

Enhancing Video-LLM Reasoning via Agent-of-Thoughts Distillation

Appendix

A. Limitation

Despite the advancements mentioned in the paper, several limitations remain and we leave them as future work: (i) similar to prior approaches [7, 10, 16], the effectiveness of our agent-based system is contingent upon the progress of the underlying visual model components. Enhancing its ability to generalize across diverse datasets is essential for broader applicability; (ii) while our primary focus has been on compositional VideoQA tasks [42], and we have demonstrated improvements across a series of benchmarks, achieving holistic enhancements will require further exploration into creating a more balanced distribution of training data; (iii) furthermore, our agent-based framework has the potential to address additional video-related tasks, such as video captioning and referring segmentation. We aim to expand our methodology to these domains, which could yield even more robust and versatile applications in the future.

B. Experimental Details

B.1. Training Details

For all models, their projection layers and language model are fine-tuned and visual encoder is frozen. We use a cosine learning rate schedule, with warm up ratio 0.03 and learning rate $4e-5$. For both Instruct and AoTD setting, we fine-tune the model with batch size 48 and totally 1 epoch. We believe that longer training will get a better performance on in-domain benchmarks but maybe a destroy on out-of-domain benchmarks.

B.2. Specialized Models Evaluation Details

In this section we will show the details about each sub-task’s evaluation from data preparation to evaluation metric.

Question decomposition. Since there may be multiple valid ways to decompose the same problem, we evaluate only the accuracy of the final output in this sub-task. Specifically, the model takes the query and instruction as input and generates an executable program. We replace all intermediate outputs within the program and focus on whether the final output matches the correct answer. If the decomposition is correct, the final output must align with the answer. Any programs that cannot be executed or that lead to an incorrect answer are considered failures.

Object detection. To evaluate the performance of detection models, we sample frames with scene graph annotations from the input video clip and provide them, along with the text query, as input to the model. The model then outputs a series of bounding boxes that exceed a confidence threshold. We select the bounding box with the highest confidence as the final output and calculate the IoU to assess accuracy.

Temporal grounding. Since scene graphs provide both the start and end frame IDs, as well as key frame IDs for each event, we use IoU and Recall as metrics to capture different aspects of model performance. The model takes the video clip and text query as input and outputs the predicted start and end frame IDs. We calculate IoU based on the alignment between the predicted and annotated start and end frame IDs, and we compute Recall using the key frame ID annotations to evaluate how well the model captures important frames.

Action recognition. For discriminative models, we provide the video clip and a list of action labels as input to complete a classification task. For generative models, we provide the video clip along with an instruction prompt, asking the model to generate five actions most relevant to the video, ranked by likelihood. We then use the top-ranked output from each model to calculate the Top-1 accuracy for both approaches.

Question answering. The evaluation of question answering follows a similar approach to previous methods. The model takes the video clip and question as input and returns an answer, from which we directly calculate the accuracy. The key difference between this sub-task and a standard QA task is that the answers are based on a series of information collected by preceding agents, allowing for a more accurate assessment of the model’s pure question-answering ability.

C. More Results

Here, we introduce some examples to show the process from query to Chain-of-Thought using our agent-based system. We can find that our system is able to decompose complex questions into easier sub-tasks and the final CoT retains step-by-step

problem-solving ideas and spatial-temporal information representing video understanding ability.



Question: Why does the man have to throw the plane first in the middle of the video?

Possible answers: (A) chase the ball (B) the toy car hit the wall (C) put on earmuffs (D) get it start flying (E) bounce baby on cushion

Answer: (D) get it start flying

Question

```
def execute_command(video_clip, query, possible_answers):
    middle_clip = trim(video_clip, start=video_clip.num_frames//4, end=video_clip.num_frames//4*3)
    throw_clip = Filter_frames_with_act(middle_clip, 'man throw the plane')
    throw_captions = Video_caption(throw_clip)
    man_bboxes = Find(throw_clip, 'man')
    throw_reason = Video_query(throw_clip, query='why does the man have to throw the plane first',
                               possible_answers=possible_answers)
    info = {'reason why the man has to throw the plane first': throw_reason, 'what happened in the target
                                                    video': throw_captions}
    answer = select_answer(query, info, possible_answers)
    return answer
```

Program

call function trim
 trimmed video from frame 8 to frame 24
 call Filter_frames_with_act
 filter action man throw the plane
 find action from frame 16 to frame 20
 call function Video_caption
 frame 16: A man and a child playing plane in a park.
 frame 17: A man and a child standing in a grassy field.
 frame 18: A man and a child standing in a grassy field.
 frame 19: A man and a child playing plane in a park.
 call function Find
 finding man
 find man at [283, 45, 330, 247] in frame 16
 find man at [215, 117, 249, 263] in frame 17
 find man at [189, 122, 220, 266] in frame 18
 find man at [201, 132, 236, 277] in frame 19
 call function Video_query
 Question: why does the man have to throw the plane first
 Answer: The man has to throw the plane first because it is a remote control airplane and needs to be launched into the air
 before the young boy can start playing with it.
 call function select_answer
 select: D
 program output: get it start flying

Execution Trace

To solve this question, we first need to find when the man is throwing the plane in the middle of the video. The man is throwing the plane from frame 16 to frame 20.
 Next, we should find where the man is to better know why he has to throw the plane first. We can do this by analyzing the man's bounding box region in each frame. The man's bounding box is [283, 45, 330, 247] in frame 16, [215, 117, 249, 263] in frame 17, [189, 122, 220, 266] in frame 18, and [201, 132, 236, 277] in frame 19.
 By analyzing the man's bounding box region, we can see that the man is playing with a remote control airplane and needs to launch it into the air before the young boy can start playing with it. This is because the airplane needs to be in the air for the boy to play with it.
 Therefore, the answer is that the man has to throw the plane first because it is a remote control airplane and needs to be launched into the air before the young boy can start playing with it. So the answer should be get it start flying.

CoT

Figure 5. Example form NExT-QA [43]

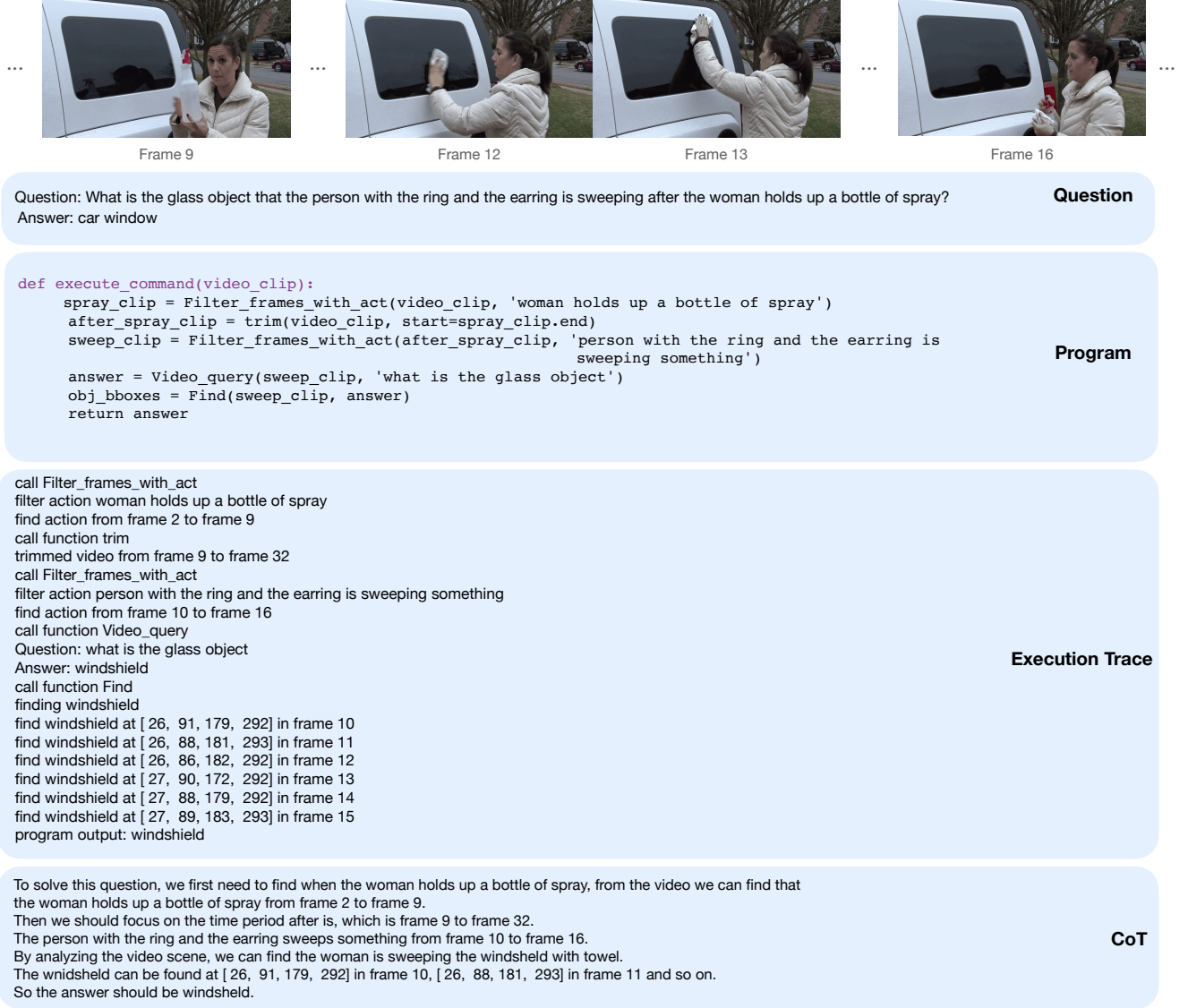


Figure 6. Example from ANetQA [51]

D. Prompts

In this section we present the prompts used in our agent-based system for generating programs, converting execution traces and filtering CoTs.

D.1. Prompt for Program Generation

For each video and query, we call a code LLM to decompose the query to a Python program under the guidance of the prompt below. We modify the ViperGPT [37] prompt to adapt to the visual agents we use.

```
1 def Query_Objs(clip, query):
2     """
3     Query the objects that appear in video clip and match the query descriptions.
4     Parameters
5     -----
6     clip:
7         a list of video frames.
8     query:
9         Description of the target object.
10    Returns
```

```

11  -----
12  a list of bounding boxes of the objects that match the query.
13  Examples
14  -----
15  #return white_objs
16  def execute_command(video_clip):
17      white_objs = Query_Objs(video_clip, "white object")
18      return white_objs
19  """
20
21  def Query_Actions(clip, obj=None):
22      """
23      Find the actions happened in the video clip, if obj is not None, query the actions related to it.
24      Parameters
25      -----
26      clip:
27          a list of the video frames.
28      obj:
29          object class which is used to query the actions related to it.
30      Returns
31      -----
32      a list of actions classes happened in the video clip.
33      Examples
34      -----
35      #return actions
36      def execute_command(video_clip, query, possible_answers):
37          actions = Query_Actions(video_clip)
38          return actions
39      """
40
41  def Filter_frames_with_act(clip, action):
42      """
43      filter a new video clip containing the time period in which the target action occurred
44      Parameters
45      -----
46      clip:
47          a list of video frames.
48      action:
49          the target action which is used to filter frames.
50      Returns
51      -----
52      a new video clip containing the time period in which the target action occurred.
53      Examples
54      -----
55      #return jump_clip
56      def execute_command(video_clip, query, possible_answers):
57          jump_clip = Filter_frames_with_act(video_clip, "person is jumping")
58          return jump_clip
59      """
60
61  def Filter_frames_with_obj(clip, obj):
62      """
63      filter a new video clip that the target object occurred.
64      Parameters
65      -----
66      clip:
67          a list of video frames.
68      obj:
69          class or description about the target object.
70      Returns
71      -----
72      a new video clip that the target object occurred in it.
73      Examples
74      -----
75      #return shoe_clip
76      def execute_command(video_clip, query, possible_answers):
77          shoe_clip = Filter_frames_with_obj(video_clip, "shoe")
78          return shoe_clip
79      """
80
81  def trim(clip, start=None, end=None):
82      """
83      Returns a new video clip containing a trimmed version of the original video at the [start, end] clip.
84      Parameters
85      -----
86      clip:
87          a list of video frames.
88      start : Union[int, None]
89          An int describing the starting frame in this video clip with respect to the original video.

```

```

90     end : Union[int, None]
91         An int describing the ending frame in this video clip with respect to the original video.
92
93     Returns
94     -----
95     a new video clip with start and end.
96     """
97 def Find(clip, obj):
98     """
99     find all bounding boxes around a certain object in the video clip,
100     and collates them into a collection of frames.
101     Parameters
102     -----
103     clip:
104         a list of video frames.
105     obj:
106         the object to look for.
107     Returns
108     -----
109     a new video clip composed of crops of the object.
110     Examples
111     -----
112     # Return the shoe_clip
113     def execute_command(video_clip, query, possible_answers):
114         shoe_clip = Find(video_clip, "shoe")
115         return shoe_clip
116     """
117
118 def select_answer(query, info, possible_answers):
119     """
120     Uses a language model to choose the option that best answers the question given the input information.
121     Parameters
122     -----
123     query:
124         the input question.
125     info:
126         Any useful information to answer the question.
127     possible_answers:
128         a list of possible answers to the question.
129     Returns
130     -----
131     one answer chosen from the possible answers.
132     Examples
133     -----
134     # Return the answer
135     def execute_command(video_clip, query, possible_answers):
136         clip_summary = Video_summary(video_clip)
137         info = {
138             "summary of the target video": clip_summary
139         }
140         answer = select_answer(query, info, possible_answers)
141         return answer
142     """
143 def exist(clip, query):
144     """
145     judge whether a object exists in the video.
146     Parameters
147     -----
148     clip:
149         a list of video frames.
150     query:
151         query to the object class.
152     Returns
153     -----
154     Return True if the object specified by query is found in the video, and False otherwise.
155     Examples
156     -----
157     # Return the flag
158     def execute_command(video_clip, query, possible_answers):
159         flag = exist(video_clip, "shoe")
160         return flag
161     """
162 def Video_summary(clip, query):
163     """
164     give a brief summary of the video clip related to the query.
165     Parameters
166     -----
167     clip:
168         a list of video frames.

```

```

169     query:
170         a question about the video.
171     Returns
172     -----
173     return a brief summary of the video clip.
174     Examples
175     -----
176     # Return the clip_summary
177     def execute_command(video_clip, query, possible_answers):
178         clip_summary = Video_summary(video_clip, query)
179         return clip_summary
180     """
181 Write a function using Python and the functions (above) that could be executed to provide an answer to the query.
182
183 Consider the following guidelines:
184 - Use base Python (comparison, sorting) for basic logical operations, start/end, math, etc.
185 - Objects with mutiple names like "phone/camera", "cup/glass/bottle" with slash, input them as a whole object name.
186 - Just use the class and function appear above except for some base python operations.
187 - Only answer with a function starting def execute_command, do not answer any extra words and symbols before and after
  the function.
188 - No text that is not related to function can appear.
189 - the answer only begins with "def execute_command" and ends with "return answer".
190
191 Here are some examples of the function you should write:
192 -----
193 question: What else is the person able to do with the door?
194 possible answers: ["Hold the door.", "Put down the door.", "Close the door.", "Open the door."]
195 def execute_command(video_clip, query, possible_answers):
196     door_clip = Filter_frames_with_obj(video_clip, "door")
197     person_clip = Find(door_clip, "person")
198     clip_summary = Video_summary(person_clip, query)
199     door_actions = Query_Actions(person_clip, "door", possible_answers=possible_answers)
200     door_actions =
201     info = {
202         "actions the person able to do with the door else": door_actions,
203         "summary of the target video": clip_summary
204     }
205     answer = select_answer(query, info, possible_answers)
206     return answer
207 -----
208 Query: INSERT_QUERY_HERE
209 possible answers: INSERT_POSSIBLE_ANSWERS_HERE

```

D.2. Prompt for Execution Trace Conversion

After getting the execution trace by running the program step by step, we use a LLM to convert the trace into a natural language CoT. The LLM takes query, execution trace, possible answers (in MC-VQA) and execution trace as input. The instruction prompt is as follow:

```

1 Given a video and a question, I wrote the function execute_command using Python, and the other functions above that
  could be executed to provide an answer to the query.
2 As shown in the code, the code will print execution traces.
3 I need you to rewrite the execution trace into a natural language rationale that leads to the answer.
4
5 Consider the following guidelines:
6 - Use all the bounding box information in the rationale, do not use words like "so on" to omit the bounding box, just
  write all of them into the rationale.
7 - Referencing the execution trace, write a reasoning chain that leads to the most common human answer. Notice that the
  output should be the same as the human answer, not necessarily the program output.
8 - If some part of the rationale lacks logic, add reasonable content to make it logical.
9
10
11 Here are some examples of the rantionale you should write:
12 -----
13 Question: What did the person do with the table?
14 def execute_command(video_clip, query, possible_answers, time_wait_between_lines, syntax):
15     table_clip = Filter_frames_with_act(video_clip, 'person interacting with table')
16     person_clip = Find(table_clip, 'person')
17     table_bboxes = Find(table_clip, 'table')
18     clip_summary = Video_summary(person_clip)
19     person_action = Query_Actions(person_clip, 'table', possible_answers=possible_answers)
20     info = {'actions the person do with the table': person_action, 'summary of the target video': clip_summary}
21     answer = select_answer(query, info, possible_answers)
22     return answer
23 Execution trace:
24 call Filter_frames_with_act
25 filter action person interacting with table
26 find action from frame 2 to frame 11

```

```

27 call function Find
28 finding person
29 find person at [139, 141, 229, 342] in frame 2
30 find person at [151, 123, 242, 349] in frame 3
31 find person at [153, 121, 242, 274] in frame 4
32 find person at [158, 123, 255, 261] in frame 5
33 find person at [163, 124, 270, 262] in frame 6
34 find person at [153, 121, 242, 351] in frame 7
35 find person at [95, 113, 196, 316] in frame 8
36 find person at [83, 113, 196, 285] in frame 9
37 find person at [112, 116, 201, 332] in frame 10
38 call function Find
39 finding table
40 find table at [183, 140, 269, 257] in frame 2
41 find table at [194, 131, 269, 255] in frame 3
42 find table at [227, 129, 269, 252] in frame 4
43 find table at [226, 165, 269, 258] in frame 5
44 find table at [233, 170, 270, 259] in frame 6
45 find table at [217, 129, 269, 256] in frame 7
46 find table at [217, 122, 270, 254] in frame 8
47 find table at [221, 123, 269, 256] in frame 9
48 find table at [225, 125, 270, 263] in frame 10
49 call function Video_summary
50 summary result: The video shows a man in a kitchen, bending over and holding an orange object, surrounded by various
  kitchen items and furniture, with a focus on his actions and the domestic setting.
51 call function Query_Actions
52 Query table
53 Answer: tidied up.
54 call function select_answer
55 the information used: - actions the person do with the table: tidied up.
56 - summary of the target video: The video shows a man in a kitchen, bending over and holding an orange object,
  surrounded by various kitchen items and furniture, with a focus on his actions and the domestic setting.
57 program output: Tidied up.
58 Rationale:
59 To solve this question, we first have to find when did the person interact with the table.
60 From the video, we can see that the person is interacting with the table from frame 2 to frame 11.
61 In this time period, we can find person at [139, 141, 229, 342] in frame 2, [151, 123, 242, 349] in frame 3, [153,
  121, 242, 274] in frame 4 and so on.
62 Table can also be found at [183, 140, 269, 257] in frame 2, [194, 131, 269, 255] in frame 3, [227, 129, 269, 252] in
  frame 4 and so on.
63 By analyzing the person and table bounding box region, we can see that the person is holding an orange object to clean
  the table in the kirchen environment.
64 So the answer should be tidied up.
65 -----
66 Now, look the question, program and execution trace, please transfer these information to a rantionale.
67 Question: INSERT_QUESTION_HERE
68 INSERT_PROGRAM_HERE
69 Execution trace:
70 INSERT_EXECUTION_TRACE_HERE
71 Rationale:

```

D.3. Prompt for CoT Filtering

In order to obtain high quality distillation data, we continue using LLM to filter CoTs. We prompt the LLM to select those CoTs that are truly helpful for solving questions and reflect the step-by-step thinking process. The prompt is as follows:

```

1 I will give you a question and a rationale to solve the question, you need to judge whether the rationale is thinking
  step by step and helpful to solve the question.
2 If yes, return True, If not, return False. no need to explain.
3 Here is the question and rationale:
4 Question: INSERT_QUESTION_HERE
5 Rationale: INSERT_RATIONALE_HERE

```

D.4. Prompt for Inference

Question: question content
 Answer in one word or phrase. / Explain the rationale to answer the question.

Question: question content

Options:

(A) option content

(B) option content

(C) option content

(D) option content

Answer with the option’s letter from the given choices directly and only give the best option. / Explain the rationale to answer the question.

E. More related works

E.1. Video-language models (Video-LLMs).

Most existing Video-LLMs [4, 44, 55] are composed of a pre-trained visual encoder (like CLIP [34] or SigLIP [53]) to encode video frames into a sequence of visual features, an adapter to transfer the visual features to tokens that can be understood by the language model, and a pretrained LLM to output the final response.

E.2. Visual Programming and Agents.

Recent work like MoReVQA [27] proposes a multi-stage system, getting a strong zero-shot VideoQA ability while is able to create interpretable intermediate outputs. VURF [25] proposes a self-refinement method to resolve the LLM hallucinations to get a more concise program based on the context cues.

F. More Ablations

F.1. Analysis on latency and computation.

We compare the agent-based system with the distilled model (LNV-AoTD) across three key metrics. As shown in Table 8, LNV-AoTD outperforms the agent-based system in both inference latency and computational efficiency, demonstrating the necessity of distillation process.

Model	Time (s) ↓	Memory (GB) ↓	TFLOPs ↓
Agent-based system	47.93	65.34	30.98
LNV-AoTD	10.58	18.21	13.53

Table 8. Latency and computation comparison.

F.2. Analysis on more models

We further experiment on VideoLLaMA2[4] and Qwen2-VL[38] using identical settings. The results confirm AoTD’s effectiveness and superior performance.

Model	Filtering	MVBench (Acc.)	STAR (Acc.)	AGQA (Acc. / Score)
LNV-AoTD	✗	53.7	73.3	59.5/3.5
LNV-AoTD	✓	55.6	74.3	60.9/3.6
VideoLLaMA2-Instruct	-	54.9	69.2	56.0/3.5
VideoLLaMA2-AoTD	✓	56.0	71.1	57.2/3.5
Qwen2-VL-Instruct	-	65.6	71.4	59.8/3.6
Qwen2-VL-AoTD	✓	66.5	73.1	61.2/3.7

Table 9. Ablation results on more models.