
Vision-Language Models Provide Promptable Representations for Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Humans can quickly learn new behaviors by leveraging background world knowl-
2 edge. In contrast, agents trained with reinforcement learning (RL) typically learn
3 behaviors from scratch. We thus propose a novel approach that uses the vast
4 amounts of general and indexable world knowledge encoded in vision-language
5 models (VLMs) pre-trained on Internet-scale data for embodied RL. We initialize
6 policies with VLMs by using them as promptable representations: embeddings
7 that encode semantic features of visual observations based on the VLM’s internal
8 knowledge and reasoning capabilities, as elicited through prompts that provide task
9 context and auxiliary information. We evaluate our approach on visually-complex,
10 long horizon RL tasks in Minecraft and robot navigation in Habitat. We find that
11 our policies trained on embeddings from off-the-shelf, general-purpose VLMs out-
12 perform equivalent policies trained on generic, non-promptable image embeddings.
13 We also find our approach outperforms instruction-following methods and performs
14 comparably to domain-specific embeddings. Finally, we show that our approach
15 can use chain-of-thought prompting to produce representations of common-sense
16 semantic reasoning, improving policy performance in novel scenes by 1.5 times.

17 1 Introduction

18 Embodied decision-making often requires representations informed by world knowledge for per-
19 ceptual grounding, planning, and control. Humans rapidly learn to perform sensorimotor tasks by
20 drawing on prior knowledge, which might be high-level and abstract (“If I’m cooking something
21 that needs milk, the milk is probably in the refrigerator”) or grounded and low-level (e.g., what
22 refrigerators and milk look like). These capabilities would be highly beneficial for reinforcement
23 learning (RL) too: we aim for our agents to interpret tasks in terms of concepts that can be reasoned
24 about with relevant prior knowledge and grounded with previously-learned representations, thus
25 enabling more efficient learning. However, doing so requires a condensed source of vast amounts of
26 general-purpose world knowledge, captured in a form that allows us to specifically index into and
27 access *task-relevant* information. Therefore, we need representations that are contextual, such that
28 agents can use a concise task context to draw out relevant background knowledge, abstractions, and
29 grounded features that aid it in acquiring a new behavior.

30 An approach to facilitate this involves integrating RL agents with the prior knowledge and reasoning
31 abilities of pre-trained foundation models. Transformer-based language models (LMs) and vision-
32 language models (VLMs) are trained on Internet-scale data to enable generalization in downstream
33 tasks requiring facts or common sense. Moreover, in-context learning [8], chain-of-thought reason-
34 ing (CoT) [71], and instruction fine-tuning [50] have provided better ways to index into (V)LMs’
35 knowledge and steer their capabilities based on user needs. These successes have seen some transfer
36 to embodied control, with (V)LMs being used to reason about goals to produce executable plans
37 [2] or as encoders of useful information (like instructions [40] or feedback [61]) that the control

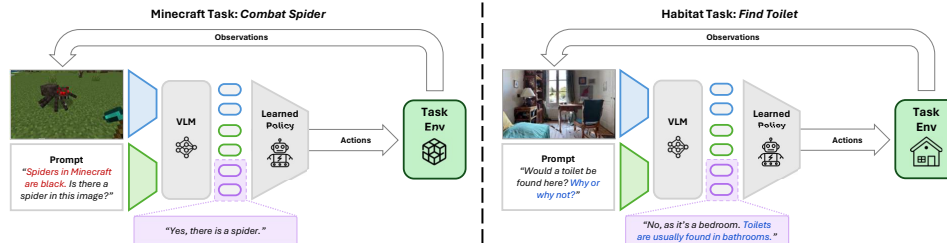


Figure 1: **Example instantiations of PR2L for tasks in Minecraft and Habitat.** We query a VLM with a *task-relevant prompt* about observations to produce *promptable representations*, which we train a policy on via RL. Rather than directly asking for actions or specifying the task, the prompt enables indexing into the VLM’s prior world knowledge to access task-relevant information. This prompt also allows us to **inject auxiliary information** and **elicit chain-of-thought reasoning**.

38 policy utilizes. Both these paradigms have major limitations: actions generated by LMs are often not
 39 appropriately grounded, unless the tasks and scenes are amenable to being expressed or captioned in
 40 language. Even then, (V)LMs are often only suited to producing subtask plans, not low-level control
 41 signals. On the other hand, using (V)LMs to simply encode inputs under-utilizes their knowledge and
 42 reasoning abilities, instead focusing on producing embeddings that reflect the compositionality of
 43 language (e.g., so an instruction-following policy may generalize). This motivates the development
 44 of an algorithm for learning to produce low-level actions that are grounded and leverage (V)LMs’
 45 knowledge and reasoning.

46 To this end, we introduce **Promptable Representations for Reinforcement Learning (PR2L)**: a flexible
 47 framework for steering VLMs into producing *semantic features*, which **(i)** integrate observations
 48 with prior task knowledge and **(ii)** are grounded into actions via RL (see Figure 1). Specifically,
 49 we ask a VLM questions about observations that are related to the given control task, priming it to
 50 attend to task-relevant features in the image based on both its internal world knowledge, reasoning
 51 capabilities, and any supplemental information injected via prompting. The VLM then encodes this
 52 information in decoded text, which is discarded, and associated embeddings, which serve as inputs
 53 to a learned policy. In contrast to the standard approach of using pre-trained image encoders to
 54 convert visual inputs into *generic* features for downstream learning, our method yields *task-specific*
 55 features capturing information particularly conducive to learning a considered task. Thus, the VLM
 56 does not just produce an un-grounded encoding of instructions, but embeddings containing semantic
 57 information relevant to the task, that is both grounded and informed by the VLM’s prior knowledge.

58 To the best our knowledge, we introduce the first approach for initializing RL policies with generative
 59 VLM representations. We demonstrate our approach on tasks in Minecraft [19] and Habitat [58], as
 60 they present semantically-rich problems representative of many practical, realistic, and challenging
 61 applications of RL. We find that PR2L outperforms equivalent policies trained on vision-only
 62 embeddings or with instruction-conditioning, popular ways of using pre-trained image models and
 63 VLMs respectively for control. We also show that promptable representations extracted from general-
 64 purpose VLMs are competitive with domain-specific representations. Our results highlight how
 65 visually-complex control tasks can benefit from accessing the knowledge captured within VLMs via
 66 prompting in both online and offline RL settings.

67 2 Related Works

68 **Vision-language models.** In this work, we utilize *generative VLMs* (like [33, 34, 14, 29]): models
 69 that generate language in response to an image and a text prompt passed as input. This is in contrast to
 70 other designs of combining vision and language that either generate images or segmentation [57, 30]
 71 and contrastive representations [52]. Formally, the VLM enables sampling from $p(x_{1:K}|I, c)$, where
 72 $x_{1:K}$ represents the K tokens of the output, I is the input image(s), c is the prompt, and p is the
 73 distribution over natural language responses produced by the VLM on those inputs. Typically, the
 74 VLM is pre-trained on tasks that require building association between vision and language such as
 75 captioning. All these tasks require learning to attend to certain semantic features of input images
 76 depending on the given prompt. For auto-regressive generative VLMs, this distribution is factorized as
 77 $\prod_t p(x_t|I, c, x_{1:t-1})$. Typical architectures parameterize these distributions using weights that define
 78 a representation $\phi_t(I, c, x_{1:t-1})$, which depends on the image I , the prompt c , and the previously
 79 emitted tokens, and a decoder $p(x_t|\phi_t(I, c, x_{1:t-1}))$, which defines a distribution over the next token.

80 **Embodied (V)LM reasoning.** Many recent works have leveraged (V)LMs as priors over effective
81 plans for a given goal. These works use the model’s language modeling and auto-regressive generation
82 capabilities to extract such priors as textual subtask sequences [2, 24, 60] or code [36, 64, 77, 68],
83 thereby using the (V)LM to decompose long-horizon tasks into executable parts. These systems
84 often need grounding mechanisms to ensure plan feasibility (e.g., affordance estimators [2], scene
85 captioners [77], or trajectory labelers [51]). They also often assume access to low-level policies
86 that can execute these subtasks, such as robot pick-and-place skills [2, 36], which is often a strong
87 assumption. These methods generally do not address how such policies can be acquired, nor how
88 these low-level skills can themselves benefit from the prior knowledge in (V)LMs. Even works in
89 this area that use RL still use (V)LMs as state-dependent priors over reasonable high-level goals to
90 learn [17]. This is a key difference from our work: instead of considering priors on plans/goals, we
91 rely on VLM’s implicit knowledge *of the world* to extract representations which encode task-relevant
92 information. We train a policy to convert these features into low-level actions via standard RL,
93 meaning the VLM does not need to know how to take actions for a task.

94 **Embodied (V)LM pre-training.** Other works use (V)LMs to embed useful information like instruc-
95 tions [40, 45, 42, 44, 49], feedback [61, 9], reward specifications [19], and data for world modeling
96 [39, 47]. These works use (V)LMs as *encoders* of the compositional semantic structure of input text
97 and images, which aids in generalization: an instruction-conditioned model may never have learned
98 to grasp apples (but can grasp other objects), but by interacting with them in other ways and receiving
99 associated language descriptions, the model might still be able to grasp them zero-shot. In contrast,
100 our method produces embeddings that are informed by world knowledge and reasoning, both from
101 prompting and pre-training. Rather than just specifying that the task is to acquire an apple, we ask a
102 VLM to parse observations into task-relevant features, like whether there is an apple in the image or
103 if the observed location likely contains apples – information that is useful even in single-task RL.
104 Thus, we use VLMs to help RL solve new tasks, not just to follow instructions.

105 These two categories are not mutually exclusive: Brohan et al. [6] use VLMs to understand instruc-
106 tions, but also reasoning (e.g., figuring out the “correct bowl” for a strawberry is one that contains
107 fruits); Palo et al. [51] use a LM to reason about goal subtasks and a VLM to know when a trajectory
108 matches a subtask, automating the demonstration collection/labeling of Ahn et al. [2], while Adeniji
109 et al. [1] use a similar approach to pretrain a language-conditioned RL policy that is transferable to
110 learning other tasks; and Shridhar et al. [63] use CLIP to merge vision and text instructions directly
111 into a form that a Transporter [76] policy can operationalize. Nevertheless, these works primarily
112 focus on instruction-following for robot manipulation. Our approach instead prompts a VLM to
113 supplement RL with representations of world knowledge, not instructions. In addition, except for
114 Adeniji et al. [1], these works focus on behavior cloning (BC), assuming access to demonstrations for
115 policy learning, whereas our framework can be used for both online RL and offline RL/BC.

116 3 PR2L: Promptable Representations for Reinforcement Learning

117 We adopt the standard framework of partially-observed Markov decision process in deep RL, wherein
118 the objective is to find a policy mapping states to actions that maximizes the expected returns. Our
119 goal is to supplement RL with task-relevant information extracted from VLMs containing general-
120 purpose knowledge. One way to index into this information is by prompting the model to get it
121 to produce semantic information relevant to a given control task. Therefore, our approach, PR2L,
122 queries a VLM with a task-relevant prompt for each visual observation received by the agent, and
123 receives both the decoded text and, critically, the intermediate representations, which we refer to
124 as *promptable representations*. Even though the decoded text might often not be correct or directly
125 usable for choosing the action, our key insight is that these VLM embeddings can still provide
126 useful semantic features for training control policies via RL. This recipe enables us to incorporate
127 semantic information without the need of re-training or fine-tuning a VLM to directly output actions,
128 as proposed by Brohan et al. [6]. Note that our method is *not* an instruction-following method, and
129 it does **not** require a task instruction to perform well. Instead, our approach still learns control via
130 RL, while benefiting from the incorporation of *background context*. In this section, we will describe
131 various components of our approach, accompanied by practical design choices and considerations.

132 3.1 Promptable Representations

133 In principle, one can directly query a VLM to produce actions for a task given a visual observation.
134 While this may work when high-level goals or subtasks are sufficient, VLMs are empirically poor at
135 yielding the low-level actions used commonly in RL [23]. As VLMs are trained to follow instructions

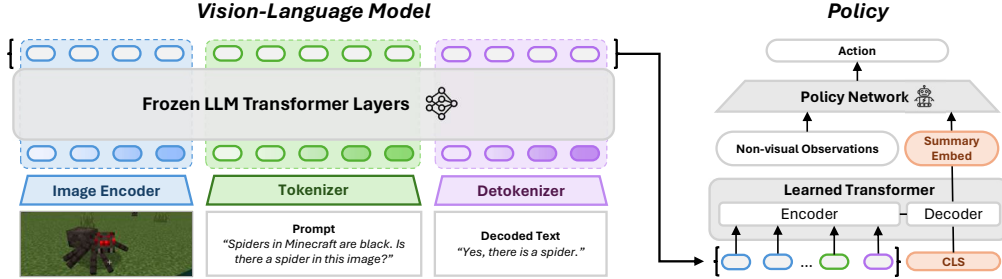


Figure 2: **Schematic of how we extract task-relevant features from the VLM and use them in a policy trained with RL.** These representations can incorporate task context from the prompt, while generic [image embeddings](#) cannot. As generative VLM’s embeddings can be variable length, the policy has a Transformer layer that takes in these embeddings and a “CLS” token, thereby condensing all inputs into a single summary vector.

136 and answer questions about images, it is more appropriate to use these models to extract and reason
 137 about *semantic features* about observations that are conducive to being linked to actions. We thus
 138 elicit features that are useful for the downstream task by querying these VLMs with *task-relevant*
 139 *prompts* that provide contextual task information, thereby causing the VLM to attend to and interpret
 140 appropriate parts of observed images. Extracting these features naïvely by only using the VLM’s
 141 *decoded text* has its own challenges: such models often suffer from hallucinations [26] and an inability
 142 to report what they “know” in language, even when their embeddings contain such information
 143 [27, 21]. However, even when the text is bad, the underlying *representations* still contain valuable
 144 granular world information that is potentially lost in the projection to language [32, 72, 22, 35]. Thus,
 145 we disregard the generated text and instead provide our policy the embeddings produced by the VLM
 146 in response to prompts asking about relevant semantic features in observations instead.

147 **Which parts of the network can be used as promptable representations?** The VLMs we consider
 148 are all based on the Transformer architecture [67], which treats the prompt, input image(s), and
 149 decoded text as token sequences. This architecture provides a source of learned representations by
 150 computing embeddings for each token at every layer based on the previous layer’s token embeddings.
 151 In terms of the generative VLM formalism introduced prior, a Transformer-based VLM’s repre-
 152 sentations $\phi_t(I, c, x_{1:t-1})$ consist of N embeddings per token (the outputs of the N self-attention
 153 layers) in the input image I , prompt c , and decoded text $x_{1:t-1}$. The decoder $p(x_t|\phi_t)$ extracts the
 154 final layer’s embedding of the most recent token x_{t-1} , projecting it to a distribution over the token
 155 vocabulary and allowing for it to be sampled. When given a visual observation and task prompt,
 156 the tokens representing the prompt, image, and answer consequently encode task-relevant semantic
 157 information. Thus, for each observation, we use the VLM to sample a response to the task prompt
 158 $x_{1:K} \sim p(x_{1:K}|I, c)$. We then use some or all of these token embeddings $\phi_K(I, c, x_{1:t-1})$ as our
 159 promptable representations and feed them, along with any non-visual observation information, as a
 160 state representation into our neural policy trained with RL.

161 In summary, our approach involves creating a task-relevant prompt that provides context and auxiliary
 162 information. This prompt, alongside the current visual observation from the environment, is fed
 163 to into the VLM to generate tokens. While these tokens are used for decoding, they are ultimately
 164 discarded. Instead, we utilize the *representations* produced by the VLM (associated with the image,
 165 prompt, and decoded text) as input for our policy, which is trained via an off-the-shelf online RL
 166 algorithm to produce appropriate actions. A schematic of our approach is depicted in Figure 2 and a
 167 code snippet example is presented in Appendix I.

168 3.2 Design Choices for PR2L

169 To instantiate this idea, we need to make some concrete design choices in practice. First, the
 170 representations of the VLM’s decoded text depend on the chosen decoding scheme: greedy decoding
 171 is fast and deterministic, but may yield low-probability decoded tokens; beam search improves on this
 172 by considering multiple “branches” of decoded text, at the cost of requiring more compute time (for
 173 potentially small improvements); lastly, sampling-based decoding can quickly yield estimates of the
 174 maximum likelihood answer, but at the cost of introducing stochasticity, which may increase variance.
 175 Given the inherent high-variance of our tasks (due to sparse rewards and partial observability) and
 176 the expense of VLM decoding, we opt for greedy decoding or fixed-seed sampling.

177 Second, one must choose which VLM layers’ embeddings to utilize in the policy. While theoretically,
 178 all layers of the VLM could be used, pre-trained Transformer models tend to encode valuable high-

179 level semantic information in their later layers [66, 25]. Thus, we opt to only feed the final few
180 layers’ representations into our policy. As these representation sequences are of variable length, we
181 incorporate an encoder-decoder Transformer layer in the policy. At each time step in a trajectory,
182 this layer receives variable-length VLM representations, which are attended to and converted into a
183 fixed-length summarization by the embeddings of a learned “CLS” token [15] in the decoder (green
184 in Figure 2). We also note that this policy can receive the observed image directly (e.g., after being
185 embedded by the image encoder), so as to not lose any visual information from being processed
186 by the VLM. However, we do not do this in our experiments in order to more clearly isolate and
187 demonstrate the usefulness of the VLM’s representations in particular.

188 Finally, while it is possible to fine-tune the VLM for RL end-to-end with the policy [6], this incurs
189 substantial compute, memory, and time overhead, particularly with larger VLMs. Nonetheless, we
190 find that our approach performs better than not using the language and prompting components of the
191 VLM. This holds true even when the VLM is frozen, and only the policy is trained via RL, or when
192 the decoded text occasionally fails to answer the task-specific prompt correctly.

193 3.3 Task-Relevant Prompt Design

194 **How do we design good prompts to elicit useful representations from VLMs?** As we aim to
195 extract good state representations from the VLM for a downstream policy, we do not use instructions
196 or task descriptions, but task-relevant prompts: questions that make the VLM attend to and encode
197 semantic features in the image that are useful for the RL policy learning to solve the task [5]. For
198 instance, if the task is to find a toilet within a house, appropriate prompts include “What room is this?”
199 and “Would a toilet be found here?” Intuitively, the answers to these questions help determine good
200 actions (e.g., look around the room or explore elsewhere), making the corresponding representations
201 good for representing the state for a policy. Answering the questions will require the VLM to attend to
202 task-relevant features in the scene, relying on the model’s internal conception of what things look like
203 and common-sense semantic relations. One can also prompt the VLM to use chain of thought [71]
204 to explain its generated text, often requiring it to reason about task-relevant features in the image,
205 resulting in further enrichment of the state representations. Finally, prompts can provide helpful
206 auxiliary information: e.g., one can describe what certain entities of interest look like, aiding the
207 VLM in detecting them even if they were not commonly found in the model’s pre-training data.

208 Note that prompts based on instructions or task descriptions do not enjoy the above properties: while
209 the goal of those prior methods is to be able to directly query the VLM for the optimal action, the
210 goal of task-relevant prompts is to produce a useful state representation, such that running RL with
211 them can accelerate learning an optimal policy. While the former is not possible without task-specific
212 training data for the VLM in the control task, the latter proves beneficial with off-the-shelf VLMs.

213 **Evaluating and designing prompts for RL.** Since the specific representations elicited from the VLM
214 are determined by the prompt, we want to design prompts that produce promptable representations
215 that maximize performance on the downstream task. The brute-force approach would involve running
216 RL with each candidate prompt to measure its efficacy, but this would be computationally very
217 expensive. In lieu of this, we evaluate candidate prompts on a small dataset of observations labeled
218 with semantic features of interest for the considered task. Example features include whether task-
219 relevant entities are in the image, the relative position of said entities, or even actions (if expert
220 demonstrations are available). We test prompts by querying the VLM and checking how well the
221 resulting decoded text for each image matches ground truth labels. As this is only practical for
222 small, discrete spaces that are easily expressed in words, we see how well a small model can fit the
223 VLM’s embeddings to the labels (akin to probing in self-supervised learning [62, 4]). While this does
224 not directly optimize for task performance, it does act as a proxy that ensures a prompt’s resulting
225 representations encode certain semantic features which are helpful for the task.

226 4 Experimental Setups

227 Our experiments analyze whether promptable representations from VLMs provide benefits to down-
228 stream control, thus providing an effective vehicle for transferring Internet-scale knowledge to RL.
229 We aim to show that PR2L is a good source of state representations, even with our current VLMs
230 that are bad at reasoning about actions – as such models become more performant, we expect such
231 representations to be even better. We thus design experiments to answer the following: (1) Can
232 promptable representations obtained via task-specific prompts enable more performant and sample-
233 efficient learning than those of non-promptable image encoders pre-trained for vision or control? (2)
234 How does PR2L compare to approaches that directly “ask” the VLM to generate good actions for a

235 task specified in the prompt? (3) How does PR2L fare against other popular learning approaches or
236 purely visual features in our domains of interest?

237 4.1 Domain 1: Minecraft

238 We first conduct experiments in Minecraft, which provides control tasks that require associating
239 visual observations with rich semantic information to succeed. Moreover, since these observations
240 are distinct from the images in the the pre-training dataset of the VLM, succeeding on these tasks
241 relies crucially on the efficacy of the task-specific prompt in meaningfully affecting the learned
242 representation, enabling us to stress-test our method. E.g., while spiders in Minecraft somewhat
243 resemble real-life spiders, they exhibit stylistic exaggerations such as bright red eyes and a large black
244 body. If the task-specific prompt is indeed effective in informing the VLM of these facts, it would
245 produce a representation that is more conducive to policy learning and this would be reflected in
246 task performance. For this domain, we use the half-precision Vicuna-7B version of the InstructBLIP
247 instruction-tuned generative VLM [14, 12] to produce promptable representations.

248 **Minecraft tasks.** We consider all programmatic Minecraft tasks evaluated by Fan et al. [19]: *combat*
249 *spider*, *milk cow*, *shear sheep*, *combat zombie*, *combat enderman*, and *combat pigman*¹. The
250 remaining tasks considered by Fan et al. [19] are creative tasks, which do not have programmatic
251 reward functions or success detectors, so we cannot directly train RL agents on them. We follow the
252 MineDojo definitions of observation/action spaces and reward function structures for these tasks: at
253 each time step, the policy observes an egocentric RGB image, its pose, and its previously action;
254 the policy can choose a discrete action to turn the agent by changing the agent’s pitch and/or yaw in
255 discrete increments, move, attack, or use a held item. These tasks are long horizon, with a maximum
256 episode length of 500 - 1000 and taking roughly 200 steps for a learned policy to complete them. See
257 Figure 3 for example observations and Appendix B.1 for more details.

258 **Comparisons.** We compare PR2L to five performant classes of approaches for RL in Minecraft: (a)
259 Methods using non-promptable representations of visual observations. This does not use prompting
260 altogether, instead using task-agnostic embeddings from the VLM’s image encoder (specifically, the
261 ViT-g/14 from InstructBLIP – blue in Figure 2). While these representations are still pre-trained, PR2L
262 utilizes prompting to produce *task-specific* representations. For a fair comparison, we use the *exact*
263 *same* policy architecture and hyperparameters for this baseline as in PR2L, ensuring that performance
264 differences come from prompting for better representations from the VLM. (b) Methods that directly
265 “asks” the VLM to output actions to execute on the agent. This adapts the approach of Brohan et al.
266 [6] to our setting and directly outputs the action from the VLM. While Brohan et al. [6] also fine-tune
267 the VLM backbone, we are unable to do so using our compute resources. To compensate, we do not
268 just execute the action from the VLM, but train an RL policy to map this decoded action to a better
269 one. Note that if the VLM already decodes good action texts, simply copying over this action via RL
270 should be easy. (c) Methods for efficient RL from pixels via model-based approaches. We choose
271 Dreamer v3, since it has proven to be successful at learning Minecraft tasks from scratch [20]. (d)
272 Methods leveraging pretrained representations specifically useful for embodied control, though which
273 are non-promptable and non-Minecraft specific. We choose VC-1 and R3M [43, 46]. (e) Methods
274 using models pre-trained on large-scale Minecraft data. These serve as “oracle” comparisons, as
275 these representations are explicitly fine-tuned on Minecraft YouTube videos, whereas our pre-trained
276 VLM is both frozen and not trained on any Minecraft video data. We choose MineCLIP, VPT, and
277 STEVE-1 as our sources of Minecraft-specific representations [19, 3, 37].

278 We use PPO [59] as our base RL algorithm for all non-Dreamer Minecraft policies. We also note that
279 we do *not* compare against non-RL methods, such as Voyager (which uses LLMs to write high-level
280 code skills, abstracting away low-level control to hand-written APIs that use oracle information). See
281 Appendix B.2 for training details and E.1 for further discussion of such non-learned systems.

282 4.2 Domain 2: Habitat

283 A major advantage of VLMs pre-trained on Internet-scale data is their reasoning and generalization
284 capabilities. To evaluate this, we run offline BC and RL experiments in the Habitat household
285 simulator. In contrast to Minecraft, tasks in this domain require connecting *naturalistic* images
286 with real-world common sense about the structure and contents of typical home environments. Our
287 experiments evaluate (1) whether PR2L confers the generalization properties of VLMs to our policies,

¹ Fan et al. [19] also consider *hunt cow/sheep*. However, we omit them as we were unable to replicate their results on those tasks; all approaches failed to learn them.

	PR2L Prompt	RT-2-style Baseline Prompt	Change Auxiliary Text Ablation Prompt
<i>Combat Spider</i>	Spiders in Minecraft are black. Is there a spider in this image?	I want to fight a spider. I can attack, move, or turn. What should I do?	Is there a spider in this image?
<i>Milk Cow</i>	Is there a cow in this image?	I want to milk a cow. I can use my bucket, move, or turn. What should I do?	Cows in Minecraft are black and white. Is there a cow in this image?
<i>Shear Sheep</i>	Is there a sheep in this image?	I want to shear a sheep. I can use my shears, move, or turn. What should I do?	Sheep in Minecraft are usually white. Is there a sheep in this image?
<i>Other Combat Tasks</i>	Is there a [target entity] in this image?	I want to fight a [target entity]. I can attack, move, or turn. What should I do?	-

Table 1: Prompts used in Minecraft for querying the VLM with PR2L, comparison (b), and the change auxiliary text ablation. For the last column, we remove the **auxiliary text** for *combat spider*, and add it in for the other two.

288 (2) whether PR2L-based policies can leverage the semantic reasoning capabilities of the underlying
289 VLM (e.g., via chain-of-thought [71]), and (3) whether PR2L can learn entirely from stale, offline
290 data sources. We use a Llama2-7B Prismatic VLM for the Habitat experiments [29].

291 **Habitat tasks.** We consider the ObjectNav task suite in 3D scanned household scenes from the
292 HM3D dataset [58, 73, 54]. These tasks involve a simulated robot traversing a home environment to
293 find an instance of a specified object (toilet, bed, sofa, television, plant, or chair) in the shortest path
294 possible. The full benchmark consists of 80 household scenes intended to train the agent and 20 for
295 validation. We change the observation space to consist of just RGB vision, previous action, pose,
296 and target object class, omitting depth images to ensure that observed performance differences come
297 from the quality of promptable representations vs. unpromptable ones. Like with MineDojo, these
298 tasks are long horizon, taking 80 steps for a privileged shortest path follower to succeed and 150+
299 for humans. See Figure 3 for example observations and Appendix C for more details.

300 **Comparisons.** To see if PR2L can leverage VLM reasoning capabilities, we train two PR2L policies,
301 one with and one without chain-of-thought prompting (see Section 4.3). We also train a policy
302 on Prismatic VLM image encoder embeddings (equivalent to Minecraft approach (a), but with
303 Dino+SigLIP [11, 78]) on a human demonstration dataset collected from the ObjectNav training
304 scenes collected with Habitat-Web [55] and used by past works on large-scale BC on pre-trained
305 visual representations [56, 74, 43]. As it previously achieved state-of-the-art performance among
306 those works, we also compare against two policies using VC-1 as an encoder [43], either using just
307 its summarizing CLS token or using a learned Transformer layer to condense its patch embeddings.
308 We adopt the same LSTM-based recurrent architecture used by that work, but replace the image
309 embeddings with a learned Transformer layer that condenses our input token embeddings (from the
310 VLM, VLM image encoder, or VC-1) into a single summary embedding, as done with Minecraft.

311 Due to computational constraints, we train all policies on just under a tenth of the full dataset of
312 77k trajectories/12M steps. In contrast, other works using this dataset train on the entire dataset.
313 Nevertheless, we evaluate on the unseen validation scenes, thereby testing how well PR2L generalizes.

314 4.3 Designing Task-Specific Prompts for Minecraft and Habitat

315 We now discuss how to design prompts for PR2L. As noted in Section 3.3, these are not instructions
316 or task descriptions, but prompts that force the VLM to encode semantic information useful for the
317 task in its representation. The simplest relevant feature for our Minecraft tasks is the presence of the
318 target entity in an observation. Thus, we choose “Is there a [target entity] in this image?” as the base
319 of our chosen prompt. We also pick two alternate prompts per task that prepend different amounts of
320 auxiliary information about the target entity. E.g., for *combat spider*, one candidate is “Spiders in
321 Minecraft are black.” To choose between these candidates, we measure how well the VLM is able
322 to decode a correct answer to the prompt question of whether or not the target entity is present in
323 the image on a small annotated dataset. Full details of this prompt evaluation scheme for the first
324 three Minecraft tasks are presented in Appendix A and Table 5. We find that auxiliary text only helps
325 with detecting spiders while systematically and significantly degrading the detection of sheep and
326 cows. Our ablations show that this detection success rate metric correlates with performance of the
327 RL policy. Additionally, the prompts used for comparison (b) follow the prompt structure prescribed
328 by Brohan et al. [6], which motivated this comparison. In these prompts, we also provide a list of
329 actions that the VLM can choose from to the policy. All chosen prompts are presented in Table 1.

330 For Habitat, we choose the prompt “Would a [target object] be found here? Why or why not?” As
331 opposed to the Minecraft prompts, this does not just identify the presence of a target object in the
332 image, but draws on general knowledge from the VLM to determine if the observed location would
333 contain the target object, even if said object is not in view. The second part of the prompt then leads
334 the VLM to provide a chain of thought (CoT) [71] rationale for its final answer. This CoT draws out

Task	PR2L (Ours)	VLM Image Encoder	Baselines				Oracles		
			RT-2-style	Dreamer	VC-1	R3M	MineCLIP	VPT	STEVE-1
<i>Combat Spider</i>	97.6 ± 14.9	51.2 ± 9.3	71.5 ± 9.7	5.4 ± 1.1	72.2 ± 9.3	72.9 ± 8.7	<i>176.9 ± 19.8</i>	<i>137.2 ± 19.2</i>	88.8 ± 14.0
<i>Milk Cow</i>	223.4 ± 35.4	95.2 ± 18.7	128.6 ± 28.9	24.0 ± 1.2	96.6 ± 16.3	100.0 ± 14.1	194.4 ± 33.3	85.5 ± 14.5	75.2 ± 15.4
<i>Shear Sheep</i>	37.0 ± 4.4	23.0 ± 3.6	26.2 ± 3.2	20.9 ± 1.2	26.5 ± 4.0	17.5 ± 2.4	23.1 ± 3.7	24.1 ± 2.9	18.2 ± 2.5
<i>Combat Zombie</i>	24.6 ± 1.6	14.8 ± 2.0	18.2 ± 2.1	1.8 ± 0.2	5.6 ± 1.0	5.8 ± 1.4	<i>56.6 ± 8.3</i>	<i>31.2 ± 3.2</i>	23.6 ± 3.4
<i>Combat Enderman</i>	52.2 ± 5.6	51.9 ± 6.8	44.6 ± 5.8	1.6 ± 0.5	27.2 ± 2.4	33.8 ± 3.8	<i>72.1 ± 7.1</i>	<i>74.4 ± 13.2</i>	59.3 ± 6.7
<i>Combat Pigman</i>	46.4 ± 3.3	36.8 ± 3.7	35.1 ± 2.5	5.8 ± 1.5	33.7 ± 4.9	31.4 ± 4.2	<i>189.0 ± 7.9</i>	<i>169.0 ± 7.8</i>	98.3 ± 8.4

Table 2: **Performance of PR2L, baseline, and oracle approaches in Minecraft tasks.** Values reported are IQM successes and standard errors. PR2L universally outperforms all baselines. As they are trained on Minecraft-specific data, the oracles outperform PR2L in half the comparisons (italicized).

Task (# Episodes)	PR2L (Ours)		VLM Image Encoder	VC-1 + CLS		VC-1 + Patch Embeds	
	With CoT	Without CoT		40 Epochs	120 Epochs	40 Epochs	120 Epochs
<i>Average (2000)</i>	41.9%	27.8%	11.6%	6.8%	8.9%	13.6%	15.8%
<i>Toilet (398)</i>	37.2%	22.9%	8.8%	2.8%	2.0%	7.0%	9.3%
<i>Bed (433)</i>	45.0%	28.9%	12.9%	6.7%	9.9%	14.8%	19.2%
<i>Sofa (376)</i>	48.1%	34.3%	11.7%	9.8%	14.4%	17.0%	19.4%
<i>Chair (428)</i>	51.2%	40.9%	17.5%	11.7%	15.0%	22.4%	23.8%
<i>Television (281)</i>	26.7%	10.3%	5.0%	2.8%	3.2%	4.6%	4.6%
<i>Plant (84)</i>	23.8%	8.3%	9.1%	1.2%	1.2%	9.5%	9.5%

Table 3: **Performance of PR2L and baselines on Habitat ObjectNav tasks.** Following prior works, values reported are average success rates in unseen validation scenes. PR2L (with or without CoT) does better than all other approaches. PR2L with CoT does the best, universally achieving more than double the performance of all non-PR2L approaches and 14.7% higher average performance than PR2L without CoT. Note that PR2L and image encoder policies were trained for 40 epochs, but VC-1 policies’ performance saturated at 120, so we report their performance at both times.

335 task-relevant VLM world knowledge by explicitly reasoning about visual semantic concepts, that are
336 useful to learning a policy (see Table 4; ObjectNav). To investigate if PR2L enables embodied agents
337 to benefit from these VLM common-sense reasoning capabilities (even if they do not directly reason
338 about actions), we train PR2L policies both with and without the second part of the prompt.

339 5 Results

340 **Minecraft results.** We report the interquartile mean (IQM) and standard error number of successes
341 over 16 seeds for all Minecraft tasks in Table 2. PR2L uniformly outperforms the non-oracle
342 approaches of (a) using non-promptable image embeddings, (b) directly asking the VLM for actions,
343 (c) learning from scratch Dreamer, and (d) using non-promptable control-specific embeddings.

344 PR2L outperforms (a) **the VLM image encoder baseline**, even though both approaches receive
345 the same visual features, with PR2L simply transforming those features via prompting an LLM
346 (with no additional information from the environment), thus supporting that prompting does shape
347 representations in a beneficial way for learning control tasks. We provide an analysis of why PR2L
348 states are better than (b) **RT-2-style** ones in Appendix H.1. We observe that PR2L embeddings are
349 bimodally distributed, with transitions leading to high reward clustered at one mode. This structure
350 likely enables more efficient learning, thereby showing how control tasks can benefit from extracting
351 prior knowledge encoded in VLMs by prompting them with task context, even when the VLM does
352 not know how to act. For (c) **the model-based comparisons**, we find that Dreamer is not as conducive
353 at learning our Minecraft tasks. We hypothesize this is because our tasks are comparatively shorter
354 than the ones considered by Hafner et al. [20], so learning a model is less beneficial (while PR2L
355 provides immediately-useful representations). Additionally, we note that all our approaches involve
356 interacting with partially-observable, non-stationary entities, which the Dreamer model may have a
357 hard time learning. See Appendix E.2 for further discussion. Finally, (e) **the oracles** outperform PR2L
358 in *combat enderman/pigman*, all but STEVE-1 do better in *combat spider/zombie*, and none do better
359 in *shear sheep/milk cow*. We hypothesize this is because endermen and pigmen are Minecraft-specific
360 entities, giving rise to comparatively poor representations in the VLM (which is trained exclusively
361 on natural images). In contrast, Minecraft zombies/spiders are heavily stylized, but still somewhat
362 resemble other depictions of such creatures, while Minecraft cows and sheep are the closest to their
363 naturalistic counterparts, making PR2L more effective. Even though our VLM is not trained on
364 Minecraft data, its representations yield better policies in half the oracle comparisons.

365 We provide ablations in Table 8 and Appendix F. We find that (1) PR2L performs worse when it is
366 unprompted or does not decode text, (2) our prompt evaluation scheme successfully identified cases
367 where auxiliary text improves/degrades performance, and (3) a policy with oracle entity detection
368 does worse than PR2L, suggesting our prompt is not just eliciting that feature from the VLM.

Target	Example VLM Responses to "Would a [target] be found here? Why or why not?"
Toilet	"No, a toilet is not found in this room as it is a bedroom. Toilets are typically found in bathrooms ."
	The image shows a large bed with a wooden headboard, nightstands, and lamps ."
	"Yes, a toilet would be found in this room because it is a bathroom. The presence of a toilet in the image is expected, as bathrooms are designed to accommodate toilets for personal hygiene ."
Bed	"No, a bed would not be found in this room because it is a dining room. The presence of a dining table and chairs suggests that this room is designed for eating meals and socializing, rather than sleeping ."
	"Yes, a bed would be found here because the room is described as a bedroom. It has a bed, a television, a lamp, and a dresser, all of which are common furnishings in a bedroom ."
	"No, a sofa is not found here because it is a kitchen area with appliances and cabinets . Sofas are typically found in living rooms or other common areas ."
Sofa	" Yes, there is a black leather sofa in the living room. It has a red pillow on it. It is a large sectional couch ."

Table 4: Example VLM responses to the Habitat prompt for various images. Beyond just **detecting the target**, prompting the VLM for CoT elicits **relevant common sense**, which it semantically relates to **other useful visual features**. By using the underlying VLM embeddings as a state representation, the policy thus integrates the VLM’s knowledge and reasoning into its decision-making.

369 **Habitat results.** Following prior works, we report success rates on the ObjectNav validation
370 episodes in Table 3. PR2L with CoT outperforms all other policies on all tasks, including an almost
371 $4\times$ performance increase over the VLM image encoder baselines – again, suggesting that using
372 promptable representations for control improves over the base purely-visual embeddings. While
373 PR2L without CoT still does better than all baselines, we find CoT prompting improves policy
374 performance (by $1.5\times$, from 27.8% success rate to 41.9%), likely because it provides the policy with
375 useful generalizable features: e.g., even if the agent comes across an unfamiliar room while searching
376 for a toilet, it still knows to look elsewhere if the VLM reasons that, due to the presence of a bed, the
377 room is likely a bedroom (which is unlikely to contain toilets). Thus, even if the VLM cannot reason
378 about actions, *our results indicate that PR2L provides a promising way of using its ability to reason*
379 *about image semantics and common sense for control*. See Table 4 for CoT examples.

380 While we do not beat VC-1’s reported SOTA BC performance (60.3% success rate when VC-1 is
381 frozen [43]), we note that said performance is achieved with (1) over ten times more training data and
382 gradient steps and (2) image augmentations to prevent overfitting. Our VC-1 policies were trained on
383 the same amount of data as our PR2L agent and for $1\text{-}3\times$ as many gradient steps, but perform far
384 worse, suggesting that PR2L is significantly more sample- and compute-efficient than VC-1 policies.
385 Additionally, PR2L does not use any explicit countermeasures to overfitting, yet still generalizes well
386 to unseen ObjectNav scenes (aided by the VLM’s representations of reasoning).

387 Finally, we analyze policies trained with offline RL in a simplified Habitat setting in Appendices D,
388 H, where we find that VLM representations align well with the returns of an optimal policy.

389 6 Conclusion

390 We propose Promptable Representations for Reinforcement Learning, a method for extracting se-
391 mantic features from images by prompting VLMs with task context to leverage their extensive
392 general-purpose prior knowledge. We demonstrate PR2L in Minecraft and Habitat, domains that
393 benefit from interpreting observations in terms of semantic concepts that can be related to task context.
394 This framework for using VLMs for control opens new directions. For example, other types of
395 foundation models pre-trained with more sophisticated methods could also be used for PR2L: e.g.,
396 ones trained on physical interactions might yield features which encode physics or action knowledge,
397 rather than just common-sense visual semantics. Developing and using such models with PR2L offers
398 an exciting way to transfer diverse prior knowledge to a broad range of control applications.

399 A limitation of PR2L is that prompts are currently hand-crafted based on the user’s conception of
400 useful task features. While coming up with good prompts for our tasks was not hard, the process of
401 evaluating and improving them could be automated, which we leave to future works. We also find that
402 the quality of representations largely depends on the VLM – e.g., InstructBLIP could not reason well
403 about Habitat scenes, but the more recent Prismatic VLMs are more capable in that regard, enabling
404 our CoT experiments. Thus, as VLM capabilities are expected to increase, we expect the quality of
405 their representations to also improve. Lastly, the size and speed of VLMs can limit their applicability.
406 Our policies typically achieve 3-5 Hz inference speeds, comparable to those of robot policies built on
407 large models [7, 6, 49]. Likewise, our VLM sizes are comparable to models used for policies in prior
408 works [6, 65]. While their inference speeds may hinder online policy learning, we find that offline
409 approaches (which can parallelize training and data generation) we used for Habitat help remedy this.

410 References

- 411 [1] A. Adeniji, A. Xie, C. Sferrazza, Y. Seo, S. James, and P. Abbeel. Language reward modulation
412 for pretraining reinforcement learning, 2023.
- 413 [2] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakr-
414 ishnan, K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano,
415 K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine,
416 Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet,
417 N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, M. Yan, and A. Zeng.
418 Do as i can and not as i say: Grounding language in robotic affordances. 2022.
- 419 [3] B. Baker, I. Akkaya, P. Zhokhov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro,
420 and J. Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos,
421 2022.
- 422 [4] Y. Belinkov and J. Glass. Analysis methods in neural language processing: A survey, 2019.
- 423 [5] J. Borja-Diaz, O. Mees, G. Kalweit, L. Hermann, J. Boedecker, and W. Burgard. Affordance
424 learning from play for sample-efficient policy learning. In *Proceedings of the IEEE International*
425 *Conference on Robotics and Automation (ICRA)*, Philadelphia, USA, 2022.
- 426 [6] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
427 A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman,
428 A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal,
429 L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao,
430 K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soriccut,
431 H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao,
432 P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web
433 knowledge to robotic control, 2023.
- 434 [7] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
435 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi,
436 R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manju-
437 nath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao,
438 M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran,
439 V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-1:
440 Robotics transformer for real-world control at scale, 2023.
- 441 [8] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan,
442 P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child,
443 A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray,
444 B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Lan-
445 guage models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and
446 H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–
447 1901. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc94967418bfb8ac142f64a-Paper.pdf)
448 [files/paper/2020/file/1457c0d6bfc94967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc94967418bfb8ac142f64a-Paper.pdf).
- 449 [9] A. Bucker, L. Figueredo, S. Haddadin, A. Kapoor, S. Ma, S. Vemprala, and R. Bonatti. Latte:
450 Language trajectory transformer, 2022.
- 451 [10] S. Cai, Z. Wang, X. Ma, A. Liu, and Y. Liang. Open-world multi-task control through goal-aware
452 representation learning and adaptive horizon prediction, 2023.
- 453 [11] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging
454 properties in self-supervised vision transformers, 2021.
- 455 [12] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E.
456 Gonzalez, I. Stoica, and E. P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%*
457 chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- 458 [13] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos. Distributional reinforcement learning
459 with quantile regression, 2017.

- 460 [14] W. Dai, J. Li, D. Li, A. M. H. Tiong, J. Zhao, W. Wang, B. Li, P. Fung, and S. Hoi. Instructblip:
461 Towards general-purpose vision-language models with instruction tuning, 2023.
- 462 [15] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
463 transformers for language understanding, 2019.
- 464 [16] Z. Ding, H. Luo, K. Li, J. Yue, T. Huang, and Z. Lu. Clip4mc: An rl-friendly vision-language
465 model for minecraft, 2023.
- 466 [17] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas. Guiding
467 pretraining in reinforcement learning with large language models, 2023.
- 468 [18] K. Ehsani, T. Gupta, R. Hendrix, J. Salvador, L. Weihs, K.-H. Zeng, K. P. Singh, Y. Kim,
469 W. Han, A. Herrasti, R. Krishna, D. Schwenk, E. VanderBilt, and A. Kembhavi. Imitating
470 shortest paths in simulation enables effective navigation and manipulation in the real world,
471 2023.
- 472 [19] L. Fan, G. Wang, Y. Jiang, A. Mandlekar, Y. Yang, H. Zhu, A. Tang, D.-A. Huang, Y. Zhu,
473 and A. Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale
474 knowledge. In *Neural Information Processing Systems, 2022*, 2022.
- 475 [20] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world
476 models, 2023.
- 477 [21] J. Hu and R. Levy. Prompt-based methods may underestimate large language models’ linguistic
478 generalizations, 2023.
- 479 [22] C. Huang, O. Mees, A. Zeng, and W. Burgard. Visual language maps for robot navigation.
480 In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*,
481 London, UK, 2023.
- 482 [23] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:
483 Extracting actionable knowledge for embodied agents, 2022.
- 484 [24] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch,
485 Y. Chebotar, P. Sermanet, N. Brown, T. Jackson, L. Luu, S. Levine, K. Hausman, and B. Ichter.
486 Inner monologue: Embodied reasoning through planning with language models, 2022.
- 487 [25] G. Jawahar, B. Sagot, and D. Seddah. What does BERT learn about the structure of language?
488 In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,
489 Florence, Italy, 2019. Association for Computational Linguistics.
- 490 [26] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung.
491 Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38,
492 mar 2023.
- 493 [27] S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. Hatfield-
494 Dodds, N. DasSarma, E. Tran-Johnson, S. Johnston, S. El-Showk, A. Jones, N. Elhage, T. Hume,
495 A. Chen, Y. Bai, S. Bowman, S. Fort, D. Ganguli, D. Hernandez, J. Jacobson, J. Kernion,
496 S. Kravec, L. Lovitt, K. Ndousse, C. Olsson, S. Ringer, D. Amodei, T. Brown, J. Clark,
497 N. Joseph, B. Mann, S. McCandlish, C. Olah, and J. Kaplan. Language models (mostly) know
498 what they know, 2022.
- 499 [28] A. Kanervisto, S. Milani, K. Ramanaukas, N. Topin, Z. Lin, J. Li, J. Shi, D. Ye, Q. Fu, W. Yang,
500 W. Hong, Z. Huang, H. Chen, G. Zeng, Y. Lin, V. Micheli, E. Alonso, F. Fleuret, A. Nikulin,
501 Y. Belousov, O. Svidchenko, and A. Shpilman. Miner1 diamond 2021 competition: Overview,
502 results, and lessons learned, 2022.
- 503 [29] S. Karamcheti, S. Nair, A. Balakrishna, P. Liang, T. Kollar, and D. Sadigh. Prismatic vlms:
504 Investigating the design space of visually-conditioned language models, 2024.
- 505 [30] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C.
506 Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything, 2023.

- 507 [31] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement
508 learning, 2020.
- 509 [32] B. Z. Li, M. Nye, and J. Andreas. Implicit representations of meaning in neural language
510 models, 2021.
- 511 [33] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified
512 vision-language understanding and generation, 2022.
- 513 [34] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with
514 frozen image encoders and large language models, 2023.
- 515 [35] K. Li, A. K. Hopkins, D. Bau, F. Viégas, H. Pfister, and M. Wattenberg. Emergent world
516 representations: Exploring a sequence model trained on a synthetic task, 2023.
- 517 [36] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as
518 policies: Language model programs for embodied control, 2023.
- 519 [37] S. Lifshitz, K. Paster, H. Chan, J. Ba, and S. McIlraith. Steve-1: A generative model for
520 text-to-behavior in minecraft, 2023.
- 521 [38] H. Lin, Z. Wang, J. Ma, and Y. Liang. Mcu: A task-centric framework for open-ended agent
522 evaluation in minecraft, 2023.
- 523 [39] J. Lin, Y. Du, O. Watkins, D. Hafner, P. Abbeel, D. Klein, and A. Dragan. Learning to model
524 the world with language. 2023.
- 525 [40] H. Liu, L. Lee, K. Lee, and P. Abbeel. Instruction-following agents with multimodal transformer,
526 2023.
- 527 [41] H. Luo, A. Yue, Z.-W. Hong, and P. Agrawal. Stubborn: A strong baseline for indoor object
528 navigation, 2022.
- 529 [42] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data,
530 2021.
- 531 [43] A. Majumdar, K. Yadav, S. Arnaud, Y. J. Ma, C. Chen, S. Silwal, A. Jain, V.-P. Berges, P. Abbeel,
532 J. Malik, D. Batra, Y. Lin, O. Maksymets, A. Rajeswaran, and F. Meier. Where are we in the
533 search for an artificial visual cortex for embodied intelligence?, 2023.
- 534 [44] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over
535 unstructured data. In *Proceedings of the IEEE International Conference on Robotics and*
536 *Automation (ICRA)*, London, UK, 2023.
- 537 [45] V. Myers, A. He, K. Fang, H. Walke, P. Hansen-Estruch, C.-A. Cheng, M. Jalobeanu, A. Kolobov,
538 A. Dragan, and S. Levine. Goal representations for instruction following: A semi-supervised
539 language interface to control, 2023.
- 540 [46] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation
541 for robot manipulation, 2022.
- 542 [47] K. Narasimhan, R. Barzilay, and T. Jaakkola. Grounding language for transfer in deep reinforce-
543 ment learning, 2018.
- 544 [48] K. Nottingham, P. Ammanabrolu, A. Suhr, Y. Choi, H. Hajishirzi, S. Singh, and R. Fox. Do
545 embodied agents dream of pixelated sheep: Embodied decision making using language guided
546 world modelling, 2023.
- 547 [49] O.M.T., D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo,
548 T. Kreiman, Y. Tan, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot
549 policy. <https://octo-models.github.io>, 2023.
- 550 [50] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
551 K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder,
552 P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with
553 human feedback, 2022.

- 554 [51] N. D. Palo, A. Byravan, L. Hasenclever, M. Wulfmeier, N. Heess, and M. Riedmiller. Towards
555 a unified agent with foundation models, 2023.
- 556 [52] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
557 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from
558 natural language supervision, 2021.
- 559 [53] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-baselines3:
560 Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22
561 (268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- 562 [54] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. Turner, E. Un-
563 dersander, W. Galuba, A. Westbury, A. X. Chang, M. Savva, Y. Zhao, and D. Batra. Habitat-
564 matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai, 2021.
- 565 [55] R. Ramrakhya, E. Undersander, D. Batra, and A. Das. Habitat-web: Learning embodied
566 object-search strategies from human demonstrations at scale, 2022.
- 567 [56] R. Ramrakhya, D. Batra, E. Wijmans, and A. Das. Pirlnav: Pretraining with imitation and rl
568 finetuning for objectnav, 2023.
- 569 [57] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
570 with latent diffusion models, 2022.
- 571 [58] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun,
572 J. Malik, D. Parikh, and D. Batra. Habitat: A Platform for Embodied AI Research. In
573 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- 574 [59] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
575 algorithms, 2017.
- 576 [60] P. Sharma, A. Torralba, and J. Andreas. Skill induction and planning with latent language, 2022.
- 577 [61] P. Sharma, B. Sundaralingam, V. Blukis, C. Paxton, T. Hermans, A. Torralba, J. Andreas, and
578 D. Fox. Correcting robot plans with natural language feedback. In *Robotics: Science and*
579 *Systems, 2022*, 2023.
- 580 [62] X. Shi, I. Padhi, and K. Knight. Does string-based neural MT learn source syntax? In
581 *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*,
582 pages 1526–1534, Nov. 2016.
- 583 [63] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipula-
584 tion. In *Proceedings of the 5th Conference on Robot Learning (CoRL)*, 2021.
- 585 [64] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and
586 A. Garg. Progprompt: Generating situated robot task plans using large language models, 2022.
- 587 [65] A. Szot, M. Schwarzer, H. Agrawal, B. Mazouze, W. Talbott, K. Metcalf, N. Mackraz, D. Hjelm,
588 and A. Toshev. Large language models as generalizable policies for embodied tasks, 2024.
- 589 [66] I. Tenney, D. Das, and E. Pavlick. Bert rediscovers the classical nlp pipeline, 2019.
- 590 [67] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and
591 I. Polosukhin. Attention is all you need, 2017.
- 592 [68] S. Vemprala, R. Bonatti, A. Bucker, and A. Kapoor. Chatgpt for robotics: Design principles and
593 model abilities. Technical report, Microsoft, 2023.
- 594 [69] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. Voyager:
595 An open-ended embodied agent with large language models, 2023.
- 596 [70] Z. Wang, S. Cai, G. Chen, A. Liu, X. Ma, and Y. Liang. Describe, explain, plan and select:
597 Interactive planning with large language models enables open-world multi-task agents, 2023.

- 598 [71] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou.
599 Chain-of-thought prompting elicits reasoning in large language models, 2023.
- 600 [72] G. Wiedemann, S. Remus, A. Chawla, and C. Biemann. Does bert make any sense? interpretable
601 word sense disambiguation with contextualized embeddings, 2019.
- 602 [73] K. Yadav, J. Krantz, R. Ramrakhya, S. K. Ramakrishnan, J. Yang, A. Wang, J. Turner,
603 A. Gokaslan, V.-P. Berges, R. Mootaghi, O. Maksymets, A. X. Chang, M. Savva, A. Clegg,
604 D. S. Chaplot, and D. Batra. Habitat challenge 2023. [https://aihabitat.org/challenge/
605 2023/](https://aihabitat.org/challenge/2023/), 2023.
- 606 [74] K. Yadav, A. Majumdar, R. Ramrakhya, N. Yokoyama, A. Baevski, Z. Kira, O. Maksymets, and
607 D. Batra. Ovrl-v2: A simple state-of-art baseline for imagenav and objectnav, 2023.
- 608 [75] H. Yuan, C. Zhang, H. Wang, F. Xie, P. Cai, H. Dong, and Z. Lu. Plan4mc: Skill reinforcement
609 learning and planning for open-world minecraft tasks, 2023.
- 610 [76] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin,
611 D. Duong, V. Sindhwani, and J. Lee. Transporter networks: Rearranging the visual world for
612 robotic manipulation. *Conference on Robot Learning (CoRL)*, 2020.
- 613 [77] A. Zeng, M. Attarian, B. Ichter, K. Choromanski, A. Wong, S. Welker, F. Tombari, A. Purohit,
614 M. Ryoo, V. Sindhwani, J. Lee, V. Vanhoucke, and P. Florence. Socratic models: Composing
615 zero-shot multimodal reasoning with language, 2022.
- 616 [78] X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer. Sigmoid loss for language image pre-training,
617 2023.
- 618 [79] B. Zhou, K. Li, J. Jiang, and Z. Lu. Learning from visual observation via offline pretrained
619 state-to-go transformer, 2023.
- 620 [80] M. Zhu, Y. Li, and T. Kong. Integrating map-based method with end-to-end learning, 2022.
621 URL <https://www.youtube.com/watch?v=N-wW3TwEqbU>.
- 622 [81] X. Zhu, Y. Chen, H. Tian, C. Tao, W. Su, C. Yang, G. Huang, B. Li, L. Lu, X. Wang, Y. Qiao,
623 Z. Zhang, and J. Dai. Ghost in the minecraft: Generally capable agents for open-world
624 environments via large language models with text-based knowledge and memory, 2023.

Target Entity	Prompt	True Positive Rate	True Negative Rate
Spider	“Is there a spider in this image?”	22.27%	100.00%
	“Spiders in Minecraft are black. Is there a spider in this image?”	73.42%	94.54%
	“Spiders in Minecraft are black and have red eyes and long, thin legs. Is there a spider in this image?”	50.50%	99.85%
Cow	“Is there a cow in this image?”	71.00%	45.41%
	“Cows in Minecraft are black and white. Is there a cow in this image?”	98.22%	2.00%
	“Cows in Minecraft are black and white and have four legs. Is there a cow in this image?”	96.67%	7.35%
Sheep	“Is there a sheep in this image?”	88.00%	59.83%
	“Sheep in Minecraft are white. Is there a sheep in this image?”	100.00%	0.00%
	“Sheep in Minecraft are white and have four legs. Is there a sheep in this image?”	100.00%	0.00%

Table 5: InstructBLIP’s performance at decoding text indicating that it detected the presence of a target entity when given different prompts. We use this as a proxy metric for prompt engineering for RL, allowing us to determine which prompt to use for PR2L.

625 A Prompt Evaluation for RL in Minecraft

626 We discuss how to evaluate prompts to use with PR2L, by showcasing an example for a Minecraft
627 task. We start by noting that the presence and relative location of the entity of interest for each task
628 (i.e., spiders, sheep, or cows) are good features for the policy to have. To evaluate if a prompt elicits
629 these features from the VLM, we collect a small dataset of videos in which each Minecraft entity
630 of interest is on the left, right, middle, or not on screen for the entirety of the clip. Each video is
631 collected by a human player screen recording visual observations from Minecraft of the entity from
632 different angles for around 30 seconds at 30 frames per second (with the exception of the video where
633 the entity is not present, which is a minute long).

634 We propose prompts that target each of the two features we labeled. First, we evaluate prompts that
635 ask “Is there a(n) [entity] in this image?” As the answers to these questions are just yes/no, we see
636 how well the VLM can directly generate the correct answer for each frame in the collected videos.
637 The VLM should answer “yes” for frames in the three videos where the target entity is on the left,
638 right, or middle of the screen and “no” for the final video. Second, we evaluate if our prompts can
639 extract the entity’s relative position (left, right, or middle) in the videos where it is present. We
640 note that the prompts we tried could not extract this feature in the decoded text (e.g., asking “Is the
641 [entity] on the left, right, or middle of the screen?” will always cause the VLM to decode the same
642 text). Thus, we try to see if this feature can be extracted from the decoded texts’ representations. We
643 measure this by fitting a three-category linear classifier of the entity’s position given the *token-wise*
644 *mean* of the decoded tokens’ final embeddings. This is an unsophisticated and unexpressive classifier,
645 i.e., we do not have to worry about the model potentially memorizing the data, which means that
646 good classification performance corresponds to an easy extractability of said feature.

647 We evaluate three types of prompts per task entity for the first feature: one simply asking if the
648 entity is present in the image (e.g., “Is there a spider in this image?”) and two others adding varying
649 amounts of auxiliary information about visual characteristics of the entity (e.g., “Spiders in Minecraft
650 are black. Is there a spider in this image?” and “Spiders in Minecraft are black and have red eyes
651 and long, thin legs. Is there a spider in this image?”). We present evaluations of all such prompts in
652 Table 5. We find that the VLM benefits greatly from auxiliary information for the spider case only,
653 likely because spiders in Minecraft are the most dissimilar to the ones present in natural images of
654 real spiders, whereas cows and sheep are still comparatively similar, especially in terms of scale and
655 color. However, adding too much auxiliary information degrades performance, perhaps because the
656 input prompt becomes too long, and therefore is out-of-distribution for the types of prompts that
657 the VLM was pre-trained on. This same argument may explain why auxiliary information degrades

Tasks in MineDojo



Tasks in Habitat

ObjectNav: Find [toilet / sofa / bed].

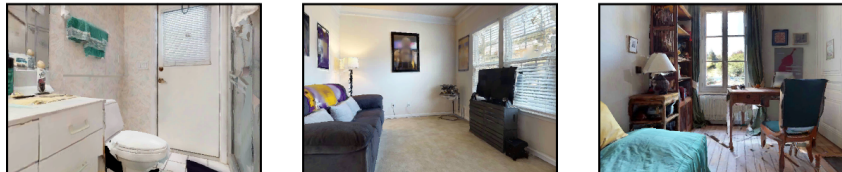


Figure 3: Example tasks, observations, and task-relevant prompts from MineDojo and Habitat.

658 performance for the other two target entities as well, causing them to almost always answer that
659 said entities are present, even when they are not. Once more, considering that these targets exhibit a
660 higher degree of visual resemblance to their real counterparts compared to Minecraft spiders, it is
661 reasonable to infer that the VLM would not benefit from auxiliary information. Furthermore, taking
662 into account that the auxiliary information we gave is more common-sense than the information given
663 for the spider, it could imply that the prompts are also more likely to be out-of-distribution (given
664 that “sheep are white” is so obvious that people would not bother expressing it in language), causing
665 the systematic performance degradation.

666 For the probing evaluation, we find that all three prompts reach similar final linear classifiabilities for
667 each of their target entities, as shown in Figure 4. While this does not aid in choosing one prompt
668 over another, it does confirm that the VLM’s decoded embeddings for each prompt still contains this
669 valuable and granular position information about the target entity, *even though the input prompt did*
670 *not ask for it.*

671 B MineDojo Details

672 B.1 Environment Details

673 **Spaces.** The observation space for the Minecraft tasks consists of the following:

- 674 1. **RGB:** Egocentric RGB images from the agent. (160, 256, 3)-size tensor of integers \in
675 $\{0, 1, \dots, 255\}$.
- 676 2. **Position:** Cartesian coordinates of agent in world frame. 3-element vector of floats.
- 677 3. **Pitch, Yaw:** Orientation of agent in world frame in degrees. Note that we limit the pitch
678 to 15° above the horizon to 75° below for *combat spider*, which makes learning easier (as
679 the agent otherwise often spends a significant amount of time looking straight up or down).
680 Two 1-element vectors of floats.
- 681 4. **Previous Action:** The previous action taken by the agent. Set to no operation at the start
682 of each episode. One-hot vector of size $|\mathcal{A}| = 53$ for *combat spider* and 89 otherwise (see
683 below).

684 This differs from the simplified observation space used in [19] in that we do not use any nearby voxel
685 label information and impose pitch limits for *combat spider*. This observation space is the same for
686 all Minecraft experiments.

687 The action space is discrete, consisting of 53 or
688 89 different actions:

- 689 1. **Turn:** Change the yaw and pitch of
690 the agent. The yaw and pitch can be
691 changed up to $\pm 90^\circ$ in multiples of
692 15° . As they can both be changed at
693 the same time, there are $9 \times 9 = 81$ total
694 different turning actions. The turning
695 action where the yaw and pitch changes
696 are both 0° is the no operation
697 action. Note that, since we impose pitch
698 limits for the spider task, we also limit
699 the change in pitch to $\pm 30^\circ$, meaning
700 there are only 45 turning actions in that
701 case.
- 702 2. **Move:** Move forward, backward, left,
703 right, jump up, or jump forward for 6
704 actions total.
- 705 3. **Attack:** Swing the held item at what-
706 ever is targeted at the center of the
707 agent’s view.
- 708 4. **Use Item:** Use the held item on what-
709 ever is targeted at the center of the
710 agent’s view. This is used to milk cows
711 or shear sheep (with an empty bucket
712 or shears respectively). If holding a
713 sword and shield, this action will block
714 attacks with the latter.

715 This non-*combat spider* action space is the same
716 as the simplified one in [19]. All experiments
717 for a given task share the same action space.

718 **World specifications.** MineDojo implements
719 a fast reset functionality that we use. Instead
720 of generating an entirely new world for each
721 episode, fast reset simply respawns the player
722 and all specified entities in the same world instance,
723 but with fully restored items, health points,
724 and other relevant task quantities. This lowers the
725 time overhead of resets significantly, but also
726 means that some changes to the world (like block
727 destruction) are persistent. However, as breaking
728 blocks generally takes multiple time steps of taking
729 the same action (and does not directly lead to any
730 reward), the agent empirically does not break many
731 blocks aside from tall grass (which is destroyed
732 with a single strike from any held item). We keep
733 all reset parameters (like the agent respawn radius,
734 how far away entities can spawn from the agent,
735 etc) at their default values provided by MineDojo.

729 We stage all tasks in the same area of the same
730 programmatically-generated world: namely, a
731 sunflower plains biome in the world with seed 123.
732 This is the default location for the implementation
733 of the spider combat task in MineDojo. We choose
734 this specific world/location as it represents a
735 prototypical Minecraft scene with relatively easily-
736 traversable terrain (thus making learning faster
737 and easier).

734 **Additional task details and reward functions.**
735 We provide additional notes about our Minecraft
736 tasks.

736 *Combat spider:* Upon detecting the agent, the
737 spider approaches and attacks; if the agent’s
738 health is depleted, then the episode terminates in
739 failure. The agent receives +1 reward for striking
740 any entity and +10 for defeating the spider. We
741 also include several distractor animals (a cow, pig,
742 chicken, and sheep) that passively wander the task
743 space; the agent can reward game by striking these
744 animals, making credit assignment of success rewards
745 and the overall task harder.

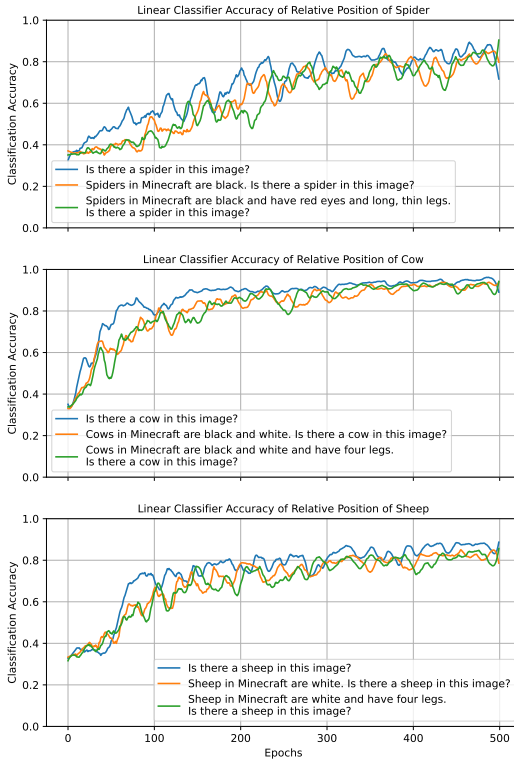


Figure 4: We train a linear classifier to predict the relative position of the target entity (left/right/middle) based on the average VLM embeddings decoded in response to each associated candidate prompt. We find that all three candidate prompts per task elicit embeddings that are similarly highly conducive to this classification scheme.

Hyperparameter	Task					
	<i>Combat Spider</i>	<i>Milk Cow</i>	<i>Shear Sheep</i>	<i>Combat Zombie</i>	<i>Combat Enderman</i>	<i>Combat Pigman</i>
Total Train Steps	150000			100000		
Rollout Steps				2048		
Action Entropy Coefficient				5e-3		
Value Function Coefficient				0.5		
Max LR	5e-5	1e-4	1e-4	5e-5	1e-4	5e-5
Min LR	5e-6	1e-4	1e-4	5e-6	1e-4	5e-6
Batch Size				64		
Update Epochs				10		
γ				0.99		
GAE λ				0.95		
Clip Range				0.2		
Max Gradient Norm				0.5		
Normalize Advantage				True		

Table 6: PPO hyperparameters for Minecraft tasks, shared by the baselines, our method, and ablations.

Policy Transformer Hyperparameters	
Transformer Token Size	512 / 128
Transformer Feedforward Dim	512 / 128
Transformer Number Heads	2
Transformer Number Decoder Layers	1
Transformer Number Encoder Layers	1
Transformer Output Dim	128
Transformer Dropout	0.1
Transformer Nonlinearity	ReLU
Policy MLP Hyperparameters	
Number Hidden Layers	1
Hidden Layer Size	128
Activation Function	tanh
VLM Generation Hyperparameters	
Max Tokens Generated	6
Min Tokens Generated	6
Decoding Scheme	Greedy

Table 7: All policy hyperparameters for all Minecraft tasks. Smaller token sizes and feedforward dimensions are used for *combat [zombie/enderman/pigman]*.

741 *Milk cow*: The agent also holds wheat in its off hand, which causes the cow to approach the agent
742 when detected and sufficiently nearby. For each episode, we track the minimum visually-observed
743 distance between the agent and the cow at each time step. The agent receives $+0.1|\Delta d_{\min}|$ reward for
744 decreasing this minimum distance (where $\Delta d_{\min} \leq 0$ is the change in this minimum distance at a
745 given time step) and $+10$ for successfully milking the cow.

746 *Shear sheep*: As with *milk cow*, the agent holds wheat in its off hand to cause the sheep to approach
747 it. The reward function also has the same structure as that task, albeit with different coefficients:
748 $+|\Delta d_{\min}|$ for decreasing the minimum distance to the sheep and $+10$ for shearing it.

749 *Combat zombie*: Same as *combat spider*, but the enemy is a zombie. We increase the episode length
750 to 1000, as the zombie has more health points than the spider.

751 *Combat enderman*: Same as *combat spider*, but the enemy is an Enderman. As with *combat zombie*,
752 we increase the episode length to 1000. Note that Endermen are non-hostile (until directly looked at
753 for sufficiently long or attacked) and have significantly more health points than other enemies. We
754 thus enchant the agent’s sword to deal more damage and decrease the initial spawn distance of the
755 enderman from the agent.

756 *Combat pigman*: Same as *combat spider*, but the enemy is a hostile zombie pigman. As with *combat*
757 *zombie*, we increase the episode length to 1000.

758 B.2 Policy and Training Details

759 For our actual RL algorithm, we use the Stable-Baselines3 (version 2.0.0) implementation of clipping-
760 based PPO [53], with hyperparameters presented in Table 6. Many of these parameters are the same
761 as the ones presented by [19]. For the spider trials, we use a cosine learning rate schedule:

$$\text{LR}(\text{current train step}) = \text{Min LR} + (\text{Max LR} - \text{Min LR}) \left(\frac{1 + \cos\left(\pi \frac{\text{current train step}}{\text{total train steps}}\right)}{2} \right) \quad (1)$$

762 We also present the policy and VLM hyperparameters in Table 7. The hyperparameters and architec-
763 ture of the MLP part of the policy are primarily defined by the default values and structure defined by
764 the Stable-Baselines3 ActorCriticPolicy class. Note that the no generation ablation, VLM image
765 encoder baseline, and MineCLIP trials do not generate text with the VLM, and so all do not use the
766 associated process’s hyperparameters. The MineCLIP trials also do not use a Transformer layer in
767 the policy, due to not receiving token sequence embeddings. It instead just uses a MLP, but with two
768 hidden layers (to supplement the lowered policy expressivity due to the lack of a Transformer layer).

769 Additionally, InstructBLIP’s token embeddings are larger than ViT-g/14’s (used in the VLM image
770 encoder baseline), and so may carry more information. However, the VLM does not receive any
771 privileged information over the image encoder *from the task environment* – any additional information
772 in the VLM’s representations is therefore purely from the model’s prompted internal knowledge. Still,
773 to ensure consistent policy expressivity, we include a learned linear layer projecting all representations
774 for this baseline and our approach to the same size (512 dimensions) so that the rest of the policy is
775 the same for both.

776 Minecraft training runs were run on 16 A5000 GPUs (to accommodate the 16 seeds).

777 C Habitat ObjectNav Details

778 C.1 Environment Details

779 The spaces and agent/task specifications are largely the same as the defaults provided by Habitat, as
780 specified in the HM3D ObjectNav configuration file [58].

781 **Spaces.** The observation space for Habitat consists of the following:

- 782 1. **RGB:** Egocentric RGB images from the agent. (480, 640, 3)-size tensor of integers \in
783 $\{0, 1, \dots, 255\}$. By default, agents also receive depth images, but we remove them to ensure
784 that state representations are grounded primarily in visual observations.
- 785 2. **Position:** Horizontal Cartesian coordinates of agent. 2-element vector of floats.
- 786 3. **Compass:** Yaw of the agent. Single floats.
- 787 4. **Previous Action:** The previous action taken by the agent. Set to no operation at the start of
788 each episode. One-hot vector of size $|\mathcal{A}| = 4$.
- 789 5. **Object Goal:** Which object the agent is aiming to find. One-hot vector of size 3.

790 The action space is the standard Habitat-Lab action space, though we remove the pitch-changing
791 actions, leaving only four:

- 792 1. **Turn:** Turn left or right, changing the yaw by 30° .
- 793 2. **Move Forward:** Move forward a fixed amount or until the agent collides with something.
- 794 3. **Stop:** Ends the episode, indicating that the agent believes it has found the goal object.

795 All observations, actions, and associated dynamics are deterministic.

796 **World specifications.** In ObjectNav, an agent is spawned in a household environment and must find
797 and navigate to an instance of a specified target object in as efficient a path as possible. Doing so
798 effectively requires a common-sense understanding of where household objects are often found and
799 the structure of standard homes.

800 Habitat provides a standardized train-validation split, consisting of 80 household scenes for training
801 (from which one can run online RL or collect data for offline RL or BC) and 20 novel scenes
802 for validation, thereby testing policies’ generalization capabilities. These scenes come from the
803 Habitat-Matterport 3D v1 dataset [54].

804 C.2 Policy and Training Details

805 In line with previous work [56, 74, 43], we train our policies with behavior cloning (BC) on the
806 Habitat-Web human demonstration dataset of 77k trajectories (12M steps) [55]. We adopt many of
807 the same design choices provided by said prior works, but with a few critical differences:

- 808 1. Due to compute limitations, we were unable to train on the full dataset (as those original
809 works used 512 parallel environments to roll out demo trajectories and collect data). Instead,
810 we used a subset of the dataset, built by dividing the dataset by both target object and scene,
811 then sampling every tenth demo. This would ensure that our training data still contained
812 examples from every training scene + target object combination that existed. In total, our
813 subsampled dataset contains approximately 1.1M steps over 7550 trajectories.
- 814 2. We adopt the same optimizer, scheduler, and associated hyperparameters as Majumdar et al.
815 [43], but find a learning rate of $1e - 4$ to be more effective than their $1e - 3$.
- 816 3. Rather than sampling partial trajectory rollouts from 512 parallel environments as done by
817 Majumdar et al. [43], our batches contain full trajectories, though with the same total number
818 of transitions per batch as in that work. This means that our batches potentially contain less
819 diverse data (due to observations from fewer different total scenes being present), but allow
820 us to compute up-to-date full trajectory hidden states for the RNN portion of our policy. We
821 use gradient accumulation to achieve this, once again due to compute limitations.
- 822 4. While Majumdar et al. [43] trains for 24k gradient steps (observing approximately 400M
823 transitions.), we find using only approximately a tenth of that (40 epochs through our smaller
824 dataset, so around 40M transitions) to reach peak performance for our policy. The scheduler
825 still assumes the full training run will last for 400M transitions, so our LR decays at the
826 same rate as with VC-1. Furthermore, for fairness, we leave our VC-1 baseline policies
827 (trained on our subsampled datasets) training beyond 40 epochs, and report their validation
828 performance at both 40 and 120 epochs (when its performance saturates).
- 829 5. For policies that receive visual observations as a sequence of tokens (PR2L, VC-1 with
830 patch embeddings), we apply 2D average pooling with kernel sizes of 4×4 to reduce down
831 to 16 tokens. Then, we pass those tokens through a learned Transformer layer, instead of the
832 learned compression layer used by Majumdar et al. [43]. We do this to ensure that policy
833 performance differences are due to representation quality, not architecture.
- 834 6. We employ inflection upweighting during training, as done by Ramrakhya et al. [56], Yadav
835 et al. [74], Majumdar et al. [43]. However, we also categorically upweight the cross entropy
836 loss of stopping and turning by 1.5 (due to them being uncommon but important), as we
837 observe this increases learning speed for all policies.
- 838 7. We do not employ any image augmentation or loss regularization to prevent overfitting.
839 However, we find our policy exhibits strong generalization performance in unseen validation
840 scenes nonetheless.

841 For PR2L-specific design choices:

- 842 1. Our chosen VLM is the Prismatic VLM [29] with Dino+SigLIP as a vision backbone and
843 Llama2-7B-pure as the language backbone. We use the 224px version, which maps images
844 to 256 visual tokens (which, as described above, get compressed into 16 via pooling).
- 845 2. To reduce the size of VLM representations for PR2L, we embed one observation (sampled
846 uniformly at random) from each trajectory in our subsampled dataset with our VLM, then
847 compute all resulting tokens’ principle component vectors. We then use said vectors to
848 lower all tokens’ dimensionality down from 4096 to 1024 (i.e., corresponding approximately
849 to their first 1024 principle components).
- 850 3. Like with the Minecraft experiments, we take the VLM’s last two layers’ embeddings and
851 treat them as our promptable representations. However, unlike with Minecraft, we stack

852 each VLM token’s two embeddings (forming new embeddings of size 2048), rather than
853 concatenate all of them.

854 4. For generating text in response to our task-relevant prompt, we use sample-based decoding
855 with fixed random seed prior to the decoding with temperature 0.4 and 32 – 48 new tokens
856 generated.

857 5. The learned Transformer layer of our policy is the same as the one used in the Minecraft
858 experiments, but with token embedding sizes of 1024.

859 All Habitat training was done on an A100 GPU server. Generation of data and evaluations were done
860 on 16 A5000 GPUs for parallelization.

861 **D Simplified Habitat Offline RL Experiments**

862 While our primary Habitat experiments use behavior cloning to stay consistent with past works, we
863 also run offline RL experiments on a simplified version of ObjectNav to better explore how VLM
864 representations aid action learning. We discuss the details of said setting now.

865 **D.1 Environment Details**

866 We pick 32 reconstructed 3D home environments with at least one instance of each of the three target
867 objects (toilet, bed, and sofa) and an annotator quality score of at least 4 out of 5. We choose to
868 remove *plants* and *televisions* from the goal object set due to finding numerous unlabeled instances
869 of them. Additionally, we remove chairs, as they are significantly more common than other goal
870 objects and thus usually can be found in much shorter episodes. This simplified problem formulation
871 enables us to remove many of the “tricks” that aid ObjectNav, such as using omnidirectional views or
872 policies with history; our agent makes action decisions purely based on its current visual observation
873 and pose, allowing us to do “vanilla” RL to better isolate the effect of PR2L.

874 To generate data, we use Habitat’s built-in greedy shortest geodesic path follower. Imitating such
875 demonstrations allows policies to learn unintuitively emergent and performant navigation behaviors
876 [18] at scale. For each defined starting location in our considered households, we autonomously
877 collect data by using the path follower to navigate to each reachable instance of the corresponding
878 goal object. This yields high quality, near-optimal data. We then supplement our dataset by generating
879 lower-quality data. Specifically, for each computed near-optimal path from a starting location to a
880 goal object instance, we choose to inject action noise partway through the trajectory (uniformly at
881 random from 0 – 90% of the way through). At that point, all subsequent actions have a 0 – 50%
882 probability (again chosen uniformly at random) of being a random action other than the one specified
883 by the path follower. To ensure that paths are sufficiently long, we choose to make the probability of
884 choosing the stop action 10% and the other two movement actions 45%. In total, we collect 107518
885 observations over 2364 trajectories.

886 **Reward functions.** The ObjectNav challenge evaluates agents based on the average "success
887 weighted by path length" (SPL) metric [73]: if an agent succeeds at taking the *stop* action while close
888 to an instance of the goal object, it gets $SPL(p, l) = \frac{l}{\max(l, p)}$ points, where l is the actual shortest
889 path from the starting point to an instance of the goal object and p is the length of the path that the
890 agent actually took during that particular episode. If the agent stops while not close to the target
891 object, the SPL is 0. Thus, taking the most efficient path to the nearest goal object and stopping yields
892 a maximum SPL of 1.

893 We use this to design our reward function. Specifically, when the agent stops, it receives a reward
894 of $+10SPL(p, l)$. Additionally, we add a shaping reward of the change in geodesic distance to the
895 nearest goal object instance each time the agent moves (where lowering that distance yields a positive
896 reward).

897 **D.2 Policy and Training Details**

898 For our offline RL experiments in Habitat, we use Conservative Q-Learning (CQL) on top of
899 the Stable-Baselines3 Contrib codebase’s implementation of Quantile Regression DQN (QR-DQN)
900 [31, 13]. We choose to multiply the QR-DQN component of the CQL loss by 0.2. Using the notation

901 proposed by Kumar et al. [31], this is equivalent to $\alpha = 5$, which said work also uses. Other
 902 hyperparameters are $\tau = 1$, $\gamma = 0.99$, fixed learning rate of $1e - 4$, 100 epochs, and 50 quantiles (no
 903 exploration hyperparameters are specified, since we do not generate any new online data).

904 The policy architecture used for Habitat experiments are the same as those used for PPO, though the
 905 final network outputs quantile Q-values for each action (rather than just a distribution over actions).
 906 The action with the highest mean quantile value is chosen at evaluation time.

907 During training, we shuffle the data and load full offline trajectories until the buffer has at least
 908 $32 \times 1024 = 32768$ transitions or all trajectories have been loaded once that epoch. We then uniformly
 909 sample and train on batches of size 512 transitions from the buffer until each transition has been
 910 trained on once in expectation (e.g., $\sim \frac{\text{number of transitions in the buffer}}{512}$ batches). Each batch is used for 8
 911 gradient steps before the next is sampled. We choose this data loading scheme to fit the training
 912 infrastructure provided by Stable-Baselines3 while not using up too much memory at once.

913 D.3 Experiments and Results

914 Our primary comparison is once again between our promptable representations and general-purpose
 915 non-promptable ones. We thus repeat the baseline described previously for Minecraft in Section 4.1,
 916 training a single agent for all three ObjectNav tasks using both PR2L and the VLM image encoder
 917 representations. We empirically note that longer visual embedding sequences tend to perform better in
 918 Habitat. To control for this, we opt to use InstructBLIP’s Q-Former unprompted embeddings instead
 919 of the ViT embeddings directly (which are much longer than PR2L’s embedding sequences). As
 920 InstructBLIP uses the former representations to extract visual features to be projected into language
 921 embedding space, this serves to close the gap in embedding sequence length between our two
 922 conditions while still providing us with general visual features that the VLM processes via prompting.
 923 In this case, we use the same InstructBLIP model as the Minecraft experiments and choose “What
 924 room is this?” as our task-relevant prompt.

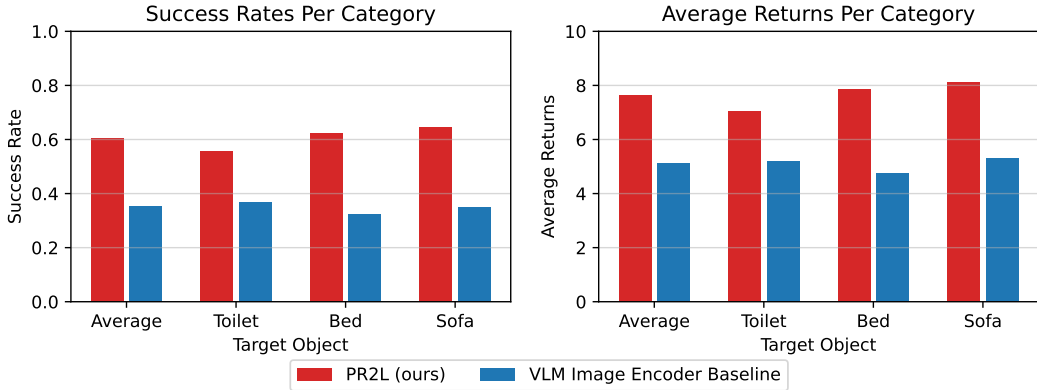


Figure 5: **Offline RL performance of PR2L and baselines in Habitat ObjectNav.** Plots show final evaluation success rates and average returns per target object and overall. PR2L outperforms the baseline in all cases.

925 We report evaluation success rates and average returns for the simplified Habitat ObjectNav setting in
 926 Figure 5. PR2L achieves nearly double the average success rate of the baseline (60.4% vs. 35.2%),
 927 supporting the hypothesis that PR2L works especially well when exploration is not needed. Lastly, in
 928 Appendix H.2, we find that PR2L causes the VLM to produce highly structured representations that
 929 correlate with an expert policy’s value function: high-value states are typically labeled by the VLM
 930 as being from a room where one would expect to find the target object.

931 E Extended Discussion of Tasks and Results

932 E.1 Notes on Task-specific Systems

933 We designed experiments to specifically investigate the use of VLM embeddings as task-specific
 934 promptable representations for downstream sensorimotor policy learning. As such, we compare with

Task	PR2L (Ours)	VLM Image Encoder	Ablations			
			No Prompt	No Generation	Change Aux. Text	Oracle Detector
<i>Combat Spider</i>	97.6 ± 14.9	51.2 ± 9.3	72.6 ± 14.2	66.6 ± 11.8	80.1 ± 12.6	58.0 ± 13.4
<i>Milk Cow</i>	223.4 ± 35.4	95.2 ± 18.7	116.6 ± 25.9	160.2 ± 23.6	80.5 ± 17.8	178.4 ± 42.5
<i>Shear Sheep</i>	37.0 ± 4.4	23.0 ± 3.6	23.8 ± 3.2	26.1 ± 4.5	27.8 ± 4.6	27.4 ± 9.3

Table 8: Minecraft ablations, VLM image encoder baseline, and our full approach. All achieve worse performance than PR2L. Values are final IQM success counts and intervals are the standard error.

935 other works that propose or evaluate either learning from scratch or from pre-trained representations,
936 but *not* to systems in Minecraft and Habitat that require domain-specific engineered systems beyond
937 just policy learning (such as Luo et al. [41], Zhu et al. [80]) or which target learning or producing
938 higher-level plans or abstractions (such as Wang et al. [70]).

939 Such comparisons are not made as these works either aim to investigate other problems in control or
940 are aiming to develop highly specialized and task-specific systems (whereas we present a general
941 approach for policy learning). For instance, Voyager shows how an LLM can reason about and
942 compose high-level hand-crafted control primitives [69]. Voyager’s ability to complete harder tasks
943 comes from its access to powerful hand-crafted high-level primitives that extensively leverage oracle
944 information, which are composed into skills by GPT-4 (which does not handle any low-level control).
945 Said hand-coded control primitives used in Voyager are very advanced and do much of the heavy-
946 lifting. In particular, Voyager gives GPT-4 access to a dedicated `killMob(<entity name>)` control
947 primitive function. This function calls a separate `bot.pvp.attack(<entity>)` (hand-written)
948 function, which calls a hard-coded oracle pathfinder, aiming controller, and attack function to
949 repeatedly approach and attack the specified entity until it is defeated. Thus, for Voyager, the skill for
950 hunting sheep simply fills in the powerful `killMob()` primitive function with “sheep” as the target,
951 abstracting away all low-level control via the oracle hand-written controllers.

952 Vitaly, unlike PR2L, Voyager does not investigate how to use (V)LMs to learn these primitives. It
953 thus cannot be applied to settings that lack such primitives (e.g., because oracle path planners are
954 not available, like in Habitat). This makes PR2L complementary: we directly learn a policy to link
955 observations to low-level actions (turning, moving, attacking, etc) via RL with no oracle information,
956 while Voyager aims to compose pre-existing primitives into skills via LLMs.

957 E.2 Notes on Dreamer v3

958 We note that PR2L just proposes to use VLMs as a source of task-specific representations for RL
959 tasks; it does not prescribe which learning algorithm to use. Therefore, in principle, one could
960 use Dreamer in conjunction with PR2L and gain benefits from both the VLM representation and
961 the choice of a strong model-based RL algorithm. However, while we leave this to future works,
962 our Minecraft comparison (c) measures how well the approach does on our Minecraft tasks (as the
963 original paper focuses more on the component subtasks involved in the *find diamond* task, all of
964 which do not involve interacting with moving entities).

965 We find that Dreamer v3 is unable to learn our six tasks given the same number of environment interac-
966 tions that PR2L+PPO was trained on. We hypothesize that this is due to its visual reconstruction-based
967 world model not being suited for tasks requiring interaction with partially-observable, non-stationary
968 autonomous entities (which all our tasks involve). We note that the last two rows of the figure
969 visualizing model reconstructions in the original Dreamer v3 paper shows that its world model
970 fails to reconstruct an observed pig [20], supporting our hypothesis. This highlights the need for
971 robust representations that are conducive to world model learning, with PR2L’s capabilities to elicit
972 task-relevant visual semantic features via prompting being one possibility for doing so.

973 F Ablations

974 We run four ablations on *combat spider*, *milk cow*, and *shear sheep* to isolate and understand the
975 importance of various components of PR2L. First, we run PR2L with *no prompt* to see if prompting
976 with task context actually tailors the VLM’s generated representations favorably towards the target
977 task, improving over an unprompted VLM. Note that this is not the same as just using the image
978 encoder (comparison (a)), as this ablation still decodes through the VLM, just with an empty prompt.
979 Second, we run PR2L with our chosen prompt, but *no generation* of text – i.e., the policy only

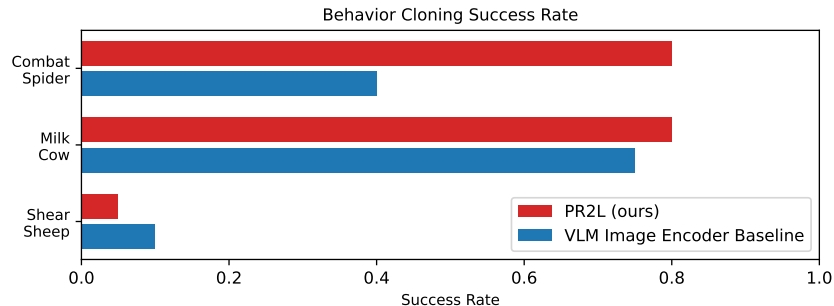


Figure 6: Success rates for BC on either PR2L or VLM image encoder baseline representations for all original tasks. PR2L excels at *combat spider*, even after the policy is trained for a single epoch.

980 receives the embeddings associated with the image and prompt (the left and middle red groupings
 981 of tokens in Figure 2, but not the right-most group). This tests the hypothesis that representations
 982 of generated text might make certain task-relevant features more salient: e.g., the embeddings for
 983 “Is there a cow in this image?”, might not encode the presence of a cow as clearly as if the VLM
 984 generates “Yes” in response, impacting downstream performance. Third, to check if our prompt
 985 evaluation strategy provides a good proxy for downstream task performance while tuning prompts
 986 for P2RL, we run PR2L with alternative prompts that were not predicted to be the best, as per our
 987 criterion in Appendix A. We thus remove the auxiliary text from the prompt for *combat spider* and
 988 add it for *milk cow* and *shear sheep*. Lastly, to see if PR2L embeddings are just better due to them
 989 encoding entity detection, we train a VLM image encoder policy with an additional ground truth
 990 oracle target entity detector as a feature.

991 Results from these additional experiments are presented in Table 8. In general, all ablations perform
 992 worse than PR2L. For *milk cow*, we note the most performant ablation is no generation, perhaps
 993 because the generated text is often wrong; among the chosen prompts, it yields the lowest true
 994 positive and negative rates for classifying the presence of its corresponding target entity (see Table 5
 995 in Appendix A), though adding auxiliary text makes it even worse, perhaps explaining why *milk cow*
 996 experienced the largest performance decrease from adding it back in. Based on these overall trends,
 997 we conclude that (i) the *promptable* and *generative* aspects of VLM representations are important for
 998 extracting good features for control tasks and (ii) our simple evaluation scheme is an effective proxy
 999 measure of how good a prompt is for PR2L.

1000 G Minecraft Behavior Cloning Experiments

1001 We collected expert policy data by training a policy on MineCLIP embeddings to completion on all
 1002 of our original tasks and saving all transitions to create an offline dataset. We then embedded all
 1003 transitions with either PR2L or the VLM image encoder. Finally, we train policies with behavior
 1004 cloning (BC) on successful trajectories under a specified length (300 for *combat spider*, 250 for *milk*
 1005 *cow*, and 500 for *shear sheep*) from either set of embeddings for all three tasks, then evaluate their
 1006 task success rates.

1007 Results are presented in Figure 6. We first note that, since the expert data was collected from a policy
 1008 trained on MineCLIP embeddings, the *shear sheep* policy is not very effective (as we found in Table
 1009 2). Both resulting *shear sheep* BC policies are likewise not very performant. We find that *combat*
 1010 *spider* in particular shows a very large gap in performance: the PR2L agent achieves approximately
 1011 twice the success rate of the VLM image encoder agent *after training for just a single epoch*. The
 1012 comparatively small amount of training and data necessary to achieve near-expert performance for
 1013 this task supports our hypothesis that promptable representations from general-purpose VLMs do
 1014 not help with exploration (they work better in offline cases, where exploration is not a problem), but
 1015 instead are particularly conducive to being linked to appropriate actions even though the VLM is not
 1016 producing actions itself. Further investigation of this hypothesis is presented in Appendix H.

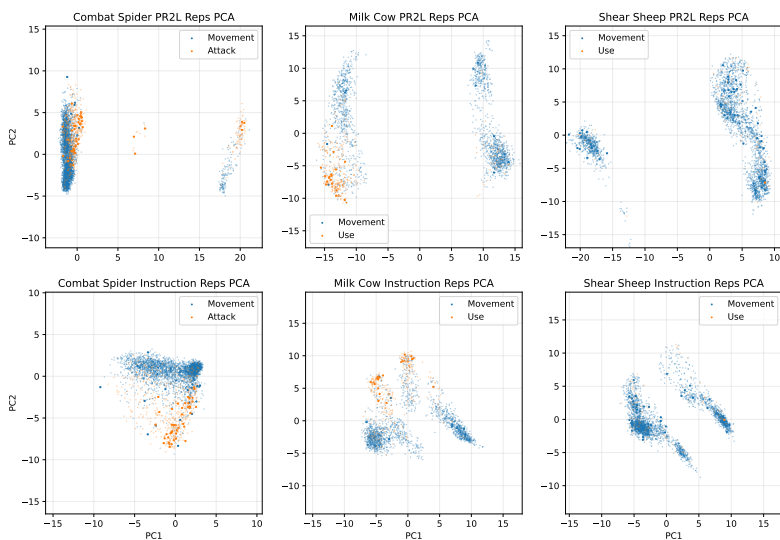


Figure 7: **PCA of PR2L representations of observations from twenty episode rollouts of expert policies in all three Minecraft tasks.** Larger points correspond to transitions where the expert received > 0.1 reward. We vary the prompt to be either our task-relevant prompt or the RT-2-style baseline instruction prompt. Our prompt’s representations are bi-modal, with the clusters on the left corresponding to the VLM outputting “yes” (the entity is in view). We find that most functional actions (orange points) that yielded rewards are located in said clusters. Note that, since these expert policies are trained on top of MineCLIP embeddings, the *shear sheep* policy is not very performant, as seen in Table 2.

1017 H Representation Analysis

1018 Why do our prompts yield higher performance than one asking for actions or instruction-following?
 1019 Intuitively, despite appropriate responses to our task-relevant prompts not directly encoding actions,
 1020 there should be a strong correlation: e.g., when fighting a spider, if the spider is in view and the
 1021 VLM detects this, then a good policy should know to attack to get rewards. We therefore wish to
 1022 investigate if our representations are conducive to easily deciding when certain rewarding actions
 1023 would be appropriate for a given task – if it is, then such a policy may be more easily learned by RL,
 1024 which would explain PR2L’s improved performance over the baselines.

1025 H.1 Minecraft Analysis

1026 To investigate this, we use the embeddings of our offline data from the BC experiments (collected
 1027 by training a MineCLIP encoder policy to high performance on all of our original three tasks, as
 1028 discussed in Appendix G). We specifically look at the embeddings produced by a VLM when given
 1029 our standard task-relevant prompts and when given the instructions used for our RT-2-style baseline.
 1030 We then perform principal component analysis (PCA) on the tokenwise average of all embeddings
 1031 for each observation, thereby projecting the embeddings to a 2D space with maximum variance.

1032 We visualize these low-dimensional space in Figure 7 for the final 20 successful observations from
 1033 each task, with the point colors of orange and blue respectively indicating whether the observation
 1034 results in a functional action (attack or use item) or movement (translation or rotation) by the expert
 1035 policy. Additionally, we enlarge points corresponding to when the agent received rewards in order to
 1036 recognize which actions aided in or achieved the task objective.

1037 We find that our considered prompts resulted in a bimodal distribution over representations, wherein
 1038 the left-side cluster corresponds to the VLM outputting “yes (the entity is in view)” and the right-side
 1039 one corresponds to “no.” Additionally, observations resulting in functional actions that received

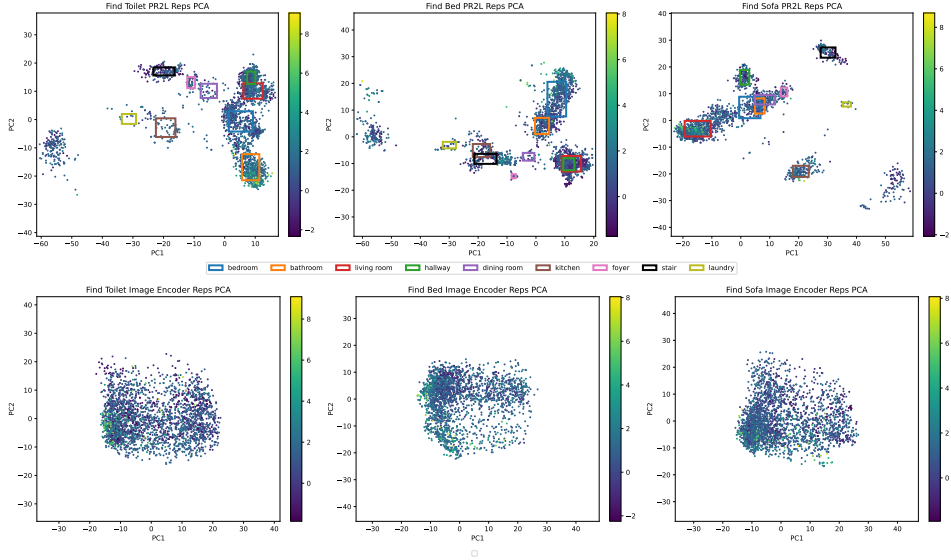


Figure 8: **PCA of PR2L (above) and image encoder (below) representations of observations from thirty episode rollouts of expert policies in all Habitat tasks.** The points’ colors correspond to their value under Habitat’s built-in oracle shortest path follower (a near-optimal policy). More yellow is better. Boxes correspond to points the VLM has labeled as a given household room, in response to the task prompt of “What room is this?” This analysis aligns with intuition: for *find toilet*, high value observations tend to be labeled as bathrooms (orange box), *find bed*’s tend to be labeled as bedrooms (blue), and *find sofa*’s are labeled as living rooms (red).

1040 rewards (large orange points in Figure 7) tend to be on the left-side (“yes”) cluster for representations
 1041 elicited by our prompt, but are more widely distributed in the instruction prompt case, in agreement
 1042 with intuition. This is especially clear in the *milk cow* plot, wherein nearly all rewarding functional
 1043 actions (using the bucket on the cow to successfully collect milk) are in the lower left corner.

1044 This analysis supports that the representations yielded by InstructBLIP in response to our chosen style
 1045 of prompts are more structured than representations from instructions. Such structure is useful in
 1046 identifying and learning rewarding actions, even when said actions were taken from an expert policy
 1047 trained on unrelated embeddings. This suggests that such representations may similarly be more
 1048 conducive to being mapped to good actions via RL, which we observe empirically (as our prompt’s
 1049 representations yield more performant policies than the instructions for the RT-2-style baseline).

1050 H.2 Habitat Analysis

1051 Likewise, we conduct a similar analysis on the Habitat data from our simplified setting. Specifically,
 1052 we wish to see if PR2L produces representations that are conducive to extracting the *value function*
 1053 of a good policy. Since the chosen Habitat ObjectNav prompt is “What room is this?” we expect the
 1054 state representations to be clustered based on room categories. Intuitively, states corresponding to the
 1055 room one is likely to find the target object should have the highest values.

1056 As shown in Figure 8, we thus used PCA to project expert trajectories’ PR2L and general image
 1057 encoder state representations (generated with Habitat’s geodesic shortest path follower) to two
 1058 dimensions, then colored each one based on their value under said near-optimal policy. We also
 1059 plotted the mean and standard deviation of all points labeled as each room, visualizing them as
 1060 axis-aligned bounding boxes. Note that each upper subplot in Figure 8 has a cluster of points far from
 1061 all boxes. These correspond to the VLM generating nothing or garbage data with no room label.

1062 This visualization qualitatively agrees with intuition. High value states tend to be grouped with the
 1063 room the corresponding target object is often found in: *find toilet* corresponds to bathrooms, *find bed*
 1064 to bedrooms, and *find sofa* to living rooms. Comparatively, the general image encoder features do
 1065 not have such semantically meaningful groupings; all observations are clustered together and, within

1066 that single grouping, high-value observations are more spread out. This all supports the idea that
1067 prompting allows representations to take on structures that correlate well to value functions of good
1068 policies.

1069 I Code Snippets

1070 We provide some code snippets showcasing instantiations of PR2L.

```
1071 class Policy(torch.nn.Module):
1072     def __init__(self, num_actions, tf_embed_dim=4096):
1073         """Policy that accepts promptable reps as input"""
1074         super().__init__()
1075         # Project down VLM embed dimensions
1076         self.embed_fc = torch.nn.Linear(tf_embed_dim, 1024)
1077         # Predict actions
1078         self.action_fc = torch.nn.Linear(1024, num_actions)
1079         # Transformer layer to condense promptable reps to 1 token
1080         self.transformer = torch.nn.Transformer(
1081             1024,
1082             1,
1083             num_encoder_layers=1,
1084             num_decoder_layers=1,
1085             dim_feedforward=1024,
1086             batch_first=True,
1087         )
1088         self.cls = torch.nn.Embedding(1, 1024) # cls tokens
1089
1090     def forward(self, x):
1091         seq, mask = x
1092         bs, traj_len, num_tokens, _ = seq.shape
1093
1094         # [batch*traj_len, num tokens, token size]
1095         seq = seq.reshape(bs * traj_len, num_tokens, -1)
1096         # [batch*traj_len, num tokens]
1097         mask = mask.reshape(bs * traj_len, num_tokens)
1098
1099         # Project down
1100         # [batch*traj_len, num tokens, tf dim]
1101         seq = self.embed_fc(seq)
1102
1103         # Get CLS embedding
1104         cls = self.cls(torch.zeros([bs * traj_len, 1],
1105                                   device=seq.device, dtype=int))
1106
1107         # Get summary embedding
1108         # [batch*traj_len, 1, tf dim]
1109         cls_embed = self.transformer(
1110             seq, # Encoder input
1111             cls, # Decoder input
1112             # Apply mask
1113             src_key_padding_mask=mask,
1114             memory_key_padding_mask=mask,
1115         )
1116
1117         # [batch, traj_len, d_model]
1118         cls_embed = cls_embed.reshape(bs, traj_len, -1)
1119
1120         # Predict actions
1121         # [batch, traj_len, actions]
1122         return self.action_fc(cls_embed)
```

Listing 1: Example policy for PR2L.

```
1123 def process_obs(model, processor, image, prompt, device, last_n=2):
```

```

1124     inputs = processor(images=image, text=prompt, return_tensors="pt")
1125     .to(device)
1126
1127     # Generate text in response to prompt and extract embeddings
1128     outputs = model.generate(
1129         **inputs,
1130         output_hidden_states=True,
1131         return_dict_in_generate=True,
1132         # Any other generation parameters (min/max tokens, temp, etc)
1133     )
1134     hs = outputs["hidden_states"]
1135
1136     # Get image and prompt token embeds
1137     # Any additional processing should happen here (eg pooling of
1138     visual tokens)
1139     # [last_n, num img + prompt tokens, tf_embed_dim]
1140     image_and_prompt_embs = torch.cat(hs[0], dim=0)[-last_n:]
1141
1142     # Get decoded token embeds
1143     # [last_n, num decoded tokens, tf_embed_dim]
1144     dec_embs = []
1145     for dec_hs in hs[1:]:
1146         # [last_n, 1, tf_embed_dim]
1147         dec_hs = torch.cat(dec_hs, dim=0)[-last_n:]
1148         dec_embs.append(dec_hs)
1149     # [last_n, num decoded tokens, tf_embed_dim]
1150     dec_embs = torch.cat(dec_embs, dim=1)
1151
1152     # [last_n, num total tokens]
1153     seq_embs = torch.cat([image_and_prompt_embs, dec_embs], dim=1)
1154     tf_embed_dim = seq_embs.shape[-1]
1155
1156     # [bs=1, seq_len=1, last_n*num total tokens, tf_embed_dim]
1157     seq_embs = seq_embs.reshape(1, 1, -1, tf_embed_dim)
1158
1159     mask = torch.zeros(seq_embs[:-1], type=int)
1160
1161     return seq_embs, mask

```

Listing 2: Example code for extracting promptable representations from a VLM.

```

1162 # Create VLM and processor (InstructBLIP, for example)
1163 model = InstructBlipForConditionalGeneration.from_pretrained(
1164     "Salesforce/instructblip-vicuna-7b"
1165 )
1166 processor = InstructBlipProcessor.from_pretrained("Salesforce/
1167     instructblip-vicuna-7b")
1168
1169 # Set device, can also change dtype if desired
1170 device = "cuda:0"
1171 model = model.to(device)
1172
1173 # Create env
1174 env = ...
1175
1176 # Create policy. This can be trained via RL or BC as needed.
1177 policy = Policy(env.num_actions).to(device)
1178
1179 # Define task-relevant prompt
1180 prompt = "Would a toilet be found here? Why or why not?"
1181
1182 # To predict an action, get an RGB obs from the env and process it
1183 with the VLM
1184 obs = env.reset()
1185 seq, mask = process_obs(model, processor, obs, prompt, device)

```

```
1186
1187 # Then, pass it through the policy to get action logits and step env
1188 act_logits = policy.forward((seq, mask)).reshape(env.num_actions)
1189 action = torch.argmax(act_logits)
1190 obs, _, _, _ = env.step(action)
```

Listing 3: Example usage of the above function and policy.

1191 J Extended Literature Review

1192 **Learning in Minecraft.** We now consider some current approaches for creating autonomous learning
1193 systems for tasks in Minecraft. Such works highlight some of the difficulties prevalent in tasks
1194 in said environment. For instance, since Minecraft tasks take place in a dynamic open world, it
1195 can be difficult for an agent to determine what goal it is attempting to reach and how close it is
1196 to reaching that goal. [10] tackles these issues by introducing and integrating a training scheme
1197 for self-supervised goal-conditioned representations and a horizon predictor. [79] learns a model
1198 from visual observations to discriminate between expert state sequences and non-expert ones, which
1199 provides a source of intrinsic rewards for downstream RL tasks (as it pushes the policy to learn
1200 to match the expert state distribution, which tend to be “good” states for accomplishing tasks in
1201 Minecraft).

1202 **Foundation Models and Minecraft.** Likewise, there has been much interest in applying foundation
1203 models – especially (V)LMs – to Minecraft tasks. [3] pretrains on large scale videos, which enabled
1204 the first agent that could learn to acquire diamond tools (thereby completing a longstanding challenge
1205 in the MineRL competition [28]). LMs have subsequently also been used to produce graphs of
1206 proposed skills to learn or technology tree advancements to make in the form of structured language
1207 [48, 81, 75, 70]. Other works propose to use the LLM to generate actions or code submodules
1208 given textual descriptions of observations or agent states [69]. Finally, VLMs have been used
1209 largely for language-conditioned reward shaping [19, 16]. In contrast, we use VLMs as a source
1210 of representations for learning of atomic tasks (as defined by [38]) that have pre-defined reward
1211 functions; the latter works can thus be used in conjunction with our proposed approach for tasks
1212 where these vision-language reward functions are appropriate.