethard. 2021. Do pretrained transformers infer telicity like humans? In <i>Proceedings of the 25th Conference on Computational Natural Language Learning</i> , pages 72–81, Online. Association for Computational Linguistics.	1360
1361	Yaoyao Zhong, Wei Ji, Junbin Xiao, Yicong Li, Wei Deng, and Tat seng Chua. 2022. Video question answering: Datasets, algorithms and challenges. <i>ArXiv</i> , abs/2203.01225.	1362
1363	A Video and Action Classes Statistics	1364
1365	Tables 3 and 4 provide details on videos and action classes in Perfect Times, respectively.	1366
1367		1368
1368		1369
1369		1370
1370		1371
1371		1372
1372		1373
1373		1374
1374		1375
1375		1376
1376		1377
1377		1378
1378		1379
1379		1380
1380		1381
1381		1382
1382		1383
1383		1384
1384		1385
1385		1386
1386		1387
1387		1388
1388		1389
1389		1390
1390		1391
1391		1392
1392		1393
1393		1394
1394		1395
1395		1396
1396		1397
1397		1398
1398		1399
1399		1400
1400		1401
1401		1402
1402		1403
1403		1404
1404		1405
1405		1406
1406		1407
1407		1408
1408		1409
1409		1410
1410		1411
1411		1412
1412		1413
1413		1414
1414		1415

A Video and Action Classes Statistics

Tables 3 and 4 provide details on videos and action classes in Perfect Times, respectively.

Table 3: General Statistics of Perfect Times

Parameter	Value
Total videos	400
Total videos duration (sec)	11386.65
Total videos duration (min)	189.78
Average video duration (sec)	29.20
Minimum video duration (sec)	7.21
Maximum video duration (sec)	54.12

Table 4: Action Annotations

Parameter	Annotation
Total actions	3115
Minimum actions per video	2
Maximum actions per video	30
Average actions per video	7.99
Average time per action (sec)	8.35
Actions with ≤ 1 s duration	1149 (36.89%)
Actions with ≤ 2 s duration	1437 (46.13%)
Actions with ≤ 5 s duration	1850 (59.39%)
Actions with ≤ 10 s duration	2259 (72.52%)

B Dataset Statistics

The breakdown by answer options with the correct answers and distractor types of the total of 3739 qa pairs is given in Table 5

Table 5: QA pair statistics and distractor types.

Parameter	Value
Total QA pairs	3739
Correct a0	953
Correct a1	966
Correct a2	960
Correct a3	860
a0_distractor_type	953
a1_distractor_type	916
a2_distractor_type	927
a3_distractor_type	943

C Templates

This section presents templates in all languages, taking into account the parsed verb phrases from video annotations. In most cases, the classes include one verb, `verb_1`. In the tables below for several verbs in one class in coordinating constructions, `conj`, the coordinating conjunction and all elements with `_2` represent the second conjunct, as in *working or playing on the laptop*. The second

conjunct in the verb phrase is optional.

1432

D Dataset Examples

1433

The aggregated examples in all languages are presented in Figures 8 - 13.

1434

1435

E Annotator’s Statistics

1436

The statistics on each annotator’s performance is given in Table 10, where each language represents the annotator of this language.

1437

1438

1439

F Models

1440

Table 11 gives details on the open-source model’s components.

1441

1442

G Statistics on Conjunctions

1443

Figures 14-17 show preference of the correct answers with *while* over the ones with *when* across all the languages in numbers.

1444

1445

1446

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

Template	MCA	DCA	Condition	n_persons	Question
1a	t	a	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	1	What did the person in the video do after v_dea_ing_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
1b	t	a	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	2	What did the person in the video do after the other person person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
2a1	t	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	1	What did the person in the video do after they v_dea_past_perf_1 the obj_1 adjunct conj v_dea_past_perf_2 the obj_2?
2a2	t	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	1	What did the person in the video do when they v_dea_past_perf_1 the obj_1 adjunct conj v_dea_past_perf_2 the obj_2?
2b1	t	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	2	What did the person in the video do when the other person v_dea_past_perf_1 the obj_1 adjunct conj v_dea_past_perf_2 the obj_2?
2b2	t	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	2	What did the person in the video do when the other person v_dea_past_perf_1 the obj_1 adjunct conj v_dea_past_perf_2 the obj_2?
3a	t	a	et_mca <= st_dca	1	What had the person in the video done before v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
3b	t	a	et_mca <= st_dca	2	What had the person in the video done before the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
4a	t	t	et_mca <= st_dca	1	What had the person in the video done before they v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
4b	t	t	et_mca <= st_dca	2	What had the person in the video done before the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
5a1	t	a	et_mca > st_dca & et_mca < et_dca	1	What did the person in the video do v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
5a2	t	a	et_mca > st_dca & et_mca < et_dca	1	What did the person in the video do while v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
5b1	t	a	et_mca > st_dca & et_mca < et_dca	2	What did the person in the video do while the other person was v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
5b2	t	a	et_mca > st_dca & et_mca < et_dca	2	What did the person in the video do when the other person was v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
6a	a	a	et_mca <= st_dca	1	What had the person in the video been doing before v_dea_ing_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
6b	a	a	et_mca <= st_dca	2	What had the person in the video been doing before the other person was v_dea_ing_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
7a	a	t	et_mca <= st_dca	1	What had the person in the video been doing before v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
7b	a	t	et_mca <= st_dca	2	What had the person in the video been doing while the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
8a1	a	a	st_mca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	1	What was the person in the video doing while v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
8a2	a	a	st_mca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	1	What was the person in the video doing when v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
8b1	a	a	st_mca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	2	What was the person in the video doing while the other person was v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
8b2	a	a	st_mca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	2	What was the person in the video doing when the other person was v_dea_ing_1 the obj_1 adjunct conj v_dea_ing_2 the obj_2?
9a	a	t	st_mca < et_dca & et_mca > et_dca	1	What was the person in the video doing when the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
9b	a	t	st_mca < et_dca & et_mca > et_dca	2	What was the person in the video doing after v_dea_ing_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
10a	a	a	st_mca > et_dca	1	What was the person in the video doing after the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
10b	a	a	st_mca > et_dca	2	What was the person in the video doing after the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
11a	a	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	1	What was the person in the video doing after v_dea_ing_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
11b	a	t	st_mca >= et_dca; st_mca < et_dca & et_mca > et_dca	2	What was the person in the video doing after the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
12a1	t	t	et_mca == et_dca	1	What did the person in the video do when they v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
12a2	t	t	et_mca == et_dca	1	What did the person in the video do the moment they v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
12b1	t	t	et_mca == et_dca	2	What did the person in the video do when the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?
12b2	t	t	et_mca == et_dca	2	What did the person in the video do the moment the other person v_dea_past_1 the obj_1 adjunct conj v_dea_past_2 the obj_2?

Table 6: English templates: `_past` is the past simple form, `_perf` is the past perfect form, `_ing` is the -ing form including the past continuous form.

Template	MCA	DCA	Condition	n_persons	Question
la	t	a	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	1	Cosa ha fatto la persona nel video dopo v_inf_pass_pross_1 obj_1 adjunct conj v_inf_pass_pross_2 obj_2 ?
1b	t	a	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	2	Cosa ha fatto la persona nel video dopo che l'altra persona v_trapass_pross_1 obj_1 adjunct conj v_trapass_pross_2 obj_2 ?
2al	t	t	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	1	Cosa fece la persona nel video dopo v_inf_pass_pross_1 obj_1 adjunct conj v_inf_pass_pross_2 obj_2 ?
2a2	t	t	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	1	Cosa fece la persona nel video quando v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
2bl	t	t	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	2	Cosa fece la persona nel video dopo che l'altra persona v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
2b2	t	t	st_meca >= et_dca; st_meca < et_dca & et_mca > et_dca	2	Cosa fece la persona nel video quando l'altra persona v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
3a	t	a	et_meca <= st_dca	1	Cosa aveva fatto la persona nel video prima di v_inf_1 obj_1 adjunct conj v_inf_2 obj_2 ?
3b	t	a	et_meca <= st_dca	2	Cosa aveva fatto la persona nel video prima che l'altra persona v_imp_1 obj_1 adjunct conj v_imp_2 obj_2 ?
4a	t	t	et_meca <= st_dca	1	Cosa aveva fatto la persona nel video prima di v_inf_1 obj_1 adjunct conj v_inf_2 obj_2 ?
4b	t	t	et_meca <= st_dca	2	Cosa aveva fatto la persona nel video prima che l'altra persona v_trapass_cong_2 obj_2 ?
5al	t	a	et_meca > st_dca & et_mea < et_dca	1	Cosa ha fatto la persona nel video prima mentre v_imp_1 obj_1 adjunct conj v_imp_2 obj_2 ?
5a2	t	a	et_meca > st_dca & et_mea < et_dca	1	Cosa ha fatto la persona nel video mentre l'altra persona v_imp_1 obj_1 adjunct conj v_imp_2 obj_2 ?
5b1	t	a	et_meca > st_dca & et_mea < et_dca	2	Cosa ha fatto la persona nel video quando l'altra persona v_imp_1 obj_1 adjunct conj v_imp_2 obj_2 ?
5b2	t	a	et_meca > st_dca & et_mea < et_dca	2	Cosa stava facendo la persona nel video prima di v_inf_1 obj_1 adjunct conj v_inf_2 obj_2 ?
6a	a	a	et_meca <= st_dca	1	Cosa stava facendo la persona nel video prima che l'altra persona v_imp_cong_1 obj_1 adjunct conj v_imp_cong_2 obj_2 ?
6b	a	a	et_meca <= st_dca	2	Cosa stava facendo la persona nel video prima che l'altra persona v_trapass_pross_1 obj_1 adjunct conj v_trapass_pross_2 obj_2 ?
7a	a	t	et_meca <= st_dca	1	Cosa stava facendo la persona nel video prima che l'altra persona v_imp_cong_1 obj_1 adjunct conj v_inf_2 obj_2 ?
7b	a	t	et_meca <= st_dca	2	Cosa faceva la persona nel video mentre v_imp_prog_1 obj_1 adjunct conj v_imp_prog_2 obj_2 ?
8al	a	a	st_meca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	1	Cosa faceva la persona nel video quando v_imp_prog_1 obj_1 adjunct conj v_imp_prog_2 obj_2 ?
8a2	a	a	st_meca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	1	Cosa faceva la persona nel video mentre v_imp_prog_1 obj_1 adjunct conj v_imp_prog_2 obj_2 ?
8b1	a	a	st_meca > st_dca & st_mca < et_dca; st_mca == st_dca; st_mca < st_dca & et_mca > st_dca	2	Cosa faceva la persona nel video quando l'altra persona v_trapass_pross_1 obj_1 adjunct conj v_trapass_pross_2 obj_2 ?
8b2	a	a	st_meca > st_dca & st_mca < st_dca & et_mca > st_dca	2	Cosa faceva la persona nel video quando v_trapass_pross_1 obj_1 adjunct conj v_trapass_pross_2 obj_2 ?
9a	a	t	st_meca < et_dca & et_mea > et_dca	1	Cosa faceva la persona nel video quando v_imp_1 obj_1 adjunct conj v_imp_2 obj_2 ?
9b	a	t	st_meca < et_dca & et_mca > et_dca	2	Cosa faceva la persona nel video dopo v_inf_pass_pross_1 obj_1 adjunct conj v_inf_pass_pross_2 obj_2 ?
10a	a	a	st_meca >= et_dca	1	Cosa faceva la persona nel video dopo v_trapass_pross_1 obj_1 adjunct conj v_trapass_pross_2 obj_2 ?
10b	a	a	st_meca >= et_dca	2	Cosa faceva la persona nel video quando v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
12al	t	t	et_meca == et_dca	1	Cosa fece la persona nel video nel momento in cui v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
12a2	t	t	et_meca == et_dca	1	Cosa fece la persona nel video quando l'altra persona v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
12b1	t	t	et_meca == et_dca	2	Cosa fece la persona nel video nel momento in cui l'altra persona v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?
12b2	t	t	et_meca == et_dca	2	Cosa fece la persona nel video quando l'altra persona v_pass_rem_1 obj_1 adjunct conj v_pass_rem_2 obj_2 ?

Table 7: Italian templates: `_inf_` is the infinitive form, `_pass_pross` is the passato prossimo (past perfect) form, `_pass_rem_` is the passato remoto (simple past) form, `_imp_` is the imperfetto (imperfect) form, `_gerund` (gerund) form, `_imp_cong_` is the imperfetto congiuntivo (imperfect subjunctive) form, `_imp_prog_` is the imperfetto progressivo (past continuous) form. Conj is a coordinating conjunction in case of several verbs in one class, as in *lavorando o giocando al computer*.

Template	MCA	DCA	Condition	n_persons	Question
1a	t	a	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	1	Что сделал человек на видео после того, как v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
1b	t	a	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	2	Что сделал человек на видео после того, как другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
2a1	t	t	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	1	Что сделал человек на видео после того, как v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
2a2	t	t	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	1	Что сделал человек на видео после того, когда v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
2b1	t	t	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	2	Что сделал человек на видео после того, как другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
2b2	t	t	st_mca >= et_dea; st_mca < et_dea & et_mca > et_dea	2	Что сделал человек на видео, когда другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
3a	t	a	et_mca <= st_dca	1	Что сделал человек на видео до того, как v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
3b	t	a	et_mca <= st_dca	2	Что сделал человек на видео перед тем, как другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
4a	t	t	et_mca <= st_dca	1	Что сделал человек на видео до того, как v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
4b	t	t	et_mca <= st_dca	2	Что сделал человек перед тем, как другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
5a1	t	a	et_mca > st_dca & et_mca < et_dca	1	Что сделал человек, пока он v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
5a2	t	a	et_mca > st_dca & et_mca < et_dca	1	Что сделал человек, когда v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
5b1	t	a	et_mca > st_dca & et_mca < et_dca	2	Что сделал человек, пока другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
5b2	t	a	et_mca > st_dca & et_mca < et_dca	2	Что сделал человек, когда другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
6a	a	a	et_mca <= st_dca	1	Что делал человек до того, как v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
6b	a	a	et_mca <= st_dca	2	Что делал человек до того, как другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
7a	a	t	et_mca <= st_dca	1	Что делал человек до того, как он v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
7b	a	t	et_mca <= st_dca	2	Что делал человек, как другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
8a1	a	a	st_mca > st_dea & st_mca < et_dea; st_mca == st_dca; st_mea > st_dca	1	Что делал человек, пока v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
8a2	a	a	st_mca > st_dea & st_mca < et_dea; st_mea == st_dca; st_mea > st_dca	1	Что делал человек, когда v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
8b1	a	a	st_mca > st_dea & st_mca < et_dea; st_mea == st_dca; st_mea > st_dca	2	Что делал человек, пока другой человек v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
8b2	a	a	st_mca > st_dea & st_mca < et_dea; st_mea == st_dca; st_mea > st_dca	2	Что делал человек, когда другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
9a	a	t	st_mca < et_dea & et_mca > et_dea	1	Что делал человек, когда v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
9b	a	t	st_mca < et_dea & et_mca > et_dea	2	Что делал человек после того, как он v_imp_1 obj_1 adjunct conj v_imp_2 obj_2?
10a	a	a	st_mca > et_dea	1	Что делал человек после того, как другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
10b	a	a	st_mca > et_dea	2	Что делал человек после того, как он v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
11a	a	t	st_mca > et_dea; st_mea < et_dea & et_mea > et_dca	1	Что делал человек после того, как он v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
11b	a	t	st_mca > et_dea; st_mea < et_dea & et_mea > et_dca	2	Что делал человек после того, как другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
12a1	t	t	et_mca == et_dca	1	Что сделал человек на видео, когда v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
12a2	t	t	et_mca == et_dca	1	Что сделал человек на видео в тот момент, когда v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
12b1	t	t	et_mca == et_dca	2	Что сделал человек на видео, когда другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?
12b2	t	t	et_mca == et_dca	2	Что сделал человек на видео в тот момент, когда другой человек v_perf_1 obj_1 adjunct conj v_perf_2 obj_2?

Table 8: Rissain Templates: `_imp_` is imperfective past, `_perf_` — perfective past. Conj is a coordinating conjunction in case of several verbs in one class, as in читал или писал.

Template	MCA	DCA	Condition	n_persons	Question
la	t	a	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	1	ビデオに写っている人はadj objをta_form後、何をしましたか？
lb	t	a	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	2	ビデオに写っている人は、他の\かadj objをta_form後、何をしましたか？
2al	t	t	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	1	ビデオに写っている人はadj objをta_form後、何をしましたか？
2a2	t	t	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	1	ビデオに写っている人は、他の\かadj objをta_form時、何をしましたか？
2b1	t	t	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	2	ビデオに写っている人は、他の\かadj objをta_form後、何をしましたか？
2b2	t	t	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	2	ビデオに写っている人はadj objをinf前に、何をしましたか？
3a	t	a	et_mca <= st_dca	1	ビデオに写っている人は、他の\かadj objをinf前に、何をしましたか？
3b	t	a	et_mca <= st_dca	2	ビデオに写っている人は、他の\かadj objをinf前に、何をしましたか？
4a	t	t	et_mca <= st_dca	1	ビデオに写っている人はadj objをinf前に、何をしましたか？
4b	t	t	et_mca <= st_dca	2	ビデオに写っている人は、他の\かadj objをinf前に、何をしましたか？
5al	t	a	et_mca > st_dea & et_mca < et_dca	1	ビデオに写っている人はadj objをte_form、何をしましたか？
5a2	t	a	et_mca > st_dea & et_mca < et_dca	1	ビデオに写っている人はadj objをte_formながら、何をしましたか？
5b1	t	a	et_mca > st_dea & et_mca < et_dca	2	ビデオに写っている人は、他の\かadj objをinf_ inp間に、何をしましたか？
5b2	t	a	et_mca > st_dea & et_mca < et_dca	2	ビデオに写っている人は、他の\かadj objをta_form時、何をしましたか？
6a	a	a	et_mca <= st_dca	1	ビデオに写っている人はadj objをinf前に、何をいましたか？
6b	a	a	et_mca <= st_dca	2	ビデオに写っている人はadj objをinf前に、何をいましたか？
7a	a	t	et_mca <= st_dca	1	ビデオに写っている人はadj objをinf前に、何をしましたか？
7b	a	t	et_mca <= st_dca	2	ビデオに写っている人は、他の\かadj objをinf前に、何をしましたか？
8al	a	a	st_mca > st_dea & st_mca < et_dea; st_mca = st_dea; st_mca < st_dca & et_mca > st_dca	1	ビデオに写っている人はadj objをti_formながら、何をしていましたか？
8a2	a	a	st_mca > st_dea & st_mca < et_dea; st_mca = st_dea; st_mca < st_dca & et_mca > st_dca	1	ビデオに写っている人はadj objをti_form時、何をしていましたか？
8b1	a	a	st_mca > st_dea & st_mca < et_dea; st_mca = st_dea; st_mca < st_dca & et_mca > st_dca	2	ビデオに写っている人は、他の\かadj objをinf_ inp間に、何をしていましたか？
8b2	a	a	st_mca > st_dea & st_mca < et_dea; st_mca = st_dea; st_mca < st_dca & et_mca > st_dca	2	ビデオに写っている人は、他の\かadj objをta_form時、何をしていましたか？
9a	a	t	st_mca < et_dea & et_mca > et_dca	1	ビデオに写っている人はadj objをta_form時、何をしましたか？
9b	a	t	st_mca < et_dea & et_mca > et_dca	2	ビデオに写っている人はadj objをimp後、何をしましたか？
10a	a	a	st_mca >= et_dca	1	ビデオに写っている人はadj objをimp後、何をしましたか？
10b	a	a	st_mca >= et_dca	2	ビデオに写っている人は、他の\かadj objをimp後、何をしましたか？
11a	a	a	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	1	ビデオに写っている人はadj objをta_form後、何をしましたか？
11b	a	t	st_mca >= et_dca; st_mca < et_dea & et_mca > et_dca	2	ビデオに写っている人は、他の\かadj objをta_form時、何をしましたか？
12al	t	t	et_mca = et_dca	1	ビデオに写っている人はadj objをta_form時、何をしましたか？
12a2	t	t	et_mca = et_dca	1	ビデオに写っている人はadj objをform_ inp間、何をしましたか？
12b1	t	t	et_mca = et_dca	2	ビデオに写っている人は、他の\かadj objをta_form瞬間、何をしましたか？
12b2	t	t	et_mca = et_dca	2	ビデオに写っている人は、他の\かadj objをta_form 瞬間、何をしましたか？

Table 9: Japanese templates: ta_form is the general past form, inf is the dictionary form, te_form is the conjunction form, i_form is the present progressive, and imp is the imperfect. Conj is a coordinating conjunction in case of several verbs in one class, as in ノートパソコンで作業したり遊んだりする.



drinking from a cup or glass or bottle

close the refrigerator

Q: What did the person in the video do after drinking from a cup or glass or bottle?

- a0: The person was sitting on the floor.
a1: The person awakened in bed.

- a2: The person closed the refrigerator.**
a3: The person was closing the refrigerator.

Q: Cosa ha fatto la persona nel video dopo aver bevuto da una tazza o un bicchiere o una bottiglia?

- a0: La persona era seduta sul pavimento
a1: La persona si è svegliata nel letto

- a2: La persona ha chiuso il frigorifero**
a3: La persona chiudeva il frigorifero

Q: Что сделал человек на видео после того, как он пил из чашки или стакана или бутылки?

- a0: Он сидел на полу.
a1: Он проснулся в кровати.

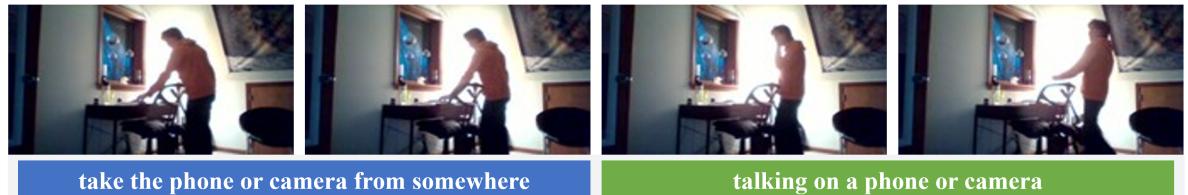
- a2: Он закрыл холодильник.**
a3: Он закрывал холодильник.

Q: ビデオに写っている人はカップかグラスか、またはボトルから飲んでいた後、何をしましたか？

- a0: その人は床に座っていた。
a1: その人はベッドで目を覚ました。

- a2: その人は冷蔵庫を閉めた。**
a3: その人は冷蔵庫を閉めていた。

Figure 8: Example from Perfect Times generated by Template 1a for all the languages.



take the phone or camera from somewhere

talking on a phone or camera

Q: What was the person in the video doing after taking the phone or camera from somewhere?

- a0: The person was washing something with a towel.
a1: The person was talking on a phone or camera.

- a2: The person talked on a phone or camera.
a3: The person put the phone or camera somewhere.

Q: Cosa faceva la persona nel video dopo aver preso il telefono oppure la fotocamera da qualche parte?

- a0: La persona lavava qualcosa con un asciugamano
a1: La persona parlava al telefono oppure alla fotocamera

- a2: La persona ha parlato al telefono oppure alla fotocamera
a3: La persona ha messo il telefono oppure la fotocamera da qualche parte

Q: Что делал человек на видео после того, как он взял телефон или камеру откуда-то?

- a0: Он протирал что-то полотенцем.
a1: Он разговаривал по телефону или камере.

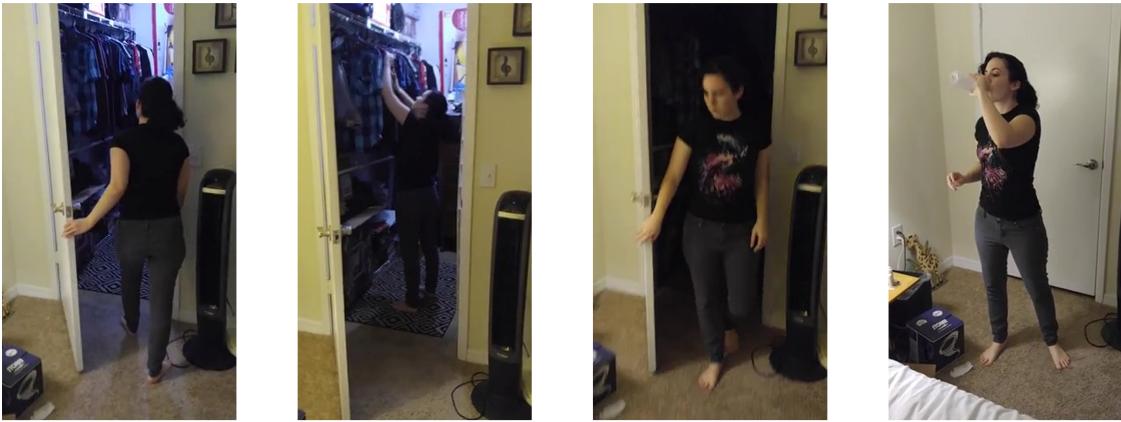
- a2: Он поговорил по телефону или камере.
a3: Он положил телефон или камеру куда-то.

Q: ビデオに写っている人はどこから携帯電話かカメラを取った後、何をしていましたか？

- a0: その人は何かをタオルで洗っていた。
a1: その人は携帯電話かカメラで話していた。

- a2: その人は携帯電話かカメラで話した。
a3: その人はどこかに携帯電話かカメラを置いた。

Figure 9: Example from Perfect Times generated by Template 11a for all the languages.



tiding up the closet

drinking from a cup or glass or bottle

Q: What had the person in the video been doing before drinking from a cup or glass or bottle?

a0: The person had been holding the laptop.

a2: The person had been tidying up the closet or cabinet.

a1: The person had tidied up the closet or cabinet.

a3: The person had closed the closet or cabinet.

Q: Cosa stava facendo la persona nel video prima di bere da una tazza o un bicchiere o una bottiglia?

a0: La persona stava tenendo il laptop

a2: La persona stava riordinando l'armadio oppure i mobili

a1: La persona aveva riordinato l'armadio oppure i mobili

a3: La persona aveva chiuso l'armadio oppure i mobili

Q: Что делал человек на видео до того, как пил из чашки или стакана или бутылки?

a0: Он держал ноутбук.

a2: Он прибирал в шкафу или шкафчике.

a1: Он убрал в шкафу или шкафчике.

a3: Он закрыл шкаф или шкафчик.

Q: ビデオに写っている人はカップかグラスか、またはボトルから飲んでいる前に、何をしていましたか？

a0: その人はノートパソコンを持っていた

a2: その人はクローゼットかキャビネットを整理していた。

a1: その人はクローゼットかキャビネットを整理した。a3: その人はクローゼットかキャビネットを閉めた。

Figure 10: Example from Perfect Times generated by Template 6a for all the languages.



dress

holding the bag

Q: What was the person in the video doing when the other person was dressing?

a0: The person was holding the bag.

a2: The person was holding the food.

a1: The person held the bag.

a3: The person took the clothes from somewhere.

Q: Cosa faceva la persona nel video quando l'altra persona si stava vestendo?

a0: La persona teneva la borsa

a2: La persona teneva il cibo

a1: La persona ha tenuto la borsa

a3: La persona ha preso i vestiti da qualche parte

Q: Что делал человек на видео, когда другой человек одевался?

a0: Он держал сумку.

a2: Он держал еду.

a1: Он подержал сумку.

a3: Он взял одежду откуда-то.

Q: ビデオに写っている人は、他の人が服を着た時、何をしていましたか？

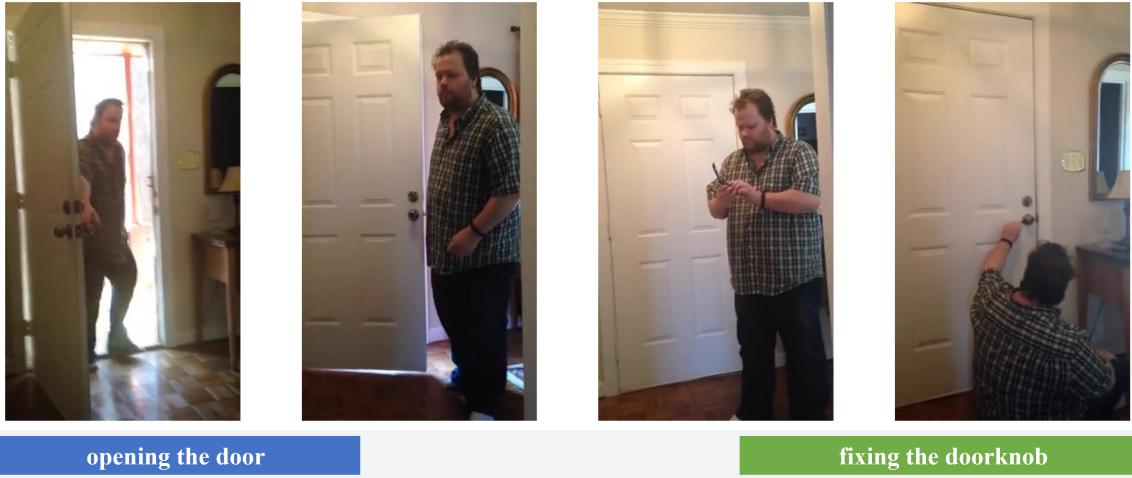
a0: その人はバッグを持っていた。

a2: その人は食べ物を持っていた。

a1: その人はバッグを持った。

a3: その人はどこから服を取った。

Figure 11: Example from Perfect Times generated by Template 8b2 for all the languages.



opening the door

fixing the doorknob

Q: What did the person in the video do after they had opened the door?

a0: The person fixed the doorknob.

a1: The person tidied the towel or towels.

a2: The person was opening the door.

a3: The person was fixing the doorknob.

Q: Cosa fece la persona nel video dopo aver aperto la porta?

a0: La persona riparò la maniglia della porta

a1: La persona riordinò l'asciugamano oppure gli asciugamani

a2: La persona stava aprendo la porta

a3: La persona stava riparando la maniglia della porta

Q: Что сделал человек на видео после того, как открыл дверь?

a0: Он починил дверную ручку.

a1: Он убрал полотенце или полотенца.

a2: Он открывал дверь.

a3: Он чинил дверную ручку.

Q: ビデオに写っている人はくしゃみをしながら、何をしていましたか？

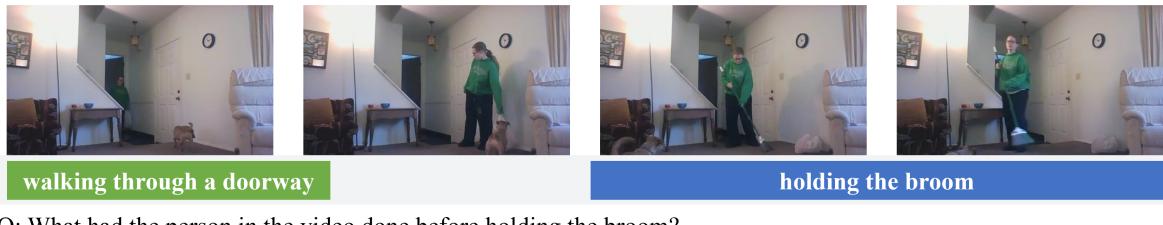
a0: その人は服を持った。

a1: その人は紙かノートで作業した。

a2: その人は紙かノートで作業していた。

a3: その人はどこかにほうきを置いていた。

Figure 12: Example from Perfect Times generated by Template 2a1 for all the languages.



walking through a doorway

holding the broom

Q: What had the person in the video done before holding the broom?

a0: The person had taken off the shoes.

a2: The person had been walking through a doorway.

a1: The person had been putting the groceries somewhere. **a3: The person had walked through a doorway.**

Q: Cosa aveva fatto la persona nel video prima di tenere la scopa?

a0: La persona si era tolta le scarpe

a2: La persona stava passando attraverso un'apertura

a1: La persona stava mettendo la spesa da qualche parte

a3: La persona era passata attraverso un'apertura

Q: Что сделал человек на видео до того, как он держал швабру?

a0: Он снял обувь.

a2: Он шел через дверной проем.

a1: Он клал продукты куда-то.

a3: Он прошел через дверной проем.

Q: ビデオに写っている人はほうきを持っている前に、何をしましたか？

a0: その人は靴を脱いだ。

a2: その人はドアを通っていた。

a1: その人はどこかに食料品を置いていた。

a3: その人はドアを通った。

Figure 13: Example from Perfect Times generated by Template 3a for all the languages.

Table 10: Annotator accuracy and distractor type distribution.

Annotator	Accuracy	Distractor Type 1	Distractor Type 2	Distractor Type 3
English	93.40	45.62	50.88	3.50
Italian	97.84	20.83	12.50	66.67
Russian	83.83	34.57	46.91	18.52
Japanese	98.36	20.00	40.00	40.00

Table 11: Model configurations and references.

Model	Vision Encoder	Language Model	Parameters
LLaVA-NeXT-Video-7B	SigLIP-400M (Zhai et al., 2023)	Qwen 1.5 (Bai et al., 2023)	7B
MiniCPM-V-2_6	SigLIP-400M (Zhai et al., 2023)	Qwen2 (Bai et al., 2023)	7B
Qwen2-VL	Qwen2-VL (Wang et al., 2024)	Qwen2 (Bai et al., 2023)	7B
InternVL2	InternViT-300M-448px (Chen et al., 2024)	internlm2_5-7b-chat (Chen et al., 2024)	8B

```

"en": {
  "chatgpt": {
    "most_true": "8a1",
    "most_true_count": 243,
    "least_true": "12b1",
    "least_true_count": 0,
    "most_false": "8a2",
    "most_false_count": 266,
    "least_false": "2b1",
    "least_false_count": 0
  },
  "llava_next": {
    "most_true": "8a1",
    "most_true_count": 191,
    "least_true": "3b",
    "least_true_count": 0,
    "most_false": "8a2",
    "most_false_count": 336,
    "least_false": "2b1",
    "least_false_count": 1
  }
},
"it": {
  "internvl8B": {
    "most_true": "8a1",
    "most_true_count": 179,
    "least_true": "6b",
    "least_true_count": 0,
    "most_false": "8a2",
    "most_false_count": 347,
    "least_false": "12b1",
    "least_false_count": 0
  },
  "minicpm": {
    "most_true": "8a1",
    "most_true_count": 194,
    "least_true": "4b",
    "least_true_count": 0,
    "most_false": "8a2",
    "most_false_count": 329,
    "least_false": "12b1",
    "least_false_count": 0
  }
}

```

Figure 14: Statistics on while (8a1) vs. when (8a2) for English in GPT-4o and LLaVA-NeXT-Video.

Figure 15: Statistics on *mentre* (8a1) vs. *quando* (8a2) for Italian in InternVL2 and MiniCPM-V.

```

"jp": {
    "gemini": {
        "most_true": "8a1",
        "most_true_count": 231,
        "least_true": "12b1",
        "least_true_count": 0,
        "most_false": "8a2",
        "most_false_count": 304,
        "least_false": "12b2",
        "least_false_count": 0
    },
}

```

Figure 16: Statistics on ながら (8a1) vs. 時 (8a2) for Japanese in Gemini-2.0-Flash-Lite.

```

"ru": {
    "minicpm": {
        "most_true": "8a1",
        "most_true_count": 203,
        "least_true": "10b",
        "least_true_count": 0,
        "most_false": "8a2",
        "most_false_count": 340,
        "least_false": "12b1",
        "least_false_count": 0
    },
    "gemini": {
        "most_true": "8a1",
        "most_true_count": 286,
        "least_true": "12b1",
        "least_true_count": 0,
        "most_false": "8a2",
        "most_false_count": 245,
        "least_false": "12b2",
        "least_false_count": 0
    },
}

```

Figure 17: Statistics on пока (8a1) vs. когда (8a2) for Russian in MiniCPM-V and Gemini-2.0-Flash-Lite.