

Vision-Language Models Provide Promptable Representations for Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

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

1 Introduction

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

An approach to facilitate this involves integrating RL agents with the prior knowledge and reasoning abilities of pre-trained foundation models. Transformer-based language models (LMs) and vision-language models (VLMs) are trained on Internet-scale data to enable generalization in downstream tasks requiring facts or common sense. Moreover, in-context learning (Brown et al., 2020), chain-of-thought reasoning (CoT) (Wei et al., 2023), and instruction fine-tuning (Ouyang et al., 2022) have provided better ways to index into (V)LMs’ knowledge and steer their capabilities based on user needs. These successes have seen some transfer to embodied control, with (V)LMs being used to reason about goals to produce executable plans (Ahn et al., 2022) or as encoders of useful information (like instructions (Liu et al., 2023) or feedback (Sharma et al., 2023)) that the control policy utilizes. Both these paradigms have major limitations: actions generated by LMs are often not appropriately grounded, unless the tasks and scenes are amenable to being expressed or captioned in language. Even then, (V)LMs are often only suited to producing subtask plans, not low-level control signals. On the other hand, using (V)LMs to simply encode inputs under-utilizes their knowledge and

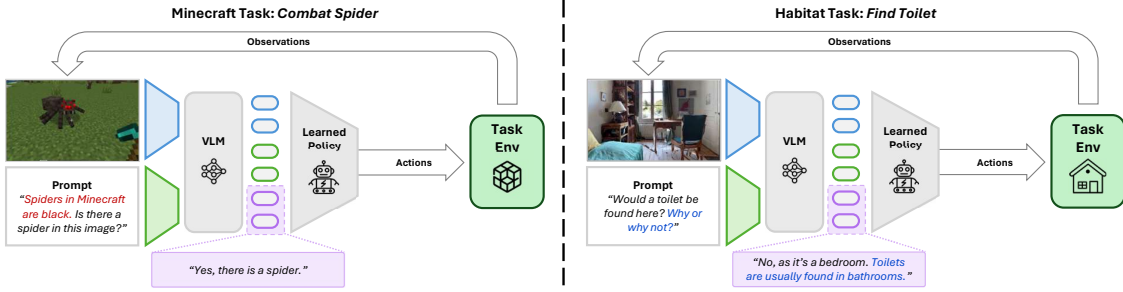


Figure 1: **Example instantiations of PR2L for tasks in Minecraft and Habitat.** We query a VLM with a *task-relevant prompt* about observations to produce *promptable representations*, which we train a policy on via RL. Rather than directly asking for actions or specifying the task, the prompt enables 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*.

reasoning abilities, instead focusing on producing embeddings that reflect the compositionality of language (e.g., so an instruction-following policy may generalize). This motivates the development of an algorithm for learning to produce low-level actions that are grounded and leverage (V)LMs’ knowledge and reasoning.

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

To the best of our knowledge, we introduce the first approach for initializing RL policies with generative VLM representations. We demonstrate our approach on tasks in Minecraft (Fan et al., 2022) and Habitat (Savva et al., 2019), as they present semantically-rich problems representative of many practical, realistic, and challenging applications of RL. We find that PR2L outperforms equivalent policies trained on vision-only embeddings or with instruction-conditioning, popular ways of using pre-trained image models and VLMs respectively for control. We also show that promptable representations extracted from general-purpose VLMs are competitive with domain-specific representations. Our results highlight how visually-complex control tasks can benefit from accessing the knowledge captured within VLMs via prompting in both online and offline RL settings.

2 Related Works

Vision-language models. In this work, we utilize *generative VLMs* (like Li et al. (2022; 2023a); Dai et al. (2023); Karamcheti et al. (2024)): models that generate language in response to an image and a text prompt passed as input. This is in contrast to other designs of combining vision and language that either generate images or segmentation (Rombach et al., 2022; Kirillov et al., 2023) and contrastive representations (Radford et al., 2021). Formally, the VLM enables sampling from $p(x_{1:K}|I, c)$, where $x_{1:K}$ represents the K tokens of the output, I is the input image(s), c is the prompt, and p is the distribution over natural language responses produced by the VLM on those inputs. Typically, the VLM is pre-trained on tasks that require building association between vision and language such as captioning. All these tasks require learning to attend to certain semantic features of input images depending on the given prompt. For auto-regressive generative VLMs, this distribution is factorized as $\prod_t p(x_t|I, c, x_{1:t-1})$. Typical architectures parameterize these distributions using weights that define a representation $\phi_t(I, c, x_{1:t-1})$, which depends on the image

I , the prompt c , and the previously emitted tokens, and a decoder $p(x_t|\phi_t(I, c, x_{1:t-1}))$, which defines a distribution over the next token.

Embodied (V)LM reasoning. Many recent works have leveraged (V)LMs as priors over effective plans for a given goal. These works use the model’s language modeling and auto-regressive generation capabilities to extract such priors as textual subtask sequences (Ahn et al., 2022; Huang et al., 2022b; Sharma et al., 2022) or code (Liang et al., 2023; Singh et al., 2022; Zeng et al., 2022; Vemprala et al., 2023), thereby using the (V)LM to decompose long-horizon tasks into executable parts. These systems often need grounding mechanisms to ensure plan feasibility (e.g., affordance estimators (Ahn et al., 2022), scene captioners (Zeng et al., 2022), or trajectory labelers (Palo et al., 2023)). They also often assume access to low-level policies that can execute these subtasks, such as robot pick-and-place skills (Ahn et al., 2022; Liang et al., 2023), which is often a strong assumption. These methods generally do not address how such policies can be acquired, nor how these low-level skills can themselves benefit from the prior knowledge in (V)LMs. Even works in this area that use RL still use (V)LMs as state-dependent priors over reasonable high-level goals to learn (Du et al., 2023). This is a key difference from our work: instead of considering priors on plans/goals, we rely on VLM’s implicit knowledge *of the world* to extract representations which encode task-relevant information. We train a policy to convert these features into low-level actions via standard RL, meaning the VLM does not need to know how to take actions for a task.

Embodied (V)LM pre-training. Other works use (V)LMs to embed useful information like instructions (Liu et al., 2023; Myers et al., 2023; Lynch & Sermanet, 2021; Mees et al., 2023; O.M.T. et al., 2023), feedback (Sharma et al., 2023; Bucker et al., 2022), reward specifications (Fan et al., 2022), and data for world modeling (Lin et al., 2023b; Narasimhan et al., 2018). These works use (V)LMs as *encoders* of the compositional semantic structure of input text and images, which aids in generalization: an instruction-conditioned model may never have learned to grasp apples (but can grasp other objects), but by interacting with them in other ways and receiving associated language descriptions, the model might still be able to grasp them zero-shot. In contrast, our method produces embeddings that are informed by world knowledge and reasoning, both from prompting and pre-training. Rather than just specifying that the task is to acquire an apple, we ask a VLM to parse observations into task-relevant features, like whether there is an apple in the image or if the observed location likely contains apples – information that is useful even in single-task RL. Thus, we use VLMs to help RL solve new tasks, not just to follow instructions.

These two categories are not mutually exclusive: Brohan et al. (2023a) use VLMs to understand instructions, but also reasoning (e.g., figuring out the “correct bowl” for a strawberry is one that contains fruits); Palo et al. (2023) use a LM to reason about goal subtasks and a VLM to know when a trajectory matches a subtask, automating the demonstration collection/labeling of Ahn et al. (2022), while Adeniji et al. (2023) use a similar approach to pretrain a language-conditioned RL policy that is transferable to learning other tasks; and Shridhar et al. (2021) use CLIP to merge vision and text instructions directly into a form that a Transporter (Zeng et al., 2020) policy can operationalize. Nevertheless, these works primarily focus on instruction-following for robot manipulation. Our approach instead prompts a VLM to supplement RL with representations of world knowledge, not instructions. In addition, except for Adeniji et al. (2023), these works focus on behavior cloning (BC), assuming access to demonstrations for policy learning, whereas our framework can be used for both online RL and offline RL/BC.

3 PR2L: Promptable Representations for Reinforcement Learning

We adopt the standard framework of partially-observed Markov decision process in deep RL, wherein the objective is to find a policy mapping states to actions that maximizes the expected returns. Our goal is to supplement RL with task-relevant information extracted from VLMs containing general-purpose knowledge. One way to index into this information is by prompting the model to get it to produce semantic information relevant to a given control task. Therefore, our approach, PR2L, queries a VLM with a task-relevant prompt for each visual observation received by the agent, and receives both the decoded text and, critically, the intermediate representations, which we refer to as *promptable representations*. Even though the decoded text might often not be correct or directly usable for choosing the action, our key insight is that these VLM embeddings can still provide useful semantic features for training control policies via RL. This recipe enables

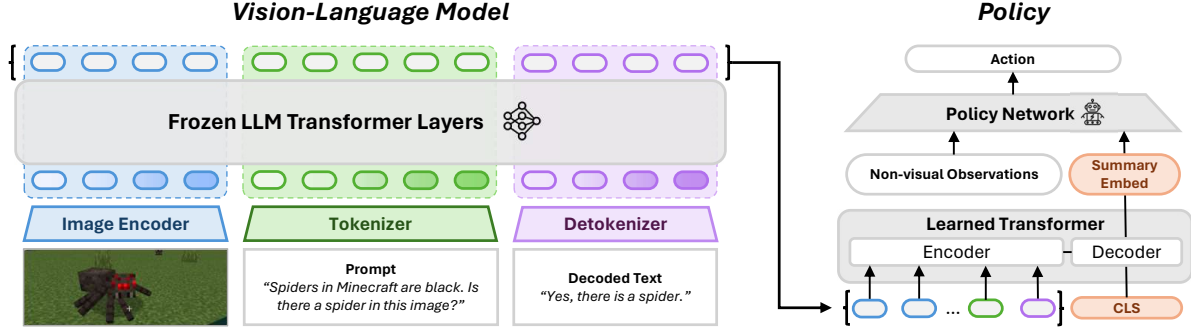


Figure 2: **Schematic of how we extract task-relevant features from the VLM and use them in a policy trained with RL.** These representations can incorporate task context from the prompt, while generic image embeddings cannot. As generative VLM’s embeddings can be variable length, the policy has a Transformer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

us to incorporate semantic information without the need of re-training or fine-tuning a VLM to directly output actions, as proposed by Brohan et al. (2023a). Note that our method is *not* an instruction-following method, and it does **not** require a task instruction to perform well. Instead, our approach still learns control via RL, while benefiting from the incorporation of *background context*. In this section, we will describe various components of our approach, accompanied by practical design choices and considerations.

3.1 Promptable Representations

In principle, one can directly query a VLM to produce actions for a task given a visual observation. While this may work when high-level goals or subtasks are sufficient, VLMs are empirically poor at yielding the low-level actions used commonly in RL (Huang et al., 2022a). As VLMs are trained to follow instructions and answer questions about images, it is more appropriate to use these models to extract and reason about *semantic features* about observations that are conducive to being linked to actions. We thus elicit features that are useful for the downstream task by querying these VLMs with *task-relevant prompts* that provide contextual task information, thereby causing the VLM to attend to and interpret appropriate parts of observed images. Extracting these features naïvely by only using the VLM’s *decoded text* has its own challenges: such models often suffer from hallucinations (Ji et al., 2023) and an inability to report what they “know” in language, even when their embeddings contain such information (Kadavath et al., 2022; Hu & Levy, 2023). However, even when the text is bad, the underlying *representations* still contain valuable granular world information that is potentially lost in the projection to language (Li et al., 2021; Wiedemann et al., 2019; Huang et al., 2023; Li et al., 2023b). Thus, we disregard the generated text and instead provide our policy the embeddings produced by the VLM in response to prompts asking about relevant semantic features in observations instead.

Which parts of the network can be used as promptable representations? The VLMs we consider are all based on the Transformer architecture (Vaswani et al., 2017), which treats the prompt, input image(s), and decoded text as token sequences. This architecture provides a source of learned representations by computing embeddings for each token at every layer based on the previous layer’s token embeddings. In terms of the generative VLM formalism introduced prior, a Transformer-based VLM’s representations $\phi_t(I, c, x_{1:t-1})$ consist of N embeddings per token (the outputs of the N self-attention layers) in the input image I , prompt c , and decoded text $x_{1:t-1}$. The decoder $p(x_t|\phi_t)$ extracts the final layer’s embedding of the most recent token x_{t-1} , projecting it to a distribution over the token vocabulary and allowing for it to be sampled. When given a visual observation and task prompt, the tokens representing the prompt, image, and answer consequently encode task-relevant semantic information. Thus, for each observation, we use the VLM to sample a response to the task prompt $x_{1:K} \sim p(x_{1:K}|I, c)$. We then use some or all of these token embeddings $\phi_K(I, c, x_{1:t-1})$ as our promptable representations and feed them, along with any non-visual observation information, as a state representation into our neural policy trained with RL.

In summary, our approach involves creating a task-relevant prompt that provides context and auxiliary information. This prompt, alongside the current visual observation from the environment, is fed into the

VLM to generate tokens. While these tokens are used for decoding, they are ultimately discarded. Instead, we utilize the *representations* produced by the VLM (associated with the image, prompt, and decoded text) as input for our policy, which is trained via an off-the-shelf online RL algorithm to produce appropriate actions. A schematic of our approach is depicted in Figure 2 and a code snippet example is presented in Appendix I.

3.2 Design Choices for PR2L

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

Second, one must choose which VLM layers’ embeddings to utilize in the policy. While theoretically, all layers of the VLM could be used, pre-trained Transformer models tend to encode valuable high-level semantic information in their later layers (Tenney et al., 2019; Jawahar et al., 2019). Thus, we opt to only feed the final few layers’ representations into our policy. As these representation sequences are of variable length, we incorporate an encoder-decoder Transformer layer in the policy. At each time step in a trajectory, this layer receives variable-length VLM representations, which are attended to and converted into a fixed-length summarization by the embeddings of a learned “CLS” token (Devlin et al., 2019) in the decoder (green in Figure 2). We also note that this policy can receive the observed image directly (e.g., after being embedded by the image encoder), so as to not lose any visual information from being processed by the VLM. However, we do not do this in our experiments in order to more clearly isolate and demonstrate the usefulness of the VLM’s representations in particular.

Finally, while it is possible to fine-tune the VLM for RL end-to-end with the policy (Brohan et al., 2023a), this incurs substantial compute, memory, and time overhead, particularly with larger VLMs. Nonetheless, we find that our approach performs better than not using the language and prompting components of the VLM. This holds true even when the VLM is frozen, and only the policy is trained via RL, or when the decoded text occasionally fails to answer the task-specific prompt correctly.

3.3 Task-Relevant Prompt Design

How do we design good prompts to elicit useful representations from VLMs? As we aim to extract good state representations from the VLM for a downstream policy, we do not use instructions or task descriptions, but task-relevant prompts: questions that make the VLM attend to and encode semantic features in the image that are useful for the RL policy learning to solve the task (Borja-Diaz et al., 2022). For instance, if the task is to find a toilet within a house, appropriate prompts include “What room is this?” and “Would a toilet be found here?” Intuitively, the answers to these questions help determine good actions (e.g., look around the room or explore elsewhere), making the corresponding representations good for representing the state for a policy. Answering the questions will require the VLM to attend to task-relevant features in the scene, relying on the model’s internal conception of what things look like and common-sense semantic relations. One can also prompt the VLM to use chain of thought (Wei et al., 2023) to explain its generated text, often requiring it to reason about task-relevant features in the image, resulting in further enrichment of the state representations. Finally, prompts can provide helpful auxiliary information: e.g., one can describe what certain entities of interest look like, aiding the VLM in detecting them even if they were not commonly found in the model’s pre-training data.

Note that prompts based on instructions or task descriptions do not enjoy the above properties: while the goal of those prior methods is to be able to directly query the VLM for the optimal action, the goal of task-relevant prompts is to produce a useful state representation, such that running RL with them can accelerate

learning an optimal policy. While the former is not possible without task-specific training data for the VLM in the control task, the latter proves beneficial with off-the-shelf VLMs.

Evaluating and designing prompts for RL. Since the specific representations elicited from the VLM are determined by the prompt, we want to design prompts that produce promptable representations that maximize performance on the downstream task. The brute-force approach would involve running RL with each candidate prompt to measure its efficacy, but this would be computationally very expensive. In lieu of this, we evaluate candidate prompts on a small dataset of observations labeled with semantic features of interest for the considered task. Example features include whether task-relevant entities are in the image, the relative position of said entities, or even actions (if expert demonstrations are available). We test prompts by querying the VLM and checking how well the resulting decoded text for each image matches ground truth labels. As this is only practical for small, discrete spaces that are easily expressed in words, we see how well a small model can fit the VLM’s embeddings to the labels (akin to probing in self-supervised learning (Shi et al., 2016; Belinkov & Glass, 2019)). While this does not directly optimize for task performance, it does act as a proxy that ensures a prompt’s resulting representations encode certain semantic features which are helpful for the task.

4 Experimental Setups

Our experiments analyze whether promptable representations from VLMs provide benefits to downstream control, thus providing an effective vehicle for transferring Internet-scale knowledge to RL. We aim to show that PR2L is a good source of state representations, even with our current VLMs that are bad at reasoning about actions – as such models become more performant, we expect such representations to be even better. We thus design experiments to answer the following: **(1)** Can promptable representations obtained via task-specific prompts enable more performant and sample-efficient learning than those of non-promptable image encoders pre-trained for vision or control? **(2)** How does PR2L compare to approaches that directly “ask” the VLM to generate good actions for a task specified in the prompt? **(3)** How does PR2L fare against other popular learning approaches or purely visual features in our domains of interest?

4.1 Domain 1: Minecraft

We first conduct experiments in Minecraft, which provides control tasks that require associating visual observations with rich semantic information to succeed. Moreover, since these observations are distinct from the images in the the pre-training dataset of the VLM, succeeding on these tasks relies crucially on the efficacy of the task-specific prompt in meaningfully affecting the learned representation, enabling us to stress-test our method. E.g., while spiders in Minecraft somewhat resemble real-life spiders, they exhibit stylistic exaggerations such as bright red eyes and a large black body. If the task-specific prompt is indeed effective in informing the VLM of these facts, it would produce a representation that is more conducive to policy learning and this would be reflected in task performance. For this domain, we use the half-precision Vicuna-7B version of the InstructBLIP instruction-tuned generative VLM (Dai et al., 2023; Chiang et al., 2023) to produce promptable representations.

Minecraft tasks. We consider all programmatic Minecraft tasks evaluated by Fan et al. (2022): *combat spider*, *milk cow*, *shear sheep*, *combat zombie*, *combat enderman*, and *combat pigman*¹. The remaining tasks considered by Fan et al. (2022) are creative tasks, which do not have programmatic reward functions or success detectors, so we cannot directly train RL agents on them. We follow the MineDojo definitions of observation/action spaces and reward function structures for these tasks: at each time step, the policy observes an egocentric RGB image, its pose, and its previously action; the policy can choose a discrete action to turn the agent by changing the agent’s pitch and/or yaw in discrete increments, move, attack, or use a held item. These tasks are long horizon, with a maximum episode length of 500 - 1000 and taking roughly 200 steps for a learned policy to complete them. See Figure 3 for example observations and Appendix B.1 for more details.

¹ Fan et al. (2022) 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.

Comparisons. We compare PR2L to five performant classes of approaches for RL in Minecraft: **(a)** Methods using non-promptable representations of visual observations. This does not use prompting altogether, instead using task-agnostic embeddings from the VLM’s image encoder (specifically, the ViT-g/14 from InstructBLIP – blue in Figure 2). While these representations are still pre-trained, PR2L utilizes prompting to produce *task-specific* representations. For a fair comparison, we use the *exact same* policy architecture and hyperparameters for this baseline as in PR2L, ensuring that performance differences come from prompting for better representations from the VLM. **(b)** Methods that directly “asks” the VLM to output actions to execute on the agent. This adapts the approach of Brohan et al. (2023a) to our setting and directly outputs the action from the VLM. While Brohan et al. (2023a) also fine-tune the VLM backbone, we are unable to do so using our compute resources. To compensate, we do not just execute the action from the VLM, but train an RL policy to map this decoded action to a better one. Note that if the VLM already decodes good action texts, simply copying over this action via RL should be easy. **(c)** Methods for efficient RL from pixels via model-based approaches. We choose Dreamer v3, since it has proven to be successful at learning Minecraft tasks from scratch Hafner et al. (2023). **(d)** Methods leveraging pretrained representations specifically useful for embodied control, though which are non-promptable and non-Minecraft specific. We choose VC-1 and R3M Majumdar et al. (2023); Nair et al. (2022). **(e)** Methods using models pre-trained on large-scale Minecraft data. These serve as “oracle” comparisons, as these representations are explicitly fine-tuned on Minecraft YouTube videos, whereas our pre-trained VLM is both frozen and not trained on any Minecraft video data. We choose MineCLIP, VPT, and STEVE-1 as our sources of Minecraft-specific representations Fan et al. (2022); Baker et al. (2022); Lifshitz et al. (2023).

We use PPO (Schulman et al., 2017) as our base RL algorithm for all non-Dreamer Minecraft policies. We also note that we do *not* compare against non-RL methods, such as Voyager (which uses LLMs to write high-level code skills, abstracting away low-level control to hand-written APIs that use oracle information). See Appendix B.2 for training details and E.1 for further discussion of such non-learned systems.

4.2 Domain 2: Habitat

A major advantage of VLMs pre-trained on Internet-scale data is their reasoning and generalization capabilities. To evaluate this, we run offline BC and RL experiments in the Habitat household simulator. In contrast to Minecraft, tasks in this domain require connecting *naturalistic* images with real-world common sense about the structure and contents of typical home environments. Our experiments evaluate **(1)** whether PR2L confers the generalization properties of VLMs to our policies, **(2)** whether PR2L-based policies can leverage the semantic reasoning capabilities of the underlying VLM (e.g., via chain-of-thought Wei et al. (2023)), and **(3)** whether PR2L can learn entirely from stale, offline data sources. We use a Llama2-7B Prismatic VLM for the Habitat experiments Karamcheti et al. (2024).

Habitat tasks. We consider the ObjectNav task suite in 3D scanned household scenes from the HM3D dataset (Savva et al., 2019; Yadav et al., 2023a; Ramakrishnan et al., 2021). These tasks involve a simulated robot traversing a home environment to find an instance of a specified object (toilet, bed, sofa, television, plant, or chair) in the shortest path possible. The full benchmark consists of 80 household scenes intended to train the agent and 20 for validation. We change the observation space to consist of just RGB vision, previous action, pose, and target object class, omitting depth images to ensure that observed performance differences come from the quality of promptable representations vs. unpromptable ones. Like with MineDojo, these tasks are long horizon, taking 80 steps for a privileged shortest path follower to succeed and 150+ for humans. See Figure 3 for example observations and Appendix C for more details.

Comparisons. To see if PR2L can leverage VLM reasoning capabilities, we train two PR2L policies, one with and one without chain-of-thought prompting (see Section 4.3). We also train a policy on Prismatic VLM image encoder embeddings (equivalent to Minecraft approach (a), but with Dino+SigLIP Caron et al. (2021); Zhai et al. (2023)) on a human demonstration dataset collected from the ObjectNav training scenes collected with Habitat-Web Ramakrishna et al. (2022) and used by past works on large-scale BC on pre-trained visual representations Ramakrishna et al. (2023); Yadav et al. (2023b); Majumdar et al. (2023). As it previously achieved state-of-the-art performance among those works, we also compare against two policies using VC-1 as an encoder (Majumdar et al., 2023), either using just its summarizing CLS token or using a learned Transformer layer to condense its patch embeddings. We adopt the same LSTM-based recurrent

	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.

architecture used by that work, but replace the image embeddings with a learned Transformer layer that condenses our input token embeddings (from the VLM, VLM image encoder, or VC-1) into a single summary embedding, as done with Minecraft.

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

4.3 Designing Task-Specific Prompts for Minecraft and Habitat

We now discuss how to design prompts for PR2L. As noted in Section 3.3, these are not instructions or task descriptions, but prompts that force the VLM to encode semantic information useful for the task in its representation. The simplest relevant feature for our Minecraft tasks is the presence of the target entity in an observation. Thus, we choose “Is there a [target entity] in this image?” as the base of our chosen prompt. We also pick two alternate prompts per task that prepend different amounts of auxiliary information about the target entity. E.g., for *combat spider*, one candidate is “Spiders in Minecraft are black.” To choose between these candidates, we measure how well the VLM is able to decode a correct answer to the prompt question of whether or not the target entity is present in the image on a small annotated dataset. Full details of this prompt evaluation scheme for the first three Minecraft tasks are presented in Appendix A and Table 5. We find that auxiliary text only helps with detecting spiders while systematically and significantly degrading the detection of sheep and cows. Our ablations show that this detection success rate metric correlates with performance of the RL policy. Additionally, the prompts used for comparison (b) follow the prompt structure prescribed by Brohan et al. (2023a), which motivated this comparison. In these prompts, we also provide a list of actions that the VLM can choose from to the policy. All chosen prompts are presented in Table 1.

For Habitat, we choose the prompt “Would a [target object] be found here? Why or why not?” As opposed to the Minecraft prompts, this does not just identify the presence of a target object in the image, but draws on general knowledge from the VLM to determine if the observed location would contain the target object, even if said object is not in view. The second part of the prompt then leads the VLM to provide a chain of thought (CoT) (Wei et al., 2023) rationale for its final answer. This CoT draws out task-relevant VLM world knowledge by explicitly reasoning about visual semantic concepts, that are useful to learning a policy (see Table 4). To investigate if PR2L enables embodied agents to benefit from these VLM common-sense reasoning capabilities (even if they do not directly reason about actions), we train PR2L policies both with and without the second part of the prompt.

5 Results

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

PR2L outperforms (a) **the VLM image encoder baseline**, even though both approaches receive the same visual features, with PR2L simply transforming those features via prompting an LLM (with no additional

Task	PR2L (Ours)	Baselines					Oracles		
		VLM Image Encoder	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>	<i>59.3 ± 6.7</i>
<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>	<i>98.3 ± 8.4</i>

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.

information from the environment), thus supporting that prompting does shape representations in a beneficial way for learning control tasks. We provide an analysis of why PR2L states are better than **(b) RT-2-style** ones in Appendix [H.1](#). We observe that PR2L embeddings are bimodally distributed, with transitions leading to high reward clustered at one mode. This structure likely enables more efficient learning, thereby showing how control tasks can benefit from extracting prior knowledge encoded in VLMs by prompting them with task context, even when the VLM does not know how to act. For **(c) the model-based comparisons**, we find that Dreamer is not as conducive at learning our Minecraft tasks. We hypothesize this is because our tasks are comparatively shorter than the ones considered by [Hafner et al. \(2023\)](#), so learning a model is less beneficial (while PR2L provides immediately-useful representations). Additionally, we note that all our approaches involve interacting with partially-observable, non-stationary entities, which the Dreamer model may have a hard time learning. See Appendix [E.2](#) for further discussion. Finally, **(e) the oracles** outperform PR2L in *combat enderman/pigman*, all but STEVE-1 do better in *combat spider/zombie*, and none do better in *shear sheep/milk cow*. We hypothesize this is because endermen and pigmen are Minecraft-specific entities, giving rise to comparatively poor representations in the VLM (which is trained exclusively on natural images). In contrast, Minecraft zombies/spiders are heavily stylized, but still somewhat resemble other depictions of such creatures, while Minecraft cows and sheep are the closest to their naturalistic counterparts, making PR2L more effective. Even though our VLM is not trained on Minecraft data, its representations yield better policies in half the oracle comparisons.

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

Habitat results. Following prior works, we report success rates on the ObjectNav validation episodes in Table [3](#). PR2L with CoT outperforms all other policies on all tasks, including an almost $4\times$ performance increase over the VLM image encoder baselines – again, suggesting that using promptable representations for control improves over the base purely-visual embeddings. While PR2L without CoT still does better than all baselines, we find CoT prompting improves policy performance (by $1.5\times$, from 27.8% success rate to 41.9%), likely because it provides the policy with useful generalizable features: e.g., even if the agent comes across an unfamiliar room while searching for a toilet, it still knows to look elsewhere if the VLM reasons that, due

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.

to the presence of a bed, the room is likely a bedroom (which is unlikely to contain toilets). Thus, even if the VLM cannot reason about actions, *our results indicate that PR2L provides a promising way of using its ability to reason about image semantics and common sense for control*. See Table 4 for CoT examples.

While we do not beat VC-1’s reported SOTA BC performance (60.3% success rate when VC-1 is frozen Majumdar et al. (2023)), we note that said performance is achieved with (1) over ten times more training data and gradient steps and (2) image augmentations to prevent overfitting. Our VC-1 policies were trained on the same amount of data as our PR2L agent and for 1-3× as many gradient steps, but perform far worse, suggesting that PR2L is significantly more sample- and compute-efficient than VC-1 policies. Additionally, PR2L does not use any explicit countermeasures to overfitting, yet still generalizes well to unseen ObjectNav scenes (aided by the VLM’s representations of reasoning).

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

6 Conclusion

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

A limitation of PR2L is that prompts are currently hand-crafted based on the user’s conception of useful task features. While coming up with good prompts for our tasks was not hard, the process of evaluating and improving them could be automated, which we leave to future works. We also find that the quality of representations largely depends on the VLM – e.g., InstructBLIP could not reason well about Habitat scenes, but the more recent Prismatic VLMs are more capable in that regard, enabling our CoT experiments. Thus, as VLM capabilities are expected to increase, we expect the quality of their representations to also improve. Lastly, the size and speed of VLMs can limit their applicability. Our policies typically achieve 3-5 Hz inference speeds, comparable to those of robot policies built on large models Brohan et al. (2023b); O.M.T. et al. (2023). Likewise, our VLM sizes are comparable to models used for policies in prior works (Brohan et al. 2023a; Szot et al. 2024). While their inference speeds may hinder online policy learning, we find that offline approaches (which can parallelize training and data generation) we used for Habitat help remedy this.

References

- Ademi Adeniji, Amber Xie, Carmelo Sferrazza, Younggyo Seo, Stephen James, and Pieter Abbeel. Language reward modulation for pretraining reinforcement learning, 2023.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as i can and not as i say: Grounding language in robotic affordances. 2022.
- Bowen Baker, Ilge Akkaya, Peter Zhokhov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos, 2022.
- Yonatan Belinkov and James Glass. Analysis methods in neural language processing: A survey, 2019.
- Jessica Borja-Diaz, Oier Mees, Gabriel Kalweit, Lukas Hermann, Joschka Boedecker, and Wolfram Burgard. Affordance learning from play for sample-efficient policy learning. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Philadelphia, USA, 2022.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023a.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2023b.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- Arthur Buckner, Luis Figueredo, Sami Haddadin, Ashish Kapoor, Shuang Ma, Sai Vemprala, and Rogerio Bonatti. Latte: Language trajectory transformer, 2022.

- Shaofei Cai, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Open-world multi-task control through goal-aware representation learning and adaptive horizon prediction, 2023.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Will Dabney, Mark Rowland, Marc G. Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression, 2017.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- Ziluo Ding, Hao Luo, Ke Li, Junpeng Yue, Tiejun Huang, and Zongqing Lu. Clip4mc: An rl-friendly vision-language model for minecraft, 2023.
- Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models, 2023.
- Kiana Ehsani, Tanmay Gupta, Rose Hendrix, Jordi Salvador, Luca Weihs, Kuo-Hao Zeng, Kunal Pratap Singh, Yejin Kim, Winson Han, Alvaro Herrasti, Ranjay Krishna, Dustin Schwenk, Eli VanderBilt, and Aniruddha Kembhavi. Imitating shortest paths in simulation enables effective navigation and manipulation in the real world, 2023.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Neural Information Processing Systems, 2022*, 2022.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models, 2023.
- Jennifer Hu and Roger Levy. Prompt-based methods may underestimate large language models’ linguistic generalizations, 2023.
- Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents, 2022a.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models, 2022b.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure of language? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, 2019. Association for Computational Linguistics.

- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, mar 2023.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislaw Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. Language models (mostly) know what they know, 2022.
- Anssi Kanervisto, Stephanie Milani, Karolis Ramanauskas, Nicholay Topin, Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, Wei Yang, Weijun Hong, Zhongyue Huang, Haicheng Chen, Guangjun Zeng, Yue Lin, Vincent Micheli, Eloi Alonso, François Fleuret, Alexander Nikulin, Yury Belousov, Oleg Svidchenko, and Aleksei Shpilman. Minerl diamond 2021 competition: Overview, results, and lessons learned, 2022.
- Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models, 2024.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything, 2023.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning, 2020.
- Belinda Z. Li, Maxwell Nye, and Jacob Andreas. Implicit representations of meaning in neural language models, 2021.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023a.
- Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task, 2023b.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control, 2023.
- Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft, 2023.
- Haowei Lin, Zihao Wang, Jianzhu Ma, and Yitao Liang. Mcu: A task-centric framework for open-ended agent evaluation in minecraft, 2023a.
- Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca Dragan. Learning to model the world with language. 2023b.
- Hao Liu, Lisa Lee, Kimin Lee, and Pieter Abbeel. Instruction-following agents with multimodal transformer, 2023.
- Haokuan Luo, Albert Yue, Zhang-Wei Hong, and Pulkit Agrawal. Stubborn: A strong baseline for indoor object navigation, 2022.
- Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data, 2021.

- Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik, Dhruv Batra, Yixin Lin, Oleksandr Maksymets, Aravind Rajeswaran, and Franziska Meier. Where are we in the search for an artificial visual cortex for embodied intelligence?, 2023.
- Oier Mees, Jessica Borja-Diaz, and Wolfram Burgard. Grounding language with visual affordances over unstructured data. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023.
- Vivek Myers, Andre He, Kuan Fang, Homer Walke, Philippe Hansen-Estruch, Ching-An Cheng, Mihai Jalobeanu, Andrey Kolobov, Anca Dragan, and Sergey Levine. Goal representations for instruction following: A semi-supervised language interface to control, 2023.
- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation, 2022.
- Karthik Narasimhan, Regina Barzilay, and Tommi Jaakkola. Grounding language for transfer in deep reinforcement learning, 2018.
- Kolby Nottingham, Prithviraj Ammanabrolu, Alane Suhr, Yejin Choi, Hannaneh Hajishirzi, Sameer Singh, and Roy Fox. Do embodied agents dream of pixelated sheep: Embodied decision making using language guided world modelling, 2023.
- O.M.T., Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. <https://octo-models.github.io>, 2023.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- Norman Di Palo, Arunkumar Byravan, Leonard Hasenclever, Markus Wulfmeier, Nicolas Heess, and Martin Riedmiller. Towards a unified agent with foundation models, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- Santhosh K. Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alex Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X. Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai, 2021.
- Ram Ramrakhya, Eric Undersander, Dhruv Batra, and Abhishek Das. Habitat-web: Learning embodied object-search strategies from human demonstrations at scale, 2022.
- Ram Ramrakhya, Dhruv Batra, Erik Wijmans, and Abhishek Das. Pirlnav: Pretraining with imitation and rl finetuning for objectnav, 2023.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- Pratyusha Sharma, Antonio Torralba, and Jacob Andreas. Skill induction and planning with latent language, 2022.
- Pratyusha Sharma, Balakumar Sundaralingam, Valts Blukis, Chris Paxton, Tucker Hermans, Antonio Torralba, Jacob Andreas, and Dieter Fox. Correcting robot plans with natural language feedback. In *Robotics: Science and Systems, 2022*, 2023.
- Xing Shi, Inkit Padhi, and Kevin Knight. Does string-based neural MT learn source syntax? In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 1526–1534, November 2016.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Proceedings of the 5th Conference on Robot Learning (CoRL)*, 2021.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models, 2022.
- Andrew Szot, Max Schwarzer, Harsh Agrawal, Bogdan Mazoure, Walter Talbott, Katherine Metcalf, Natalie Mackraz, Devon Hjelm, and Alexander Toshev. Large language models as generalizable policies for embodied tasks, 2024.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- Sai Vemprala, Rogerio Bonatti, Arthur Buckner, and Ashish Kapoor. Chatgpt for robotics: Design principles and model abilities. Technical report, Microsoft, 2023.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, 2023a.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents, 2023b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
- Gregor Wiedemann, Steffen Remus, Avi Chawla, and Chris Biemann. Does bert make any sense? interpretable word sense disambiguation with contextualized embeddings, 2019.
- Karmesh Yadav, Jacob Krantz, Ram Ramrakhyia, Santhosh Kumar Ramakrishnan, Jimmy Yang, Austin Wang, John Turner, Aaron Gokaslan, Vincent-Pierre Berges, Roozbeh Mootaghi, Oleksandr Maksymets, Angel X Chang, Manolis Savva, Alexander Clegg, Devendra Singh Chaplot, and Dhruv Batra. Habitat challenge 2023. <https://aihabitat.org/challenge/2023/>, 2023a.
- Karmesh Yadav, Arjun Majumdar, Ram Ramrakhyia, Naoki Yokoyama, Alexei Baevski, Zsolt Kira, Oleksandr Maksymets, and Dhruv Batra. Ovrl-v2: A simple state-of-art baseline for imagenav and objectnav, 2023b.
- Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks, 2023.
- Andy Zeng, Pete Florence, Jonathan Tompson, Stefan Welker, Jonathan Chien, Maria Attarian, Travis Armstrong, Ivan Krasin, Dan Duong, Vikas Sindhwani, and Johnny Lee. Transporter networks: Rearranging the visual world for robotic manipulation. *Conference on Robot Learning (CoRL)*, 2020.

Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. Socratic models: Composing zero-shot multimodal reasoning with language, 2022.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023.

Bohan Zhou, Ke Li, Jiechuan Jiang, and Zongqing Lu. Learning from visual observation via offline pretrained state-to-go transformer, 2023.

Minzhao Zhu, Yifeng Li, and Tao Kong. Integrating map-based method with end-to-end learning, 2022. URL <https://www.youtube.com/watch?v=N-wW3TwEqbU>

Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, Yu Qiao, Zhaoxiang Zhang, and Jifeng Dai. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory, 2023.