

# Video PreTraining (VPT): Learning Watching Unlabeled Online Videos

Bowen Baker<sup>\*†</sup>  
bowen@openai.com

Jie Tang<sup>\*†</sup>  
jietang@openai.com

Ilge Akkaya<sup>\*†</sup>  
ilge@openai.com

Adrien Ecoffet<sup>\*†</sup>  
adrien@openai.com

Peter Zhokhov<sup>\*†</sup>  
peterz@openai.com

Brandon Houghton<sup>\*†</sup>  
brandon@openai.com

Jeff Clune<sup>\*‡</sup>  
jclune@gmail.com

Joost Huizinga<sup>\*†</sup>  
joost@openai.com

Raul Sampedro<sup>\*†</sup>  
raulsamg@gmail.com

## Abstract

Pretraining on noisy, internet-scale datasets has been heavily studied as a technique for training models with broad, general capabilities for text, images, and other modalities.<sup>1–6</sup> However, for many sequential decision domains such as robotics, video games, and computer use, publicly available data does not contain the labels required to train behavioral priors in the same way. We extend the internet-scale pretraining paradigm to sequential decision domains through semi-supervised imitation learning wherein agents learn to act by watching online unlabeled videos. Specifically, we show that with a small amount of labeled data we can train an inverse dynamics model accurate enough to label a huge unlabeled source of online data – here, online videos of people playing Minecraft – from which we can then train a general behavioral prior. Despite using the native human interface (mouse and keyboard at 20Hz), we show that this behavioral prior has nontrivial zero-shot capabilities and that it can be fine-tuned, with both imitation learning and reinforcement learning, to hard-exploration tasks that are impossible to learn from scratch via reinforcement, and we are the first to report computer agents that can craft diamond tools, which can take proficient humans upwards of 20 minutes (24,000 environment actions) of gameplay to accomplish.

## 1 Introduction

Work in recent years has demonstrated the efficacy of pretraining large and general found models<sup>7</sup> on noisy internet-scale datasets for use in downstream tasks in natural language computer vision.<sup>5,6</sup> For sequential decision domains (e.g. robotics, game playing, and c usage) where agents must repeatedly act within an environment, a wealth of data also exists on the web; however, most of this data is in the form of *unlabeled* video (i.e. without also existing at each frame), making it much less straightforward to train a behavioral prior in the environment than it is in e.g. natural language. In a few rare settings, such as Chess, Go, and Shogi, this was a large effort by a dedicated team. Each author made huge contributions to the original VPT project team and were full time on the project for over six months. BB, IA, PZ, and Joost were thus involved for even longer (over a year). It was also randomized between IA and PZ.

<sup>\*</sup>This was a large effort by a dedicated team. All members were full time on the project for over six months. BB, IA, PZ, and Joost were thus involved for even longer (over a year). It was also randomized between IA and PZ.

<sup>†</sup>OpenAI

<sup>‡</sup>University of British Columbia

already exist large datasets with action labels from various domains used for imitation learning.<sup>9,10</sup> When large labeled datasets do not exist, training capable agents is reinforcement learning (RL),<sup>11</sup> which can be very hard to learn with RL and do not have general-purpose foundation models to sequential decision domains by utilizing freely available sources of labeled data.<sup>12-18</sup> Many virtual tasks, e.g. navigation, Photoshop, booking flights, etc., can be very hard to learn with RL and do not have large, general-purpose foundation datasets with a simple semi-supervised imitation learning method Video PreTraining (VPT) and demonstrate its efficacy in the domain of Minecraft.

Existing semi-supervised imitation learning methods aim to learn with few or no explicit action labels; however, they generally rely on the policy’s ability to explore the environment throughout training, making them susceptible to exploration bottlenecks.<sup>21-25</sup> Furthermore, most prior semi-supervised imitation learning work was tested in the relatively low data regime; because we experiment with far more data ( $\sim 70k$  hours of unlabeled video), we hypothesize that we can achieve good performance with a much simpler method, a trend that has proven true for pretraining in other modalities such as text.<sup>1</sup> In particular, given a large but unlabeled dataset, we propose generating pseudo-labels from the action taken at each timestep in a video. Behavioral cloning (BC) can require a large amount of data because the model must learn to infer intent and the distribution over future behaviors from only past observations. In contrast, the inverse dynamics model (IDM) that predicts environment mechanics are far simpler than the breadth of human behavior that can take place within the environment. Using pseudo-labels generated from the IDM, we then train a model to mimic the distribution of behavior in the previously unlabeled dataset with standard behavioral cloning at scale, which does not require any model rollouts and thus does not suffer from any potential exploration bottlenecks.

We chose to test our method in Minecraft because (a) it is one of the most actively played games in the world<sup>26</sup> and thus has a wealth of commonly available video data online,<sup>(b)</sup> it is a fairly open-ended sandbox game with an extremely wide variety of potential things to do, build, and collect, making our results more applicable to real-world applications such as computer usage, which also tends to be varied and open-ended, and (c) it has already garnered interest by the RL community as a research domain due to its complexity and corresponding difficult exploration challenges.<sup>27-31</sup> In this work we use the native human interface for Minecraft so that we can (1) most accurately experience the game without modification, (2) make data collection easier per second and must use a mouse and keyboard interface to interact with the environment, and (3) eliminate the need to hand-engineer a custom interface for models to interact with the environment. This choice means that our models play at 20 frames per second, trading, etc., including dragging items to specific slots or navigating the native human interface to interact with human contractors to craft items, construct crafting and attacking macros,<sup>30,32-34</sup> using the native GUIs for crafting video data and the environment and correspondingly difficult exploration tasks near impossible facing a naive target.

In Section 4 we show that the VPT foundation model has nontrivial zero-shot generalization tasks impossible to learn with RL alone, such as crafting planks and perch required a human proficient in Minecraft a median of 50 seconds or  $\sim 970$  distributions, our agent is able to push even further into the technology space.

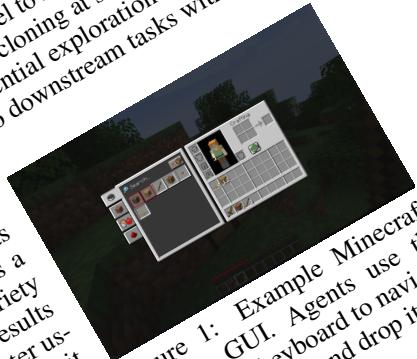


Figure 1: Example Minecraft crafting GUI. Agents use the mouse and keyboard to navigate menus and drag and drop items.

(taking a human a median of 2.3 minutes or  $\sim 2790$  actions),  
 the most dramatic improvements: our agent is able to craft diamond  
 in Minecraft made even more challenging by using the native human.  
 of this work are (1) we are the first to show promising results applying semi-supervised learning to extremely large, noisy, and freely available video datasets for sequential decision making, (2) we show that such pretraining plus fine-tuning enables agents to solve tasks that were impossible to learn, (3) we show that labeled contractor data is far more efficiently used than the VPT method than it would be by directly training a foundation model from it and (4) we open source our contractor data, trained model weights, and Minecraft environment for future research into learning to act via semi-supervised imitation learning at scale.

## 2 Preliminaries and Related Work

Imitation learning methods<sup>35–38</sup> seek to construct a policy that accurately models the distribution of behavior in some dataset  $D = \{(o_i, a_i)\}, i \in \{1\dots N\}$  of action-observation pairs. In order to roll out these policies in an environment, they must be causal, meaning they condition on observations from the current timestep  $t$  and past timesteps only, i.e.  $\pi \sim p(a_t | o_1 \dots o_t)$ . Imitation learning is simplest when demonstrations are labeled with corresponding actions. Imitating labeled trajectories is GAIL<sup>23</sup> constructs an adversarial objective incentivizing the trained policy to exhibit behaviors seen success in aerial vehicles,<sup>39,40</sup> self-driving cars,<sup>41,42</sup> board games,<sup>9,43</sup> and video games.<sup>10,44</sup> When labeled demonstrations are not available, standard behavioral cloning will not work; however, there is a large body of work in imitating behavior from unlabeled demonstrations and then train RL agents to match these waypoints; however, they construct waypoints that are embeddings from unsupervised feature learning models. Pathak et al.<sup>49</sup> and Nair et al.<sup>50</sup> train a small amount of environment interactions and then map the learned latent actions to real actions with agent positions in videos and then train RL agents to match these waypoints. Similarly, Behbahani et al.<sup>47</sup> and Aytar et al.<sup>48</sup> task a RL agent to take actions that advance the IDM on trajectories of past and future timesteps given observations of past and future timesteps, e.g.  $p_{IDM}(a_t | o_t, o_{t+1})$ , and simultaneously train (1) an inverse dynamics model (IDM),<sup>51</sup> which aims to uncover the underlying states expressed as high dimensional visual waypoints. Most similar to our own work, Torabi et al.<sup>24</sup> train the IDM is collected by rolling out the BC model in the target environment such that both models improve in tandem. However, at any point in training if there are sequences in the dataset that the IDM performs poorly on, it requires that the BC model perform on trajectories of observations labeled with the IDM fixed throughout BC training. Compared to most previous work in semi-supervised imitation learning, we experiment more complex and open-ended environment of Minecraft. Minecraft is a voxel-based game that, due its popularity and wide variety of mechanics, has attracted a vast base of research.<sup>27,28,30–34,52–60</sup> A large body of work focuses on small, custom-made environments with tasks such as navigation,<sup>53,60</sup> block placing,<sup>54,55</sup> instruction following,<sup>58,59</sup> others.<sup>28,31,57</sup> Work operating in the massive, randomly generated curriculum learning<sup>30</sup> and, most closely related to this work, automated diamond mining<sup>27,32–34</sup> However, to the best of our knowledge, there is no published work that operates in the full, unmodified human action space, drag-and-drop inventory management and item crafting.

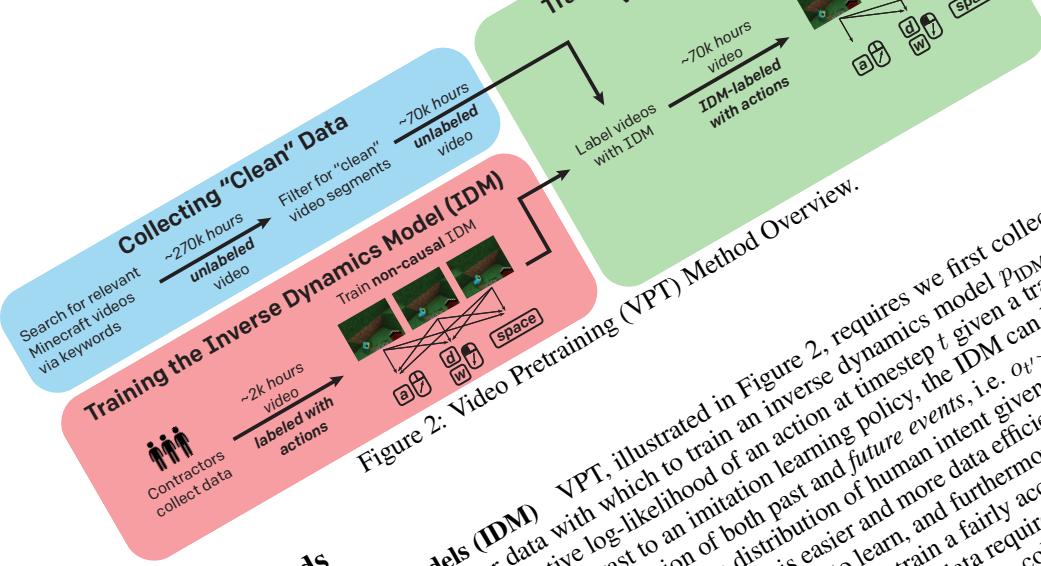


Figure 2: Video Pretraining (VPT) Method Overview.

### 3 Methods

#### Inverse Dynamics Models (IDM)

VPT, illustrated in Figure 2, requires we first collect a small amount of labeled contractor data with which to train an inverse dynamics model  $p_{\text{IDM}}(a_t|o_1 \dots o_T)$ , which seeks to minimize the negative log-likelihood of an action at timestep  $t$  given a trajectory of  $T$  observations  $o_t : t \in [1 \dots T]$ . In contrast to an imitation learning policy, the IDM can be non-causal, meaning its prediction for  $a_t$  can be a function of both past and future events, i.e.  $o'_t > t$ . Compared to the behavioral cloning objective of modeling the distribution of both past and future events, the IDM is much easier to learn, and more data frames only, we hypothesize that inverting environment dynamics is easier and more data intent given past frames only, Sec. 4.1 will show that the IDM objective is much easier to learn, and furthermore Sec. 4.6 will show that with very little labeled data (as few as 100 hours) we can train a fairly accurate IDM. This IDM can be used to label online videos, providing the large amount of data required for the harder task of behavioral cloning. See appendices D and B for IDM training and data collection details.

**Data Filtering** We gather a large dataset of Minecraft videos by searching the web for related keywords (Appendix A). Online videos often (1) include overlaid artifacts, such as a video feed of the player’s face, channel logos, watermarks, etc., (2) are collected from platforms other than a computer with different game modes, e.g. in Minecraft other than

We call data “survival mode” if it does not contain visual artifacts and is from scratch and must gather or craft all their items. We only want “survival mode” where players start from scratch and must gather or craft all their items. We call data “unclean.” With enough data, a large enough model, and enough training compute, a BC model trained on both unclean and clean videos would likely still perform well in a clean Minecraft environment. However, for simplicity and training compute efficiency, we choose to filter out unclean segments of video (note that a video may contain both clean and unclean segments). We do this by training a model to filter out unclean segments using a small dataset (800) of images sampled from online videos labeled by contractors as clean or unclean (Appendix A.2).

**VPT Foundation Model** We train a foundation model with standard behavioral cloning, i.e. minimizing the negative log-likelihood of actions predicted by the IDM on clean data. For a particular trajectory of length  $T$  we minimize

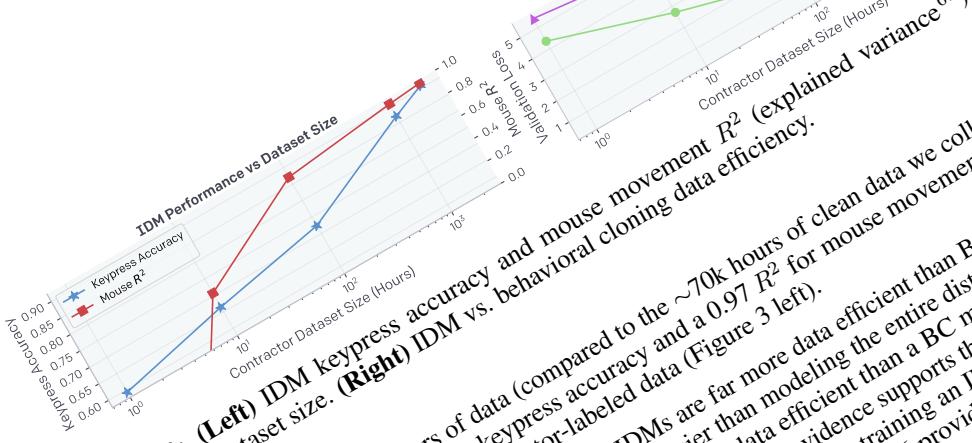
$$\min_{\theta} \sum_{t \in [1 \dots T]} -\log \pi_{\theta}(a_t|o_1, \dots, o_t)$$

As we will see in the following sections, this model exhibits nontrivial zero-shot behavior and fine-tuned with both imitation learning and RL to perform even more complex skills.

## 4 Results

### 4.1 Performance of the Inverse Dynamics Model

The IDM architecture is comprised primarily of a temporal convolution layer, processing stack, and residual unmasked attention layers, from which the IDM the IDM predicts keypresses and mouse movements (see Appendix D for IDM architecture details). A key hypothesis behind our work is that IDMs can be trained with a large amount of labeled data. While more data improves both mouse movement and keypress prediction, we find that the IDM performs best with a small amount of labeled data (Sec. 4.6). This is likely due to the fact that the IDM is a non-causal model, and thus needs to consider future events to make predictions. As a result, it is more difficult to train the IDM with a large amount of labeled data, as the model needs to learn to ignore irrelevant information. In contrast, the VPT model is a causal model, and thus only needs to consider past events to make predictions. As a result, it is easier to train the VPT model with a large amount of labeled data.



**Figure 3: (Left) IDM Performance vs Dataset Size**

Contractor Dataset Size (Hours)	Keypress Accuracy	Mouse R <sup>2</sup>
10 <sup>1</sup>	~0.65	~0.55
10 <sup>2</sup>	~0.85	~0.75
10 <sup>3</sup>	~0.95	~0.90

**(Right) IDM vs. behavioral cloning data efficiency.**

Contractor Dataset Size (Hours)	Mouse movement R <sup>2</sup>	Explained variance (%)
10 <sup>1</sup>	~0.35	~10
10 <sup>2</sup>	~0.45	~20
10 <sup>3</sup>	~0.55	~30

Figure 3: (Left) IDM keypress accuracy and mouse movement  $R^2$  (explained variance) as a function of dataset size. (Right) IDM vs. behavioral cloning data efficiency.

Figure 3 (left) shows IDM keypress accuracy (blue line with circles) and mouse movement  $R^2$  (green line with circles) plotted against Contractor Dataset Size (Hours) on a log scale. The x-axis ranges from  $10^0$  to  $10^3$ , and the y-axis ranges from 0.0 to 1.0. Both metrics increase monotonically with dataset size, with mouse movement  $R^2$  reaching approximately 0.9 at 1000 hours.

Figure 3 (right) shows IDM vs. behavioral cloning data efficiency. The x-axis is Contractor Dataset Size (Hours) on a log scale from  $10^0$  to  $10^2$ , and the y-axis is Validation Loss from 0.0 to 1.0. A red line with squares represents IDM, which starts at ~0.7 loss at 10 hours and decreases to ~0.2 loss at 1000 hours. A blue line with triangles represents behavioral cloning, which starts at ~0.8 loss at 10 hours and decreases to ~0.1 loss at 1000 hours.

Figure 4: (Left) Training and validation loss on the main IDM contractor dataset with IDM pseudo-labels, and loss on the web\_clean internet dataset with ground-truth labels but is out-of-distribution (see text). (Right) Amount a given item was collected per episode averaged over 2500 60-minute survival episodes as a function of training epoch, shaded with the standard error of the mean. Basic mining refers to collection of dirt, gravel, or sand (all materials that can be gathered without tools). Logs are obtained by repeatedly hitting trees for three seconds, a difficult feat for the RL agent to achieve as we show in Sec. 4.4. Planks are crafted from logs, and crafting tables are crafted from planks. Crafting requires using in-game crafting actions to make a crafting table.

Figure 4 (left) is a line plot showing Loss (Y-axis, 0.0 to 4.0) versus Training Progress (Epochs) (X-axis, 0 to 30). It compares Train loss (blue circles), Validation loss (green dashed line), and Contractor Data loss (red dashed line). All losses decrease rapidly initially and then level off.

Figure 4 (right) is a scatter plot showing Crafting Count (Y-axis, 0 to 30) versus Training Progress (Epochs) (X-axis, 0 to 25). It tracks the collection of various items: basic mining (blue circles), logs (green circles), planks (orange circles), crafting tables (purple circles), and total crafting (yellow diamonds). The total crafting count increases steadily over time, reaching approximately 30 by epoch 25.

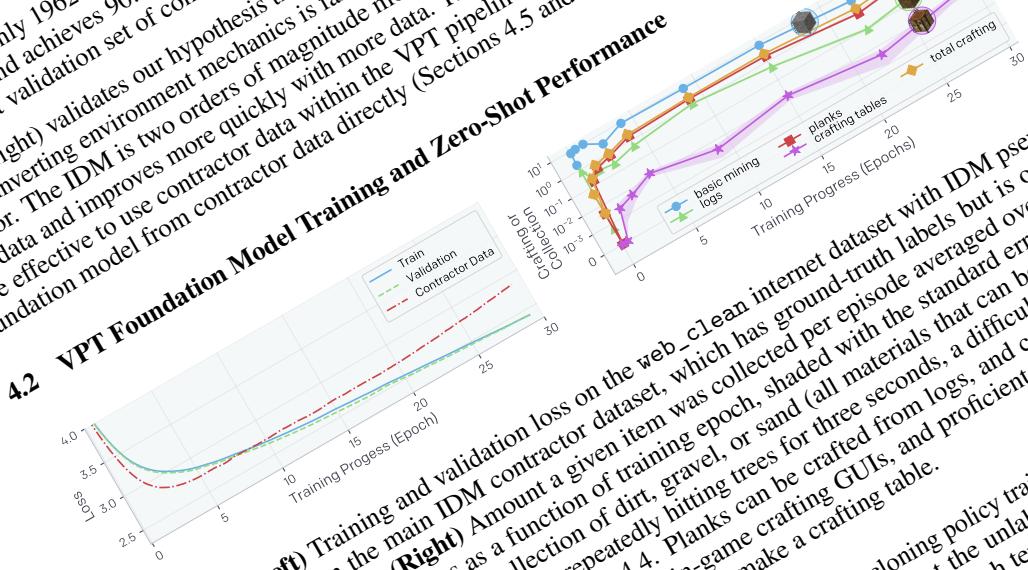


Figure 4: **(Left)** Training and validation loss on the main IDM contractor database distribution (see text). **(Right)** Amount a given item can be harvested as a function of training episodes as a function of collection size. Logs are obtained by repeated mining, as we show in Sec. 4.4. Mining requires to collection of repeated items.

**(Left)** Scatter plot showing the relationship between Contractor Dataset Size (Hours) on a logarithmic x-axis (from 10<sup>0</sup> to 10<sup>2</sup>) and the explained variance ( $R^2$ ) on a logarithmic y-axis (from 10<sup>-6</sup> to 10<sup>0</sup>). The data points show a positive correlation, with a green line representing a linear fit.

**(Right)** Line graph titled "Training and Zero-Shot Performance" showing Crating Collected (y-axis, log scale from 10<sup>-3</sup> to 10<sup>1</sup>) versus Training Progress (Epochs) (x-axis, 0 to 25). The graph includes four series: Train (blue circles), Validation (green triangles), Contractor Data (red dashed line), and Logs (purple stars). A shaded gray area represents the standard error of the mean. The validation and contractor data lines are very close, starting at ~10<sup>-2</sup> and rising to ~10<sup>0</sup>. The logs series starts at ~10<sup>-3</sup>, rises to ~10<sup>-2</sup> by epoch 10, and then continues to rise more slowly, reaching ~10<sup>0</sup> by epoch 25. A purple arrow points to the logs series with the label "basic mining".

The figure contains two line plots. The left plot shows validation loss on a logarithmic y-axis (from  $10^{-3}$  to  $10^0$ ) against Training Progress (Epochs) from 0 to 30. The right plot shows the Crafting rate of Contractor Data (from 0 to 30) against the same x-axis. Both plots include five data series: logs (green circles), basic mining (orange squares), planks (red triangles), crafting tables (purple diamonds), and total crafting (blue stars). Error bars represent standard deviation. In the left plot, all series show a downward trend as training progresses. In the right plot, the total crafting series reaches the highest crafting rate, followed by crafting tables, planks, basic mining, and logs.

to craft a crafting table, which are required to unlock most other items. Human proficient in Minecraft approximately 50 seconds (970 actions) to craft these behaviors at a rate far below that of our proficient contractors, e.g. on average 5.44 crafting tables in 60 minutes of play versus 0.19 for the foundation model. The contractor also collects various berries and mushrooms from them; kills zombies that appear during the night; hunts wild flowers and crafts dyes from chests. The model also learned to navigate uneven terrain, swing, and pillar jump, which involves the agent repeatedly jumping and quickly placing a block below itself such that it climbs upward by making a pillar.<sup>(iv)</sup>

While training and validation loss decrease healthily over training (Fig. 4, left), loss on our contractor dataset (which the VPT model does not train on) begins increasing after 7 epochs. Contractor data could be out-of-distribution because our contractors may have a different distribution of play or because there is some impactful visual domain shift compared to videos from the web. While one could have expected this would be predictive of declining evaluation performance, we do not see notable game statistics from the VPT foundation model rollouts (Figure 4, right) decrease over training, and in the next section we show that BC fine-tuning performance continually improves as the VPT foundation model trains. We provide more insight into this curious phenomenon in Appendix H.

### 4.3 Fine-Tuning with Behavioral Cloning

Foundation models are designed to have a broad behavior profile and be generally capable across a wide variety of tasks. To incorporate new knowledge or allow them to specialize on a narrower task distribution, it is common practice to fine-tune these models to smaller, more specific datasets.<sup>1</sup> The VPT foundation model trained on the broad `web-clean` dataset yet unable to improve the VPT foundation model’s ability to collect and craft these “early game” items by fine-tuning, we attempt to go past this in the technology tree and craft these `contractor_house`, contractors have 10 minutes to build a basic house from scratch using primarily wood, sand, and dirt. Collecting contractor data can be difficult and expensive, so we also construct a dataset `earlygame_keyword` by searching for videos online with descriptions that match keywords such as “new world”, “let’s play episode 1”, etc.; this is a subset of `web_clean` and is labeled with the IDM. See Appendix B.4 and A.3 for full descriptions of both datasets.

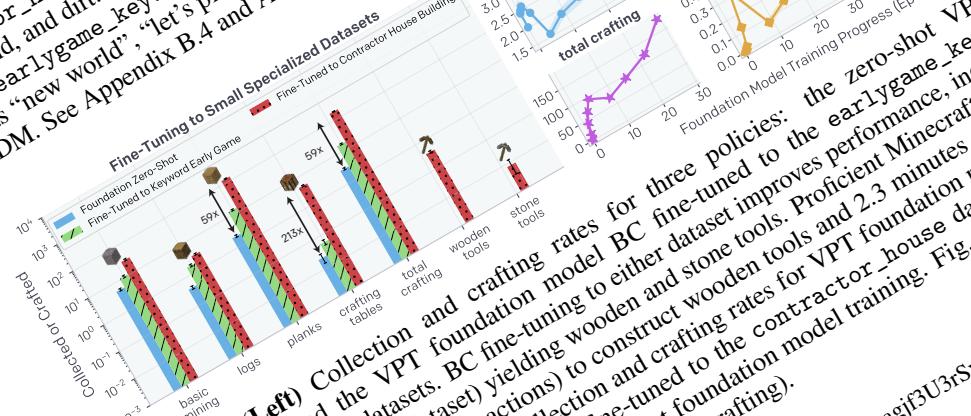


Figure 5: (Left) Collection and crafting rates for three policies: the zero-shot VPT foundation model, and the VPT foundation model BC fine-tuned to either dataset improves the `contractor_house` dataset. BC fine-tuning model yielding wooden and stone tools. Proficient Minecraft, a median of 1.2 minutes (1390 actions) to construct wooden tools and 2.3 minutes to construct stone tools. (Right) Collection and crafting rates for VPT foundation throughout training after they are BC fine-tuned to the `contractor_house` dataset. Proficient Minecraft, crafting-related behaviors increase throughout foundation model training. Fig. 6

(iv) Sample videos: [https://www.youtube.com/playlist?list=PLNAOb\\_agj3U3rS](https://www.youtube.com/playlist?list=PLNAOb_agj3U3rS)

**Fine-tuning with Reinforcement Learning**

Action	Actions	Time	Success Rate
Wooden Pickaxe	200	1.7 minutes	98% 10 min
Cobblestone	2700	2.3 minutes	98% 10 min
Stone Pickaxe	6540	5.4 minutes	84% 10 min
Iron Ore	6540	6.1 minutes	78% 10 min
Furnace	7320	6.1 minutes	60% 10 min
Iron Ingot	10970	9 minutes	60% 10 min
Iron Pickaxe	11181	9.3 minutes	57% 10 min
Diamond Pickaxe	23975	20 minutes	56% 10 min

Despite the foundation model's zero-shot rollout performance plateauing 1/3 into training (Fig. 4, right), fine-tuning performance does continue to increase throughout foundation model training (Fig. 5, right). Additionally, there is a stark difference in performance when training from scratch vs. fine-tuning from the VPT foundation model (Fig. 5 right, comparing the left and rightmost points).

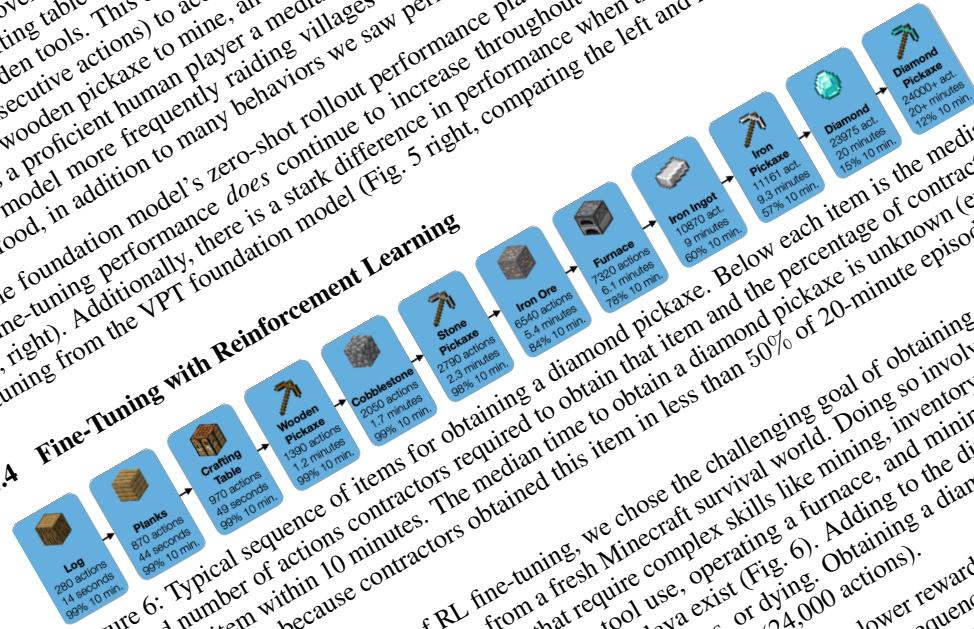
Fine-tuning to earlygame\_keyword results in a large boost (Fig. 5). However, when fine-tuning to this dataset we did not see only a refinement of existing skills. We saw an even bigger improvement and 59x more crafting tables. This entire sequence takes a proficient human player a median of 23 minutes (1390 consecutive actions) to accomplish. The model goes further and collects cobblestone, which requires a wooden pickaxe to mine, and crafts stone tools, requiring it to reveal a new crafting interface; also saw this model more frequently raiding villages that randomly spawn in the game, hunting animals for food, in addition to many behaviors we saw performed by the foundation model.<sup>(v)</sup>

model: 2.5x more crafting tables, 6.1x more planks, 4.3x more logs and 59x more crafting over all. In addition, we saw the emergence of crafting wooden table; this takes a proficient human player a median of 2.3 minutes (213x more crafting tools) to accomplish. The model goes further and collects cobblestone, which requires a wooden pickaxe to mine, and crafts stone tools, requiring it to reveal a new crafting interface; also saw this model more frequently raiding villages that randomly spawn in the game, hunting animals for food, in addition to many behaviors we saw performed by the foundation model.<sup>(v)</sup>

The slide features a large title '4.4 Fine-Tuning with Reinforcement Learning' at the top. Below it is a flowchart with four blue rounded rectangular boxes arranged in a sequence. Each box contains an icon of a game item (Crafting Table, Wooden Pickaxe, Cobblestone, and Stone Pickaxe) and some text. Arrows point from one box to the next in the sequence.

- Crafting Table  
970 actions  
190 seconds  
99% 10 min.
- Wooden Pickaxe  
1390 actions  
12 minutes  
99% 10 min.
- Cobblestone  
2080 actions  
1.7 minutes  
99% 10 min.
- Stone Pickaxe

At the bottom right, there is additional text: 'of items for ob...', 'tractors The...', and a small portion of the previous slide's text: 'is food, in', 'the foundation mo...', 'fine-tuning performance...', 'e-tuning from the VPT foundation mo...', 'e-tuning (...', '5, right). Additionally, there is ...'.



**Re-Tuning with Reinforcement Learning**

Item	Actions	Seconds	Minutes
Planks	970 actions	44 seconds	99% 10 min
Crating Table	970 actions	49 seconds	99% 10 min
Wooden Pickaxe	1380 actions	12 minutes	99% 10 min
Cobblestone	2050 actions	17 minutes	99% 10 min

**Figure 6: Typical sequence of actions for items that got the item within 10 minutes. That it is > 20m) because contractor**

To demonstrate the efficacy of RL fine-tuning, we chose the challenging goal of obtaining a diamond pickaxe within 10 minutes starting from a fresh Minecraft survival world. Doing so involves a diamagnetic sequence of difficult-to-obtain items that require complex skills like mining, inventory management, crafting with and without a crafting table, tool use, operating a furnace, and mining at the depths, where many hazards like enemies and lava exist (Fig. 6). Adding to the difficulty is that items are easily lost by dropping items, destroying items, or dying. Obtaining a diamond pickaxe often than not takes a proficient human over 20 minutes (24,000 actions).

RL experiments use a  $\sim 248$  million parameter VPT policy gradient<sup>64</sup> RL algorithm for  $\sim 1.3$  million episodes (roughly 2 minutes. See Appendix G.1 for reward function and RL training details. RL training with RL is catastrophic forgetting<sup>65,66</sup> because it is realized. For instance, while our VPT required to smelt iron zero-shot combat the catastrophic latent ability throughout RL training.

6: Typical sequence of items for obtaining a diamond pickaxe. Below each item is the name, and number of actions contractors required to obtain that item and the percentage of contractors obtained this item in less than 50% of 20-minute episodes.

Agents can be collected with the phasic policy. Episodes last for 10 minutes and computational constraints limit the number of skills that can be learned. A major problem when fine-tuning the policy is that the sequence of behaviors exhibited by the players smelting with furnaces exhibits the entire sequence of behaviors learned by the agent. It therefore may be necessary to do so have been performed. To overcome this, the Kullback-Leibler (KL) divergence loss between the RL model and the hard an exploration challenge initialized policy fails to achieve almost (Fig. 7a). The model never learns the diamond pickaxe task is for RL in the a diamond pickaxe (Fig. 7b). RL fine-tuning from the VPT foundation model better (Fig. 7a), learning everything up to mining iron ore and crafting furnaces this agent fails at smelting an iron ingot, the next item required to get further information.

(iv) Sample Videos: <https://www.youtube.com>

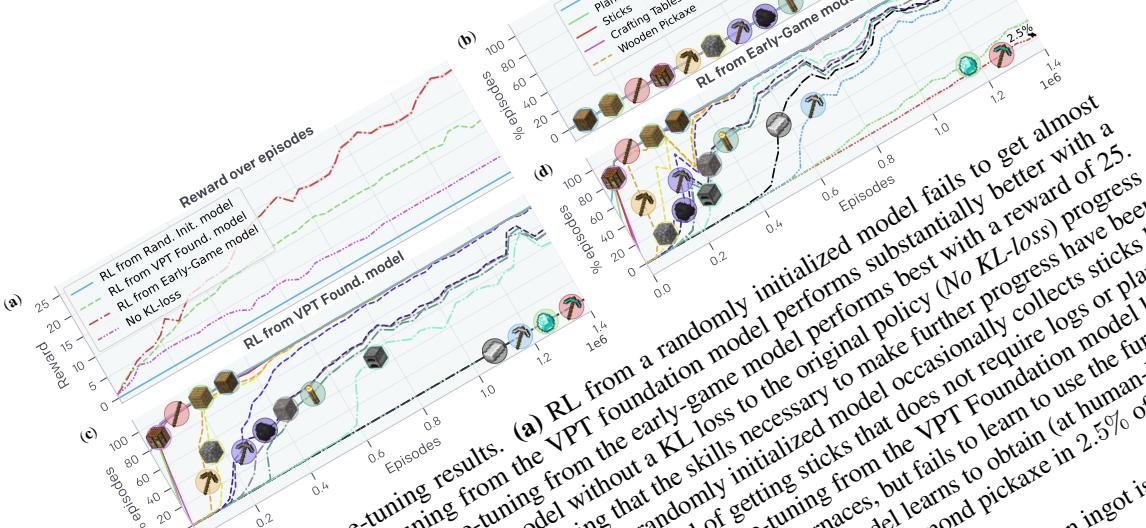


Figure 7: RL Fine-tuning results. (a) RL from a randomly initialized model fails to get almost any reward, RL fine-tuning from the VPT foundation model performs substantially better with a reward near 13, and RL fine-tuning from the early-game model without a KL loss to the original policy (No KL-loss) progresses over 100,000 episodes, suggesting that the skills necessary to make further progress have been catastrophically forgotten. (b) RL from a randomly initialized model occasionally collects sticks by breaking leaves (an easy but inefficient method of getting sticks that does not require logs or planks) and never learns to reliably collect logs. (c) RL fine-tuning from the VPT Foundation model learns everything in the curriculum up to iron ore and making furnaces, but fails to learn to use the furnace to smelt iron ingots. (d) RL fine-tuning from the early-game model learns to obtain (at human-level) all items in the sequence towards a diamond pickaxe and crafts a diamond pickaxe in 2.5% of episodes.

because the zero-shot probability that the VPT foundation model smelts an iron ingot is too low, even when given the prerequisite materials.

Results further improve by first BC fine-tuning the VPT Foundation Model to the `earlygame_keyword` dataset (the *early-game model*, Sec. 4.3) and then fine-tuning with RL (Fig. 7a), which in preliminary experiments we found to perform better than first fine-tuning to `contractor_house` followed by BC fine-tuning, and then RL fine-tuning with RL (Appendix G.2). The three-phase training to (pretraining, BC fine-tuning, and then RL fine-tuning) succeeds in learning extremely difficult tasks: it achieves over 80% reliability on obtaining a diamond pickaxe (Fig. 7d). For comparison, human players give 2.5% reliability on obtaining a diamond pickaxe collect these items in 57%, 15%, and 12% of episodes, respectively, meaning our model is human-level for collecting diamonds and simplifying action space designed to ease exploration. To the best of our knowledge, **we are the first to report non-zero success rates on crafting a diamond pickaxe**. Qualitatively, the model developed useful skills for diamond mining, such as efficient mining patterns, cave exploration, returning to previously placed objects like crafting tables, and advanced techniques like using wooden pickaxes as fuel when moving on to iron tools.<sup>(vi)</sup>

The treatment without the importance of the KL loss to the pretrained model during RL fine-tuning (Fig. 7a) limits its reward (Fig. 7a). This failure to progress further into the sequence (logs, planks, sticks, and crafting tables) likely because, while the initial skills of chopping logs and crafting planks are lost due to catastrophic forgetting,

#### 4.5 Data Scaling Properties of the Foundation Model

In this section we validate a core hypothesis behind this work: that it is far more efficient to labeled contractor data to train an IDM within the VPT method than it is to directly collect contractor data from that same small contractor dataset. If we could cheaply collect a contractor dataset of a similar order of magnitude as `web_clean`, then this would however, collecting that scale of data would have cost millions of dollars. The `web_clean` dataset contains foundation models trained on increasing orders of magnitude of data from 1 hour to 1k hours<sup>(vii)</sup>.

<sup>(vi)</sup> Videos found at [https://www.youtube.com/playlist?list=PLNAOlb\\_agjf3e\\_KE](https://www.youtube.com/playlist?list=PLNAOlb_agjf3e_KE)

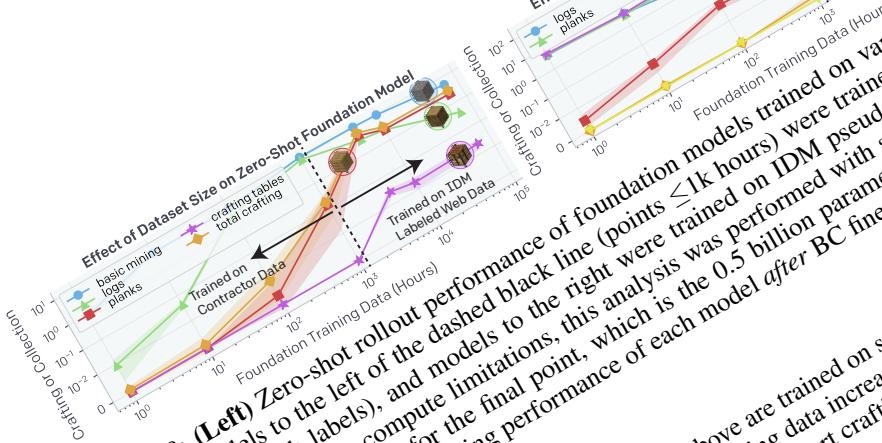


Figure 8: **(Left)** Zero-shot rollout performance of foundation models trained on varying amounts of data. Models to the left of the dashed black line (points  $\leq 1k$  hours) were trained on contractor data (ground-truth labels), and models to the right were trained on IDM pseudo-labeled subsets of `web_clean`. Due to compute limitations, this analysis was performed with smaller (71 million parameter) models except for the final point, which is the 0.5 billion parameter VPT foundation model. **(Right)** The corresponding performance of each model after BC fine-tuning each model to the `contractor_house` dataset.

contractor data, and those trained on 5k hours and above are trained on subsets of `web_clean`, which does not contain any IDM contractor data. Scaling training data increases log collection, mining, and crafting capabilities. The zero-shot model only begins to start crafting crafting tables at over 5000 hours of training data. When fine-tuning each foundation model to `contractor_house`, we see that entire  $\sim 70k$  hour `web_clean` dataset. We furthermore only see the emergence of crafting stone tools at the largest data scale.

#### 4.6 Effect of Inverse Dynamics Model Quality on Behavioral Cloning

This section investigates how downstream BC performance is affected by IDM quality. We train IDMs on increasingly larger datasets and use each to independently label the `earlygame_keyword` dataset (this smaller dataset was chosen due to a limited compute budget). We then train a BC model from scratch on each dataset and report game statistics for each model as a function of IDM contractor dataset size (Fig. 9).

IDMs trained on at least 10 hours of data are required for any crafting, and the crafting rate increases quickly up until 100 hours of data, likely due to noise. Similarly, crafting tables are again gains plateau after 100 hours. While in all previous experiments we use our best IDM trained on 1962 hours of data, these results suggest we could reduce that number to as low as 100 hours.

## 5 Discussion and Conclusion

The results presented in this paper help pave the path to utilizing the wealth of unlabeled data on the web for sequential decision domains. Compared to generative video modeling or contrastive methods that would only yield representational priors, VPT offers the exciting possibility of directly *to act* during pretraining and using these learned behavioral priors as extremely effective downstream task priors for RL. VPT could even be a better general representation learning method than this present in features trained to correctly predict the distribution over future human actions in this intriguing direction to future work.



Figure 9: Zero-shot performance of BC models trained from scratch on the `earlygame_keyword` dataset labeled with IDMs that were trained on increasing amounts of contractor data. In all cases, gains are few to no gains and differences are small, with the best BC models trained on 100 hours of data.

## References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[3] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

[4] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.

[5] Dhruv Mahajan, Ross Girshick, Vignesh Ramamathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pages 181–196, 2018.

[6] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.

[7] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Emma Brunskill, Sydne Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Simran Arora, and Lucas Beyer. Scaling vision opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[8] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Julian Schrittwieser. Taming the game of go with deep neural networks and tree search. *Nature*, 572(7771):492–496, 2019.

[9] David Silver, Aja Huang, Chris J Maddison, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Julian Schrittwieser, Ioannis Antonoglou, Arthur Guez, Laurent Sifre, George van den Driessche, and Lucas Beyer. Scaling vision representations. *CoRR*, abs/2106.04560, 2021. URL <https://arxiv.org/abs/2106.04560>.

[10] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Timo Ewalds, Petko Georgiev, and Yuhuai Wu. Young chung, David H Choi, Richard Powell, and Tom Schaul. Starcraft ii using multi-agent reinforcement learning. *Nature*, 572(7771):497–501, 2019.

- [11] Richard S Sutton and Andrew G Barto. *Reinforcement learning*. MIT press, 2018.
- [12] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, PZ Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [13] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob M. Munos. Emergent tool use from multi-agent autocurricula. *arXiv preprint arXiv:1909.07528*, 2019.
- [14] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marrs, Guy Lever, Antonio Garcia Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-level performance in 3d multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019.
- [15] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapurowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human benchmark. In *International Conference on Machine Learning*, pages 507–517. PMLR, 2020.
- [16] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- [17] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- [18] Adrien Ecoffet, Joost Huijzinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then explore. *Nature*, 590(7847):580–586, 2021.
- [19] Peter C Humphreys, David Raposo, Toby Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair Muldal, Josh Abramson, Petko Georgiev, Alex Goldin, Adam Santoro, et al. A data-driven approach for learning to control computers. *arXiv preprint arXiv:2202.08137*, 2022.
- [20] Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In *Doina Precup and Yee Whye Teh, editors, Proceedings of the 34th International Conference on Machine Learning*, pages 3135–3144. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/shi17a.html>.
- [21] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *Icmi*, volume 1, page 2, 2000.
- [22] Faraz Torabi, Garrett Warnell, and Peter Stone. Recent advances in imitation learning from observation. *arXiv preprint arXiv:1905.13566*, 2019.
- [23] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.
- [24] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.
- [25] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation via context translation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.
- [26] Twinfinite Staff. Most played games in 2020-ranked-by-peaks. Twinfinite. URL <https://twinfinite.net/2020-ranked-by-peaks/>.
- [27] William H Guss, Brandon Houghton, Nicholas Topin, Phillip Wang, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of multiagent reinforcement learning environments. *arXiv preprint arXiv:1907.13440*, 2019.

- [28] Chen Tessler, Shahar Givony, Tom Zahavy, Daniel Mankowitz, and Yair Feinerman. Hierarchical approach to lifelong learning in minecraft. In *Proc. of Artificial Intelligence*, volume 31, 2017.
- [29] Christian Scheller, Yanick Schraner, and Manfred Vogel. Sample efficient learning through learning from demonstrations in minecraft. In *NeurIPS 2019 Conference Demonstration Track*, pages 67–76. PMLR, 2020.
- [30] Ingmar Kanitscheider, Joost Huizinga, David Farhi, William Heben Guss, Brandon Houghteling, Raul Sampedro, Peter Zhokhov, Bowen Baker, Adrien Ecoffet, Jie Tang, et al. Multi-task curriculum learning in a complex, visual, hard-exploration domain: Minecraft. *arXiv preprint arXiv:2106.14876*, 2021.
- [31] Junhyuk Oh, Valliappa Chockalingam, Honglak Lee, et al. Control of memory, active perception, and action in minecraft. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2016.
- [32] Vihang P Patil, Markus Hofmarcher, Marius-Constantin Dinu, Matthias Dorfer, Patrick M Blies, Johannes Brandstetter, Jose A Arjona-Medina, and Sepp Hochreiter. Align-rudder: Learning from few demonstrations by reward redistribution. *arXiv preprint arXiv:2009.14108*, 2020.
- [33] Alexey Skrynnik, Aleksey Staroverov, Ermek Aitykulov, Kirill Aksenov, Vasili Davydov, and Aleksandr I Panov. Forgetful experience replay in hierarchical reinforcement learning. *arXiv preprint arXiv:2006.09939*, 2020.
- [34] Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, and Wei Yang. Juewu-mc: Playing minecraft with sample-efficient hierarchical reinforcement learning. *arXiv preprint arXiv:2112.04907*, 2021.
- [35] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1, 1988.
- [36] Stefan Schaal. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3(6):233–242, 1999.
- [37] Breanna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- [38] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Christina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- [39] Claude Sammut, Scott Hurst, Dana Kedzier, and Donald Michie. Learning to fly. In *Machine Learning Proceedings 1992*, pages 385–393. Elsevier, 1992.
- [40] Alessandro Giusti, Jérôme Guzzi, Dan Cireşan, Fang-Lin He, Juan P Rodríguez, Flavio Fontana, Matthias Faessler, Christian Forster, Jürgen Schmidhuber, Gianni Di Caro, et al. A machine learning approach to visual perception of forest trails for mobile robots. *IEEE Robotics and Automation Letters*, 1(2):661–667, 2015.
- [41] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Urs Müller, Jiaakai Zhang, et al. End-to-end driving via conditional imitation learning. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 4693–4700. IEEE, 2018.
- [42] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. Computing “elo ratings” of move patterns in the game of go. In *Proc. of Artificial Intelligence*, volume 31, 2017.
- [43] Rémi Coulom. Computing “elo ratings” of move patterns in the game of go. In *Proc. of Artificial Intelligence*, volume 31, 2017.
- [44] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, et al. Deep q-learning from demonstrations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.

- [45] Ashley Edwards, Himanshu Sahni, Yannick Schroeck, and Jitendra Malik. 2019. Learning policies from observation. In *International conference on machine learning*. PMLR, 1–10.

[46] Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, and Jitendra Malik. 2018. Reinforcement learning of physical skills from videos. *ACM Transactions On Graphics* 37(6):1–14.

[47] Feryal Behbahani, Kyriacos Shiarlis, Xi Chen, Vitaly Kurin, Sudhanshu Kasewa, Ciprian Stirbu, Joao Gomes, Supratik Paul, Frans A Oliehoek, Joao Messias, et al. Learning from demonstration in the wild. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 775–781. IEEE, 2019.

[48] Yusuf Aytar, Tobias Pfaff, David Budden, Thomas Paine, Ziyu Wang, and Nando De Freitas. Playing hard exploration games by watching youtube. *Advances in neural information processing systems*, 31, 2018.

[49] Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Dian Chen, Yide Shen, Evan Shelhamer, Jitendra Malik, Alexei A. Efros, and Trevor Darrell. Zero-shot visual imitation. In *ICLR*, 2018.

[50] Ashvin Nair, Dian Chen, Pulkit Agrawal, Philip Isola, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Combining self-supervised learning and imitation for vision-based rope manipulation. In *European symposium on artificial neural networks, number CONF, pages 2146–2153*, 05 2017. doi: 10.1109/ICRA.2017.7989247.

[51] Duy Nguyen-Tuong, Jan Peters, Matthias Seeger, and Bernhard Schölkopf. Learning inverse dynamics: a comparison. In *European symposium on artificial neural networks, number CONF, pages 2146–2153*, 05 2017. doi: 10.1109/ICRA.2017.7989247.

[52] David Abel, Alekh Agarwal, Fernando Diaz, Akshay Krishnamurthy, and Sergey Levine. Exploratory gradient boosting for reinforcement learning in complex domains. *arXiv preprint arXiv:1603.04119*, 2016.

[53] Dilip Arumugam, Jun Ki Lee, Sophie Saskin, and Michael L Littman. Deep reinforcement learning from policy-dependent human feedback. *arXiv preprint arXiv:1902.04257*, 2019.

[54] Alexander Trott, Stephan Zheng, Caiming Xiong, and Richard Socher. Keeping your distance: Solving sparse reward tasks using self-balancing shaped rewards. *Advances in Neural Information Processing Systems*, 32, 2019.

[55] Stephan Alanić. Deep reinforcement learning with model learning and monte carlo tree search in minecraft. *arXiv preprint arXiv:1803.08456*, 2018.

[56] Hiroto Udagawa, Tarun Narasimhan, and Shim-Young Lee. Fighting zombies in minecraft with deep reinforcement learning. Technical report, Technical report, Technical report, Technical report, Stanford University, 2016.

[57] Tianmin Shu, Caiming Xiong, and Richard Socher. Hierarchical and interpretable skill generation in multi-task reinforcement learning. *arXiv preprint arXiv:1712.07294*, 2017.

[58] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. *arXiv preprint arXiv:1712.07294*, 2017.

[59] Zhengxiang Shi, Yue Feng, and Aldo Lipani. Learning to execute or ask clarifications. *IEEE transactions on neural networks and learning systems*, 31(1), 2020.

[60] Tambe Matiisen, Avital Oliver, Taco Cohen, and John Schulman. Teacher-student learning. *IEEE transactions on neural networks and learning systems*, 31(1), 2020.

[61] Robert George Douglas Steel, James Hiram Torrie, et al. *Principles and procedures of statistics*, 1960.

- [62] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [63] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [64] Karl W Cobbe, Jacob Hilton, Oleg Klimov, and John Schulman. Phasic policy gradient. In *Machine Learning and Tong Zhang, editors, Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 2020–2027. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/cobbe21a.html>*.
- [65] Dhireesha Kudithipudi, Mario Aguilera-Simon, Jonathan Babb, Maxim Bazhenov, Douglas Blackiston, Josh Bongard, Andrew P Brna, Suraj Chakravarthi Raja, Nick Cheney, Jeff Clune, et al. Biological underpinnings for lifelong learning machines. *Nature Machine Intelligence*, 4(3):196–210, 2022.
- [66] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [67] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big???. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 2021.
- [68] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [70] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [71] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backpropagation for neural networks: Tricks of the trade, pages 9–48. Springer, 2012.
- [72] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [73] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurti, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance.pdf>.
- [74] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional convex optimization. *Advances in neural information processing systems*, 2016.
- [75] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features? *Advances in neural information processing systems*, 2013.

- [76] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Transformer-xl: Attentive language models beyond a fixed arXiv:1901.02860, 2019.
- [77] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [78] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. Dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [79] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- [80] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- [81] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *Advances in neural information processing systems*, 30, 2017.
- [82] Tom Schaul, Daniel Horgan, Karol Gregor, and Wojciech Zaremba. Hindsight experience proximators. In *International conference on machine learning*, pages 1312–1320. PMLR, 2015.
- [83] Open Ended Learning Team, Adam Stooke, Anuj Mahajan, Catarina Barros, Charlie Deck, Jakob Bauer, Jakub Sygnowski, Maja Trelacz, Max Jaderberg, Michael Mathieu, et al. Open-ended learning leads to generally capable agents. *arXiv preprint arXiv:2107.12808*, 2021.
- [84] Jelena Luketina, Nantas Nardelli, Gregory Farquhar, Jakob Foerster, Jacob Andreas, Edward Grefenstette, Shimon Whiteson, and Tim Rocktäschel. A survey of reinforcement learning informed by natural language. In *Proceedings of the International Joint Conference on Artificial Intelligence*, IJCAI-19, pages 6309–6317. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/880. URL <https://doi.org/10.24963/ijcai.2019/880>.
- [85] DeepMind Interactive Agents Team, Josh Abramson, Arun Ahuja, Arthur Brussee, Federico Carnevale, Mary Cassin, Felix Fischer, Petko Georgiev, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical multimodal interactive agents with imitation and self-supervised learning. *arXiv preprint arXiv:2112.03763*, 2021.
- [86] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [87] Dong Yu and Li Deng. Automatic speech recognition, volume 1. Springer, 2016.
- [88] Daulet Nurmanbetov. punct, May 25 2021. URL <https://github.com/Felflare/rpulse>. accessed 2022-04-22.
- [89] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Zhu, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*, 2022.

## Acknowledgements

We thank the following people for helpful discussions and support: Bob McGrew, Joel Lehman, Ilya Sutskever, Wojciech Zaremba, Ingmar Kanitscheider, David Gordon, and Christopher Berner.

## Supplementary Information

# A Collecting Internet Data

## Collecting Internet Data

### A.1 Initial Unclean Dataset Curation

**Collecting Internet Data**

**1.1 Initial Unclean Dataset Curation**

Our goal was to curate a video dataset of Minecraft gameplay from the survival game mode. Additionally, we prefer the data come from game modes as close as possible to our evaluation environment, meaning preferably coming from Minecraft version 1.16, being on a computer (which uses a mouse and keyboard vs. video game controllers with keypads and other buttons), being single- (vs. multi-) player, and having the default look of the game (vs. modifications that alter that style, such as to make it look realistic). To try to accomplish these goals, we collect a dataset by performing keyword searches of publicly available videos on the internet. A list of search queries we used are given in Table 1.

Table 1.

Supplementary Information

minecraft survival longplay  
minecraft gameplay no webcam  
minecraft survival tutorial mode  
minecraft survival let's play  
minecraft survival guide  
craft survival for beginners  
beginners guide  
raft starter guide  
guide 1.16  
new survival world

minecraft  
rollers with kc  
book of the game (vs.  
o accomplish these goals, w  
e videos on the internet. A list of  
minecraft survival longplay  
minecraft gameplay no webcam  
minecraft survival tutorial  
minecraft survival guide  
minecraft survival let's play  
minecraft survival for beginners  
minecraft beginners guide  
ultimate minecraft starter guide  
minecraft how to start a new survival world  
minecraft survival fresh start  
minecraft survival let's play episode 1  
let's play minecraft episode 1  
minecraft survival 101  
minecraft survival learning to play  
how to play minecraft  
how to play minecraft basic  
minecraft survival for noobs  
minecraft survival for dummies  
how to play minecraft tutorial for beginners  
minecraft survival new world  
minecraft survival a new beginning  
minecraft survival episodio 1  
minecraft survival 1. bőlüm  
i made a new minecraft survival world  
arch terms used for generating the initial  
e perform an additional  
tribution. In this  
videos th

Table 1: Search terms used for generating the initial web dataset.

gplay  
no webcam  
y survival mode  
l tutorial  
al guide  
ival let's play  
ival for beginners  
rvival for beginners  
beginners guide  
inecraft starter guide  
raft survival guide 1.16  
raft how to start a new survival world  
raft survival let's play episode 1  
craft survival let's play episode 1  
craft survival episode 1  
t's play minecraft 101  
inecraft survival learning to play  
minecraft survival basic  
how to play minecraft  
minecraft survival for noobs  
minecraft survival for dummies  
minecraft survival tutorial for beginners  
how to play minecraft new world  
minecraft survival a new beginning  
minecraft survival episodio 1  
minecraft survival europa 1  
minecraft survival 1. bolum  
i made a new minecraft survival world

Table 1: Search terms used for generating the initial web dataset.

For videos that have metadata available, we perform an additional step of metadata-based filtering to eliminate videos that do not fit our target distribution. In this step, we look for a list of keywords in the video title and description and reject videos that contain these terms. The keywords we use are: {ps3, ps4, ps5, xbox 360, playstation, timelapse, multiplayer, The pocket edition, skyblock, realistic minecraft, how to install, how to download, realmcr... This process yielded us ~270k hours of unlabeled data, which we filter down to only as described in the next section.

### Training a Model to Filter out Unclean Video Segments

The scope of this work to the Minecraft Survival game mode includes clips that are obtained from this mode that are relatively

at have metadata in videos that do not fit on the video title and descriptions we use are: {ps3, ps4, ps5, xbox one, skyblock, realistic minecraft, how to process yielded us ~270k hours of unlabeled data, described in the next section.

## 4.2 Training a Model to Filter out Unclean Video Segments

We restrict the scope of this work to the Minecraft Survival game mode training dataset to clips that are obtained from this mode that are relatively

To do so, we asked contractors to label a set of random video frames (N=8800). These images were from a random subset of the videos of the project (Section A.1).

### A.2.1 Label Collection

We asked 5 workers on Amazon Mechanical Turk (mTurk) that we selected with a sample query task to label random screen capture images to be used in training the classifier. A sample interface that the workers saw on mTurk is given in Figure 10.

We asked workers to label videos as being in one of the following three categories (see Figure 11 for visual examples of each class):

1. Minecraft Survival Mode - No Artifacts: Video frames (images) that correspond to the Minecraft Survival game mode that do not contain any non-game visual artifacts (e.g. subscribe buttons, channel logos, advertisements, picture-in-picture of the narrator, etc.).
2. Minecraft Survival Mode - with Artifacts: Video frames (images) that include such visual artifacts.
3. None of the Above: Video frames (images) that are not from the Minecraft survival game mode, including those from other Minecraft game modes such as creative mode or even other games/topics entirely.

The full set of instructions workers received are as follows (note that we also included multiple image examples from each category in the worker instructions, similar to the sample subset provided in Figure 11):

Please help us identify screenshots that belong only to the survival mode in Minecraft. Everything else (Minecraft creative mode, other games, music videos, etc.) should be marked as None of the Above. Survival mode is identified by the info at the bottom of the screen:

- a health bar (row of hearts)
- a hunger bar (row of chicken drumsticks)
- a bar showing items held

#### Survival Mode Creative Mode

Creative mode only has an item hotbar and should be classified as None of the Above.

#### Label Descriptions

- Minecraft Survival Mode - No Artifacts: These images will be clean screenshots from the Minecraft survival mode gameplay without any noticeable artifacts.
- Minecraft Survival Mode - with Artifacts: These images will be valid screenshots mode overlays (a logo/brand), text annotations, a picture-in-picture of the player, etc.
- None of the Above: Use this category when the image is not a valid Minecraft screenshot. It may be a non-Minecraft frame or from a different game mode. In no game modes such as the creative mode, the health/hunger bars will be missing.

In total, we spent \$319.96 on human labeling experiments on mTurk, of which \$159.96 was paid to workers. The remaining amount was spent towards Amazon platform fees received \$0.01 per labeled image, at an hourly compensation of \$7.20 (based on an average time of 5 seconds/image – in our internal sample run of the same task, we found an average time to be < 3 seconds).

Since we perform rigorous keyword and metadata based filtering of videos (against which we served sample images to be labeled, serving offensive content to the user),

low risk and no such images were detected during our main experiment, and the workers were fully anonymized via personally identifiable information (PII) was collected.

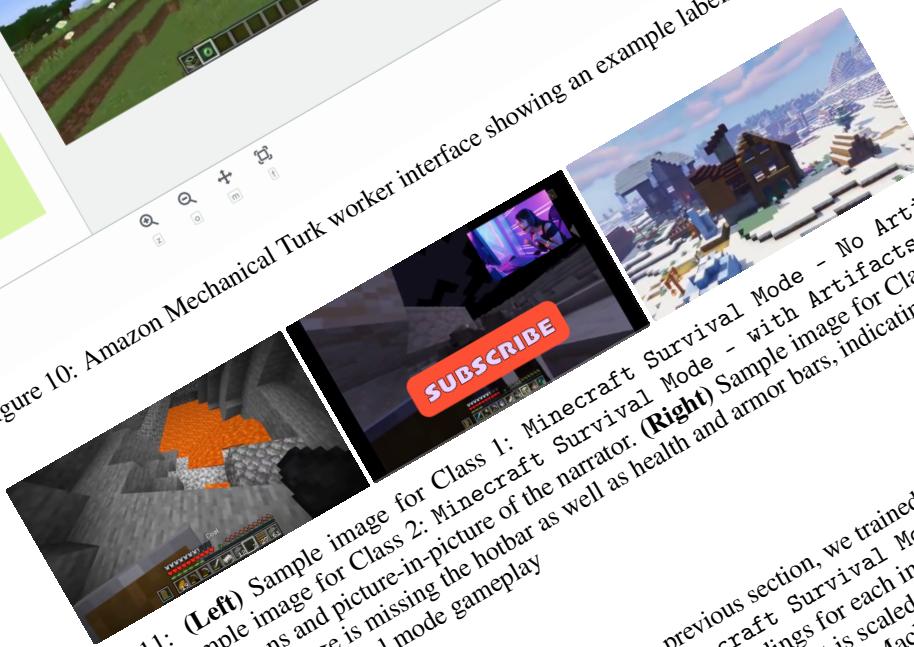
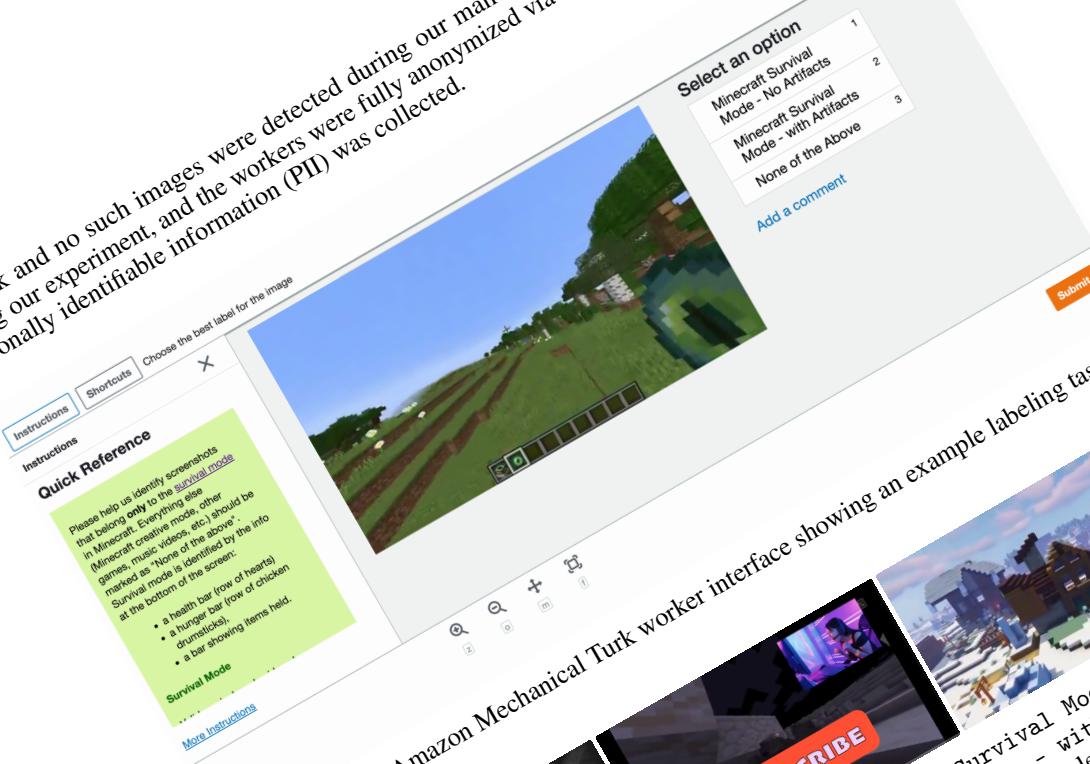


Figure 11: **(Left)** Sample image for Class 1: Minecraft Survival Mode - No Artifacts. **(Middle)** Sample image for Class 2: Minecraft Survival Mode - with Artifacts - Image contains annotations and picture-in-picture of the narrator. **(Right)** Sample image for Class 3: None not captured during survival mode gameplay

VM Training

Labels collected as described in the previous section, we trained a classifier to extract frames from the Minecraft Survival Mode - No Artifacts labeled images, we obtain embeddings for each image using the RN50x64 ResNet-based CLIP model that is scaled up to have approximately 16 times the resolution of the original image. We then train a Support Vector Machine (SVM) using the Scikit-learn<sup>68</sup> SVM implementation with the parameters: C=1, gamma=0.001, kernel='rbf', class\_weight='balanced'. The training sequences at a rate of 3 frames/second. We then consider not clean sequences at a rate of 3 frames/second. Finally, we segment videos by frame. The final web-based

## A.2.2 SVM Training

The image is a collage of several screenshots from the video game Minecraft. One screenshot shows a first-person view of a snowy landscape with a wooden house in the background. Another shows a character standing in a dark, stone-walled room. A third shows a character in a creative mode environment with floating blocks. A fourth screenshot is a close-up of a computer interface showing a frame classifier with various labels like 'Survival Mode - No Artifacts', 'Survival Mode - with Artifacts - 1', and 'Survival Mode - with Artifacts - 2'. A large red 'SUBSCRIBE' button is overlaid on one of the screenshots.

Finally, we apply the classification for videos that consist of a "clean" segments that we apply a median filter with Artifacts Mode – with Artifacts Mode. Given a set of frame segments, we can compute a frame classifier configuration given in Table 2.

- apply a median "clean" segments that

4.3 **early-game Dataset**  
The early-game dataset is a ~3000 hour subset of web\_clean and determine, where players start in a fresh world, i.e. instances where the videos in web\_clean and determine match. The behavior, i.e. instances where players start in a fresh world, i.e. instances where the videos in web\_clean and determine match, is accompanied by text that accompanies the expressions match:

	Kernel	Kernel Size	Number of Epochs	Batch Size	Learning Rate	Optimizer	Model Type	Model Path
CLIP Model Specification	Ridge	448x448	1000	1024	0.001	SGD	SVM	models/clip_svm.pkl
CLIP Input Image Resolution	Kernel	1024	2000	C	1.0	SGD	SVM	models/clip_svm.pkl
CLIP Embedding Feature Length	Gamma	20	2000	Gamma	0.001	SGD	SVM	models/clip_svm.pkl
SVM Parameters	Class 1	2000	2000	Class 2	0.001	SGD	SVM	models/clip_svm.pkl
Sample Size	Class 3	2000	4400	Class 3	0.001	SGD	SVM	models/clip_svm.pkl

Section Details and SVM Configuration. The parameters are as follows:

- Kernel: Ridge
- Kernel Size: 448x448
- Number of Epochs: 1000
- Batch Size: 1024
- Learning Rate: 0.001
- Optimizer: SGD
- Model Type: SVM
- Model Path: models/clip\_svm.pkl

Table 2: Feature Extraction Details and SVM Classification Results for the Scikit-learn implementation in Scikit-learn<sup>68</sup>.

- (episodes|decks|series|part|partner)\*(.1#1|1/.01#01|0)
  - start
  - beginning
  - (new|fresh|clean).\*(world|game|play)
  - from scratch

of videos, we take only the first

- start
- beginning
- (new|fresh|clean):\*(world|game|play)
- from scratch

From this set of videos, we take only the first 5 minutes of each video.

## Contractor Data

## Contractor Play

Minecraft recorder that we built to implement that using the MCE modpacks. It records all the player movements and events in a database.

**B Contractor Data**

## **Contractor Data**

### **B.1 Recording Contractor Play**

Our contractors use a custom Minecraft recorder that we built that records their actions and game video feeds as they play. The recorder is implemented using the MCP-Reborn (github.com/Hexception/MCP-Reborn) modding package. To ensure that the recorder environment is as close as possible to the Minecraft underlying game engine for both. The recorder is a Java app that runs in a window mode, with constant resolution of 1280x760. Brightness is set to 0 (the "gloomy" setting in Minecraft), which is the default setting. Other graphics settings (field of view, GUI scale) are fixed to the values used in the Minecraft environment (C.1); we explicitly prevented users from changing graphics settings. Unlike the environment, the recorder allows all keyboard key presses and continuous (as opposed to binned) mouse actions. On every game step (or “tick”) the frame buffer used to display the game window is downsized to 640x360 and written into a video file. In-game actions are recorded in a separate JSONL file (a text file where each line is a JSON-formatted string). All recordings are chunked into 5 minute clips: after each 5 minute segment of contractor state is recorded, the recorder automatically uploads the video file, the JSONL file with actions, as well as a Minecraft state file to individual cloud bucket, as well as with credentials giving write access only to that bucket. Credentials also included adjective-adjective-noun names (e.g. grumpy-amethyst-chipmunk), generated by the namegenerator python package to ensure contractor anonymity when we publish the data.

## Contractor Data

### 1 Recording Contractor Play

#### Contractor Contract

Contractors by posting the following offer on the UpWork freelancing platform can earn data for training AI models in Minecraft. You'll need to provide us with a modified version of Minecraft (that collects and stores contractor data) and a draft survival mode! Paid per hour of game time played. We do not collect any data that is necessary. We do not collect any data that is unnecessary.

## B.2 Contractor Contract

**2.2 Contractor Contract**

We recruited contractors by posting the following offer on the UpWork freelancing site:

“We are collecting data for training AI models in Minecraft. You’ll need to download the modified version of Minecraft (that collects and uploads data), and play Minecraft survival mode! Paid per hour of game time in Minecraft not necessary. We do not collect any data from your computer.”

We had the applications open for a day, and then randomly selected contractors. Later in the project, as we needed more data and as some working contractors, we added more applicants from the original pool as well as applicable taxes). All of the contractors were paid \$20 per hour (minus Upwork fees). Over the course of the project, we collected some data based on about \$90,000. Bugs in the recorder and for some ideas of human play that was not used for training, for contractor compensation over the course of the project, we did not use any bugs in the recorder. However, as total, we spent about \$40,000. Could likely obtain most of our results with an IDM trained using only \$2000 worth of data, i.e. the foundation VPT model, BC fine-tuning to the earlygame keyword dataset, and the RL fine-tuning results. Collecting the contractor data, the actual cost about \$8000. Because we used the IDM trained on about 2000 hours of contractor data, the cost of contractor data for those results was around \$40,000.

In early stages of the project, we were planning to use contractor data solely for the purpose of training the IDM. As such, no specific tasks were given, other than “play the survival mode of Minecraft like you normally would.” Later in the project, we requested that contractors perform specific tasks in Minecraft, such as:

- Collect as many units of wood as possible, using only wooden or stone tools (treechop)
- Start a new world every 30 minutes of game play
- Build a basic house in 10 minutes using only dirt, wood, sand, and either wooden or stone tools (contractor\_house, more details below in Appendix B.4).
- Starting from a new world and an empty inventory, find resources and craft a diamond pickaxe in 20 minutes (`obtain_diamond_pickaxe`). This dataset was used to obtain statistics for how long it takes humans on average to complete this task (and the subtasks required to complete it) when obtaining a diamond pickaxe is their goal.

Since we only recorded in-game events and videos, the data does not include personally identifiable information. That being said, the contractors could theoretically use Minecraft’s open-world property blocks to write their name or offensive content and/or offensive content (e.g. by using Minecraft blocks to write their name or offensive content, then finding a spot from which the message would be visible). In practice, we have not seen any attempts to do so in the contractor videos that we watched. Of course, our BC models on videos from the internet of people playing Minecraft, such behavior is rare enough that our model could also potentially learn it, although we expect

### B.3 Data for the Inverse Dynamics Model.

Since the IDM’s task is to infer actions given the video, any labelled data is appropriate for IDM training. In practice, we included general gameplay as well as the treechop task data described in the previous section, which amounted to a total of 1962 hours. Due to collecting datasets like contractor\_house only at late stages of the project, they were not included in IDM training.

### B.4 contractor\_house.

The contractor\_house contains about 420 hours of data. We asked contractors to build a house in 10 minutes, using only basic dirt, wood, and sand, blocks. Each trajectory starts in generated world and a timer forcibly ends a trajectory after a 20 minute limit. For this reason, contractors chose to begin their trajectories by crafting basic tools and building blocks, it was common for the first 2 minutes to be spent crafting a wooden pickaxe and then for an assortment of stone tools before gathering more building blocks and beginning a structure.

## C Minecraft environment details

Our Minecraft training environment is a hybrid between MineRL<sup>27</sup> (github.com/Hexception/MCP-Reborn) Minecraft modding package. Unlike



Figure 12: **(Left)** Sample of a Minecraft frame in the original resolution (640x360) with an in-game GUI open. The mouse cursor can be seen in the center of the image. This particular GUI shows the player's inventory and can be used to craft very basic items. **(Middle)** We downsample images to 128x128 for computational reasons. Shown is a downsampled observation (400x360) with an in-game GUI. The health, hunger, hotbar overlays, and agent hand can be seen in the lower part of the image.

game, in which the server (or the "world") always runs at 20Hz and the client runs as fast as rendering can complete (typically at 60-100Hz), in our version the client and server run in the same thread at the same frequency. This allows us to run the environment slower or faster than real time, while avoiding artifacts like missing chunks of the world. The action and observation spaces are similar to those of MineRL environments and are described in more detail in the following subsections. The environment also returns diagnostic information, such as in-game stats, contents of the agent's inventory, whether any in-game GUI is open, etc., which we use for tracking and recording but not as inputs to the models. The episode length is 10 minutes for RL experiments and 60 minutes for BC model evaluations. The agent can "die" in a number of ways, such as staying under water for too long and drowning, being killed by hostile mobs, or falling from a tall structure. We do not terminate the episode on agent "death". Instead, just as for humans in the regular Minecraft game, the agent drops all its items when it dies and respawns at a random spot close to the initial spawning spot in the same Minecraft world. The policy state is not masked on death, so the model can remember the fact that it has died and act accordingly.

### C.1 Observation space

The environment observations are simply the raw pixels from the Minecraft game that a human would see. Unlike MineRL, we do not remove overlays like the hotbar, health indicators, and the animation of a moving hand shown in response to the attack or "use" actions. The field of view is 70 degrees, which corresponds to the Minecraft default. GUI scale (a parameter controlling the size of the in-game GUI) is set to 2, and brightness is set to 2 (which is not a Minecraft default, but as very frequently used in online videos). The rendering resolution is 640x360, which is downsampled to 128x128 before being input to the models. We empirically found 128x128 to be the smallest resolution for which in-game GUI elements are still discernible, and then chose that to minimize compute costs. Whenever an in-game GUI is open, we additionally render an image of a mouse cursor at the appropriate mouse position to match what a human player's operating system does (Figure 12).

### C.2 Action space

Our action space includes almost all actions directly available to human players, such as keyboard movements, and clicks. The specific binary actions we include are shown in Table 12. One difference between the human action space and our agent's is that we disallow recipe letters, which is only useful for entering text into the search bar of the crafting recipe book. However, because we do allow the agent to press letters that are also short-circuited outside of the GUI, the "W" key triggers the agent to press letters that are also short-circuited within the GUI (W, A, S, D, E, Q) that produce letters if the forward action agents have not seen agents attempt to search the recipe book with these letters. Instead, the recipe book with the mouse or craft by dragging items around the crafting interface.

Action	Human action	Description
forward	W key	Move forward.
back	S key	Move backward.
left	A key	Strafe left.
right	D key	Strafe right.
jump	space key	Jump.
inventory	E key	Open or close inventory and the 2x2 crafting grid.
sneak	shift key	Move carefully in current direction of motion. In the GUI it acts as a modifier key: when used with attack it moves item from/to the inventory to/from the hotbar, and when used with craft it crafts the maximum number of items possible instead of just 1.
sprint	ctrl key	Move fast in the current direction of motion.
attack	left mouse button	Attack; In GUI, pick up the stack of items or place the stack of items in a GUI cell; when used as a double click (attack - no attack - attack sequence), collect all items of the same kind present in inventory as a single stack.
use	right mouse button	Place the item currently held or use the block the player is looking at. In GUI, pick up the stack of items or place a single item from a stack held by mouse.
drop	Q key	Drop a single item from the stack held by mouse. If the player presses ctrl-Q then it drops the entire stack. In the GUI, the same thing happens except to the item the mouse is hovering over.
hotbar . [1-9]	keys 1 - 9	Switch active item to the one in a given hotbar cell.

Table 3: Binary actions included in the action space. Controls has more detailed descriptions of each action.

In addition to the binary (on/off) keypress actions, our action space also includes mouse movements. As with human gameplay, when in-game GUIs are not open, mouse X and Y actions change the agent’s yaw and pitch, respectively. When a GUI is open, camera actions move the mouse cursor and thus their effect depends on the current position. Inventory interaction in Minecraft requires fine-grained mouse movements to achieve both mining and navigating the world can be achieved with coarser mouse movements such as crafting and smelting, while they move the mouse or camera relative to the current position. To be able to achieve both with the same action space, we implemented mouse movements to achieve tasks such as experiments we found to improve crafting performance.

## D Inverse Dynamics Model Training Details

### D.1 IDM Architecture

The IDM model has approximately 0.5 billion trainable weights. The input to the model consists of consecutive image frames (128 frames of video), each of which has dimensions 128 by 255.0 such that they lie within the range [0, 1]. The first layer of the IDM value function has 128 learnable filters with a temporal kernel width of 5 and spatial kernel width of 3, meaning that embeddings at time index  $t$  are function of times  $t - 2, t - 1, t, t + 1$ , and  $t + 2$ . We found this layer to be extremely fun-



categorical for both the discretized horizontal and vertical movement details on the action space).  
Each dense layer or convolutional layer in the network is preceded by a layer normalization layer. Weights are initialized with Fan-In initialization<sup>71</sup> and biases are initialized with zero.

## D.2 IDM Training

The total loss for the network is the sum of each independent action prediction loss (one for each key and one for both mouse directions). Each independent loss is the negative log-likelihood of the correct action. We use the ADAM<sup>72</sup> optimizer with a learning rate of 0.003, a batch size of 128 (where each item in the batch is a video sequence), and a weight decay of 0.01. Hyperparameters were tuned in preliminary experiments. The IDM is trained on our contractor collected dataset for 20 epochs. This took 4 days on 32 A100 GPUs. We add data augmentation to each video segment; augmentations are randomly sampled once per segment such they are temporally consistent. Using the Pytorch<sup>73</sup> transforms library, we adjust the hue by a random factor between -0.2 and 0.2, saturation between 0.8 and 1.2, brightness between 0.8 and 1.2, and contrast between 0.8 and 0.2. We also randomly rotate the image between -2 and 2 degrees, scale it by a random factor between 0.98 and 1.02, shear it between 0.8 and 1.2, brightness between 0.8 and 1.2, and contrast between -0.2 and 0.2. We also randomly rotate the image between -2 and 2 degrees, translate it between -2 and 2 pixels in both the  $x$  and  $y$  dimensions.

Due to the large computational cost of running all of the experiments (for IDM, BC, and RL training), this non-ideal situation is mitigated because deep learning training tends to be low variance<sup>74,75</sup> and because we often have data points from sweeps (e.g. on dataset size) that suggest overall trends.

## D.3 Generating Pseudo Labels with the IDM

Section 4.1 shows that inverse dynamics modeling is a much easier task than behavioral cloning because IDMs can be non-causal. The IDM is trained to simultaneously predict all 128 actions for each video sequence, so the IDM will effectively be causal for frames at the end of the video clip because future frames are not included in the sequence. For this reason, we apply the IDM over a video using a sliding window with stride 64 frames and only use the pseudo-label prediction over a frames 32 to 96 (the center 64 frames). By doing this, the IDM prediction at the boundary of the video clip is never used except for the first and last frames of a full video.

## E Foundation Model Behavioral Cloning

### E.1 Foundation Model Architecture

The behavioral cloning model architecture is the same as the IDM architecture described in Appendix D.1 except that we modify the architecture so that it is causal (i.e. cannot see the future when making predictions). This means the BC architecture does not have the initial non-causal convolution layers that the IDM has (this layer is omitted completely). Furthermore, the residual transformer layers are now causally masked (as is standard in language modeling) and we do Transformer-XL-style<sup>76</sup> training where frames can attend to keys and values from past batches within the same video. We also use Transformer-XL-style relative attention position embedding.

### E.2 Null Action Filtering

The most common action humans take is the null action (no keypresses or mouse movements). It accounts for 35% of all actions they take. Among other reasons, a player may take the null action to wait for something in the game to finish, to pause between actions, or to take a break. Early on in the project we found that the BC model would take a much higher percentage of null actions, often upwards of 95%. In order to prevent this behavior we have been 1, 3, or 21 frames of consecutive null actions, and include a treatment: we filter any null filtering. Null action filtering generally helps, increasing a treatment's effectiveness.

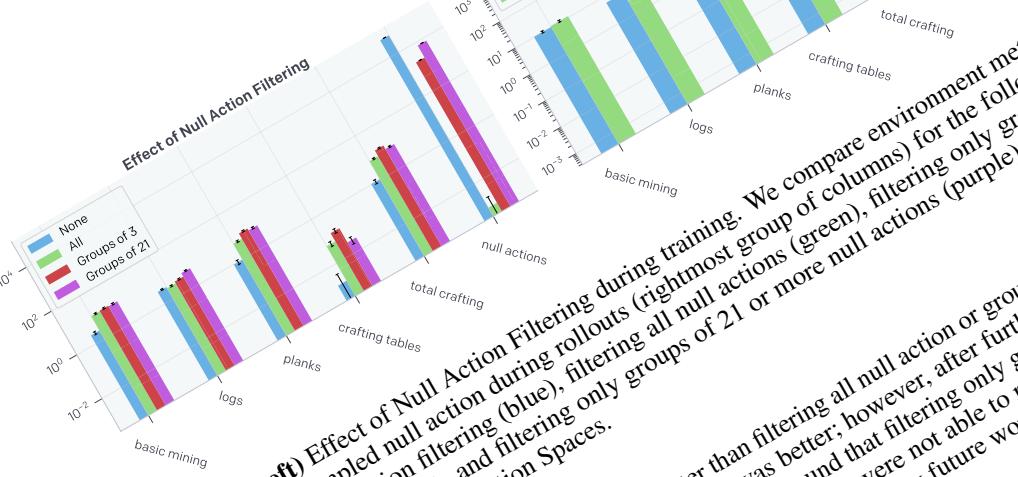


Figure 15: (**Left**) Effect of Null Action Filtering during training. We compare environment metrics and number of sampled null action during rollouts (rightmost group of columns) for the following treatments: no null action filtering (blue), filtering all null actions (green), filtering only groups of 3 or more null actions (red), and filtering only groups of 21 or more null actions (purple). (**Right**) Hierarchical versus Factored Action Spaces.

Filtering only groups of 3 performed slightly better than filtering all null actions was better; however, after further model tuning and after we had already trained our largest models, we found that filtering only groups of 3 or more null actions performed best. Due to compute constraints we were not able to redo all experiments with this setting, but doing so would be a reasonable choice for any future work.

### E.3 Joint Hierarchical Action Space

We originally worked with a factored action space, where each keypress could be independently on or off, and this choice was independent of whether the mouse was being moved. This could cause issues for modeling the human behavior distribution exactly. Say for a given state, humans either cause 50% probability (a) move forward and attack or with 50% probability (b) move left and drop their item. The best a factored distribution can do is to assign 50% probability to each of the 4 constituent actions because it chooses to press each button simultaneously and independently. See Appendix C.2 for details on the entire action space.

For this reason, we implemented a joint distribution over actions; however, the full joint distribution over 20 binary buttons and two mouse movement dimensions discretized into 11 bins each would result in  $2^{20} \times 11^2 \approx 1.2 \times 10^8$  possible combinations. This is far too large for many reasons, e.g. the final layer from the transformer stack with a dimension of 4096 would need to be mapped to each combination resulting in  $4096 \times 1.2 \times 10^8 \approx 5.2 \times 10^{11}$  parameters for this final layer alone. In order to reduce this we noted that many buttons in Minecraft have no effect when simultaneously pressed; for example, if a player tries to move forward and backward at the same time, they remain in place. Below we list the the sets of mutually exclusive actions. Furthermore, the inventory button is exclusive with all other buttons and mouse movement.

<b>Mutually Exclusive Actions</b>
forward, back
left, right
sprint, sneak
hotbar. [1-9]

Even reducing the joint action space to reflect these mutually exclusive combinations still huge action space when combined with the discretized results from  $3^3$  for the 3 sets of 2 mutually exclusive movements, i.e.  $3^3 \times 10^{11}$  taking neither in the set is an option,  $\times 10$  for the 9 hotbar keys or no hotbar key remaining binary 4 keys: use, drop, attack, and jump,  $\times 12$  for mouse movement,  $\times 12$  for camera movement,  $\times 1$  for the inventory button which is mutually exclusive with all other actions +1 for the inventory button which is mutually exclusive with all other actions this action is on, then there is a secondary discrete action head with 121 classes of mouse movements because each discretized mouse direction has 11 bins

to move the mouse. If the hierarchical action is off, then the secondary mouse movement action is masked during training and need not be sampled during evaluations. While this no longer models whether or not to move the mouse in the factored action space since dependencies between camera movement has dimension  $3^3 \times 10 \times 2^4 \times 2 + 1 = 8461$  (the  $11^2$  dimensions for whether or not to move the mouse) with an additional 121-dimension head for the joint camera movements. In the future it would be interesting to implement sequential conditional action spaces to more completely model the joint distribution.

In Figure 15 (right) we compare environment rollout performance between BC models with the hierarchical joint action space and with the factored action space. Environment statistics are comparable; however, we see that the factored action space failing to correctly model the distribution in the dataset because, due to null action filtering, there are 0 null actions in the dataset these models key is not conditioned on other keypresses.

## E.4 Foundation Model Training

The foundation model training is similar to the IDM training, with the exception of labels being IDM-generated pseudo labels. The hyperparameters used for foundation model training are listed in Table 4.

Hyperparameter	Value
Learning rate	0.002147
Weight decay	0.0625
Epochs	30
Batch size	880

Table 4: Hyperparameters for foundation model training

## F Behavioral Cloning Fine-Tuning

Behavior cloning fine-tuning is similar to the foundation model training, except we either use a focused subset of all the videos (`early_game` dataset, described in A.3) with pseudo labels, or contractor data (`contractor_house` dataset, described in B.4) with ground-truth labels. The hyperparameters used for behavior cloning fine-tuning are listed in Table 5. We used 16 A100 GPUs for about 6 hours when fine-tuning on `contractor_house` dataset, and 16 A100 GPUs for about 2 days when fine-tuning on `early_game` dataset.

Hyperparameter	Value
Learning rate	0.000181
Weight decay	0.039428
Epochs	2
Batch size	16

Table 5: Hyperparameters for behavior cloning fine-tuning

## G Reinforcement Learning Fine-Tuning

### G.1 Reinforcement Learning Fine-Tuning Training Details

RL experiments were performed with the phasic policy gradient (PPG) algorithm<sup>77</sup> that increases based on the proximal policy optimization (PPO) algorithm<sup>77</sup> performing additional passes over the collected data to optimize the value function.

$L_{klpt} = \rho KL(\pi_{pt}, \pi_\theta)$

In the fine-tuning experiments, this KL divergence loss replaces the common entropy maximization loss, which is often added to RL experiments to encourage exploration.<sup>79,80</sup> The idea behind entropy maximization is that, when all actions appear to have equal value, such as when the agent has learned about the next reward. Blindly exploring by maximizing its entropy to increase the chance that it discovers new states and action spaces are sufficiently small or the reward is sufficiently dense, but becomes infeasible when the state and action spaces are large and rewards are sparse, which is the case in the diametric task. Instead of blindly exploring through uniform-random actions, we assume that a policy has an action distribution that is much more likely to take sequences of actions, it should mimic the action-distribution of the pretrained policy instead of a uniform-random action distribution. In experiments with a coefficient of 0.01, which has been an effective work.<sup>30</sup> Empirically, we found that a high coefficient  $\rho$  for this KL divergence loss with a coefficient of 0.01, which has been an effective work.<sup>30</sup> Empirically, we found that a high coefficient  $\rho$  for this KL divergence loss while a low coefficient  $\rho$  for the auxiliary value function optimization.

Hyperparameter	Value
Learning rate:	40
Weight decay:	48
Batch size:	128
Batches per iteration:	0.999
Context length:	0.95
Discount factor ( $\gamma$ ):	0.2
GAE $\lambda$ :	5
PPO clip:	2
Max Grad norm:	2
Max Staleness:	0.5
PPG sleep cycles:	0.5
PPG sleep auxiliary value-function coefficient:	1.0
PPG sleep max Sample Reuse:	6
KL divergence coefficient $\rho$ :	0.2
Coefficient $\rho$ decay:	0.9995

For RL experiments. These are the hyperparameters for fine-tuning from the early-game model without self-play, loss being set to 0, the learning rate was set to  $2 \times 10^{-5}$ , and the agent did not require regularization steps that change the function head. The learning rate was set to  $2 \times 10^{-5}$  and the learning rate was set to  $2 \times 10^{-5}$ , as we found that performance needed to be lowered to achieve better performance.

Batch size:	100
Context length:	10
Discount factor ( $\gamma$ ):	0.9995
GAE $\lambda$ :	0.95
PPO clip:	0.2
Max Grad norm:	6
Max Staleness:	1.0
PPG sleep cycles:	0.5
PPG sleep value-function coefficient:	0.5
PPG sleep auxiliary value-function coefficient:	0.5
PPG sleep KL coefficient:	0.5
PPG sleep max Sample Reuse:	0.5
KL divergence coefficient $\rho$ :	0.2
Coefficient $\rho$ decay:	0.9995

Table 6: Hyperparameters for RL experiments. These are the hyperparameters for all treatments with two exceptions. First, when fine-tuning from the early-game model without a KL divergence loss, setting out of a sweep over 5 different learning rates, as we found that performance was substantially lower with the standard learning rate of  $2 \times 10^{-5}$  and the agent did not even learn to collect logs. We suspect that the reason that the learning rate needed to be lowered when fine-tuning without a KL loss is that the KL loss prevents making optimization steps that change the policy too much in a single step, especially in early iterations when the value function has not been optimized yet, and the KL loss thus makes it possible to optimize with a higher learning rate. Second, when running RL from a randomly initialized policy there is no KL divergence loss or KL divergence decay, but instead we use an entropy bonus of 0.01, which reportedly worked well in previous work.<sup>30</sup>

Item	Quantity rewarded	Total quantity rewarded	
		Per item	Total
Log	8	8	8
Planks	20	20	20
Stick	16	16	16
Crafting table	1	1	1
Wooden pickaxe	11	11	11
Cobblestone	1	1	1
Stone pickaxe	5	5	5
Furnace	16	16	16
Coal	3	3	3
Torch	3	3	3
Iron ore	1	1	1
Iron ingot	inf	inf	inf
Iron pickaxe	inf	inf	inf
Diamond	4/3	4/3	4/3
Diamond pickaxe	8/3	8/3	8

Table 7: Reward per item and total quantity rewarded.

at the agent does not over-value items. The agent gets 8 reward for up to 20 planks but only 1 reward for creating cobblestone. The expense of creating cobblestone is 1 reward. The reward for a diamond pickaxe is 3 and the reward for a diamond is 4/3.

Table 7: Reward per item and total quantity rewarded.

**G.2 Reinforcement Learning Fine-Tuning Additional Data**

Figure G.2 shows additional figures that are helpful for understanding the main results presented in this section. First, we show the items-over-time loss for the model without a KL loss (Fig. 16). When fine-tuning the model, we can either fine-tune all the four items that the early-game model learned or only the first item. Second, we show the results of experiments that compare fine-tuning from the house and crafting tables. Third, we show the efficiency difference between fine-tuning with and without a KL loss.

3  
1  
inf  
inf

and total quantity rewarded.

The agent does not over-value items that are supposed to be rewarded for up to 20 planks but only up to 1 crafting table, at the expense of creating a crafting table), we divide the total quantity that the agent gets rewarded for (for the purpose of quantity for diamonds is 3 and the total quantity for diamond is not put a limit on the number of these items being rewarded). For 3 iron ore, which has a base reward of 4 for being in the iron tier rewarded, thus the reward per block of iron ore is  $4/3$ . The quantity is still presented in Table 7.

ence towards a diamond pickaxe is rewarded, the reward function is still even deceptive. The sparsity comes from the fact that it can take thousands reward, even after the agent has acquired all the necessary prerequisites even take more than 10,000 actions to find a diamond after crafting an iron function can be deceptive when the most efficient method for getting one item difficult to get the next item. For example, a good strategy for the agent to quickly is to mine (i.e. spend a few seconds to pick up) its crafting table as soon as pickaxe, such that the agent has immediate access to a crafting table as soon as it through cobblestone. However, the fastest way to get a reward for gathering cobblestone immediately after crafting a wooden pickaxe, while leaving the crafting table behind. the optimal strategy for gathering cobblestone makes it more difficult to learn to craft axe.

ents ran for approximately 6 days (144 hours) on 80 GPUs (for policy optimization) and CPUs (mostly for collecting rollouts from Minecraft). In this time the algorithm episodes consist of 4,000 optimization iterations and collected roughly 1.4 million Minecraft episodes consisting of 2,000 frames each, for a total of 16.8 billion frames.

## G.2 Reinforcement Learning Fine-Tuning Additional Data

Additional figures that are helpful for understanding the main results of the RL fine-tuning are presented in this section. First, we show the items-over-training figure when RL fine-tunes the early-game model without a KL loss (Fig. 16). When training without a KL loss, the model only learns to obtain the four items that the early-game model is capable of getting: logs, planks, sticks, and crafting tables. Second, we present preliminary experiments directly compare RL fine-tuning from the house-building model and RL fine-tuning the early-game model (Fig. 17). These experiments differ from the main experiments shown here, the KL loss coefficient was set to 0.4, the learning rate did not have an increased reward. While RL fine-tuning from the house-

worked better than RL fine-tuning from the early-game model, worked better after 800,000 episodes and showed signs of smelting a crafting table. In contrast to the early-game model can craft zero-shot, which are logs, planks, model thus does not initially see any reward, are catastrophically forgotten while the first four items are learned.

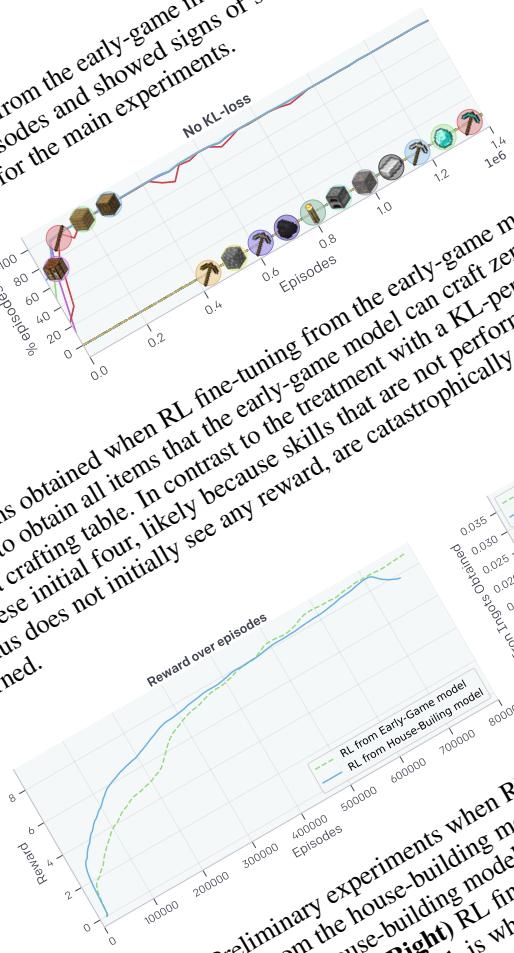


Figure 16: Items obtained when RL fine-tuning from the early-game model learns to obtain all items that the early-game model can craft beyond these initial four, likely because skills that are not performed zero-shot, it does not learn any items are learned.

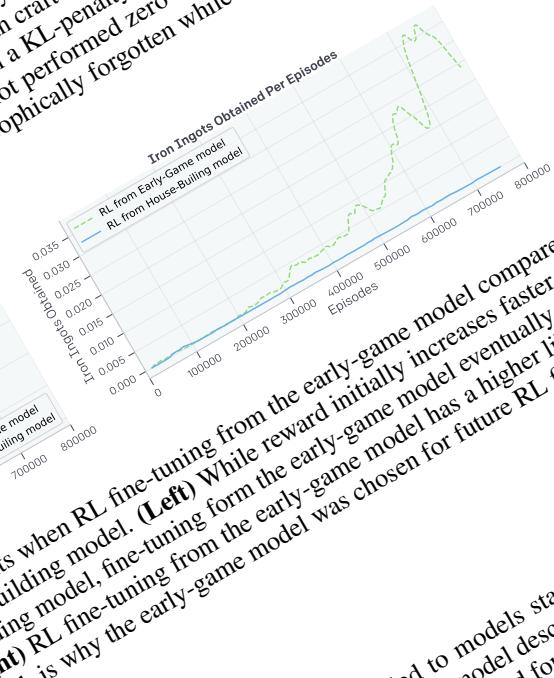


Figure 17: Preliminary experiments when RL fine-tuning from the house-building model. **(Left)** While reward initially increases faster when RL fine-tuning from the house-building model, fine-tuning from the early-game model eventually obtains a slightly higher reward. **(Right)** RL fine-tuning from the early-game model has a higher likelihood of smelting an iron-ingot, which is why the early-game model was chosen for future RL fine-tuning experiments.

## H Foundation Model Scaling

In early experiments we found that increasing model size led to models staying in the efficient learning regime longer into training.<sup>63</sup> Here we compare the 0.5B model described in Section 4.2 to both a 248M and 71M parameter model. Both of these models are trained for 15 epochs as compared to the 30 epochs the 0.5B model trained for. These models have the same architecture as the 0.5B parameter model but each layer in the 248M parameter model has 1/2 the width and each layer in the 0.5B batch size of 480, and weight decay of 0.044506. The 248M model was trained with an initial learning rate of 0.001831, batch size of 640, and weight decay of 0.051376.

In Figure 18 we show validation loss on web\_clean with IDM pseudo-labels, loss on the validation environment used to train the IDM with ground truth labels collected during contractor play, and validation loss on web\_clean for the 71M, 248M, and 0.5B models. While larger models also appear to be better at crafting than the 0.5B, and also has lower contractor dataset loss (Fig. 18 bottom left) and in environment performance tell

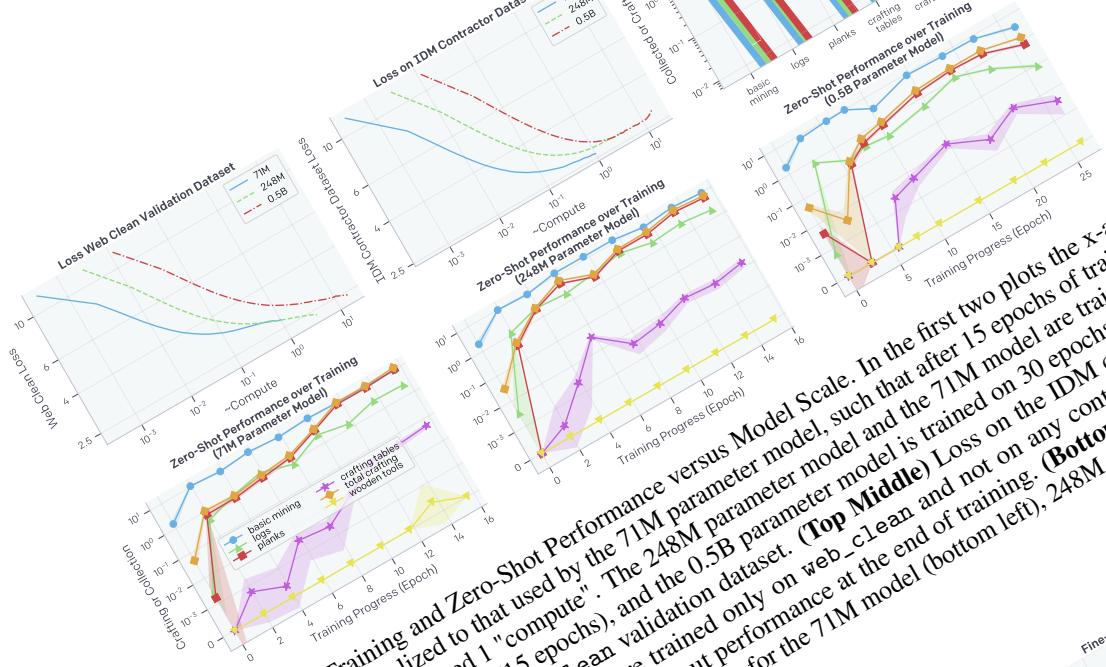


Figure 18: Training and Zero-Shot Performance versus Model Scale. In the first two plots the x-axis is compute normalized to that used by the 71M parameter model, such that after 15 epochs of training the 71M model has used 1 "compute". The 248M parameter model, and the 24M parameter model, and the 0.5B parameter model are trained on 30 epochs of data. (**Top Left**) Loss on the web\_clean validation dataset; note that these models were trained only on web\_clean and not on any contractor data. (**Top Right**) Zero-shot environment rollout performance over training for the 71M model (bottom left), 248M model (middle), and 0.5B model (bottom right).



Figure 19: contractor\_house fine-tuning performance versus model size. (**Left**) Loss on the contractor\_house validation set. (**Middle**) Loss on the full contractor dataset collected from contractor\_house validation set. (**Right**) Environment rollout performance at the end of fine-tuning.

followed by the 248M model and then the 71M model. Environment model rollouts are performed using the same game engine that we use to collect contractor data, which could be visually performed from videos taken from the web. It is plausible that the larger models overfocus on the peculiarities in web data during pretraining since they have worse contractor data loss on the middle), and this causes them to perform more poorly in the environment prior (as indicated by the validation loss in Fig. 18 top left) can quickly shift their low level features to better environment coming from our game engine, resulting in better environment rollout performance. After just a few training, which has no overlap with contractor\_house. While not conclusive, we believe this contractor\_house, all models quickly improve in loss on the contractor\_house. After just a few training, all models now performing best. While not conclusive, we believe this some intuition for future studies of model scaling for sequential decision making.

## I Text Conditioning

Goal-conditioned policies<sup>81,82</sup> make it possible for a single agent to perform goals in a single environment, which is particularly relevant in open-ended environments like Minecraft. In recent work, goal specification has increasingly taken the form of domain languages<sup>83</sup>, or even natural language<sup>84,85</sup>. The benefits of language-conditioned agents are tremendous, especially natural-language<sup>84,85</sup>. Text conditional agents, as their goal space contains a wide variety of potentially very complex tasks. Text conditional models have shown an amazing ability to perform tasks zero-shot (or learn them few-shot) including generalizing in impressive ways via the compositional and combinatorial possibilities allowed by natural language (e.g. GPT<sup>1</sup> and DALL-E<sup>2</sup><sup>86</sup>). We hypothesize that we should expect similar capabilities to emerge with natural-language-conditioned virtual agents, if they are similarly trained on enormous amounts of data (that goes from a natural language description to a sequence of actions that completes the specified goal). In this section we take preliminary steps toward that future. Our preliminary experiments provide evidence that it is possible to pretrain a natural-language-conditioned model for Minecraft using the general approach presented in this paper (VPT) plus conditioning on the speech that often accompanies videos.

In online videos, the human actor sometimes indicates their intent in their verbal commentary (e.g. “Let’s go chop some trees to make a wooden axe” or “now let’s learn how to crop photos in Photoshop”). Conditioning on this closed caption data could produce a steerable pre-trained model: i.e., it may later be possible to condition the model with text such as “I am going to craft a wooden pickaxe” or “I am going to build a house,” and have the agent perform those tasks specifically rather than simply follow typical human behavior (as was investigated in the rest of this paper). An alternate way to produce a steerable agent is via RL fine-tuning, which we could have done in prior work<sup>30</sup>. However, by adding a bit indicating the task to be completed, as has been done in prior work in Section 4.4 conditioning on natural language offers many benefits over that approach. First, it is flexible and powerful, being able to express any task. This would allow for general, capable, zero-shot agents like GPT, but be completed ahead of time. This would allow for embodied tasks such as completing tasks on computers or in simulated 3D worlds. Third, text conditioning can be used even when tasks are difficult to specify via reward functions (e.g. “Let’s build a castle surrounded by a moat”). In the limit, VPT+text could conceivably produce worlds, complete tasks on computers, and in other similar embodied sequential decision domains. “I will now build a castle zero or few shot, but in the form of agents that can act in virtual instructions, and complete tasks on computers, and in other similar embodied sequential decision domains. We do not reach those lofty goals in this work, but we began a first step towards exploring in that direction.

Many Minecraft videos feature audio commentary from the player, or could be extracted post-hoc using automated speech recognition (ASR).<sup>87</sup> Our dataset features about 17k hours of content with associated closed captions.

We fine-tuned the 220 million parameter VPT foundation model used in the RL-fine-tuning experiments (chosen vs. 0.5B for the same reason: to reduce compute costs) with an additional text-conditioning input on the subset of our data for which closed captions are available. To obtain a frame in a given chunk, we first split videos into 30 second chunks. The same text is associated with the line of text preceding and is made up of all the closed captions occurring within that chunk (95%) of our closed caption data following the chunk (if any). Because the vast majority of the library<sup>88</sup> is processed by a randomly initialized multi-layer perceptron (MLP) with two hidden size 2,048. The resulting activations are added for each frame to the pretrained mlp(textEncoder) before the transformer layers (pretransformerActivations += mlp(textEncoder)). Our model shows evidence of steerability. When conditioned on sentences like “I’m going to explore” and “I’m going to find water” the agent explores farther from its spawn point (Figure 20a). Additionally, we can steer the agent

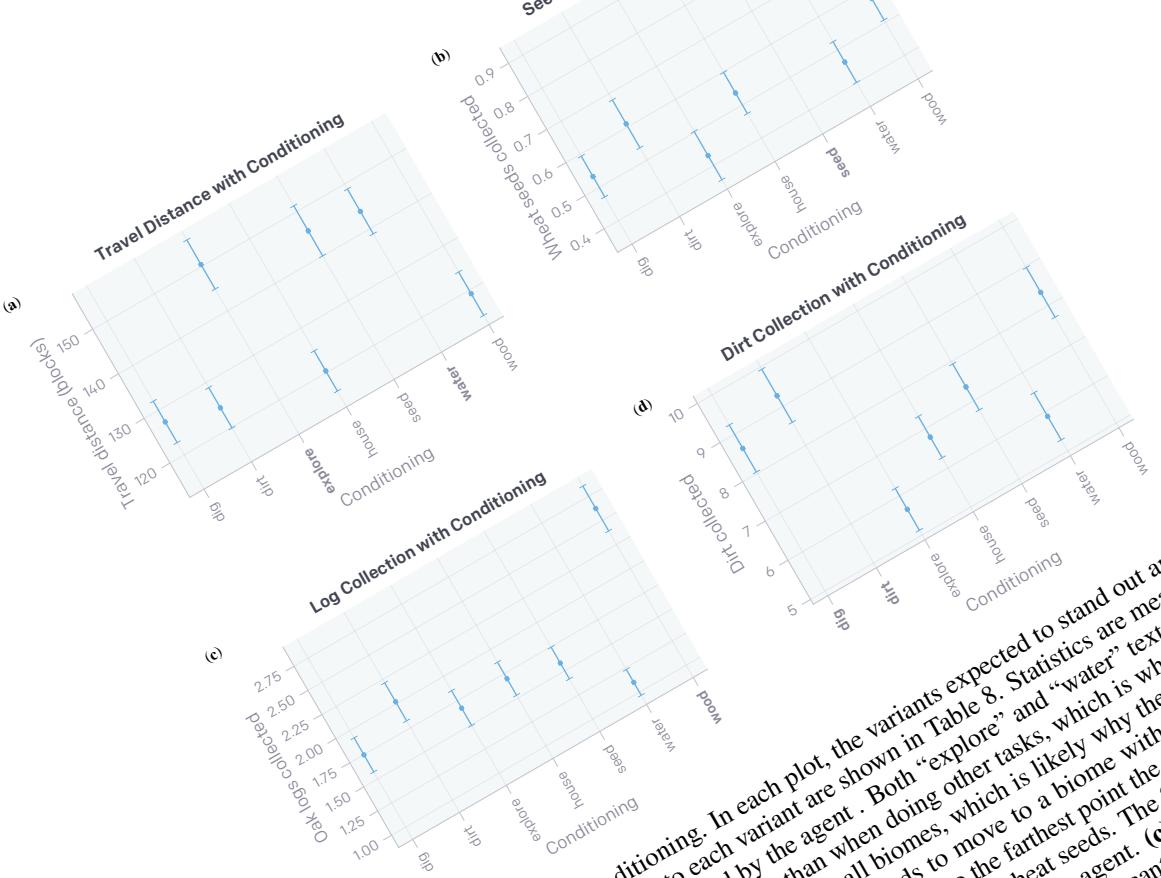


Figure 20: Evidence for conditioning. In each plot, the variants expected to stand out are shown in bold. The strings corresponding to each variant are shown in Table 8. Statistics are measured over 5 minute episodes. (a) Distance traveled by the agent is not present in all biomes, which is what occurs. (b) Both “explore” and “water” text strings are measured over 5 minute episodes. (c) Collection of wheat seeds (as the agent sometimes needs to move to a biome with grass) produces more travel (as the agent reaches the farthest point of the episode on the Euclidean distance from the spawn point to the spawn point of the “seed” condition). (d) Collection of dirt (as expected of a steerable agent). The “seed” variant reaches during the episode on the horizontal (x-z) plane substantially more than other variants. (e) Collection of oak logs (as expected of a steerable agent). The “seed” variant collects more oak logs, as is to be expected of a steerable agent (we speculate that the “wood” variant collects less because there are no trees in water). The “seed” variant collects significantly more oak logs, as is to be expected of a steerable agent. The “seed” variant collects less because there are no trees in water. The “dig” and “dirt” variants collect dirt. It is easy to mistakenly aim at the ground rather than at grass or trees when collecting seeds or wood, which likely explains the slightly higher amount of dirt collected by these variants. In all cases, the error bars are confidence intervals of the mean, over 1,000 episodes per conditioning variant. The bars in each bar plot do not overlap are statistically significantly different at a  $p < 0.05$  level.

Variant name	String
dig	I'm going to dig as far as possible
dirt	I'm going to collect dirt
explore	I'm going to explore
house	I'm going to make a house
seed	I'm going to collect seeds
water	I'm going to find water
wood	I'm going to chop wood

Table 8: Strings corresponding to each conditioning variant

early game items such as seeds, wood, and dirt by conditioning seeds/chop wood/collect dirt” (Figure 20b,c,d). While our results show some level of steerability, more work is required to successfully steer agents to gather flowers or to hunt, both in the early game, but less common (and, in the case of hunting animals, much more gathering window and various resources, and conditioned to craft a given item is present). Likewise, an experiment in which the agent is presented with a wooden axe’ failed to show that the agent was more influenced by the prior, unconditional probability of what human players would craft next given the resources available, which is not too surprising since in Minecraft, especially in the early game, there is a relatively consistent path to gathering resources in a specific order go produce more powerful tools (Fig. 6). For example, if the agent had often would make the stone pickaxe and we asked it instead to make a (weaker) wooden pickaxe, it convinces us that the “house” conditioning causes the agents to take more steps towards building a house than other variants.

Thus, our results show that it is possible to train a somewhat steerable natural-language-conditioned agent. However, its steerability is still too weak to be practically useful, and it is far from what we believe could be accomplished with more research, data, and training compute. Another exciting research direction is to have the model predict future text as well as just the next action.