CAVIAR: CRITIC-AUGMENTED VIDEO AGENTIC REASONING SUPPLEMENTAL MATERIAL

Anonymous authors

000

001

002003004

005

011 012

013

014 015

016

017

018

019

020 021

022

023

024

026 027

029

031

032 033

034

035

036 037 038

039 040

041

042

043

045

046 047

048

049

050

051

052

053

Paper under double-blind review

1 Module API

The full module API can be seen in 1. The modules included vary by the strategy being sampled.

```
@dataclasses.dataclass
class VideoSegment:
  """Class containing a segment of a video, defined by start and end
     times as MM:SS strings."""
  start: str
  end: str
def think(thought: str) -> None:
  """Tool to perform intermediate reasoning that doesn't directly look
      at the video.
   thought (str): the thought to print.
 Returns:
   None
def get_segment(start: str, end: str) -> VideoSegment:
    'Clips the segment of the input video between the times indicated
     by 'start' and 'end'.
  Args:
    start (str): the start time for the segment as an MM:SS string.
    end (str): the end time for the segment as an MM:SS string.
   A VideoSegment made from the input video from 'start' to 'end'.
def find_when(query: str, video_segment: VideoSegment | None) -> str:
  """Tool to determine timing of a query such as an event or action.
 Tries to find time ranges in the video that may correspond to the
     query. Works best for simple queries.
  The reasoning may be faulty, so you must carefully consider each
     time range and determine if it is relevant to best answering the
      question.
    query (str): the event or action to localize.
    video_segment (VideoSegment): if specified, only looks at the
       given segment of the input video.
 Returns:
   A string of time ranges identified as potentially relevant and the
        justification for each time range.
```

```
054
      def asr_understanding(question: str, answer_options: list[str] | None)
055
056
        """Tool to understand the ASR transcript of the video. Considers the
             ASR transcript of the full video and tries to answer the
057
            question.
058
        If the question cannot be answered from the transcript, returns any
            potentially relevant information or timestamps based on the
060
            transcript.
061
        Args:
062
          question (str): the question to answer.
          answer_options (list[str]): Optional. If the set of possible
063
             answers is known, they should be specified here as a list.
064
        Returns:
065
          A string containing the predicted answer to the question.
066
067
      def retrieval_qa(question: str, answer_options: list[str] | None,
068
          video_segment: VideoSegment | None) -> str:
069
        """Tool to answer a question from frames selected independently from
070
             each minute of the video.
071
        Looks at each minute of the video and tries to identify relevant
            frames.
        Then tries to answer the question using the retrieved frames.
073
074
075
          question (str): the question to answer.
076
          answer_options (list[str]): Optional. If the set of possible
077
              answers is known, they should be specified here as a list.
          video_segment (VideoSegment): if specified, only looks at the
078
              given segment of the input video.
079
080
        Returns:
081
          A string containing the predicted answer to the question.
082
083
      def finish(final_answer: str) -> str:
084
        """Print final answer and exit.
085
        After explaining your reasoning, output the final answer in the
086
            format "Final Answer: (X)" where X is the correct digit choice.
            Never say "unknown" or "unsure", or "None", instead provide your
087
             most likely guess.
090
          final_answer (str): reasoning followed by final answer.
        Returns:
092
          The parsed final answer.
094
```

Listing 1: Module API.

2 Module Implementation Details

096 097 098

099 100

101

102

103

104 105

106

107

VideoSegment. This is a convenience class for identifying segments of the input video by their start and end time.

get_segment. This module selects the appropriate segment from the input video given start and end times, returning a VideoSegment corresponding to that range.

find_when. This corresponds to coarse temporal localization with a sliding window. Using a window of 100 frames, the model is instructed to identify any time ranges that could be relevant for the given query, ultimately providing all of these ranges along with justifications for each of them as a string.

asr_understanding. This module uses the ASR transcript of all the speech in the video with timestamps aligned for the words said at each second. If the query is a question that can be answered from the transcript, it tries to provide an answer. If not, it aims to provide any relevant information along with timestamps from the ASR given.

retrieval_qa. This module tries to use the base model to retrieve frames that appear visually relevant to a query, then tries to answer the query from them. It uses a sliding window of 64 frames, returning the indices of frames with high visual relevance from each window. Those frames are then retrieved and provided back to the base model, along with 56 uniform frames from the rest of the video for additional context, to produce a response to the query.

finish. This module indicates the agent has reached a final answer and ends the agentic inference procedure, returning the final answer.

You are a video understanding agent. Your task is to answer the question provided using the available tools. These tools are already imported and accessible to you. You will receive the question and answer choices. You should use the tools iteratively, deciding which tool to use next based on the results of the previous tool call. Each turn, you should write a call to one tool, enclosed in triple backticks. Continue this process until you can confidently answer the question, selecting one of the provided answer choices. Do not import any tools; they are already available. Your final response should be one of the provided answer choices.

Step by Step Instructions

- **Analyze the question:** Carefully read the provided 'question' to understand the specific information required. Identify key entities, actions, and relationships mentioned in the question.
- 2. **Select an initial tool:** Based on your understanding of the question from step 1, choose the most appropriate tool from the available set of video interaction functions. Consider which tool will provide the most relevant information to answer the question efficiently.
- 3. **Execute the selected tool:** Call the chosen tool, providing any necessary parameters derived from the question analysis in step 1.
- 4. **Analyze the tool's output:** Examine the results returned by the tool. This output might be textual descriptions, numerical data, or other relevant information about the video.
- 5. **Evaluate progress:** Based on the tool's output from step 4, assess whether you have enough information to confidently answer the question. If yes, proceed to step 6. If not, go back to step 2 and select a different tool or use the same tool with different parameters based on the insights gained from the previous tool call. Consider how the previous tool's output informs your choice of the next tool.
- 6. **Formulate the final answer:** Using the information gathered from all tool calls, select one of the provided answer choices that best answers the 'question'. Clearly state your chosen answer choice.

Tools:

Listing 2: The preamble used for the reasoning agent.

3 STRATEGY SAMPLING DETAILS

For each dataset, 3 strategies are selected, each corresponding to a different subset of modules. These correspond to the strategies that can feasibly produce an answer for the given task, as well as the information allowed for the dataset (not including the base 'think' and 'finish' tools). For LVBench (and EgoSchema), the strategies are get_segment and retrieval_qa; retrieval_qa directly; and get_segment, retrieval_qa, and find_when (that is, all modules in use). For Neptune, as audio information is permitted, the same strategies are used with the addition of the asr_understanding module for the multi-module strategies. For the RTL task, the strategies are get_segment and find_when; find_when only; and get_segment, retrieval_qa, and find_when (again, all modules in use).

4 Critic Prompts

The critic is given the following preamble.

```
You are an expert in assessing and critiquing reasoning about videos.

Given a question and possible strategies from video reasoning agents for solving it, assess which strategy seems the most likely to result in a correct answer based on the provided examples.

If multiple strategies produce the same final answer, you may pick multiple; you may not pick multiple answers with conflicting final answers.

Observe that the reliability of the tools depends on their input; you must consider the examples to determine when a given tool is likely to hallucinate or give an accurate output, as evidenced by the ultimate result.

Produce a critique per the examples followed by "Winning Strategies:" followed by a comma separated list of the letters corresponding to those strategies.
```

Listing 3: The preamble used for the critic.

The critic is given 4 in-context examples per prompt. These examples are written by the authors with the aim of being in-distribution for the datasets shown.

```
194
         Input:
         You will be given a question about a video and four possible answer
195
            options. Question: At 12:10, what color suit is the host wearing
196
            ?Possible answer choices:
197
         (1) Navy
198
         (2) Tan
199
         (3) Black
200
         (4) White
         Video length: 40:50
201
202
         Strategy A (retrieval_qa, get_segment):
203
         input_video = get_segment(start='00:00', end='40:50')
205
206
207
208
         segment = get_segment(start='12:10', end='12:20')
209
         answer = retrieval_qa(question='What color suit is the host wearing?
             , video_segment=segment)
210
211
        The host is wearing a navy suit.
212
213
214
         final_answer = finish(f"The retrieval_qa tool determined the host
            was wearing a navy suit at 12:10. Final Answer: (1)")
215
```

```
216
        The retrieval_qa tool determined the host was wearing a navy suit at
217
             12:10. Final Answer: (1)
218
219
         Strategy B (direct retrieval_qa):
        Here's how to solve this:
220
221
        1. **Locate the timestamp:** Go to the 12:10 mark in the video.
222
223
        2. **Observe the host's attire: ** At this point, the host (Jimmy
224
            Kimmel) is wearing a dark suit and a dark-colored tie.
225
        3. **Compare to options:** The suit appears to be black.
226
227
        Therefore, the correct answer is (3) Black.
228
        Final Answer: (3)
229
         Strategy C (retrieval_qa, get_segment, find_when):
230
231
         input_video = get_segment(start='00:00', end='40:50')
232
233
234
         . . .
235
         segment = get_segment(start='12:10', end='12:20')
236
237
         ['12:10', '12:20']
238
239
         answer = retrieval_qa(question='At 12:10, what color suit is the
240
            host wearing?', answer_options=['Navy', 'Tan', 'Black', 'White'
241
            ], video_segment=segment)
242
243
        The image at 12:10 shows the host wearing a navy suit.
244
        Final Answer: (1)
245
246
247
        final_answer = finish(final_answer='The image at 12:10 shows the
248
            host wearing a navy suit. Final Answer: (1)')
249
         The image at 12:10 shows the host wearing a navy suit. Final Answer:
250
             (1)
251
252
         Critique:
253
         A: 1. B: 3. C: 1. Strategies A and C succeed while Strategy B fails
            because Strategies A and C successfully use the get_segment tool
254
             to isolate the relevant part of the video.
255
        Strategy B fails as the retrieval_qa tool alone cannot actually
256
            locate a given timestamp in a video, resulting in a hallucinated
257
             answer.
        Winning Strategies:
259
        A, C
260
261
        Input:
262
         You will be given a question about a video and four possible answer
263
            options. Question: What does the man in charcoal suit do in the
            next scene after the scene of a photo of six people end?Possible
264
             answer choices:
265
         (1) He shows his scars
266
         (2) He takes off his shirt
267
         (3) He shows his jewelry
268
         (4) He shows his physique
        Video length: 39:59
269
```

```
270
        Strategy A (retrieval_qa, get_segment):
271
272
        input_video = get_segment(start='00:00', end='39:59')
273
274
275
276
        think("I need to identify the scene where a photo of six people is
277
            shown. Then I need to find the scene immediately following it
278
            to see what the man in the charcoal suit does.")
279
        I need to identify the scene where a photo of six people is shown.
280
            Then I need to find the scene immediately following it to see
281
            what the man in the charcoal suit does.
282
283
        segment_with_photo = get_segment(start='00:00', end='01:00')
284
        answer = retrieval_qa(question="What does the man in charcoal suit
285
            do in the scene immediately after the scene with a photo of six
286
            people?", video_segment=segment_with_photo)
287
288
        After the scene with the photo of six people, the man in the
            charcoal suit walks down a hallway and is greeted by a group of
289
            people who clap and cheer for him.
290
291
292
        think ("The provided text doesn't describe any of the options (1) He
            shows his scars, (2) He takes off his shirt, (3) He shows his
293
            jewelry, (4) He shows his physique. The man is celebrated. I
294
            will guess the closest option, which is none of them, but I will
             choose the least violent option.")
296
297
        The provided text doesn't describe any of the options (1) He shows
            his scars, (2) He takes off his shirt, (3) He shows his jewelry,
             (4) He shows his physique. The man is celebrated. I will
            guess the closest option, which is none of them, but I will
300
            choose the least violent option.
301
302
        final_answer = finish("Final Answer: (1)")
303
304
        Final Answer: (1)
305
306
        Strategy B (direct retrieval_qa):
307
        Here's how to solve this step-by-step:
308
        1. **Identify the relevant scene: ** The question refers to the scene
309
             after a photo of six people. Looking through the timestamps,
310
            this appears to be around 20:49.
311
312
        2. **Observe the next scene: ** The next scene shows the man in the
313
            charcoal suit (Harry Styles) pointing to and showing jewelry on
314
315
        3. **Match to the answer choices:** This action directly corresponds
316
             to answer choice (3).
317
        Therefore, the correct answer is **(3) He shows his jewelry**.
318
        Final Answer: (3)
319
320
        Strategy C (retrieval_qa, get_segment, find_when):
321
322
        input_video = get_segment(start='00:00', end='39:59')
323
```

```
324
325
        ...
326
        think("First, I need to find when the scene with the photo of six
            people ends. Then I can examine the next scene to see what the
327
            man in the charcoal suit does.")
328
329
        First, I need to find when the scene with the photo of six people
330
            ends. Then I can examine the next scene to see what the man in
331
            the charcoal suit does.
332
333
        photo_end_time = find_when(query="scene with a photo of six people",
334
             video_segment=None)
335
336
        ["00:05", "00:24"]: The query asks to find the scene with a photo of
             six people. From the provided frames, a photo of six people is
337
             visible on the wall in the background starting at 00:05 and
338
            remaining visible until 00:24. The photo is partially obscured
339
            at times by other people and objects, but it is consistently
340
            present in the background during this timeframe.
341
        ["10:39", "10:43"]: The query asks for a scene with a photo of six
342
            people. At 10:39 through 10:43 a large group photo of more than
             six people is shown.
343
        ["17:47", "17:53"]: The query asks for a scene with a photo of six
344
            people. At 17:47-17:53, a photo of approximately six people is
345
            visible behind a curtain as Harry Styles walks out onto the
346
                   The photo is partially obscured and the exact number of
            stage.
347
            people is difficult to ascertain with certainty, but it appears
            to show at least five.
348
        ["18:47", "19:59"]: The query asks for the time range of a scene
            with a photo of six people. A photo of what appears to be six
350
            people is visible on the wall behind the main set in frames
351
            starting at 18:47 and continuing through 19:59. The photo is
352
            partially obscured at times by other elements in the scene, but
            it is consistently visible enough to be identified.
353
        ["20:03", "20:04"]: The query asks for the time range of a scene
354
            with a photo of six people. At 20:03, a wide shot of the studio
355
             audience is shown, and a photo of six people is visible on a
356
            screen in the background. This photo remains visible until
            approximately 20:04.
357
         ["23:15", "23:25"]: The query asks for the timestamp of a scene with
358
             a photo of six people. In the provided video frames, a photo % \left( 1\right) =\left( 1\right) \left( 1\right) 
359
            of six people is visible on the wall in the background starting
360
            at 23:15 and remaining visible until 23:25. The photo is
361
            partially obscured at times by other elements in the frame, but
            it is consistently present during this time range.
362
        ["35:17", "35:20"]: The provided video frames show a wide shot of a
363
            studio audience at the 35:17 mark. The audience appears to
364
            consist of more than six people. While the exact number is
365
            difficult to determine from the provided frames, it is clearly a
366
             large group. Therefore, this time range satisfies the query.
        ["39:03", "39:18"]: The provided video frames show a live
367
            performance with an audience. At the end of the performance,
368
            from 39:03 to 39:18, the camera focuses on the audience, and
369
            there are enough people visible to potentially include a group
370
            of five. The resolution isn't high enough to definitively
371
            confirm the presence of a specific photo of six people, but the
            scene depicts a large group of people, making it a plausible
372
            match for the query.
373
374
375
        think ("The most likely candidate for the 'scene with a photo of six
376
            people' is the longest continuous segment: 18:47-19:59. Let's
377
            use that as the end point for the previous scene.")
```

```
378
        The most likely candidate for the 'scene with a photo of six people'
379
             is the longest continuous segment: 18:47-19:59. Let's use that
380
            as the end point for the previous scene.
381
382
        next_scene_segment = get_segment(start='19:59', end='20:05')
383
384
        ['19:59', '20:05']
385
386
        answer = retrieval_qa(question="What does the man in the charcoal
387
            suit do in the scene from 19:59 to 20:05?", answer_options=['He
388
            shows his scars', 'He takes off his shirt', 'He shows his
389
            jewelry', 'He shows his physique'], video_segment=
390
            next_scene_segment)
        ...
        The provided video frames show a talk show segment featuring three
392
            men. From 19:59 to 20:02, the man in the gray suit (Harry
393
            Styles) sits on a couch and gestures with his hands. At 20:03,
394
            the camera angle changes to show the audience, and the three men
395
             are still on the couch. There is no indication that the man in
396
             the gray suit shows his scars, takes off his shirt, shows a
            tattoo, or shows his muscles. He is simply participating in a
397
            conversation.
398
399
        Final Answer: (4)
400
401
        final_answer = finish("Based on the analysis of the video segment
402
            from 19:59 to 20:05, the man in the charcoal suit (assuming this
             refers to one of the men in the talk show segment, as there's
404
            no charcoal suit explicitly mentioned) does not perform any of
405
            the actions listed in the options. However, since an answer must
             be chosen, and the retrieval_qa tool did not provide a
406
            definitive answer, I will select the option that is most likely
407
            to be present in a talk show setting, which is showing his physique. This is a weak inference, but given the constraints,
408
409
            it's the best guess. Final Answer: (4)")
410
        Based on the analysis of the video segment from 19:59 to 20:05, the
411
            man in the charcoal suit (assuming this refers to one of the men
412
             in the talk show segment, as there's no charcoal suit
413
            explicitly mentioned) does not perform any of the actions listed
414
             in the options. However, since an answer must be chosen, and
415
            the retrieval_qa tool did not provide a definitive answer, I
            will select the option that is most likely to be present in a
416
            talk show setting, which is showing his physique. This is a
417
            weak inference, but given the constraints, it's the best guess.
418
            Final Answer: (4)
419
420
        Critique:
        A:1. B: 3. C: 4. Strategy B succeeds while Strategies A fails
421
            because Strategy A incorrectly applies the get_segment tool to
422
            the first minute of the video and only considers that, resulting
423
             in it not being able to find a correct answer and ultimately
            guessing. Strategy C applies the find_when tool, then trims a
425
            segment and does not find the relevant information, then
426
            guessing as well.
427
        Winning Strategies:
428
429
430
        Input:
```

```
432
        You will be given a question about a video and four possible answer
433
            options. Question: What is the main reason the vlogger goes out
434
            again, where she also buys makeup? Possible answer choices:
        (1) To go exercise
435
        (2) To eat chicken
436
        (3) To buy coffee mugs
437
        (4) To go to the park
438
        Video length: 30:25
439
        Strategy A (retrieval_qa, get_segment):
441
        input_video = get_segment(start='00:00', end='30:25')
443
        . . .
445
        think ("I need to find a segment of the video where the vlogger buys
446
            makeup. I'll try to identify this segment using retrieval_qa,
447
            and then I can use that information to answer the question about
448
             why she went out.")
        ""
449
450
        I need to find a segment of the video where the vlogger buys makeup.
              I'll try to identify this segment using retrieval_qa, and then
451
             I can use that information to answer the question about why she
452
             went out.
453
454
455
        answer = retrieval_qa(question="When does the vlogger buy makeup?",
            answer_options=None, video_segment=input_video)
456
        The vlogger buys makeup at 16:10-16:11 and 17:20-17:25.
458
459
        think("The vlogger buys makeup in two segments: 16:10-16:11 and
460
            17:20-17:25. I'll use retrieval_qa again, focusing on these
461
            segments, to determine the reason for her outing.")
462
463
        The vlogger buys makeup in two segments: 16:10-16:11 and
464
            17:20-17:25. I'll use retrieval_qa again, focusing on these
            segments, to determine the reason for her outing.
465
466
467
        segment1 = get_segment(start='16:10', end='16:11')
468
        segment2 = get_segment(start='17:20', end='17:25')
469
        reason = retrieval_qa(question="What is the main reason the vlogger
470
            goes out again, where she also buys makeup?", answer_options=['
            To go exercise', 'To eat chicken', 'To buy coffee mugs', 'To go
471
            to the park'], video_segment=segment1)
472
        if reason == "": #check if the first segment didn't provide enough
473
            information
474
            reason = retrieval_qa(question="What is the main reason the
475
                vlogger goes out again, where she also buys makeup?"
                answer_options=['To go exercise', 'To eat chicken', 'To buy
476
                coffee mugs', 'To go to the park'], video_segment=segment2)
477
        . . .
478
        None
479
480
        final_answer = finish("Based on the analysis of the video segments
481
            where the vlogger buys makeup, I was unable to definitively
482
            determine the main reason for her outing from the provided
483
            options. However, given the context of beauty product purchases,
484
             it's less likely she's going to the gym, eating salmon, or
            going to the park. Therefore, I will eliminate those options.
485
            Final Answer: (3)")
```

```
486
487
        Based on the analysis of the video segments where the vlogger buys
488
            makeup, I was unable to definitively determine the main reason
            for her outing from the provided options. However, given the
489
            context of beauty product purchases, it's less likely she's
490
            going to the gym, eating salmon, or going to the park. Therefore
491
            , I will eliminate those options. Final Answer: (3)
492
493
        Strategy B (direct retrieval_qa):
        Here's how to solve this step-by-step:
495
        1. **Identify the relevant section: ** The question asks about the
            vlogger's second outing, where she buys makeup. This occurs
497
            after the gym visit.
498
        2. **Review the frames:** The frames show the vlogger going to the
499
            park after the gym. She then goes shopping and buys makeup at
500
            Kicks.
501
502
        3. **Select the answer:** The main reason for her second outing is
503
            to go to the park.
504
        Therefore, the correct answer is (4) To go to the park.
505
        Final Answer: (4)
506
507
        Strategy C (retrieval_qa, get_segment, find_when):
508
509
        input_video = get_segment(start='00:00', end='30:25')
510
511
512
513
        think("I need to find when the vlogger goes out and buys makeup. I'
            ll use 'find_when' to locate that event, then use 'retrieval_qa'
514
             to determine the reason.")
515
516
        I need to find when the vlogger goes out and buys makeup. I'll use
517
            'find_when' to locate that event, then use 'retrieval_qa' to
518
            determine the reason.
519
520
        when_out_buying_makeup = find_when(query="vlogger goes out and buys
521
            makeup", video_segment=input_video)
522
523
        ["00:01", "01:59"]: The video shows a vlogger waking up in her room,
524
             doing some laundry, and then seemingly getting ready to go out.
              While there's no explicit depiction of buying makeup, the
525
            context strongly suggests this is part of her preparation. The
526
            time range encompasses her actions from waking up to preparing
527
            laundry, which could be interpreted as part of her routine
            before going out to buy makeup.
        ["02:19", "02:25"]: The video shows the vlogger in her bathroom at
529
            02:15, where she picks up a toothbrush and looks at her tablet.
530
             From 02:19-02:25, she is shown selecting items from a bathroom
531
            cabinet. This is consistent with the query.
532
        ["05:24", "05:31"]: The vlogger is seen leaving her apartment at
533
            04:02 and returning at 04:48. At 05:24-05:31, she is seen
            putting away what appears to be a package of makeup on a shelf
534
            in her bathroom. This is supported by the context of the query
535
            and the visual of a package being placed on a shelf.
536
        ["08:51", "09:16"]: The video shows a vlogger leaving their home at
537
            08:30, walking outside, and then entering a grocery store at
538
                    The vlogger is seen selecting various grocery items,
539
            including celery, grapes, blueberries, carrots, broccoli, and
            lettuce, between 08:51 and 09:16. Therefore, the time range
```

```
540
            reflects the vlogger's grocery shopping trip, which aligns with
541
            the query's description of buying products.
542
        ["14:00", "14:02"]: The query asks about the vlogger buying makeup.
             The provided video frames show the vlogger in their home,
543
            interacting with items that could be interpreted as makeup (
544
            frames 14:00-14:02) and later organizing clothes in their closet
545
             (frames 14:39-15:25). There is no visual evidence of the
546
            vlogger going out to buy anything. Therefore, only the initial
547
            interaction with what appears to be makeup is considered
            relevant.
        ["16:09", "16:18"], ["17:17", "17:25"]: The video shows the vlogger
549
            going to a store, Kicks, and buying makeup. The vlogger is seen
550
             in the store from 16:09 to 16:18. At 17:17-17:25, the vlogger
551
            unpacks the Kicks bag, revealing makeup.
552
        ["18:12", "18:17"]: The query asks for the time range when the
            vlogger goes out and buys makeup. The provided video frames
553
            show the vlogger in their room, at their vanity. At 18:12, there
554
             are shopping bags from the store Kicks on the vanity. This
555
            strongly suggests that the vlogger has recently returned from a
556
            shopping trip where they purchased makeup from Kicks. There is
557
            no footage of the actual shopping trip itself within this video
558
            chunk.
        ["20:25", "20:51"]: The video shows a person preparing coffee.
559
             20:25 to 20:51, the person is shown unpacking two ceramic brown
560
             mugs. The contents of the mugs are visible and appear to be
561
            coffee. Therefore, this time range is identified as the time
562
            when the vlogger is shown preparing coffee.
        ["29:20", "29:41"]: The video shows the vlogger in her bathroom at
563
            28:00-28:18 applying makeup and then brushing her teeth. At
564
            29:20-29:41, she is shown painting her nails. The presence of
            various makeup in both locations supports the query.
566
567
        reason = retrieval_qa(question="What is the main reason the vlogger
568
            goes out again, where she also buys makeup?", answer_options=['
569
            To go exercise', 'To eat chicken', 'To buy coffee mugs', 'To go
570
            to the park'], video_segment=get_segment(start='16:09', end='
571
            17:25'))
572
        The vlogger goes to a store and buys coffee mugs. Then she goes to
573
            another store and buys makeup. Therefore, the main reason she
574
            goes out again is to buy coffee mugs.
575
576
        Final Answer: (3)
577
        . . .
578
        final_answer = finish("The provided text from 'retrieval_qa'
579
            indicates the vlogger buys makeup at a store called Kicks.
580
            However, it also mentions buying coffee mugs. The question asks
581
            for the *main* reason for going out again. While makeup is
582
            mentioned, the provided text strongly suggests the primary
            reason for the second outing is to buy coffee mugs. Therefore,
583
            the answer is (3). Final Answer: (3)")
584
        ""
585
        The provided text from 'retrieval_qa' indicates the vlogger buys
586
            makeup at a store called Kicks. However, it also mentions
587
            buying coffee mugs. The question asks for the *main* reason for
            going out again. While makeup is mentioned, the provided text
588
            strongly suggests the primary reason for the second outing is to
589
            buy coffee mugs. Therefore, the answer is (3). Final Answer:
590
            (3)
591
592
        Critique:
        A: 3. B: 4. C: 3.
593
```

595

596

597

598

599

600 601

```
Strategy B succeeds while Strategies A fails because Strategy A
   looks at particular time ranges of the video and concludes it is
    unable to find a correct answer, thus guessing.
While the approach is generally sound, Strategy A gives up upon not
   finding an answer initially. Strategy C looks at a narrow range
   and follows an implication to guess an answer. Thus Strategy B
   is preferred.
Winning Strategy:
```

Listing 4: LVBench critic examples.

```
604
605
      Input:
606
      You will be given a question about a video and five possible answer
607
          options. Question: What are the main points regarding proper form
           and breathing during the barbell squat? Possible answer choices:
608
          Keep your back straight, core engaged, and chest up throughout
609
          the movement. Inhale as you descend and exhale forcefully as you
610
          drive up.
611
      (2) Maintain a slight arch in your back, push your hips back first,
612
          and keep your knees aligned with your toes. Breathe in during the
          eccentric phase (lowering) and breathe out during the concentric
613
          phase (rising).
614
      (3)
          Focus on pushing through your heels, keep your head down, and
615
          hold your breath at the bottom of the squat for stability. Exhale
616
          as you return to the starting position.
617
          Prioritize depth over weight, allowing your hips to descend below
           your knees. Inhale deeply at the top, hold your breath throughout
618
           the squat, and exhale at the top.
619
          Use a wide stance, keep your toes pointed outwards, and focus on
620
          squeezing your glutes at the top of the movement. Breathe
621
          naturally throughout the exercise.
622
      Video length: 01:35
623
624
      Strategy A (get_segment, retrieval_qa, asr_understanding):
625
626
      input_video = get_segment(start='00:00', end='01:35')
627
628
629
630
      think ("I need to understand the main points about form and breathing
631
          during barbell squats. I'll use the ASR understanding tool to
          analyze the video transcript. If that doesn't work, I'll try '
632
          retrieval_qa' on the video frames.")
633
634
      I need to understand the main points about form and breathing during
635
          barbell squats. I'll use the ASR understanding tool to analyze
636
          the video transcript. If that doesn't work, I'll try 'retrieval_qa
637
          ' on the video frames.
638
      "
639
      asr_understanding(question='What are the main points regarding proper
640
          form and breathing during the barbell squat?', answer_options=['
641
          (1) Keep your back straight, core engaged, and chest up
642
          throughout the movement. Inhale as you descend and exhale
          forcefully as you drive up.', '(2) Maintain a slight arch in your
643
          back, push your hips back first, and keep your knees aligned with
644
           your toes. Breathe in during the eccentric phase (lowering) and
645
          breathe out during the concentric phase (rising).', '(3) Focus on
646
          pushing through your heels, keep your head down, and hold your
          breath at the bottom of the squat for stability. Exhale as you
647
          return to the starting position.', '(4) Prioritize depth over
```

```
648
          weight, allowing your hips to descend below your knees. Inhale
649
          deeply at the top, hold your breath throughout the squat, and
650
          exhale at the top.', '(5) Use a wide stance, keep your toes
          pointed outwards, and focus on squeezing your glutes at the top of
651
           the movement. Breathe naturally throughout the exercise.'])
652
653
      The video emphasizes maintaining a slight arch in the back
654
          (00:23-00:27), pushing the hips back first (00:28-00:31), and
655
          ensuring knees track with the toes (00:32-00:35). For breathing,
656
          it advises inhaling during the descent and exhaling during the
          ascent (00:55-01:00).
657
658
659
      think ("The ASR transcript supports answer choice (2). The video
660
          highlights maintaining a slight arch in your back, initiating the
          movement with the hips, aligning knees with toes, and coordinating
661
           breathing with the phases of the squat.")
662
663
      The ASR transcript supports answer choice (2). The video highlights
664
          maintaining a slight arch in your back, initiating the movement
665
          with the hips, aligning knees with toes, and coordinating
666
          breathing with the phases of the squat.
667
668
      final_answer = finish("Final Answer: (2)")
669
670
      Final Answer: (2)
671
      Strategy B (direct retrieval_qa):
672
      The video primarily shows a demonstration of the barbell squat and
673
          highlights key points on form. It doesn't explicitly detail the
674
          breathing technique. However, general weightlifting principles
675
          suggest inhaling during the lowering phase and exhaling during the
676
           lifting phase. Additionally, proper form includes keeping the
          back straight, core engaged, and chest up.
677
678
679
      Final Answer: (1)
680
      Strategy C (get_segment, retrieval_qa find_when, asr_understanding):
681
      input_video = get_segment(start='00:00', end='01:35')
682
683
684
685
      think ("I need to find the main points regarding form and breathing in
686
          the barbell squat. I'll use 'asr_understanding' to see if the
687
          transcript mentions these details.")
688
689
      I need to find the main points regarding form and breathing in the
690
          barbell squat. I'll use 'asr_understanding' to see if the
691
          transcript mentions these details.
692
      ...
693
      asr_understanding_result = asr_understanding(question="What are the
694
          main points regarding proper form and breathing during the barbell
695
           squat?", answer_options=[
          "Keep your back straight, core engaged, and chest up throughout
696
              the movement. Inhale as you descend and exhale forcefully as
697
              you drive up.",
698
          "Maintain a slight arch in your back, push your hips back first,
699
              and keep your knees aligned with your toes. Breathe in during
700
              the eccentric phase (lowering) and breathe out during the
701
              concentric phase (rising).",
```

```
702
          "Focus on pushing through your heels, keep your head down, and
703
              hold your breath at the bottom of the squat for stability.
704
              Exhale as you return to the starting position.",
705
          "Prioritize depth over weight, allowing your hips to descend below
               your knees. Inhale deeply at the top, hold your breath
706
              throughout the squat, and exhale at the top.",
707
          "Use a wide stance, keep your toes pointed outwards, and focus on
708
              squeezing your glutes at the top of the movement. Breathe
709
              naturally throughout the exercise."
710
      ])
      . . .
711
      The video focuses on maintaining a slight arch in your back
712
          (0:23-0:27), pushing your hips back first (0:28-0:31), and keeping
713
          your knees aligned with your toes (0:32-0:35). For breathing, it
714
          advises to inhale during the descent (eccentric phase) and exhale
715
          during the ascent (concentric phase) (0:55-1:00).
716
      . . .
717
      think ("The transcript suggests answer (2), highlighting back arch, hip
718
           movement, knee alignment, and breathing coordinated with the
719
          squat phases. Let's use 'retrieval_qa' to cross-check this.")
720
      The transcript suggests answer (2), highlighting back arch, hip
721
          movement, knee alignment, and breathing coordinated with the squat
722
           phases. Let's use 'retrieval_qa' to cross-check this.
723
724
725
      retrieval_qa_result = retrieval_qa(question="What are the main points
          regarding proper form and breathing during the barbell squat?",
726
          answer_options=[
727
          "Keep your back straight, core engaged, and chest up throughout
728
              the movement. Inhale as you descend and exhale forcefully as
729
              you drive up.",
          "Maintain a slight arch in your back, push your hips back first,
730
              and keep your knees aligned with your toes. Breathe in during
731
              the eccentric phase (lowering) and breathe out during the
732
              concentric phase (rising).",
733
          "Focus on pushing through your heels, keep your head down, and
734
              hold your breath at the bottom of the squat for stability.
              Exhale as you return to the starting position.",
735
          "Prioritize depth over weight, allowing your hips to descend below
736
              your knees. Inhale deeply at the top, hold your breath
737
              throughout the squat, and exhale at the top.",
738
          "Use a wide stance, keep your toes pointed outwards, and focus on
739
              squeezing your glutes at the top of the movement. Breathe
              naturally throughout the exercise."
740
      ], video_segment=input_video)
741
742
      The video demonstrates proper form for barbell squats, emphasizing
743
          keeping the back straight, core engaged, and chest up. While
744
          breathing is not explicitly addressed, general knowledge suggests
          inhaling on the descent and exhaling forcefully on the ascent.
745
746
      Final Answer: (1)
747
748
      ""
749
      final_answer = finish("The ASR transcript suggests (2), focusing on
          back arch and hip movement. However, 'retrieval_qa' prioritizes a
750
          straight back and doesn't fully confirm the breathing technique.
751
          Considering general weightlifting knowledge and the potential for
752
          ASR misinterpretation, (1) is the more reliable answer. Final
753
          Answer: (1)")
754
      The ASR transcript suggests (2), focusing on back arch and hip
755
          movement. However, 'retrieval_qa' prioritizes a straight back and
```

```
756
          doesn't fully confirm the breathing technique. Considering general
757
           weightlifting knowledge and the potential for ASR
758
          misinterpretation, (1) is the more reliable answer. Final Answer:
759
          (1)
760
      Critique:
761
      A: 2. B: 1. C: 1. Strategies B and C consider the visual content, but
762
          deny the information can be found in the video and rely on general
763
           knowledge.
764
      Winning Strategies:
765
766
767
      Input:
768
      You will be given a question about a video and five possible answer
769
          options. Question: To prepare this dish, what are the main steps
          involved?Possible answer choices:
770
      (1) The main steps include chopping vegetables, boiling them in water,
771
           adding spices and herbs, pouring in coconut milk, simmering the
772
          mixture, and finally garnishing with cilantro.
773
      (2) The main steps include chopping vegetables, sauteing them in a
774
          skillet, adding spices and dry fruits, pouring in beef broth,
          baking the mixture in the oven, and finally garnishing with green
775
          and red bell peppers.
776
      (3) The main steps include chopping vegetables, sauteing them in a wok
777
          , adding milk and herbs, pouring in vegetable broth, simmering the
778
           mixture, and finally garnishing with parsley.
      (4) The main steps include chopping vegetables, sauteing them in a pan
779
          , adding lentils, pouring in water, adding spices and cilantro,
780
          simmering the lentils, and finally garnishing with chopped
781
          tomatoes.
782
      (5) The main steps include chopping fruits, sauteing them in a pan,
783
          adding spices and herbs, pouring in chicken broth, simmering the
          mixture, and finally garnishing with cilantro.
784
      Video length: 09:43
785
786
      Strategy A (get_segment, retrieval_qa, asr_understanding):
787
788
      input_video = get_segment(start='00:00', end='09:43')
789
790
791
792
      answer = asr_understanding(question='To prepare this dish, what are
793
          the main steps involved?', answer_options=['The main steps include
           chopping vegetables, boiling them in water, adding spices and
794
          herbs, pouring in coconut milk, simmering the mixture, and finally
795
           garnishing with cilantro.', 'The main steps include chopping
796
          vegetables, sauteing them in a skillet, adding spices and dry
797
          fruits, pouring in beef broth, baking the mixture in the oven, and
798
           finally garnishing with green and red bell peppers.', 'The main
          steps include chopping vegetables, sauteing them in a wok, adding
799
          milk and herbs, pouring in vegetable broth, simmering the mixture,
800
           and finally garnishing with parsley.', 'The main steps include
801
          chopping vegetables, sauteing them in a pan, adding lentils,
802
          pouring in water, adding spices and cilantro, simmering the
803
          lentils, and finally garnishing with chopped tomatoes.', 'The main
           steps include chopping fruits, sauteing them in a pan, adding
804
          spices and herbs, pouring in chicken broth, simmering the mixture,
805
           and finally garnishing with cilantro.'])
806
807
      The main steps are chopping onions and carrots (00:32-01:15), sauteing
808
           them in olive oil with garlic and ginger (01:23-02:58), adding
          lentils and water (03:12-04:05), adding turmeric and cumin
809
          (04:30-05:20), simmering until the lentils are tender
```

```
810
          (05:35-07:10), and finally garnishing with chopped cilantro and
811
          tomatoes (07:22-08:01).
812
      . . .
813
      think ("None of the provided options exactly match the steps in the ASR
814
           transcript. Option (1) is closest since it mentions chopping
815
          vegetables, adding spices, and simmering. However, it omits key
816
          ingredients and steps like sauteing.")
817
818
      None of the provided options exactly match the steps in the ASR
          transcript. Option (1) is closest since it mentions chopping
819
          vegetables, adding spices, and simmering. However, it omits key
820
          ingredients and steps like sauteing.
821
822
      final_answer = finish("Final Answer: (1)")
823
824
      Final Answer: (1)
825
826
      Strategy B (direct retrieval_qa):
827
      The video demonstrates a recipe for lentil soup. The chef chops
          vegetables like onions, carrots, and garlic. They then saute these
828
           vegetables in a pan with oil. Lentils and water are added,
829
          followed by spices like cumin and turmeric. The mixture simmers
830
          until the lentils are cooked through. Finally, it's garnished with
831
           chopped tomatoes and cilantro.
832
833
      The closest answer is (4).
834
      Final Answer: (4)
835
      Strategy C (get_segment, retrieval_qa find_when, asr_understanding):
836
837
      input_video = get_segment(start='00:00', end='09:43')
838
839
840
      ...
841
      answer = asr_understanding(question='To prepare this dish, what are
842
          the main steps involved?', answer_options=['The main steps include
           chopping vegetables, boiling them in water, adding spices and
843
          herbs, pouring in coconut milk, simmering the mixture, and finally
844
           garnishing with cilantro.', 'The main steps include chopping
845
          vegetables, sauteing them in a skillet, adding spices and dry
846
          fruits, pouring in beef broth, baking the mixture in the oven, and
847
           finally garnishing with green and red bell peppers.', 'The main
848
          steps include chopping vegetables, sauteing them in a wok, adding
          milk and herbs, pouring in vegetable broth, simmering the mixture,
849
           and finally garnishing with parsley.', 'The main steps include
850
          chopping vegetables, sauteing them in a pan, adding lentils,
851
          pouring in water, adding spices and cilantro, simmering the
852
          lentils, and finally garnishing with chopped tomatoes.', 'The main
853
           steps include chopping fruits, sauteing them in a pan, adding
          spices and herbs, pouring in chicken broth, simmering the mixture,
854
           and finally garnishing with cilantro.'])
855
856
      The main steps include chopping onions and carrots (0:35-1:10),
857
          sauteing them in olive oil (1:15-2:05), adding spices like cumin
          and turmeric (2:10-2:45), then adding lentils and water
858
          (2:50-3:30). The mixture simmers for some time (3:45-5:00) and is
859
          finally garnished with cilantro and chopped tomato (5:10-5:45).
860
861
862
      think ("None of the options perfectly align with the ASR transcript.
          Option (1) is somewhat similar as it mentions chopping vegetables,
863
```

```
864
           spices, and simmering. However, it misses crucial details like
865
          sauteing and the specific ingredients mentioned.")
866
867
      None of the options perfectly align with the ASR transcript. Option
          (1) is somewhat similar as it mentions chopping vegetables, spices
868
          , and simmering. However, it misses crucial details like sauteing
869
          and the specific ingredients mentioned.
870
871
872
      final_answer = finish("Final Answer: (1)")
873
      Final Answer: (1)
874
875
      Critique:
876
      A: 1. B: 4. C: 1. A and C consider the ASR, but give up after finding
877
          it doesn't directly result in an answer.
878
      Winning Strategies:
879
880
881
      Input:
882
      You will be given a question about a video and five possible answer
          options. Question: In what order does the video discuss the
883
          following topics: introduction, stock taillight operation,
884
          Morimoto XB LED taillight replacement, taillight features, and
885
          installation and warranty? Possible answer choices:
886
      (1) Introduction, stock taillight operation, taillight features,
887
          Morimoto XB LED taillight replacement, installation and warranty.
      (2) Morimoto XB LED taillight replacement, introduction, taillight
888
          features, stock taillight operation, installation and warranty.
      (3) Introduction, Stock taillight operation, Morimoto XB LED taillight
890
           replacement, taillight features, Installation and warranty.
891
      (4) Introduction, Morimoto XB LED taillight replacement, installation
892
          and warranty, taillight features, stock taillight operation.
      (5) Stock taillight operation, introduction, Morimoto XB LED taillight
893
           replacement, taillight features, installation and warranty.
894
      Video length: 01:15
895
896
      Strategy A (get_segment, retrieval_qa, asr_understanding):
897
      input_video = get_segment(start='00:00', end='01:15')
898
899
900
901
      answer = asr_understanding(question='In what order does the video
902
          discuss the following topics: introduction, stock taillight
903
          operation, Morimoto XB LED taillight replacement, taillight
904
          features, and installation and warranty?', answer_options=['
905
          Introduction, stock taillight operation, taillight features,
906
          Morimoto XB LED taillight replacement, installation and warranty.'
          , 'Morimoto XB LED taillight replacement, introduction, taillight
907
          features, stock taillight operation, installation and warranty.',
908
          'Introduction, Stock taillight operation, Morimoto XB LED
909
          taillight replacement, taillight features, Installation and
910
          warranty.', 'Introduction, Morimoto XB LED taillight replacement,
911
          installation and warranty, taillight features, stock taillight
          operation.', 'Stock taillight operation, introduction, Morimoto XB
912
           LED taillight replacement, taillight features, installation and
913
          warranty.'])
914
915
      The video first provides an introduction (00:02-00:10). Next, it
916
          discusses the stock taillight operation (00:11-00:25).
          details the Morimoto XB LED taillight replacement (00:26-00:38).
917
          Following this, the video describes the taillight features
```

```
918
          (00:39-00:55). Finally, it covers installation and warranty
919
          (00:56-01:08). Therefore, the correct order is (3).
920
      ...
921
      final_answer = finish("Final Answer: (3)")
922
923
      Final Answer: (3)
924
925
      Strategy B (direct retrieval_qa):
926
      The video begins with an introduction of Morimoto and then shows the
          stock taillight operation. After that, it shows the Morimoto XB
927
          LED taillight replacement and then discusses its features. Finally
928
           it shows the installation process and discusses its warranty.
929
          Therefore, the correct order is introduction, stock taillight
930
          operation, Morimoto XB LED taillight replacement, taillight
          features, and installation and warranty.
931
932
933
      Final Answer: (1)
934
      Strategy C (get_segment, retrieval_qa find_when, asr_understanding):
935
936
      input_video = get_segment(start='00:00', end='01:15')
937
938
939
940
      answer = asr_understanding(question='In what order does the video
941
          discuss the following topics: introduction, stock taillight
          operation, Morimoto XB LED taillight replacement, taillight
942
          features, and installation and warranty?', answer_options=['
          Introduction, stock taillight operation, taillight features,
944
          Morimoto XB LED taillight replacement, installation and warranty.'
945
          , 'Morimoto XB LED taillight replacement, introduction, taillight
946
          features, stock taillight operation, installation and warranty.',
          'Introduction, Stock taillight operation, Morimoto XB LED
          taillight replacement, taillight features, Installation and
948
          warranty.', 'Introduction, Morimoto XB LED taillight replacement,
949
          installation and warranty, taillight features, stock taillight
950
          operation.', 'Stock taillight operation, introduction, Morimoto XB
           LED taillight replacement, taillight features, installation and
951
          warranty.'])
952
953
      The video first provides an introduction (00:02-00:10). Next, it
954
          discusses the stock taillight operation (00:11-00:25).
                                                                   Then, it
955
          details the Morimoto XB LED taillight replacement (00:26-00:38).
956
          Following this, the video describes the taillight features
          (00:39-00:55). Finally, it covers installation and warranty
957
          (00:56-01:08). Therefore, the correct order is (3).
958
959
960
      final_answer = finish("Final Answer: (3)")
961
      Final Answer: (3)
962
963
      Critique:
964
      A: 3. B: 1. C: 3. The question asks about what is discussed in the
965
          video, implying speech, but B is not able to consider the ASR and
          fails.
966
967
      Winning Strategies:
968
      A, C
969
970
      Input:
971
      You will be given a question about a video and five possible answer
          options. Question: Why does Jamie decide to stay home instead of
```

```
972
          going to the library at 11:00 a.m. on a Saturday? Possible answer
973
          choices:
974
      (1) Jamie gets a call from a friend and decides to meet them for lunch
975
           instead.
       (2) Jamie realizes the library is closed on Saturdays.
976
      (3) Jamie's car breaks down, preventing them from driving to the
977
          library.
978
       (4) Jamie starts a new book at home and is too engrossed to leave.
979
      (5) Jamie remembers they have a prior engagement, a family gathering,
980
          preventing them from going to the library.
      Video length: 01:59
981
982
      Strategy A (get_segment, retrieval_qa, asr_understanding):
983
984
      input_video = get_segment(start='00:00', end='01:59')
985
986
987
      ...
988
      answer = asr_understanding(question="Why does Jamie decide to stay
989
          home instead of going to the library at 11:00 a.m. on a Saturday
          ?", answer_options=['Jamie gets a call from a friend and decides
990
          to meet them for lunch instead.', 'Jamie realizes the library is
991
          closed on Saturdays.', 'Jamie\'s car breaks down, preventing them
992
          from driving to the library.', 'Jamie starts a new book at home
993
          and is too engrossed to leave.', 'Jamie remembers they have a
994
          prior engagement, a family gathering, preventing them from going
995
          to the library.'])
996
      Jamie decides to stay home instead of going to the library because
997
          they remember they have a prior family gathering
998
          [01:25:00-01:32:00].
999
      ""
1000
      final_answer = finish("Based on the ASR transcript, Jamie decides to
1001
          stay home instead of going to the library because they remember
1002
          they have a prior family gathering. Final Answer: (5)")
1003
1004
      Based on the ASR transcript, Jamie decides to stay home instead of
          going to the library because they remember they have a prior
1005
          family gathering. Final Answer: (5)
1006
1007
      Strategy B (direct retrieval_qa):
1008
      The video shows Jamie checking the time on their watch, which reads
1009
          11:00 AM. They then pack a bag with books and prepare to leave the
           house, even putting on their shoes. However, they receive a text
1010
          message and after reading it, unpack their bag and take off their
1011
          shoes. They then proceed to sit on the couch and relax. The video
1012
          doesn't explicitly say why the plans changed, but the text message
1013
           is likely the reason. The video ends with a short scene of a
1014
          group of people, suggesting a possible social gathering.
1015
1016
      Final Answer: (1)
1017
1018
      Strategy C (get_segment, retrieval_qa find_when, asr_understanding):
1019
      input_video = get_segment(start='00:00', end='01:59')
1020
1021
1022
1023
1024
      found_times = find_when(query="Jamie's plans are interrupted",
1025
         video_segment=input_video)
```

```
["00:35", "00:45"]: The query asks for the time range when Jamie's
1027
          plans are interrupted. Based on the video and transcript, Jamie's
           plans to go to the library are interrupted when they trip over
          their shoelaces after putting on their shoes. This occurs between
1029
          00:35 and 00:42. The video shows them putting on their shoes,
1030
          tripping, and then sitting down.
1031
1032
      . . .
      segment_of_interest = get_segment(start='00:35', end='00:45')
1033
1034
      ['00:35', '00:45']
1035
1036
1037
      answer = retrieval_qa(question="Why does Jamie decide to stay home
1038
          instead of going to the library at 11:00 a.m. on a Saturday?",
          answer_options=[', Jamie gets a call from a friend and decides to
1039
          meet them for lunch instead.', 'Jamie realizes the library is
1040
          closed on Saturdays.', 'Jamie\'s car breaks down, preventing them
1041
          from driving to the library.', 'Jamie starts a new book at home
1042
          and is too engrossed to leave.', 'Jamie remembers they have a
1043
          prior engagement, a family gathering, preventing them from going
          to the library.'], video_segment=segment_of_interest)
1044
1045
      The video shows Jamie putting on their shoes and then tripping over
1046
          their shoelaces. There is no indication they got a call, realized
1047
           the library was closed, had car trouble, or started a new book.
1048
          The video also does not show a family gathering. Therefore, none
1049
          of the options 1-5 are supported by the video.
1050
      However, if we must choose the most likely option, option (3) is the
1051
          most plausible. Jamie tripping could be interpreted as an
1052
          indication of car trouble, albeit a very indirect one.
1053
      Final Answer: (3)
1054
1055
1056
      final_answer = finish(final_answer="The video shows Jamie putting on
1057
          their shoes and then tripping over their shoelaces. There is no
1058
          indication they got a call, realized the library was closed, had
          car trouble, or started a new book. The video also does not show a
1059
           family gathering. Therefore, none of the options 1-5 are
          supported by the video. However, if we must choose the most likely
1061
           option, option (3) is the most plausible. Jamie tripping could be
1062
           interpreted as an indication of car trouble, albeit a very
1063
          indirect one. Final Answer: (3)")
1064
      The video shows Jamie putting on their shoes and then tripping over
1065
          their shoelaces. There is no indication they got a call, realized
1066
           the library was closed, had car trouble, or started a new book.
          The video also does not show a family gathering. Therefore, none
1068
          of the options 1-5 are supported by the video. However, if we must
           choose the most likely option, option (3) is the most plausible.
1069
          Jamie tripping could be interpreted as an indication of car
1070
          trouble, albeit a very indirect one. Final Answer: (3)
1071
1072
      Critique:
1073
      A: 5. B: 1. C: 3. B makes a guess after not finding an explicit answer
           in the visuals. C selects an overly narrow range to consider, as
1074
          the reason for the plans being interrupted may not be exactly when
1075
           they are interrupted.
1076
1077
      Winning Strategies:
1078
1079
```

Listing 5: Neptune critic examples.

```
1080
1081
      Input:
      Question: When is the trampoline first used in the video?
1082
      Video length: 04:23
1083
1084
      Strategy A (get_segment, find_when):
1085
1086
      input_video = get_segment(start='00:00', end='04:23')
1087
1088
1089
1090
      time_ranges = find_when(query='the trampoline gets its first use',
1091
          video_segment=input_video)
1092
      ["00:45", "01:15"]: The video shows a person assembling a trampoline.
1093
          The trampoline is shown fully assembled at 00:45 and the assembly
1094
          process is shown up to 01:15.
1095
      ["02:30", "02:55"]: The video shows a person jumping on the trampoline
1096
          . This is the first time the trampoline is shown being used, from
          02:30 to 02:55.
1097
1098
      . . .
1099
      think ("The first time the trampoline is used is arguably at 00:45,
1100
          when it is assembled. However, the question asks when it gets its
1101
           *first use*, which implies actual use by a person. Therefore, the
           relevant time range is 02:30-02:55.")
1102
1103
      The first time the trampoline is used is arguably at 00:45, when it is
1104
           assembled. However, the question asks when it gets its *first
1105
          use*, which implies actual use by a person. Therefore, the
1106
          relevant time range is 02:30-02:55.
1107
      . . .
1108
      final_answer = finish("The question asks when the trampoline gets its
1109
          first use, implying actual use by a person. While the trampoline
1110
          is assembled at 00:45-01:15, its first actual use begins at 02:30
1111
          when a person jumps on it. Therefore, the relevant time range is
          02:30-02:55. Final Answer: [02:30, 02:55]")
1112
1113
      The question asks when the trampoline gets its first use, implying
1114
          actual use by a person. While the trampoline is assembled at
1115
          00:45-01:15, its first actual use begins at 02:30 when a person
1116
          jumps on it. Therefore, the relevant time range is 02:30-02:55.
          Final Answer: [02:30, 02:55]
1117
1118
      Strategy B (get_segment, find_when, retrieval_qa):
1119
1120
      input_video = get_segment(start='00:00', end='04:23')
1121
1122
1123
      ""
1124
      first_use = find_when(query='first use of the trampoline',
1125
         video_segment=input_video)
1126
      ["00:45", "01:50"]: The video shows a person assembling a trampoline.
1127
          The first appearance of the trampoline occurs at 00:45, when it's
1128
          partially assembled. The trampoline is shown in several stages of
1129
           assembly until 01:50.
1130
      ["03:10", "03:40"]: The video shows the finished trampoline being
1131
          showcased in a backyard. The trampoline is fully assembled and
          ready for use at 03:10, and is shown from various angles until
1132
          03:40.
1133
```

```
1134
1135
      think("The first time range [00:45, 01:50] seems to describe the
1136
          trampoline being assembled, but the second range [03:10, 03:40]
          shows the trampoline in its final assembled state, ready for use.
1137
           Therefore, the answer is the second range.")
1138
      . . .
1139
      The first time range [00:45, 01:50] seems to describe the trampoline
1140
          being assembled, but the second range [03:10, 03:40] shows the
1141
          trampoline in its final assembled state, ready for use. Therefore
1142
          , the answer is the second range.
1143
1144
      final_answer = finish("The first time range [00:45, 01:50] seems to
1145
          describe the trampoline being assembled, but the second range
1146
          [03:10, 03:40] shows the trampoline in its final assembled state,
          ready for use. Therefore, the answer is the second range. Final
1147
          Answer: [03:10, 03:40]")
1148
1149
      The first time range [00:45, 01:50] seems to describe the trampoline
1150
          being assembled, but the second range [03:10, 03:40] shows the
1151
          trampoline in its final assembled state, ready for use. Therefore
1152
          , the answer is the second range. Final Answer: [03:10, 03:40]
1153
      Strategy C (find_when):
1154
1155
      input_video = get_segment(start='00:00', end='03:53')
1156
1157
1158
      query = "When is the trampoline first used in the video?"
1160
      result = find_when(query, None)
1161
      ["00:45", "01:15"]: The video shows a person assembling a trampoline.
1162
          The trampoline is shown fully assembled at 00:45 and the assembly
1163
          process is shown up to 01:15.
1164
      ["02:30", "02:55"]: The video shows a person jumping on the trampoline
1165
          . This is the first time the trampoline is shown being used, from
1166
          02:30 to 02:55.
1167
1168
      think ("The first result mentions the trampoline being assembled from
1169
          00:45 to 01:15, but this is assembly, not use. The second result
1170
          shows the trampoline in use from 02:30 to 02:55. I will choose
1171
          the second range as the answer.")
1172
      The first result mentions the trampoline being assembled from 00:36 to
1173
           01:59, but this is assembly, not use. The second result shows the
1174
           trampoline in use from 02:30 to 02:55. I will choose the second
1175
          range as the answer.
1176
      . . .
1177
      final_answer = "The trampoline's first use is shown from 02:30 to
1178
          02:55. Final Answer: [02:30, 02:55]"
1179
      finish(final_answer)
1180
1181
      The trampoline's first use is shown from 02:30 to 02:55. Final Answer:
           [02:30, 02:55]
1182
1183
1184
      Critique:
1185
      Strategies A and C actually perceive a person using the trampoline and
1186
           pick the associated range. Strategy B only sees the finished
          trampoline and guesses based on that, but does not directly see
1187
          the part relevant to the query.
```

```
1188
1189
      Winning Strategies:
1190
      A, C
1191
1192
      Question: When is the woman showing a harder variation of the exercise
1193
1194
      Video length: 01:18
1195
1196
      Strategy A (get_segment, find_when):
1197
      input_video = get_segment(start='00:00', end='01:18')
1198
1199
      ...
1200
      harder_exercise_times = find_when(query='woman showing a harder
1201
          variation of the exercise', video_segment=input_video)
1202
1203
      ["00:03", "00:27"], ["00:30", "01:15"]: The video shows a woman
1204
          performing a modified push-up against a wall from 0:03 to 0:27.
          From 0:30 to 1:15, she performs standard push-ups on the floor.
1205
1206
          The query asks for a harder version of the exercise. The standard
          push-up on the floor is considered more challenging than the
1207
          modified version against the wall.
1208
1209
1210
      think("The tool 'find_when' returned two time ranges: ['00:03', '00:27
          '] and ['00:30', '01:15']. The first range shows a modified push-
1211
          up against a wall, and the second shows standard push-ups. The
1212
          standard push-up is described as the harder variation. Therefore,
1213
          the final answer should be the second range.")
1214
1215
      The tool 'find_when' returned two time ranges: ['00:03', '00:27'] and
          ['00:30', '01:15']. The first range shows a modified push-up
1216
          against a wall, and the second shows standard push-ups. The
1217
          standard push-up is described as the harder variation. Therefore,
1218
          the final answer should be the second range.
1219
1220
      final_answer = finish("The tool 'find_when' returned two time ranges:
1221
          ['00:03', '00:27'] and ['00:30', '01:15']. The first range shows a
1222
           modified push-up against a wall, and the second shows standard
1223
          push-ups. The standard push-up is described as the harder
1224
          variation. Therefore, the final answer should be the second range.
1225
           Final Answer: [00:30, 01:15]")
1226
      The tool 'find_when' returned two time ranges: ['00:03', '00:27'] and
1227
          ['00:30', '01:15']. The first range shows a modified push-up
1228
          against a wall, and the second shows standard push-ups. The
1229
          standard push-up is described as the harder variation. Therefore,
1230
          the final answer should be the second range. Final Answer: [00:30,
           01:15]
1231
1232
      Strategy B (get_segment, find_when, retrieval_qa):
1233
1234
      input_video = get_segment(start='00:00', end='01:18')
1235
1236
1237
1238
      times = find_when(query='woman showing a harder variation of the
1239
          exercise', video_segment=input_video)
1240
      ["00:03", "00:27"], ["00:30", "01:15"]: The video shows a woman
1241
          performing a modified push-up against a wall from 0:03 to 0:27.
```

```
1242
          From 0:30 to 1:15, she performs standard push-ups on the floor.
1243
          The query asks for a harder version of the exercise. The standard
1244
          push-up on the floor is considered more challenging than the
1245
          modified version against the wall.
1246
1247
      think("The tool 'find_when' returned two time ranges: 00:03-00:27 and
1248
          00:30-01:15. The description indicates that 00:30-01:15 shows a
1249
          harder variation of the exercise. I will use this time range as
1250
          the final answer.")
1251
      The tool 'find_when' returned two time ranges: 00:03-00:27 and
          00:30-01:15. The description indicates that 00:30-01:15 shows a
1253
          harder variation of the exercise. I will use this time range as
1254
          the final answer.
1255
      . . .
1256
      1257
1258
          00:30-01:15 shows a harder variation of the exercise. I will use
1259
          this time range as the final answer. Final Answer: [00:30, 01:15]"
1260
      . . .
1261
      The tool 'find_when' returned two time ranges: 00:03-00:27 and
1262
          00:30-01:15. The description indicates that 00:30-01:15 shows a
1263
          harder variation of the exercise. I will use this time range as
1264
          the final answer. Final Answer: [00:30, 01:15]
1265
      Strategy C (find_when):
1266
      input_video = get_segment(start='00:00', end='01:18')
1268
1269
1270
      ""
1271
      think("I need to find when the woman is showing a harder variation of
1272
          the exercise. I'll first try to find any mentions of 'harder' or
1273
          similar terms related to the exercise difficulty.")
1274
1275
      I need to find when the woman is showing a harder variation of the
          exercise. I'll first try to find any mentions of 'harder' or
1276
          similar terms related to the exercise difficulty.
1277
1278
1279
      query_result = find_when(query="harder exercise", video_segment=None)
1280
      ["00:03", "00:27"]: The query "harder exercise" is subjective. However
1281
          , based on the video content, the frames from 00:03 to 00:27 show
1282
          a person performing a modified exercise, which could be considered
1283
          an "easier exercise." The exercise is shown for a sustained
1284
          period. Since a harder exercise was requested, it is likely shown
          at a different time in the video.
1285
1286
1287
      final_answer = finish("Based on the tool's analysis, the woman
1288
          performs a harder version of the exercise between 00:03 and 00:27.
1289
          The tool identified this segment, but it may have mistaken it for
          an easier variation. Final Answer: [00:03, 00:27]")
1290
1291
      Based on the tool's analysis, the woman performs a harder version of
1292
          the exercise between 00:03 and 00:27. The tool identified this
1293
          segment, but it may have mistaken it for an easier variation.
1294
          Final Answer: [00:03, 00:27]
1295
      Critique:
```

```
1296
      Strategies A and B find ranges corresponding to both versions of the
1297
          exercise and compare them to find the harder variation.
1298
      Strategy C only identified one range, for a potentially easier version
          , then did not recognize the alternative to compare to.
1299
1300
      Winning Strategies:
1301
      A, B
1302
1303
      Question: When are we given a zoomed in view of the guitar strings?
1304
      Video length: 01:53
1305
      Strategy A (get_segment, find_when):
1306
1307
      input_video = get_segment(start='00:00', end='01:53')
1308
1309
1310
      . . .
1311
      times = find_when(query='zoomed in view of the guitar strings',
1312
          video_segment=input_video)
1313
      ["00:00", "01:52"]: The video shows a person playing a guitar
1314
          throughout its duration. The query asks for a zoomed in view of
1315
          the guitar strings. Therefore, the entire video shows the event.
1316
1317
1318
      final_answer = finish("The query asks for when a zoomed in view of the
           guitar strings is given. The response from 'find_when' indicates
1319
          that a guitar is played throughout the entire video. Therefore,
1320
          the entire video is relevant. Final Answer: [00:00, 01:53]")
1321
1322
      The query asks for when a zoomed in view of the guitar strings is
1323
          given. The response from 'find_when' indicates that a guitar is
          played throughout the entire video. Therefore, the entire video is
1324
           relevant. Final Answer: [00:00, 01:53]
1325
1326
      Strategy B (get_segment, find_when, retrieval_qa):
1327
1328
      input_video = get_segment(start='00:00', end='01:53')
1329
1330
1331
1332
      times = find_when(query='zoomed in view of the guitar strings',
1333
         video_segment=input_video)
1334
      ["00:00", "01:52"]: The video shows a person playing a guitar
1335
          throughout its duration. The query asks for a zoomed in view of
1336
          the guitar strings. Therefore, the entire video shows the event.
1337
1338
      final_answer = finish(final_answer="The video shows a person playing a
1339
           guitar throughout its duration. The query asks for a zoomed in
1340
          view of the guitar strings. Therefore, the entire video shows the
1341
          event. Final Answer: [00:00, 01:53]")
1342
1343
      The video shows a person playing a guitar throughout its duration. The
           query asks for a zoomed in view of the guitar strings. Therefore,
1344
           the entire video shows the event. Final Answer: [00:00, 01:53]
1345
1346
      Strategy C (find_when):
1347
1348
      input_video = get_segment(start='00:00', end='01:53')
1349
```

```
1350
1351
      think ("I need to find when the video shows a close-up of the guitar
1352
          strings. I'll use the 'find_when' tool to search for this.")
1353
      I need to find when the video shows a close-up of the guitar strings.
1354
          I'll use the 'find_when' tool to search for this.
1355
1356
1357
      find_when(query="close-up of guitar strings", video_segment=None)
1358
      ["01:10", "01:25"]: The video shows a person playing guitar. From
1359
          [01:10] to [01:25], the camera zooms in to show a close-up view of
1360
           the guitar strings as the person plays them.
1361
1362
      . . .
1363
      think("The 'find_when' tool returned a single time range: 01:10-01:25.
1364
            This seems to directly answer the question. I will use this as
1365
          my final answer.")
1366
1367
      The 'find_when' tool returned a single time range: 01:10-01:25. This
1368
          seems to directly answer the question. I will use this as my final
1369
1370
      ...
1371
      finish("The find_when tool identified the time range 01:10-01:25 as
1372
          showing a close-up of the guitar strings. Final Answer: [01:10,
1373
          01:25]")
1374
      The find_when tool identified the time range 01:10-01:25 as showing a
1375
          close-up of the guitar strings. Final Answer: [01:10, 01:25]
1376
1377
      Critique:
      Strategies A and B choose the whole video, despite the query asking
1378
          for a zoomed in view, implying a relative change in the video.
1379
      Strategy C instead identifies a particular range.
1380
1381
      Winning Strategies:
1382
      C
1383
      Input:
1384
      Question: When is the person least energetic in the video?
1385
      Video length: 00:49
1386
1387
      Strategy A (get_segment, find_when):
1388
      input_video = get_segment(start='00:00', end='00:49')
1389
1390
1391
      "
1392
      activity_times = find_when(query='person least energetic',
1393
          video_segment=input_video)
1394
1395
      ["00:11", "00:18"], ["00:19", "00:24"], ["00:33", "00:41"]: The query
1396
          asks to identify the time ranges where the person is least
1397
          energetic. Based on the provided video frames, the person appears
          to be sleeping during the periods from 00:11 to 00:18, from 00:19
1398
          to 00:24, and from 00:33 to 00:41. These segments represent
1399
          periods of minimal movement and apparent rest.
1400
1401
1402
      think("The tool returned multiple time ranges: ['00:11', '00:18'], ['
         00:19', '00:24'], ['00:33', '00:41']. Since the question asks for
1403
```

```
1404
           the time the person is LEAST energetic, and the first two ranges
1405
          are contiguous, I will combine them into a single range.")
1406
      ...
      The tool returned multiple time ranges: ['00:11', '00:18'], ['00:19',
1407
          '00:24'], ['00:33', '00:41']. Since the question asks for the
1408
          time the person is LEAST energetic, and the first two ranges are
1409
          contiguous, I will combine them into a single range.
1410
1411
1412
      final_answer = "The person is least energetic from 00:11 to 00:24, and
          from 00:33 to 00:41. Final Answer: [00:11, 00:24], [00:33, 00:41]
1413
1414
      finish(final_answer)
1415
1416
      The person is least energetic from 00:11 to 00:24, and from 00:33 to
          00:41. Final Answer: [00:11, 00:24], [00:33, 00:41]
1417
1418
      Strategy B (get_segment, find_when, retrieval_qa):
1419
1420
      input_video = get_segment(start='00:00', end='00:49')
1421
1422
1423
1424
      activity_times = find_when(query='person is least energetic',
1425
          video_segment=input_video)
1426
      ["00:11", "00:18"]: The query asks to identify when a person is least
1427
          energetic in the video. The video shows a person who appears to be
1428
           sleeping between 00:11 and 00:18. This is the period of least
1429
          visible movement and likely the time of lowest energy.
1430
      ...
1431
      \verb|think("The find_when tool suggests the person is least energetic| \\
1432
          between 00:11 and 00:18. This seems reasonable based on the
1433
          description. I will use this as my final answer.")
1434
1435
      The find_when tool suggests the person is least energetic between
1436
          00:11 and 00:18. This seems reasonable based on the description.
          I will use this as my final answer.
1437
1438
1439
      final_answer = finish("Final Answer: [00:11, 00:18]")
1440
1441
      Final Answer: [00:11, 00:18]
1442
      Strategy C (find_when):
1443
1444
      input_video = get_segment(start='00:00', end='00:49')
1445
1446
1447
      . . .
1448
      think("I need to find when the person is least energetic. I'll use
1449
          find_when to search for inactivity.")
1450
1451
      I need to find when the person is least energetic. I'll use find_when
1452
          to search for inactivity.
1453
1454
      activity_times = find_when(query="person is least energetic",
1455
          video_segment=None)
1456
      ["00:11", "00:18"]: The query asks to identify when a person is least
1457
          energetic in the video. The person appears to be sleeping between
```

Input:

```
1458
          00:11 and 00:18, indicating low energy. There may be other periods
1459
           of low activity, but this period of apparent sleep is the
1460
          clearest indication of low energy.
1461
1462
      think ("The find_when tool suggests the person is least energetic
1463
          between 00:11 and 00:18. This seems like the most obvious period
1464
          of inactivity (sleeping). I will use this as my final answer.")
1465
1466
      The find_when tool suggests the person is least energetic between
          00:11 and 00:18. This seems like the most obvious period of
1467
          inactivity (sleeping). I will use this as my final answer.
1468
1469
1470
      final_answer = finish("The find_when tool suggests the person is least
           energetic between 00:11 and 00:18. This seems like the most
1471
          obvious period of inactivity (sleeping). Therefore, I'm choosing
1472
          this as the final answer. Final Answer: [00:11, 00:18]")
1473
1474
      The find_when tool suggests the person is least energetic between
1475
          00:11 and 00:18. This seems like the most obvious period of
          inactivity (sleeping). Therefore, I'm choosing this as the final
1476
          answer. Final Answer: [00:11, 00:18]
1477
1478
      Winning Strategies:
1479
1480
```

Listing 6: ActivityNet-RTL critic examples.

```
You will be given a question about a video and five possible answer
1484
          options, where C refers to the person wearing the camera. Question
1485
          : What is the main action being performed by C in the video, and
1486
          provide specific moments that support your observation. Possible
          answer choices:
1487
      (1) The main action performed by C is applying makeup. This is evident
1488
           as C picks up a brush, looks in a mirror, and applies something
1489
          to their face.
1490
      (2) C's main action in the video is eating a meal. This is clearly
          evident because C lifts a fork to their mouth multiple times and
1491
          chews.
1492
      (3) The main action performed by C in the video is reading a book.
1493
          This is evident in the fact that C holds a book open, turns pages,
1494
          and moves their eyes across the page.
1495
      (4) In the video, the main action performed by C is talking on the
1496
          phone. This is evident because C holds a phone to their ear and
          speaks into it several times.
1497
      (5) The main action performed by {\tt C} in the video is typing on a
1498
          computer. This is clearly seen as C sits at a desk, places their
1499
          hands on a keyboard, and moves their fingers.
1500
      Video length: 01:48
1501
      Strategy A (retrieval_qa, get_segment):
1502
1503
      input_video = get_segment(start='00:00', end='01:48')
1504
1505
1506
1507
      answer = retrieval_qa(question='What is the main action being
1508
          performed by C in the video, and provide specific moments that
1509
          support your observation.', answer_options=['The main action
1510
          performed by C is applying makeup. This is evident as C picks up a
1511
          brush, looks in a mirror, and applies something to their face.',
          'C\'s main action in the video is eating a meal.
                                                             This is clearly
```

```
1512
          evident because C lifts a fork to their mouth multiple times and
1513
          chews.', 'The main action performed by C in the video is reading a
1514
           book. This is evident in the fact that C holds a book open, turns
           pages, and moves their eyes across the page.', 'In the video, the
1515
           main action performed by C is talking on the phone. This is
1516
          evident because C holds a phone to their ear and speaks into it
1517
          several times.', 'The main action performed by C in the video is
1518
          typing on a computer. This is clearly seen as C sits at a desk,
1519
          places their hands on a keyboard, and moves their fingers.'],
1520
          video_segment=input_video)
1521
      C is clearly typing on a computer throughout the video.
                                                                 They are
1522
          sitting at a desk with their hands on a keyboard and their fingers
1523
           moving. While there may be other actions present, the consistent
1524
          and primary action is typing.
1525
      Final Answer: (5)
1526
1527
1528
      final_answer = finish(final_answer='C is clearly typing on a computer
1529
          throughout the video. They are sitting at a desk with their hands
1530
           on a keyboard and their fingers moving. While there may be other
          actions present, the consistent and primary action is typing.\n
1531
          nFinal Answer: (5)')
1532
1533
      C is clearly typing on a computer throughout the video.
                                                                 They are
1534
          sitting at a desk with their hands on a keyboard and their fingers
1535
           moving. While there may be other actions present, the consistent
          and primary action is typing.
1536
1537
      Final Answer: (5)
1538
1539
      Strategy B (direct retrieval_qa):
1540
      C is clearly typing on a computer throughout the video. They are
          sitting at a desk with their hands on a keyboard and their fingers
1541
           are moving. This is the primary and most consistent action shown
1542
1543
1544
      Final Answer: (5)
1545
1546
      Strategy C (retrieval_qa, get_segment, find_when):
1547
1548
      input_video = get_segment(start='00:00', end='01:48')
1549
1550
1551
      ""
1552
      think ("The question asks about the main action C performs. I'll use '
1553
          find_when' to check for specific actions mentioned in the options,
1554
           like 'applying makeup', 'eating', 'reading', 'talking on the
1555
         phone', and 'typing'. This will help identify the main action.")
1556
      The question asks about the main action C performs. I'll use '
1557
          find_when' to check for specific actions mentioned in the options,
1558
           like 'applying makeup', 'eating', 'reading', 'talking on the
1559
          phone', and 'typing'. This will help identify the main action.
1560
1561
      makeup_times = find_when(query='C applies makeup', video_segment=
1562
          input_video)
1563
      eating_times = find_when(query='C eats a meal', video_segment=
1564
          input_video)
      reading_times = find_when(query='C reads a book', video_segment=
1565
          input_video)
```

```
1566
      phone_times = find_when(query='C talks on the phone', video_segment=
1567
          input_video)
1568
       typing_times = find_when(query='C types on a computer', video_segment=
1569
          input_video)
1570
       ["00:12", "00:18"]: At around 00:12, C appears to pick up a small
1571
          object, which is potentially a phone, and brings it to their ear.
1572
          The interaction with the object continues until about 00:18.
1573
          Although it is not definitively clear whether C is talking, their
1574
          actions suggest they may be on a phone call.
1575
       . . .
1576
      think ("'find_when' only returned a possible phone call. The limited
1577
          information isn't conclusive. Let's use 'retrieval_qa' to check if
1578
           C is focused on work, which could suggest typing.")
1579
       'find_when' only returned a possible phone call. The limited
1580
          information isn't conclusive. Let's use 'retrieval_qa' to check if
1581
           C is focused on work, which could suggest typing.
1582
       ...
1583
      work_focused = retrieval_qa(question="Is C focused on work in the
1584
          video?", answer_options=["Yes", "No"], video_segment=input_video)
1585
1586
      The video shows C sitting at a desk and using a computer. Although it'
1587
          s not shown exactly what they are doing, their posture and actions
1588
           suggest they may be working.
1589
      Final Answer: (2)
1590
1592
      think ("The response suggests work at a desk and computer use. That
1593
          supports typing. However, the short phone interaction is also a
          possibility. Since 'eating' seems unlikely in a work context, I'll
1594
           guess between phone and typing, leaning towards phone due to the
1595
           'find_when' result.")
1596
1597
      The response suggests work at a desk and computer use. That supports
1598
          typing. However, the short phone interaction is also a
          possibility. Since 'eating' seems unlikely in a work context, I'll
1599
           guess between phone and typing, leaning towards phone due to the
1600
           'find_when' result.
1601
1602
1603
      final_answer = finish(final_answer="The 'find_when' result hinted at
          phone use. 'retrieval_qa' suggested work at a desk with a computer , which supports 'typing.' Lacking definitive proof, I'll guess '
1604
1605
          talking on the phone' based on the slight 'find_when' evidence.
1606
          Final Answer: (4)")
1607
1608
      The 'find_when' result hinted at phone use. 'retrieval_qa' suggested
          work at a desk with a computer, which supports 'typing.' Lacking
1609
          definitive proof, I'll guess 'talking on the phone' based on the slight 'find_when' evidence. Final Answer: (4)
1610
1611
1612
      Critique:
1613
      A: 5. B: 5. C: 4. Strategy C attempts an overly complex strategy with
          the find_when tool and ends up guessing due to limited information
1614
1615
1616
      Winning Strategies:
1617
      A, B
1618
1619
      Input:
```

```
1620
      You will be given a question about a video and five possible answer
1621
          options, where {\tt C} refers to the person wearing the camera. Question
1622
          : Describe how C prepared the ingredients and assembled the
          sandwich in the video, noting key steps and transitions. Possible
1623
          answer choices:
1624
      (1) C prepared the ingredients by ordering them online and having them
1625
           delivered. They then assembled the sandwich by layering the
1626
          ingredients between two slices of bread. They added condiments
1627
          and cut the sandwich in half.
1628
      (2) C prepared the ingredients by taking pre-sliced ingredients from
          the refrigerator. They then assembled the sandwich by layering the
1629
           ingredients between two slices of bread. They added condiments
1630
          and cut the sandwich in half.
1631
      (3) C prepared the ingredients by chopping vegetables and slicing meat
1632
          . They then assembled the sandwich by layering the ingredients
          between two slices of bread. They added condiments and cut the
1633
          sandwich in half.
1634
      (4) C prepared the ingredients by microwaving pre-packaged meals. They
1635
           then assembled the sandwich by layering the ingredients between
1636
          two slices of bread. They added condiments and cut the sandwich
1637
          in half.
1638
      (5) C prepared the ingredients by going to a restaurant and getting
          takeout. They then assembled the sandwich by layering the
1639
          ingredients between two slices of bread. They added condiments
1640
          and cut the sandwich in half.
1641
      Video length: 02:45
1642
1643
      Strategy A (retrieval_qa, get_segment):
1644
      input_video = get_segment(start='00:00', end='02:45')
1646
1647
      ""
1648
      answer = retrieval_qa(question='Describe how C prepared the
1649
          ingredients and assembled the sandwich in the video, noting key
1650
          steps and transitions. ', answer_options=['C prepared the
1651
          ingredients by ordering them online and having them delivered.
1652
          They then assembled the sandwich by layering the ingredients
          between two slices of bread. They added condiments and cut the
1653
          sandwich in half.', 'C prepared the ingredients by taking pre-
1654
          sliced ingredients from the refrigerator. They then assembled the
1655
          sandwich by layering the ingredients between two slices of bread.
1656
           They added condiments and cut the sandwich in half.', 'C prepared
1657
           the ingredients by chopping vegetables and slicing meat. They
          then assembled the sandwich by layering the ingredients between
1658
          two slices of bread. They added condiments and cut the sandwich
1659
          in half.', 'C prepared the ingredients by microwaving pre-packaged
1660
           meals. They then assembled the sandwich by layering the
1661
          ingredients between two slices of bread. They added condiments
          and cut the sandwich in half.', 'C prepared the ingredients by
          going to a restaurant and getting takeout. They then assembled the
1663
           sandwich by layering the ingredients between two slices of bread.
1664
            They added condiments and cut the sandwich in half.'],
1665
          video_segment=input_video)
1666
      . . .
1667
      The video shows C assembling a sandwich. They take various ingredients
          , like sliced meat and vegetables, from the refrigerator.
1668
          layer these ingredients between two slices of bread, add
1669
          condiments, and then cut the sandwich in half. There is no
1670
          indication of ordering food online, chopping vegetables,
1671
          microwaving meals, or getting takeout. The ingredients appear to
```

1673

be pre-sliced.

Final Answer: (2)

```
1674
1675
      ...
1676
      final_answer = finish("The video clearly shows C using pre-sliced
          ingredients from the refrigerator. Therefore, the most accurate
1677
          description is option (2). Final Answer: (2)")
1678
1679
      The video clearly shows C using pre-sliced ingredients from the
1680
          refrigerator. Therefore, the most accurate description is option
1681
          (2). Final Answer: (2)
1682
      Strategy B (direct retrieval_qa):
1683
      The video shows C making a sandwich. They retrieve pre-sliced
1684
          ingredients, like lettuce, tomato, and meat, from the refrigerator
1685
             They then assemble the sandwich, add condiments, and cut it in
1686
          half. The ingredients were prepared by being pre-sliced and stored
           in the refrigerator.
1687
1688
      Final Answer: (3)
1689
1690
      Strategy C (retrieval_qa, get_segment, find_when):
1691
      input_video = get_segment(start='00:00', end='02:45')
1692
1693
1694
1695
1696
      think ("To understand how the ingredients were prepared, I'll use '
1697
          find_when' to look for actions like chopping or retrieving from
          the refrigerator, and then use 'retrieval_qa' to select the best
1698
          answer.")
1699
1700
      To understand how the ingredients were prepared, I'll use 'find_when'
1701
          to look for actions like chopping or retrieving from the
          refrigerator, and then use 'retrieval_qa' to select the best
1702
          answer.
1703
1704
1705
      preparation_times = find_when("C preparing the sandwich ingredients",
1706
          input_video)
1707
      ["00:05", "00:40"]: C retrieves pre-sliced sandwich ingredients from
1708
          the refrigerator.
1709
1710
1711
      answer = retrieval_qa(question="Describe how C prepared the
          ingredients and assembled the sandwich in the video, noting key
1712
          steps and transitions.", answer_options=["C prepared the
1713
          ingredients by ordering them online and having them delivered.
1714
          They then assembled the sandwich by layering the ingredients
1715
          between two slices of bread. They added condiments and cut the
1716
          sandwich in half.", "C prepared the ingredients by taking pre-
          sliced ingredients from the refrigerator. They then assembled the
1717
          sandwich by layering the ingredients between two slices of bread.
1718
           They added condiments and cut the sandwich in half.", "C prepared
1719
           the ingredients by chopping vegetables and slicing meat. They
1720
          then assembled the sandwich by layering the ingredients between
1721
          two slices of bread. They added condiments and cut the sandwich
          in half.", "C prepared the ingredients by microwaving pre-packaged
1722
           meals. They then assembled the sandwich by layering the
1723
          ingredients between two slices of bread. They added condiments
1724
          and cut the sandwich in half.", "C prepared the ingredients by
1725
          going to a restaurant and getting takeout. They then assembled the
1726
           sandwich by layering the ingredients between two slices of bread.
1727
            They added condiments and cut the sandwich in half."],
          video_segment=input_video)
```

```
1728
1729
      The video shows C making a sandwich. They retrieve pre-sliced
          ingredients, like lettuce, tomato, and meat, from the refrigerator
1730
1731
             They then assemble the sandwich, add condiments, and cut it in
          half. The ingredients were prepared by being pre-sliced and stored
1732
           in the refrigerator.
1733
1734
      Final Answer: (3)
1735
      . . .
1736
      final_answer = finish("Since C takes pre-sliced ingredients from the
1737
         refrigerator, the correct answer is (3). Final Answer: (3)")
1738
1739
      Since C takes pre-sliced ingredients from the refrigerator, the
1740
          correct answer is (3). Final Answer: (3)
1741
      Critique:
1742
      A: 2. B: 3. C: 3. Strategy A correctly identifies that the ingredients
1743
           are pre-sliced from the refrigerator but fails to choose the
1744
          matching answer.
1745
1746
      Winning Strategies:
      B, C
1747
1748
      Input:
1749
      You will be given a question about a video and five possible answer
1750
          options. Question: Compare the lighting conditions in the room at
          the start of the video to those near the end. How has the lighting
1751
           evolved over the course of the video? Possible answer choices:
1752
      (1) The lighting in the earlier segments of the video is primarily
1753
          provided by a single overhead lamp, while the lighting in the
1754
          latter segments is primarily provided by multiple spotlights.
1755
      (2) The lighting in the earlier segments of the video is primarily
          natural light from a window, while the lighting in the latter
1756
          segments is primarily provided by an overhead lamp.
1757
      (3) The lighting in the earlier segments of the video is primarily
1758
          provided by multiple spotlights, while the lighting in the latter
1759
          segments is primarily natural light from a window.
1760
      (4) The lighting in the earlier segments of the video is primarily
1761
          natural light from a window, while the lighting in the latter
          segments is dimmer, perhaps due to cloud cover or sunset.
1762
      (5) The lighting in the earlier segments of the video is dimmer,
1763
          perhaps due to cloud cover or sunset, while the lighting in the
1764
          latter segments is primarily natural light from a window.
1765
      Video length: 03:01
1766
      Strategy A (retrieval_qa, get_segment):
1767
1768
      input_video = get_segment(start='00:00', end='03:01')
1769
1770
1771
      . . .
1772
      think ("To compare the lighting in the earlier vs. latter segments, I'
1773
          ll first use retrieval_qa on the first minute and the last minute.
1774
            This should give a general idea of the differences. Then I can
1775
          refine my approach if needed.")
1776
      To compare the lighting in the earlier vs. latter segments, I'll first
1777
           use retrieval_qa on the first minute and the last minute. This
1778
          should give a general idea of the differences. Then I can refine
1779
          my approach if needed.
1780
1781
      first_minute_segment = get_segment(start='00:00', end='01:00')
```

```
1782
      first_minute_answer = retrieval_qa(question="Describe the lighting in
1783
          this segment.", video_segment=first_minute_segment)
1784
1785
      last_minute_segment = get_segment(start='02:00', end='03:00')
      last_minute_answer = retrieval_qa(question="Describe the lighting in
1786
          this segment.", video_segment=last_minute_segment)
1787
1788
      think(f"First minute: \{\{first\_minute\_answer\}\}. \ Last minute: \{\{first\_minute\_answer\}\}.
1789
          last_minute_answer}}. Now I need to compare these and select the
1790
          best matching answer choice.")
1791
      First minute: The lighting appears to be dim, likely from a weak
1792
          overhead light source.. Last minute: The room is brightly lit by
1793
          natural light coming in from a large window.. Now I need to
1794
          compare these and select the best matching answer choice.
1795
      . . .
1796
      think ("The first minute shows dim lighting, which could be interpreted
1797
           as cloud cover or sunset. The last minute shows bright natural
1798
          light. Option (4) seems close, mentioning dimmer lighting
          initially and then a change, but it doesn't specifically mention
1799
          the natural light. However, it's the closest option.")
1800
1801
      The first minute shows dim lighting, which could be interpreted as
1802
          cloud cover or sunset. The last minute shows bright natural light.
1803
            Option (4) seems close, mentioning dimmer lighting initially and
1804
           then a change, but it doesn't specifically mention the natural
          light. However, it's the closest option.
1805
1806
1807
      final_answer = finish("Reasoning: The initial dim lighting suggests
1808
          cloud cover or sunset as per option (4), even though the later
1809
          segment has bright natural light, not just a generic change. Final
           Answer: (4)")
1810
1811
      Reasoning: The initial dim lighting suggests cloud cover or sunset as
1812
          per option (4), even though the later segment has bright natural
1813
          light, not just a generic change. Final Answer: (4)
1814
1815
      Strategy B (direct retrieval_qa):
1816
      The video starts with dim lighting, likely due to cloud cover or the
1817
          time of day. As the video progresses, the lighting becomes
1818
          brighter and it appears more natural light is entering the room,
1819
          likely from a window. This suggests the sun might have come out
          from behind the clouds or the video was filmed as the day
1820
          progressed.
1821
1822
1823
      Final Answer: (5)
1824
1825
      Strategy C (retrieval_qa, get_segment, find_when):
1826
1827
      input_video = get_segment(start='00:00', end='03:01')
1828
1829
1830
      ""
1831
      think ("To compare the lighting, I'll analyze the first and last
1832
          minutes using find_when and then use retrieval_qa for a more
1833
          detailed comparison.")
1834
```

```
1836
      To compare the lighting, I'll analyze the first and last minutes using
1837
           find_when and then use retrieval_qa for a more detailed
1838
          comparison.
1839
1840
      early_segment = get_segment(start='00:00', end='01:00')
1841
1842
      ['00:00', '01:00']
1843
      "
1844
      early_lighting = find_when(query="Describe the lighting conditions.",
1845
         video_segment=early_segment)
1846
1847
      ["00:00", "00:59"]: The lighting conditions appear \dim, likely due to
1848
          a weak overhead light source or cloud cover outside.
1849
      ...
1850
      late_segment = get_segment(start='02:00', end='03:01')
1851
1852
      ['02:00', '03:01']
1853
1854
      late_lighting = find_when(query="Describe the lighting conditions.",
1855
          video_segment=late_segment)
1856
1857
      ["02:00", "03:00"]: The room is brightly lit by natural light
1858
          streaming in from a large window.
1859
1860
1861
      comparison = retrieval_qa(question="Compare the lighting conditions in
1862
           the room at the start of the video to those near the end. How has
1863
           the lighting evolved over the course of the video?",
1864
          answer_options=[
          "The lighting in the earlier segments of the video is primarily
1865
              provided by a single overhead lamp, while the lighting in the
1866
              latter segments is primarily provided by multiple spotlights."
1867
1868
          "The lighting in the earlier segments of the video is primarily
              natural light from a window, while the lighting in the latter
1869
              segments is primarily provided by an overhead lamp.",
1870
          "The lighting in the earlier segments of the video is primarily
1871
              provided by multiple spotlights, while the lighting in the
1872
              latter segments is primarily natural light from a window.",
1873
          "The lighting in the earlier segments of the video is primarily
1874
              natural light from a window, while the lighting in the latter
              segments is dimmer, perhaps due to cloud cover or sunset.",
1875
          "The lighting in the earlier segments of the video is dimmer,
1876
              perhaps due to cloud cover or sunset, while the lighting in
1877
              the latter segments is primarily natural light from a window."
1878
      ], video_segment=input_video)
1879
      The video starts with dim lighting, likely due to cloud cover or the
1880
          time of day. As the video progresses, the lighting becomes
1881
          brighter and it appears more natural light is entering the room,
1882
          likely from a window. This suggests the sun might have come out
1883
          from behind the clouds or the video was filmed as the day
1884
          progressed.
1885
      Final Answer: (5)
1886
1887
1888
      final_answer = finish("Based on the analysis, the early segment has
1889
          dim lighting (potentially due to cloud cover), while the late
```

```
1890
          segment is lit by natural light. This matches option (5). Final
1891
          Answer: (5)")
1892
      Based on the analysis, the early segment has dim lighting (potentially
1893
           due to cloud cover), while the late segment is lit by natural
1894
          light. This matches option (5). Final Answer: (5)
1895
1896
1897
      Critique: A: 4. B: 5. C: 5. Strategy A initially uses the valid
1898
          strategy of using segments at the beginning and end, but
          ultimately does not logically connect the information to the
1899
          appropriate answer choice.
1900
1901
      Winning Strategies: B, C
1902
      Input:
1903
      You will be given a question about a video and five possible answer
1904
          options, where C refers to the person wearing the camera. Question
1905
          : What is the primary action performed with the backpack at the
1906
          beginning and the main purpose of using a backpack in the context
1907
          of this video? please highlight the significance of this action in
1908
           the overall travel process. Possible answer choices:
      (1) The primary action performed with the backpack at the beginning is
1909
           to open it. the main purpose of using a backpack in the context
1910
          of this video is to store belongings. opening the backpack allows
1911
          for easy access to items needed for the trip.
1912
      (2) The primary action performed with the backpack at the beginning is
           to close it. the main purpose of using a backpack in the context
1913
          of this video is to carry belongings. closing the backpack ensures
1914
           that items are secure and do not fall out during travel.
1915
      (3) The primary action performed with the backpack at the beginning is
1916
           to put it on. the main purpose of using a backpack in the context
1917
           of this video is to carry belongings. putting the backpack on
          allows for hands-free travel and even weight distribution.
1918
      (4) The primary action performed with the backpack at the beginning is
1919
           to pack it. the main purpose of using a backpack in the context
1920
          of this video is to carry belongings. packing the backpack ensures
1921
          that all necessary items are brought on the trip.
1922
      (5) The primary action performed with the backpack at the beginning is
          to take it off. the main purpose of using a backpack in the
1923
          context of this video is to carry belongings. taking the backpack
1924
          off allows for a rest and reduces strain on the back.
1925
      Video length: 02:48
1926
1927
      Strategy A (retrieval_qa, get_segment):
1928
      input_video = get_segment(start='00:00', end='02:48')
1929
1930
1931
      ""
1932
      answer = retrieval_qa(question='What is the primary action performed
1933
          with the backpack at the beginning and the main purpose of using a
1934
           backpack in the context of this video? please highlight the
1935
          significance of this action in the overall travel process.',
1936
          answer_options=['The primary action performed with the backpack at
1937
           the beginning is to open it. the main purpose of using a backpack
           in the context of this video is to store belongings. opening the
1938
          backpack allows for easy access to items needed for the trip.', '
1939
          The primary action performed with the backpack at the beginning is
1940
          to close it. the main purpose of using a backpack in the context
1941
          of this video is to carry belongings. closing the backpack ensures
1942
          that items are secure and do not fall out during travel.', 'The
          primary action performed with the backpack at the beginning is to
1943
         put it on. the main purpose of using a backpack in the context of
```

```
1944
          this video is to carry belongings. putting the backpack on allows
1945
          for hands-free travel and even weight distribution.', 'The primary
1946
           action performed with the backpack at the beginning is to pack it
          . the main purpose of using a backpack in the context of this
1947
          video is to carry belongings. packing the backpack ensures that
1948
          all necessary items are brought on the trip.', 'The primary action
1949
           performed with the backpack at the beginning is to take it off.
1950
          the main purpose of using a backpack in the context of this video
1951
          is to carry belongings. taking the backpack off allows for a rest
1952
          and reduces strain on the back.'], video_segment=input_video)
1953
      The video shows a person packing a backpack at the beginning. They
1954
          carefully organize and place items inside. This is essential for
1955
          travel as it ensures all needed items are brought along. The main
1956
          purpose of the backpack in the video is to carry belongings.
1957
      Final Answer: (4)
1958
1959
      . . .
1960
      think ("The response aligns with option 4, showing the backpack is
1961
          initially packed, crucial for travel to ensure having necessary
1962
          items. The main purpose is indeed carrying belongings.")
1963
      The response aligns with option 4, showing the backpack is initially
1964
          packed, crucial for travel to ensure having necessary items. The
1965
          main purpose is indeed carrying belongings.
1966
1967
      final_answer = finish("Final Answer: (4)")
1968
      Final Answer: (4)
1970
1971
      Strategy B (direct retrieval_qa):
      The video shows a person packing a backpack at the beginning. They
1972
          carefully organize and place items inside. This is essential for
1973
          travel as it ensures all needed items are brought along. The main
1974
          purpose of the backpack in the video is to carry belongings.
1975
1976
      Final Answer: (4)
1977
1978
      Strategy C (retrieval_qa, get_segment, find_when):
1979
1980
      input_video = get_segment(start='00:00', end='02:48')
1981
1982
1983
      . . .
1984
      think("I'll use 'find_when' to identify the initial backpack action
1985
          and 'retrieval_qa' to determine its purpose in the video, then
          analyze the options.")
1987
      I'll use 'find_when' to identify the initial backpack action and '
1988
          retrieval_qa' to determine its purpose in the video, then analyze
1989
          the options.
1990
      ...
1991
      initial_action_times = find_when("someone interacts with a backpack",
1992
          input_video)
1993
1994
      ["00:05", "00:25"]: The video shows a person actively packing a
1995
          backpack. They are placing clothes, toiletries, and other items
1996
          inside the backpack.
1997
```

```
1998
       ["01:50", "01:55"]: The person is seen wearing the backpack and
1999
          walking through an airport terminal.
2000
2001
2002
      purpose = retrieval_qa(question="What is the main purpose of using a
2003
          backpack in the context of this video?", answer_options=[
2004
           "The main purpose of using a backpack in the context of this video
2005
                is to store documents.",
2006
           "The main purpose of using a backpack in the context of this video
               is to carry belongings.",
2007
           "The main purpose of using a backpack in the context of this video
2008
               is to organize toiletries."
2009
           "The main purpose of using a backpack in the context of this video
2010
                is to conceal valuables.",
           "The main purpose of using a backpack in the context of this video
2011
                is to transport food.'
2012
      ], video_segment=input_video)
2013
2014
      The video shows a person packing and then carrying a backpack filled
2015
          with personal belongings during their travel. Therefore, the main
2016
          purpose of the backpack is to carry belongings.
2017
      Final Answer: (2)
2018
2019
2020
      final_answer = finish("'find_when' shows the backpack is first packed,
           and 'retrieval_qa' confirms it's used to carry belongings.
2021
          Knowing packing is essential for bringing necessary items, option
2022
          (3), putting on the backpack for hands-free carrying, aligns with
2023
          the main purpose and initial action. Final Answer: (3)")
2024
2025
       'find_when' shows the backpack is first packed, and 'retrieval_qa'
2026
          confirms it's used to carry belongings. Knowing packing is
          essential for bringing necessary items, option (3), putting on the
2027
           backpack for hands-free carrying, aligns with the main purpose % \left( 1\right) =\left( 1\right) \left( 1\right) 
2028
          and initial action. Final Answer: (3)
2029
2030
      Critique:
      A: 4. B: 4. C: 2. Strategy C correctly finds that the backpack is
2031
          first packed and then used to carry belongings, and reasons that
2032
          the main purpose is to carry belongings, but ultimately
2033
          illogically concludes the purpose is to put the backpack on
2034
          despite this.
2035
      Winning Strategies:
2036
      A, B
2037
       """).strip()
2038
2039
         else:
           CRITIC_EXAMPLES = textwrap.dedent(f"""
         Input:
2041
         You will be given a question about a video and four possible answer
2042
            options. Question: At 12:10, what color suit is the host wearing
2043
            ?Possible answer choices:
2044
         (1) Navy
2045
         (2) Tan
         (3) Black
2046
         (4) White
2047
         Video length: 40:50
2048
2049
         Strategy A (retrieval_qa, get_segment):
2050
         input_video = get_segment(start='00:00', end='40:50')
2051
```

```
2052
2053
2054
2055
         segment = get_segment(start='12:10', end='12:20')
         answer = retrieval_qa(question='What color suit is the host wearing?
2056
            ', video_segment=segment)
2057
2058
         The host is wearing a navy suit.
2059
2060
         final_answer = finish(f"The retrieval_qa tool determined the host
2061
           was wearing a navy suit at 12:10. Final Answer: (1)")
2062
2063
         The retrieval_qa tool determined the host was wearing a navy suit at
2064
             12:10. Final Answer: (1)
2065
         Strategy B (direct retrieval_qa):
2066
        Here's how to solve this:
2067
2068
         1. **Locate the timestamp: ** Go to the 12:10 mark in the video.
2069
2070
         2. **Observe the host's attire: ** At this point, the host (Jimmy
            Kimmel) is wearing a dark suit and a dark-colored tie.
2071
2072
         3. **Compare to options:** The suit appears to be black.
2073
2074
         Therefore, the correct answer is (3) Black.
        Final Answer: (3)
2075
2076
         Strategy C (retrieval_qa, get_segment, find_when):
2077
2078
         input_video = get_segment(start='00:00', end='40:50')
2079
2081
2082
         segment = get_segment(start='12:10', end='12:20')
2083
2084
         ['12:10', '12:20']
2085
2086
         answer = retrieval_qa(question='At 12:10, what color suit is the
2087
            host wearing?', answer_options=['Navy', 'Tan', 'Black', 'White'
2088
            ], video_segment=segment)
2089
        The image at 12:10 shows the host wearing a navy suit.
2090
2091
        Final Answer: (1)
2092
2093
2094
         final_answer = finish(final_answer='The image at 12:10 shows the
2095
           host wearing a navy suit. Final Answer: (1)')
2096
         The image at 12:10 shows the host wearing a navy suit. Final Answer:
2097
             (1)
2098
2099
         Critique:
2100
         A: 1. B: 3. C: 1. Strategies A and C succeed while Strategy B fails
            because Strategies A and C successfully use the get_segment tool
2101
             to isolate the relevant part of the video.
2102
         Strategy B fails as the retrieval_qa tool alone cannot actually
2103
            locate a given timestamp in a video, resulting in a hallucinated
2104
             answer.
2105
        Winning Strategies:
```

```
2106
        A, C
2107
2108
        Input:
2109
         You will be given a question about a video and four possible answer
            options. Question: What does the man in charcoal suit do in the
2110
            next scene after the scene of a photo of six people end?Possible
2111
             answer choices:
2112
         (1) He shows his scars
2113
         (2) He takes off his shirt
2114
         (3) He shows his jewelry
         (4) He shows his physique
2115
        Video length: 39:59
2116
2117
         Strategy A (retrieval_qa, get_segment):
2118
         input_video = get_segment(start='00:00', end='39:59')
2119
2120
2121
2122
2123
         think ("I need to identify the scene where a photo of six people is
2124
            shown. Then I need to find the scene immediately following it
            to see what the man in the charcoal suit does.")
2125
2126
        I need to identify the scene where a photo of six people is shown.
2127
            Then I need to find the scene immediately following it to see
2128
            what the man in the charcoal suit does.
2129
2130
         segment_with_photo = get_segment(start='00:00', end='01:00')
2131
         answer = retrieval_qa(question="What does the man in charcoal suit
2132
            do in the scene immediately after the scene with a photo of six
2133
            people?", video_segment=segment_with_photo)
2134
        After the scene with the photo of six people, the man in the
2135
            charcoal suit walks down a hallway and is greeted by a group of
2136
            people who clap and cheer for him.
2137
2138
         think ("The provided text doesn't describe any of the options (1) He
2139
            shows his scars, (2) He takes off his shirt, (3) He shows his
2140
            jewelry, (4) He shows his physique. The man is celebrated. I
2141
            will guess the closest option, which is none of them, but I will
2142
             choose the least violent option.")
2143
        The provided text doesn't describe any of the options (1) He shows
2144
            his scars, (2) He takes off his shirt, (3) He shows his jewelry,
2145
             (4) He shows his physique. The man is celebrated. I will
2146
            guess the closest option, which is none of them, but I will
2147
            choose the least violent option.
2148
         . . .
2149
         final_answer = finish("Final Answer: (1)")
2150
2151
        Final Answer: (1)
2152
2153
        Strategy B (direct retrieval_qa):
        Here's how to solve this step-by-step:
2154
2155
         1. **Identify the relevant scene:** The question refers to the scene
2156
             after a photo of six people. Looking through the timestamps,
2157
            this appears to be around 20:49.
2158
```

```
2160
                    2. **Observe the next scene: ** The next scene shows the man in the
2161
                             charcoal suit (Harry Styles) pointing to and showing jewelry on
2162
                             his arm.
2163
                    3. **Match to the answer choices:** This action directly corresponds
2164
                               to answer choice (3).
2165
2166
                     Therefore, the correct answer is **(3) He shows his jewelry**.
2167
                    Final Answer: (3)
2168
                     Strategy C (retrieval_qa, get_segment, find_when):
2169
2170
                     input_video = get_segment(start='00:00', end='39:59')
2171
2172
2173
                     . . .
2174
                     think ("First, I need to find when the scene with the photo of six
2175
                             people ends. Then I can examine the next scene to see what the
2176
                            man in the charcoal suit does.")
2177
2178
                    First, I need to find when the scene with the photo of six people
                             ends. Then I can examine the next scene to see what the man in
2179
                             the charcoal suit does.
2180
2181
2182
                    photo_end_time = find_when(query="scene with a photo of six people",
2183
                               video_segment=None)
2184
                     ["00:05", "00:24"]: The query asks to find the scene with a photo of
2185
                               six people. From the provided frames, a photo of six people is
2186
                               visible on the wall in the background starting at 00:05 and
2187
                             remaining visible until 00:24. The photo is partially obscured
2188
                             at times by other people and objects, but it is consistently
                             present in the background during this timeframe.
2189
                     ["10:39", "10:43"]: The query asks for a scene with a photo of six
2190
                             people. At 10:39 through 10:43 a large group photo of more than
2191
                               six people is shown.
2192
                     ["17:47", "17:53"]: The query asks for a scene with a photo of six
                             people. At 17:47-17:53, a photo of approximately six people is
2193
                             visible behind a curtain as Harry Styles walks out onto the
2194
                             stage. The photo is partially obscured and the exact number of
2195
                             people is difficult to ascertain with certainty, but it appears % \left( 1\right) =\left( 1\right) +\left( 1\right) +\left(
2196
                             to show at least five.
2197
                     ["18:47", "19:59"]: The query asks for the time range of a scene
                             with a photo of six people. A photo of what appears to be six
2198
                             people is visible on the wall behind the main set in frames
2199
                             starting at 18:47 and continuing through 19:59. The photo is
2200
                             partially obscured at times by other elements in the scene, but
2201
                             it is consistently visible enough to be identified.
                     ["20:03", "20:04"]: The query asks for the time range of a scene
2203
                             with a photo of six people. At 20:03, a wide shot of the studio
                             audience is shown, and a photo of six people is visible on a screen in the background. This photo remains visible until
2204
2205
                             approximately 20:04.
2206
                     ["23:15", "23:25"]: The query asks for the timestamp of a scene with
2207
                               a photo of six people. In the provided video frames, a photo
                             of six people is visible on the wall in the background starting
2208
                             at 23:15 and remaining visible until 23:25. The photo is
2209
                             partially obscured at times by other elements in the frame, but
2210
                             it is consistently present during this time range.
2211
                     ["35:17", "35:20"]: The provided video frames show a wide shot of a
2212
                             studio audience at the 35:17 mark. The audience appears to
                             consist of more than six people. While the exact number is
2213
```

```
2214
            difficult to determine from the provided frames, it is clearly a
2215
             large group. Therefore, this time range satisfies the query.
2216
         ["39:03", "39:18"]: The provided video frames show a live
2217
            performance with an audience. At the end of the performance,
            from 39:03 to 39:18, the camera focuses on the audience, and
2218
            there are enough people visible to potentially include a group
2219
            of five. The resolution isn't high enough to definitively
2220
            confirm the presence of a specific photo of \operatorname{six} people, but the
2221
            scene depicts a large group of people, making it a plausible
2222
            match for the query.
2223
         . . .
2224
         think ("The most likely candidate for the 'scene with a photo of six
2225
            people' is the longest continuous segment: 18:47-19:59. Let's
2226
            use that as the end point for the previous scene.")
2227
        The most likely candidate for the 'scene with a photo of six people'
2228
             is the longest continuous segment: 18:47-19:59. Let's use that
2229
            as the end point for the previous scene.
2230
2231
2232
        next_scene_segment = get_segment(start='19:59', end='20:05')
2233
         ['19:59', '20:05']
2234
2235
2236
         answer = retrieval_qa(question="What does the man in the charcoal
            suit do in the scene from 19:59 to 20:05?", answer_options=['He shows his scars', 'He takes off his shirt', 'He shows his
2237
2238
            jewelry', 'He shows his physique'], video_segment=
2239
            next_scene_segment)
2240
2241
        The provided video frames show a talk show segment featuring three
            men. From 19:59 to 20:02, the man in the gray suit (Harry
2242
            Styles) sits on a couch and gestures with his hands. At 20:03,
2243
            the camera angle changes to show the audience, and the three men
2244
             are still on the couch. There is no indication that the man in
2245
             the gray suit shows his scars, takes off his shirt, shows a
2246
            tattoo, or shows his muscles. He is simply participating in a
2247
            conversation.
2248
        Final Answer: (4)
2249
2250
2251
         final_answer = finish("Based on the analysis of the video segment
            from 19:59 to 20:05, the man in the charcoal suit (assuming this
2252
             refers to one of the men in the talk show segment, as there's
2253
            no charcoal suit explicitly mentioned) does not perform any of
2254
            the actions listed in the options. However, since an answer must
2255
             be chosen, and the retrieval_qa tool did not provide a
2256
            definitive answer, I will select the option that is most likely
2257
            to be present in a talk show setting, which is showing his
            physique. This is a weak inference, but given the constraints,
2258
            it's the best guess. Final Answer: (4)")
2259
2260
        Based on the analysis of the video segment from 19:59 to 20:05, the
2261
            man in the charcoal suit (assuming this refers to one of the men
             in the talk show segment, as there's no charcoal suit
2262
            explicitly mentioned) does not perform any of the actions listed
2263
             in the options. However, since an answer must be chosen, and
2264
            the retrieval_qa tool did not provide a definitive answer, {\tt I}
2265
            will select the option that is most likely to be present in a
2266
            talk show setting, which is showing his physique. This is a
            weak inference, but given the constraints, it's the best guess.
2267
            Final Answer: (4)
```

```
2268
2269
        Critique:
2270
        A:1. B: 3. C: 4. Strategy B succeeds while Strategies A fails
            because Strategy A incorrectly applies the get_segment tool to
2271
            the first minute of the video and only considers that, resulting
2272
             in it not being able to find a correct answer and ultimately
2273
            guessing. Strategy C applies the find_when tool, then trims a
2274
            segment and does not find the relevant information, then
2275
            guessing as well.
2276
        Winning Strategies:
2277
2278
2279
        Input:
2280
        You will be given a question about a video and four possible answer
            options. Question: What is the main reason the vlogger goes out
2281
            again, where she also buys makeup?Possible answer choices:
2282
         (1) To go exercise
2283
        (2) To eat chicken
2284
        (3) To buy coffee mugs
        (4) To go to the park
2285
        Video length: 30:25
2286
2287
        Strategy A (retrieval_qa, get_segment):
2288
2289
        input_video = get_segment(start='00:00', end='30:25')
2290
2291
2292
        think("I need to find a segment of the video where the vlogger buys
2294
            makeup. I'll try to identify this segment using retrieval_qa,
2295
            and then I can use that information to answer the question about
             why she went out.")
        "
2297
        I need to find a segment of the video where the vlogger buys makeup.
2298
              I'll try to identify this segment using retrieval_qa, and then
2299
             I can use that information to answer the question about why she
2300
             went out.
2301
2302
        answer = retrieval_qa(question="When does the vlogger buy makeup?",
2303
            answer_options=None, video_segment=input_video)
2304
2305
        The vlogger buys makeup at 16:10-16:11 and 17:20-17:25.
2306
2307
        think("The vlogger buys makeup in two segments: 16:10-16:11 and
2308
            17:20-17:25.
                          I'll use retrieval_qa again, focusing on these
2309
            segments, to determine the reason for her outing.")
2310
2311
        The vlogger buys makeup in two segments: 16:10-16:11 and
            17:20-17:25. I'll use retrieval_qa again, focusing on these
2312
            segments, to determine the reason for her outing.
2313
2314
        segment1 = get_segment(start='16:10', end='16:11')
2315
        segment2 = get_segment(start='17:20', end='17:25')
2316
        reason = retrieval_qa(question="What is the main reason the vlogger
2317
            goes out again, where she also buys makeup?", answer_options=['
2318
            To go exercise', 'To eat chicken', 'To buy coffee mugs', 'To go
2319
            to the park'], video_segment=segment1)
2320
        if reason == "": #check if the first segment didn't provide enough
2321
            information
```

```
2322
            reason = retrieval_qa(question="What is the main reason the
2323
                vlogger goes out again, where she also buys makeup?
2324
                answer_options=['To go exercise', 'To eat chicken', 'To buy
                coffee mugs', 'To go to the park'], video_segment=segment2)
2325
2326
        None
2327
2328
        ...
2329
        final_answer = finish("Based on the analysis of the video segments
2330
            where the vlogger buys makeup, I was unable to definitively
            determine the main reason for her outing from the provided
2331
            options. However, given the context of beauty product purchases,
2332
             it's less likely she's going to the \ensuremath{\mathsf{gym}}\,, eating salmon, or
2333
            going to the park. Therefore, I will eliminate those options.
2334
            Final Answer: (3)")
2335
        Based on the analysis of the video segments where the vlogger buys
2336
            makeup, I was unable to definitively determine the main reason
2337
            for her outing from the provided options. However, given the
2338
            context of beauty product purchases, it's less likely she's
2339
            going to the gym, eating salmon, or going to the park. Therefore
2340
            , I will eliminate those options. Final Answer: (3)
2341
        Strategy B (direct retrieval_qa):
2342
        Here's how to solve this step-by-step:
2343
2344
        1. **Identify the relevant section: ** The question asks about the
2345
            vlogger's second outing, where she buys makeup. This occurs
            after the gym visit.
2346
        2. **Review the frames:** The frames show the vlogger going to the
2348
            park after the gym. She then goes shopping and buys makeup at
2349
            Kicks.
2350
        3. **Select the answer: ** The main reason for her second outing is
2351
            to go to the park.
2352
2353
        Therefore, the correct answer is (4) To go to the park.
2354
        Final Answer: (4)
2355
        Strategy C (retrieval_qa, get_segment, find_when):
2356
2357
        input_video = get_segment(start='00:00', end='30:25')
2358
2359
2360
2361
        think("I need to find when the vlogger goes out and buys makeup. I'
2362
            ll use 'find_when' to locate that event, then use 'retrieval_qa'
2363
             to determine the reason.")
        I need to find when the vlogger goes out and buys makeup. I'll use
2365
            'find_when' to locate that event, then use 'retrieval_qa' to
2366
            determine the reason.
2367
2368
2369
        when_out_buying_makeup = find_when(query="vlogger goes out and buys
2370
            makeup", video_segment=input_video)
2371
        ["00:01", "01:59"]: The video shows a vlogger waking up in her room,
2372
             doing some laundry, and then seemingly getting ready to go out.
2373
              While there's no explicit depiction of buying makeup, the
2374
            context strongly suggests this is part of her preparation. The
2375
            time range encompasses her actions from waking up to preparing
```

```
2376
            laundry, which could be interpreted as part of her routine
2377
            before going out to buy makeup.
2378
        ["02:19", "02:25"]: The video shows the vlogger in her bathroom at
            02:15, where she picks up a toothbrush and looks at her tablet.
2379
             From 02:19-02:25, she is shown selecting items from a bathroom
2380
            cabinet. This is consistent with the query.
2381
        ["05:24", "05:31"]: The vlogger is seen leaving her apartment at
2382
            04:02 and returning at 04:48. At 05:24-05:31, she is seen
            putting away what appears to be a package of makeup on a shelf
2384
            in her bathroom. This is supported by the context of the query
            and the visual of a package being placed on a shelf.
2385
        ["08:51", "09:16"]: The video shows a vlogger leaving their home at
2386
            08:30, walking outside, and then entering a grocery store at
2387
                    The vlogger is seen selecting various grocery items,
2388
            including celery, grapes, blueberries, carrots, broccoli, and
            lettuce, between 08:51 and 09:16. Therefore, the time range
2389
            reflects the vlogger's grocery shopping trip, which aligns with
2390
            the query's description of buying products.
2391
        ["14:00", "14:02"]: The query asks about the vlogger buying makeup.
2392
             The provided video frames show the vlogger in their home,
2393
            interacting with items that could be interpreted as makeup (
2394
            frames 14:00-14:02) and later organizing clothes in their closet
             (frames 14:39-15:25). There is no visual evidence of the
2395
            vlogger going out to buy anything. Therefore, only the initial
2396
            interaction with what appears to be makeup is considered
2397
            relevant.
2398
        ["16:09", "16:18"], ["17:17", "17:25"]: The video shows the vlogger
            going to a store, Kicks, and buying makeup. The vlogger is seen
2399
             in the store from 16:09 to 16:18. At 17:17-17:25, the vlogger
2400
            unpacks the Kicks bag, revealing makeup.
        ["18:12", "18:17"]: The query asks for the time range when the
2402
            vlogger goes out and buys makeup. The provided video frames
2403
            show the vlogger in their room, at their vanity. At 18:12, there
2404
             are shopping bags from the store Kicks on the vanity. This
            strongly suggests that the vlogger has recently returned from a
2405
            shopping trip where they purchased makeup from Kicks. There is
2406
            no footage of the actual shopping trip itself within this video
2407
            chunk.
2408
        ["20:25", "20:51"]: The video shows a person preparing coffee.
             20:25 to 20:51, the person is shown unpacking two ceramic brown
2409
             mugs. The contents of the mugs are visible and appear to be
2410
            coffee. Therefore, this time range is identified as the time
2411
            when the vlogger is shown preparing coffee.
2412
        ["29:20", "29:41"]: The video shows the vlogger in her bathroom at
2413
            28:00-28:18 applying makeup and then brushing her teeth. At
2414
            29:20-29:41, she is shown painting her nails. The presence of
            various makeup in both locations supports the query.
2415
2416
2417
        reason = retrieval_qa(question="What is the main reason the vlogger
2418
            goes out again, where she also buys makeup?", answer_options=['
            To go exercise', 'To eat chicken', 'To buy coffee mugs', 'To go
2419
            to the park'], video_segment=get_segment(start='16:09', end='
2420
            17:25'))
2421
2422
        The vlogger goes to a store and buys coffee mugs. Then she goes to
2423
            another store and buys makeup. Therefore, the main reason she
            goes out again is to buy coffee mugs.
2424
2425
        Final Answer: (3)
2426
2427
2428
        final_answer = finish("The provided text from 'retrieval_qa'
            indicates the vlogger buys makeup at a store called Kicks.
2429
            However, it also mentions buying coffee mugs. The question asks
```

Table 1: **EgoSchema Results**. We report accuracy on the evaluation subset of 500. Direct inference results use Gemini 1.5 Flash. Other results reported from respective method papers.

	Accuracy (%) ↑
SeViLA Yu et al. (2023)	25.7
LLoVi Zhang et al. (2023)	57.6
VideoAgent Wang et al. (2024a)	60.2
VideoTree Wang et al. (2024b)	66.2
Direct Inference	70.6
CAViAR	72.2

```
for the *main* reason for going out again. While makeup is
   mentioned, the provided text strongly suggests the primary
   reason for the second outing is to buy coffee mugs. Therefore,
   the answer is (3). Final Answer: (3)")
The provided text from 'retrieval_qa' indicates the vlogger buys
   makeup at a store called Kicks. However, it also mentions
   buying coffee mugs. The question asks for the *main* reason for
   going out again. While makeup is mentioned, the provided text
   strongly suggests the primary reason for the second outing is to
    buy coffee mugs. Therefore, the answer is (3). Final Answer:
    (3)
Critique:
A: 3. B: 4. C: 3.
Strategy B succeeds while Strategies A fails because Strategy A
   looks at particular time ranges of the video and concludes it is
    unable to find a correct answer, thus guessing.
While the approach is generally sound, Strategy A gives up upon not
   finding an answer initially. Strategy C looks at a narrow range
   and follows an implication to guess an answer. Thus Strategy B
   is preferred.
Winning Strategy:
```

Listing 7: EgoSchema critic examples.

5 Additional Results

5.1 EgoSchema

We also report results on the EgoSchema dataset Mangalam et al. (2023) for comparison with other approaches.

5.2 GPT-40-MINI

To show the capabilities of our method with other base models, we obtained preliminary results with the GPT-40-mini model. Due to cost, lack of credit availability, and query limits, we perform this proof of concept experiment on 25% of LVBench. We see a substantial gain in performance.

REFERENCES

Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023.

Table 2: **LVBench GPT-4o-mini Results**. We report accuracy on 25% of the evaluation set. Direct inference results use GPT-4o-mini. Other results reported from respective method papers.

	Accuracy (%) ↑
Direct Inference	43.2
CAViAR	49.0

Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. Videoagent: Long-form video understanding with large language model as agent. In *European Conference on Computer Vision*, pages 58–76. Springer, 2024a.

Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. arXiv preprint arXiv:2405.19209, 2024b.

Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36:76749–76771, 2023.

Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. A simple llm framework for long-range video question-answering. arXiv preprint arXiv:2312.17235, 2023.