
Supplementary material for ViPR

Anonymous Author(s)

Affiliation

Address

email

1 Running the code

To reproduce the results from the paper, follow the steps below.

1. Make sure you use Python 3.9.18.
2. Download ViCLIP-L_InternVid-FLT-10M.pth from <https://huggingface.co/OpenGVLab/ViCLIP/tree/main> and put it into .cache/encoders/viclip
3. Put your OpenAI API key into .env, under the key OPENAI_API_KEY=<your_key>
4. Create a virtual environment, install the dependencies with `pip install --editable .`
5. Download the dataset files to a local directory **TODO: from where?** and note the directory of the .../videos and .../tasks directories
6. Set the video dir and task dir in the config files to the respective directories above. For example, in configs/alfred_replication.yaml you should set task_dir: "[.../tasks/alfred]"
7. Run the evaluation on the selected environment by running `vlm evaluate configs/[alfred_replication,minecraft_replication,real_life_replication].yaml`
8. After the run is complete (which can take upwards of 12 hours), you can create plots from the logs by running `vlm plot --experiment-dir experiments --task-dir [the directory]`

Note: to reproduce Minecraft experiments, you should rename src folder to tmp_src, and minecraft_src to src, since those use slightly different codebases (the difference is mostly about the usage of GPT-4o). To experiment with alfred and real life, you should rename folders back.

2 Dataset details

The evaluation code we use works with the dataset being separated into mp4 video files, yaml task files, and json task data files, which is the structure we detail below. We also supply the dataset.json file which lists all metadata for all videos in one file, to make it easier to use the dataset with your own evaluation code.

2.1 General data structure

Each of the three environments in the dataset has two parts:

- Videos: mp4 video clips organised into sub-directories. By convention, the directory structure corresponds to the three environments we have, and within them, to the different difficulty levels and task types (foundation/base, permutation, remix). The specific paths serve as identifiers when matching videos to classes in classification problems.

32 • Classification problems: yaml and json files organised into sub-directories. For a given
 33 problem, the yaml files contains its definition (detailed below), and the json file references
 34 to the videos and their classes. The path of the yaml file serves as identifier for the problem
 35 within the evaluation code.

36 **Problem definition format** The yaml file has a top-level dictionary with the following keys:

- 37 • *label_prompts*: a dictionary, where keys are class ids (labels) and their associated values are
 38 natural language descriptions which are given to the models
- 39 • *prompt_gpt*: a 0-shot prompt describing the task, given to GPT-4o during evaluation
- 40 • *example_gpt*: a prompt containing environment-specific task hints (e.g. what holding looks
 41 like, common mistakes, etc.), given to GPT-4o during 1-shot evaluation
- 42 • *metadata*: optional, a dictionary that can contain any metadata associated with the task. The
 43 metadata are not currently used by the code in any way.

44 **Problem data format** For problem *[problem].yaml*, the json file must be named *[prob-*
 45 *lem]_data.json*. The json file is a list of dictionaries with the following two keys:

- 46 • *path*: the path to a video, relative to the *video_path* given in the experiment configuration.
- 47 • *label*: the class id (label) of the video. This has to match one of the class ids in the
 48 *label_prompts* dictionary in the yaml file.

49 2.2 Virtual home video acquisition protocols

50 This section details the design decisions and preparation protocol behind the virtual home videos.
 51 The code to generate the virtual home videos is in the `generate-virtual-home` directory.

52 In general, the videos were created by cutting and pasting frames from different ALFRED trajectories.
 53 We were able to automatize this process because ALFRED contains step-by-step descriptions of
 54 actions in each trajectory, together with a list of frames each action corresponds to.

55 **Base problems** Base problem videos were created by preprocessing videos from the ALFRED
 56 dataset, to isolate individual actions. When possible, we include a few surrounding frames for context
 57 — but, we only do this if the frames do not contain a different action (i.e., they are only included if
 58 they only contain movement). When we say we created n tasks, each with an action being performed
 59 on a different object, we selected the objects to be included at random.

60 Specifically, we have the following problem types for object recognition:

- 61 • *Object recognition*: We isolate all PickupObject actions and their associated frames and
 62 descriptions from the original ALFRED videos, for each object keeping 10 (if possible)
 63 videos where it’s being picked up from a single type of container. Finally, we group the
 64 videos based on where the object was picked up from and create classification problems
 65 from those groups.

66 The goal is to have problems with videos that all show the agent picking up an object from
 67 one type of container, so that the type of object is the only major difference between the
 68 videos within the problem.

69 We only keep problems that have more than one class, ending up with “pickup from
 70 countertop”, “pickup from dining table”, and “pickup from somewhere” (i.e. instances
 71 where the container was not specified in the ALFRED dataset).

- 72 • *Container recognition*: We perform a similar selection process to the one above. We isolate
 73 all PutObject actions and their associated frames and descriptions from the original ALFRED
 74 videos, for each container keeping 10 (if possible) videos where a certain object is being
 75 picked up from it. Finally, we group the videos them based on which object was picked up
 76 and create classification problems from those groups.

77 We end up with four problems: “put butter knife”, “put keychain”, “put mug”, and “put soap”,
 78 each problem containing videos of the respective object being picked up from different
 79 containers.

80 We also have the following ones for action understanding:

- 81 • *Cleaning*: We isolate all “clean” actions and their associated frames and descriptions from
82 the original ALFRED videos, and then drop some frames to create the following alternative
83 videos: “we clean the [object] in the sink basin under running water” (original video), “we
84 go to the sink basin, holding [object] in hand, we don’t put it in the sink basin”, “we put the
85 [object] in the sink basin and pick it back up without running water over it”.
86 We create 5 such problems, corresponding to the agent cleaning 5 different objects, with 10
87 videos (if possible) per class in each problem.
- 88 • *Cooling*: Similarly to Cleaning, we isolate all “cool” actions from ALFRED, and cut some
89 frames to end up with the following alternatives: “we cool the [object] by putting it in the
90 fridge for a while” (original video), “we open the fridge, put in the [object] and immediately
91 pick it back up without leaving it to cool in a closed fridge”, “we open the fridge and then
92 immediately close it without putting the [object] in”, and “we go to the fridge and then we
93 leave, without even opening it”.
94 Again, we create 5 such problems, corresponding to the agent cooling 5 different objects,
95 with 10 videos (if possible) per class in each problem.
- 96 • *Heating*: Similarly to Cleaning, we isolate all “heat” actions from ALFRED, and cut some
97 frames to end up with the following alternatives: “we heat the [object] in the microwave”,
98 “we go to the microwave and then immediately leave without putting the [object] in”, “we
99 open the microwave and immediately close it again, without putting the [object] in”, “we
100 open the microwave, put the [object] in, then immediately pick it back up without heating
101 it”, “we put the [object] in the microwave for a while, we do not turn the microwave on, then
102 we pick up the [object] back again”.
103 Again, we create 5 such problems, corresponding to the agent heating 5 different objects,
104 with 10 videos (if possible) per class in each problem.
- 105 • *Picking*: For 5 objects, we isolate 10 (is possible) PickupObject and PutObject videos that
106 have the agent interact with (pick from or put to) the same container. The problem then is
107 for the VLM to decide whether the object has been picked up or put down, since that is the
108 single largest difference between the videos in a given problem.
- 109 • *Slicing*: Similarly to the above, we isolate all “slice” actions from ALFRED together with
110 their frames and descriptions, and then create alternative versions of these videos by cutting
111 the last part where the slicing actually happens out. Thus, in the alternative videos, the agent
112 holds a knife and walks towards an object, preparing to slice it, but does not perform the
113 slicing.
114 Again, we create 5 such problems, corresponding to the agent (not) slicing 5 different
115 objects, with 10 videos (if possible) per class in each problem.
- 116 • *Toggling*: There is no ToggleOff action in ALFRED, so to simulate it, we reverse the frames
117 from a ToggleObjectOn. Specifically, we isolate 10 toggle on and toggle off videos (if
118 possible) for a desk lamp, floor lamp, faucet, and microwave (creating a separate problem
119 for each object), since those are the only objects that can have a turned-on state in ALFRED.

120 And finally, we have two problem types for object state recognition:

- 121 • *On v. off*: We isolate ToggleObjectOn actions from ALFRED, and cut their corresponding
122 videoclips to halves — the first has the agent approach the object, which is turned off for
123 the entire duration of the clip, and the second already shows the object when it’s been turned
124 on. The problems then ask the VLM to say whether the object we see is turned on or off.
125 As above, we isolate 10 turned on and turned off videos (if possible) for a desk lamp, floor
126 lamp, faucet, and microwave (creating a separate problem for each object), since those are
127 the only objects that can have a turned-on state in ALFRED.
- 128 • *Sliced v. whole*: We isolate PickupObject actions from ALFRED, and group together videos
129 that have the agent pick up either a slice of something, or the whole thing. We create 5
130 such problems, corresponding to the agent picking up 5 different objects, with 10 videos (if
131 possible) per class (sliced and whole) in each problem.

```

for prefix_len in base_length, ..., 0 do
  for rest in find_rest(base[:prefix_len]) do
    if we have enough videos for this prefix length then
      break
    end if
    if rest[0] has already been taken after this exact prefix then
      continue
    end if
    trajectory ← base[:prefix_len] + rest
    add_to_trajectories(trajectory)
  end for
end for

```

132 **Remix problems** Remix problems are generated iteratively step-by-step from a base trajectory:

133 The find_rest procedure performs DFS; two actions (and their corresponding frames) can come after
 134 each other if they are set in the same scene (room type), if the content of agent’s hands is the same,
 135 and if there exists a sequence of frames somewhere in ALFRED that shows the agent going from the
 136 place where the first action happened to the place of the second action. We deduce the place from
 137 action metadata in ALFRED.

138 What “having enough videos” means is different for different difficulty levels and prefix lengths;
 139 below we list the number of videos we aim to generate:

- 140 • Level 2: 4 for prefix 0, 4 for prefix 1
- 141 • Level 3: 2 for prefix 0, 3 for all larger prefixes
- 142 • Level 4: 2 for prefix 0, 2 for all larger prefixes
- 143 • Level 5: 4 for prefix 0, 1 for all larger prefixes
- 144 • Level 6: 3 for prefix 0, 1 for all larger prefixes
- 145 • Level 7: 2 for prefix 0, 1 for all larger prefixes
- 146 • Level 8: 1 for prefix 0, 1 for all larger prefixes

147 We generate 12 such problems per level. The base trajectory is selected randomly without replacement
 148 each time. If it is impossible to generate a 9-class problem for a given base trajectory, we try it with
 149 the next base trajectory.

150 **Permutation problems** Permutation problems are generated iteratively from a base trajectory:
 151 we loop through all permutations, keeping those in which all pairs of neighbouring actions can be
 152 connected by frames (coming from other ALFRED videos) that show the agent going from the place
 153 where to first action happened to the place of the second action. In the end, we only keep the first 3
 154 permutations (i.e. classes) of those we found.

155 If possible, we also prefer permutations that are “slicing-consistent”, i.e. those in which all actions
 156 (like moving, heating, ...) with a whole object O come in sequence before it is sliced, and all
 157 actions with a slice of O come after it was sliced. A permutation that is not slicing-consistent is still
 158 admissible and strictly speaking logically consistent — consider that there could be multiple copies
 159 of a given object, for example, and one can thus be sliced while the other stays whole. On some
 160 levels, we do include slicing-inconsistent trajectories because we cannot generate enough consistent
 161 ones.

162 We generate 33 problems per level. The base trajectory is selected randomly without replacement
 163 each time. If it is impossible to generate a 3-class problem for a given base trajectory, we try it with
 164 the next base trajectory (except for ones on level 2).

165 2.3 Minecraft video acquisition protocols

166 To reiterate, we identify 7 fundamental actions in Minecraft:

- 167 1. Place blocks
- 168 2. Break blocks
- 169 3. Craft
- 170 4. Combat
- 171 5. Find something
- 172 6. Pick up an item, dropped after breaking a block / defeating an enemy or animal
- 173 7. Mine stone-like blocks with a pickaxe

174 **Base problems.** For base problems, several videos were selected from BASALT Benchmark and
 175 several videos were manually recorded. Then, we cut short clips (3 or 10 seconds) containing
 176 examples of fundamental actions. There are 22 3-second and 32 10-second clips.

177 **Multi-step problems.** For those problems, only manually recorded videos were used. Again, the
 178 clips with desired patterns (like "break-place" or "craft-mine-place") were cut. The length is mostly
 179 between 30 and 60 seconds. We only collect trajectories of 2 and 3 steps, which totals to 21 unique
 180 clips.

181 In both cases, we used LossLessCut to cut the clips and OBS Studio to record the videos.

182 2.4 Used dataset licenses

183 See table 1 for an overview of datasets we used to construct our dataset and their licenses.

Dataset	License
ALFRED	MIT
BASALT Benchmark	MIT
Kinetics-700	CC 4.0

Table 1: Overview of used datasets and their licenses.

184 3 Model details

185 Models get equally-spaced frames from a video; the time elapsed between frames thus depends on
 186 video length, which is non-uniform. We make sure that the last frame the model receives is also the
 187 very last frame of the video, by duplicating the last frame of the video before subsampling it. The
 188 complete subsampling logic is in figure 1.

189 3.1 Model specification

190 We used three models in our final evaluation:

- 191 1. OpenGVLAB/ViCLIP (available on HuggingFace)
- 192 2. laion/clip-vit-big-14-laion2b-39b-b160k (available on HuggingFace)
- 193 3. GPT-4o (version from 1. 6. 2024)

194 3.2 Compute

195 We ran experiments on the following hardware. Our specific choices were based on availability only;
 196 by varying batch sizes, all of the experiments could be run on any of these.

- 197 • Compute cluster, with NVIDIA A100 GPU 80GB
- 198 • Rented cloud machines (vast.ai), with RTX 3090
- 199 • locally, on the M1 Pro processor

200 Varying the batch sizes does not influence the results.

```

def subsample(x: t.Tensor, n_frames: int) -> t.Tensor:
    total_frames, *_ = x.shape

    if (total_frames - n_frames) % (n_frames - 1) != 0:
        # Replicate the last frame to make sure it will be selected
        last_frame = einops.repeat(
            x[-1], "... -> n ...", n=(n_frames - (total_frames % (n_frames - 1)))
        )
        x = t.cat([x, last_frame])
        total_frames, *_ = x.shape
        assert (total_frames - n_frames) % (n_frames - 1) == 0
    step = (total_frames - n_frames) // (n_frames - 1) + 1
    x_subsampled = x[::step]
    assert len(x_subsampled) == n_frames
    assert (x[0] == x_subsampled[0]).all() and (x[-1] == x_subsampled[-1]).all()
    return x_subsampled

```

Figure 1: The frame subsampling logic, where x is the video in shape (time, w, h, c) and n_frames is the model-dependent number of frames to subsample to.

201 3.3 GPT-4o prompts

202 We evaluate GPT-4o in a few-shot manner, using the following interaction scheme:

- 203 1. First conversation, to get frame-by-frame descriptions:
 - 204 (a) First and only message: system prompt (figure 2) + frame-by-frame prompt (figure 3)
- 205 2. In a second conversation, separate from the first one:
 - 206 (a) First message: The same system prompt (figure 2) + class match prompt (figure 4)
 - 207 (b) Second message, same conversation, after GPT-4o replies: scoring prompt (figure 5)

208 We found this prompting setup to yield significantly better results than forcing the model to do
 209 everything (describing the frames, then matching the descriptions to classes) in a single forward pass.
 210 We extract the scores from the last response which has predictable easy-to-parse format. For the full
 211 prompts see figs. 2 to 5.

212 Figures 6, 7 and 8 mention the prompts used within Minecraft environment.

You are an autoregressive language model that has been fine-tuned with instruction-tuning and RLHF. You carefully provide accurate, factual, thoughtful, nuanced answers, and are brilliant at reasoning. Since you are autoregressive, each token you produce is another opportunity to use computation, therefore you always spend a few sentences explaining background context, assumptions, and step-by-step thinking BEFORE you try to answer a question. You always use precise, plain scientific language, without unneeded flourish. You provide details where it might help the explanation.

Figure 2: The system prompt, used in both requests.

213 4 Notes about Croissant and the hosting platform

214 We do include the croissant spec file, although (to our knowledge) it does not directly support datasets
 215 with video files, so it likely is not a good fit for our dataset.

216 In the meantime, we look for more long-term hosting solutions and more fitting standardized
 217 documentation formats; we hope the current temporary solution will be enough for the reviewers to
 218 be able to review the work.

```

You will be given {{n_frames}} frames from a first-person video. The frames are given to
you in chronological order.

{{task.prompt_gpt}}

{{task.example_gpt}}

Input frames. Describe each frame separately.

{{frames}}

```

Figure 3: The prompt used to obtain frame-by-frame descriptions. *n_frames* depends on the model setting (by default, 16), *task.prompt_gpt* and *task.example_gpt* are problem-specific, the latter only being used in 1-shot mode (which is the default). *frames* refer to the actual frames, which are base64 encoded and passed through the OpenAI API as detailed in their documentation.

```

Consider the following sequence of frame-by-frame descriptions:

'''
{{frame_descriptions}}
'''

Break down each of the following summaries into individual steps, and mention what
frames or frame ranges each step matches in parentheses after each step. Note even
partial matches, e.g. matching a kind of action (put, pick, ...) even though the
object might be incorrect. Also provide a one-sentence commentary on how well each
summary matches the sequence of the frame descriptions. In the commentary, the most
important thing is to match the kinds of actions performed and their order. For
example if put (of anything) was described before a pick (of anything) in the
frames, maintaining this order in the summary is more important than getting the
exact object right. Do not comment on the overall quality of the summaries.

Summaries, given in the format '- (label) summary':
{{class_list}}

```

Figure 4: The prompt used to make GPT-4o think step by step about how well each class description matches the frame descriptions it produced in the first request. The *frame_descriptions* are the generated frame-by-frame descriptions, *class_list* is the list of label-description pairs from the task yaml definition. The raw frames are not supplied to the model here anymore.

```

Based on your findings in the previous section, score the summaries from 0 to 5, where
0 is "likely does not describe the video" and 5 is "among the given options, this
one most likely describes the video". Make sure to score each summary individually
. At least one score should be non-zero, even if it's not a perfect match.
Whatever scores you pick, there !must be! exactly one summary with the highest
score. Follow the format below, verbatim:

'''
- (label) score
'''

```

Figure 5: The prompt used to obtain the final scores in a predictable format.

```

# TASK
You will be given several frames from a Minecraft footage. The frames are given in
chronological order.
First, describe each frame. Focus on what activities the player is performing in each
frame and what the player interacts with. Be brief and to the point.
# EXAMPLE OF FRAME DESCRIPTIONS
Input: [ten frames]

Assistant:
1. We see a flat grassy area with small houses and trees nearby. Player holds an oak
fence block in their hands.
2. Same scene from a slightly different angle.
3. Player breaks a tuft of grass using an oak fence block.
4. Player proceeds removing grass.
5. Player places a fence block next to a house.
6. Player places more fence blocks, building a straight fence line.
7. Player turns around and finds a pig.
8. Player places one more fence block in the line.
9. Player breaks another tuft of grass.
10. Player proceeds removing grass.

```

Figure 6: The prompt used to obtain frame-by-frame descriptions for Minecraft videos. *frames* refer to the actual frames, which are base64 encoded and passed through the OpenAI API as detailed in their documentation.

```

Then, given the original frames and your description, score the following potential
video descriptions from 0 to 1 based on how well they describe the video you've
seen. Feel free to use values between 0 and 1, too. There could be more than one
correct' description with score 1. The descriptions are given in the following
format:

- (label) description

Options:

{classes}

The final scores should be in the following format, verbatim:

'''
- (label) your score
'''

Be sure not to alter the label in any way, since we will use it to match your scores to
the potential descriptions we've given you.

```

Figure 7: The prompt used to obtain the final scores in a predictable format for Minecraft videos in **multilabel** setting.

Then, given the original frames and your description, score the following potential video descriptions from 0 to 1 based on how well they describe the video you've seen. Feel free to use values between 0 and 1, too. There should be exactly one 'correct' description with score 1. The descriptions are given in the following format:

- (label) description

Options:

{classes}

The final scores should be in the following format, verbatim:

```
""
- (label) your score
""
```

Be sure not to alter the label in any way, since we will use it to match your scores to the potential descriptions we've given you.

Figure 8: The prompt used to obtain the final scores in a predictable format for Minecraft videos in **multiclass** setting.