

Video PreTraining (VPT): Learning by Watching Unlabeled Online Video

Bowen Baker*†
bowen@openai.com

Ilge Akkaya*†
ilge@openai.com

Peter Zhokhov*†
peterz@openai.com

Jie Tang*†
jietang@openai.com

Adrien Ecoffet*†
adrien@openai.com

Brandon Houghton*†
brandon@openai.com

Jeff Clune*‡
jclune@gmail.com

Abstract

Pretraining on noisy, internet-scale datasets has been heavily studied for training models with broad, general capabilities for text, image, and video modalities.^{1–6} However, for many sequential decision domains such as video games, and computer use, publicly available data does not contain the behavioral priors required to train behavioral priors in the same way. We extend the pretraining paradigm to sequential decision domains through sequential imitation learning wherein agents learn to act by watching online unlabeled video. Specifically, we show that with a small amount of labeled data via a learned inverse dynamics model accurate enough to label a huge unlabeled dataset – here, online videos of people playing Minecraft – from which we can train a general behavioral prior. Despite using the native human input (mouse and keyboard at 20Hz), we show that this behavioral prior has near-human shot capabilities and that it can be fine-tuned, with both imitation and reinforcement learning, to hard-exploration tasks that are impossible to learn from scratch via reinforcement learning. For many tasks our models exceed human-level performance, and we are the first to report computer agents using diamond tools, which can take proficient humans upwards of 20 minutes to complete (e.g., environment actions) of gameplay to accomplish.

1 Introduction

Work in recent years has demonstrated the efficacy of pretraining large language models⁷ on noisy internet-scale datasets for use in downstream tasks in natural language processing, computer vision,^{5,6,8} and robotics.⁹ For sequential decision domains (e.g. robotics, game playing, and computer use) where agents must repeatedly act within an environment, a wealth of unlabeled data is available on the web; however, most of this data is in the form of *unlabeled* video (i.e. without labels indicating what action was taken at each frame), making it much less straightforward to train a behavioral prior than it is in e.g. natural language. In a few rare settings, such as Chess, Go,

^{*}This was a large effort by a dedicated team. Each author made huge contributions over several time periods. All members were full time on the project for over six months. BB, IL, and PZ were original VPT project team and were thus involved for even longer (over a year). Aside from these members, author order is random. It was also randomized between IA and PZ.

[†]OpenAI

[‡]University of British Columbia

already exist large datasets with action labels from various online platforms used for imitation learning.^{9,10} When large labeled datasets do not exist, for training capable agents is reinforcement learning (RL),¹¹ which can be expensive for hard-exploration problems.^{12–18} Many virtual tasks, e.g., navigation, Photoshop, booking flights, etc., can be very hard to learn with RL and do not have available sources of labeled data.^{19,20} In this paper, we seek to extend the large, general-purpose foundation models to sequential decision domains by internet-scale unlabeled video datasets with a simple semi-supervised imitation call this method Video PreTraining (VPT) and demonstrate its efficacy in the

Existing semi-supervised imitation learning methods aim to learn with few or no labeled data, however, they generally rely on the policy's ability to explore the environment, making them susceptible to exploration bottlenecks.^{21–25} Furthermore, most existing imitation learning work was tested in the relatively low data regime; because we can collect more data ($\sim 70k$ hours of unlabeled video), we hypothesize that we can achieve better performance with a much simpler method, a trend that has proven true for pretraining inverse dynamics models as text.¹ In particular, given a large but unlabeled dataset, we propose generating pseudo-labels by gathering a small amount of labeled data to train an inverse dynamics model to predict the action taken at each timestep in a video. Behavioral cloning (BC) can also benefit from a large amount of data because the model must learn to infer intent and the distribution over actions from only past observations. In contrast, the inverse dynamics modeling task is *non-causal*, meaning it can look at both past and future frames to infer actions. The environment mechanics are far simpler than the breadth of human behavior that exists in the environment, suggesting that non-causal IDM could require far less data than causal models. Using pseudo-labels generated from the IDM, we then train a model to predict the distribution of behavior in the previously unlabeled dataset with standard behavioral cloning. This model does not require any model rollouts and thus does not suffer from any potential issues with the environment. Finally, we show we can fine-tune this model to downstream tasks such as behavioral cloning or reinforcement learning.

We chose to test our method in Minecraft because (a) it is one of the most actively played games in the world²⁶ and thus has a wealth of commonly available video data online, (b) 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 correspondingly difficult exploration challenges.^{27–31} In this work we use the native human interface for Minecraft so that we can (1) most accurately model the human behavior distribution and reduce domain shift between video data and the environment, (2) make data collection easier by allowing our agent to play the game without modification, and (3) eliminate the need to hand-engineer a complex interface for models to interact with the environment. This choice means that our model can act up to 10 times faster than prior work and must use a mouse and keyboard interface to interact with the environment, including smelting, trading, etc., including dragging items to specific slots or navigating menus with the mouse cursor (Fig. 1). Compared to prior work in Minecraft that uses a mouse and keyboard to construct crafting and attacking macros,^{30,32–34} using the native human interface allows us to interact with the environment's exploration difficulty, making most simple tasks near impossible to learn from scratch. Even the simple task of gathering a single wooden log while already having a sword equipped and performing consecutive attack actions with the human interface, meaning the chance for success is $\frac{1}{2}^{60}$. While this paper shows results in Minecraft only, the VPT framework could be applied to any domain.

In Section 4 we show that the VPT foundation model has nontrivial zero-shot learning capabilities for tasks impossible to learn with RL alone, such as crafting planks and planks into wooden logs, requiring a human proficient in Minecraft a median of 50 seconds or ~ 90 attempts. Through fine-tuning with behavioral cloning to smaller datasets that target specific distributions, our agent is able to push even further into the technology space.



Figure 1: Crafting with a mouse and keyboard interface.

(taking a human a median of 2.3 minutes or ~ 2790 actions). Finally, fine-tuning on the most dramatic improvements: our agent is able to craft diamond tools, which in Minecraft made even more challenging by using the native human interface. Compared to a proficient human a median upwards of 20 minutes or ~ 24000 actions. The contributions 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 tasks, (2) we show that such pretraining plus fine-tuning enables agents to solve tasks that were previously impossible to learn, (3) we show that labeled contractor data is far more effective than unlabeled contractor data for learning VPT methods, and (4) we demonstrate how to source our contractor data, trained model weights, and Minecraft environment into learning to act via semi-supervised imitation learning at scale.

2 Preliminaries and Related Work

Imitation learning methods^{35–38} seek to construct a policy that accurately matches the behavior in some dataset $D = \{(o_i, a_i)\}$, $i \in \{1\dots N\}$ of action-observation pairs. To learn these policies in an environment, they must be *causal*, meaning they can only depend on the current timestep t and past timesteps only, i.e., $\pi \sim p(a_t|o_1\dots o_{t-1})$. This is the simplest when demonstrations are labeled with corresponding actions. Imitation learning has seen success in aerial vehicles,^{39,40} self-driving cars,^{41,42} board games,⁹ and robotics.

When labeled demonstrations are not available, standard behavioral cloning methods⁴³ assume there is a large body of work in imitating behavior from unlabeled demonstrations. GAIL²³ constructs an adversarial objective incentivizing the trained policy to be indistinguishable from those in the target dataset. Edwards et al.⁴⁵ propose a similar approach using unlabeled demonstrations and then map the learned latent actions to a small amount of environment interaction. Peng et al.⁴⁶ first use motion-capture data to extract agent positions in videos and then train RL agents to match these waypoints. Pathak et al.⁴⁷ and Aytar et al.⁴⁸ task a RL agent to match waypoints; however, they represent waypoints as embeddings from unsupervised feature learning models. Pathak et al.⁴⁷ also propose goal conditioned policies to take actions that advance the current state toward the goal states expressed as high dimensional visual waypoints. Most similar to our own work, we simultaneously train (1) an inverse dynamics model (IDM),⁵¹ which aims to predict the next action between timesteps given observations of past and future timesteps, e.g., $a_{t+1} \sim p(a_{t+1}|o_t, o_{t+1}, \dots)$, and (2) a behavioral cloning (BC) model on trajectories of observations labeled by contractors. We train the IDM is collected by rolling out the BC model in the target environment. The BC model improves in tandem. However, at any point in training, if there are sequences of observations where the IDM performs poorly on, it requires that the BC model perform those sequences in order for the IDM to improve and correctly label them. Therefore, if the BC model does not learn efficiently, it could severely slow down learning. In order to avoid this potential bottleneck, we propose a simpler two-stage approach: we first train an IDM on a small number of labeled demonstrations from human contractors (they play the game as would normally as we record their mouse movements). Because human contractors reach most relevant parts of the game quickly, we hold the IDM fixed throughout BC training.

Compared to most previous work in semi-supervised imitation learning, we focus on a more complex and open-ended environment of Minecraft. Minecraft is a massive, open-world game that, due to its popularity and wide variety of mechanics, has attracted a large body of research.^{27,28,30–34,52–66} A large body of work focuses on small, custom-designed environments with tasks such as navigation,^{53,60} block placing,^{54,55} instruction following,^{56,57} and others.^{28,31,57} Work operating in the massive, randomly generated environments of Minecraft has included hill climbing,⁵⁸ automated curriculum learning³⁰ and, most recently, experiments presented in Sec. 4.4, diamond mining.^{27,32–34} However, to the best of our knowledge, there is no published work that operates in the full, unmodified human action space of Minecraft, specifically drag-and-drop inventory management and item crafting.

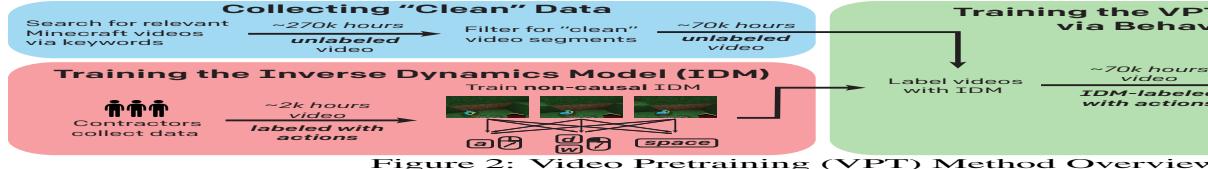


Figure 2: Video Pretraining (VPT) Method Overview

3 Methods

Inverse Dynamics Models (IDM) VPT, illustrated in Figure 2, requires a large amount of labeled contractor data with which to train an inverse dynamics model. The model, which seeks to minimize the negative log-likelihood of an action at timestep t given observations $\omega_t : t \in [1..T]$, is in contrast to an imitation learning policy, where the action a_t is a function of past observations ω_t , meaning its prediction for a_t can be a function of both past and future events. We hypothesize that the behavioral cloning objective of modeling the distribution of human intention is easier than inverting environment dynamics. In Sec. 4.1 we will show that the IDM objective is much easier to learn, and further that with very little labeled data (as few as 100 hours) we can train a fairly accurate model that can be used to label online videos, providing the large amount of data required for behavioral cloning. See appendices D and B for IDM training and data collection.

Data Filtering We gather a large dataset of Minecraft videos by searching for relevant keywords (Appendix A). Online videos often (1) include overlaid artifacts of the player’s face, channel logos, watermarks, etc., (2) are collected from a computer with different gameplay, or (3) are from different game modes. We only want “survival mode” where players start from scratch and must gather resources. We call data “clean” if it does not contain visual artifacts and is from survival mode. We call data “unclean.” With enough data, a large enough model, and enough training time, a model trained on both unclean and clean videos would likely still perform well in the 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 train a model to filter out unclean segments using a small dataset (8800) of online videos labeled by contractors as clean or unclean (Appendix A.2).

VPT Foundation Model We train a foundation model with standard behavior cloning loss, minimizing the negative log-likelihood of actions predicted by the IDM on clean trajectories of length T . We minimize

$$\min_{\theta} \sum_{t \in [1..T]} -\log \pi_{\theta}(a_t | \omega_1, \dots, \omega_t), \text{ where } a_t \sim p_{\text{IDM}}(a_t | \omega_1, \dots, \omega_t)$$

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

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 model 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 small amount of labeled data. While more data improves both mouse movement and keypress prediction accuracy, we find that the model can be trained with as few as 100 hours of labeled data.

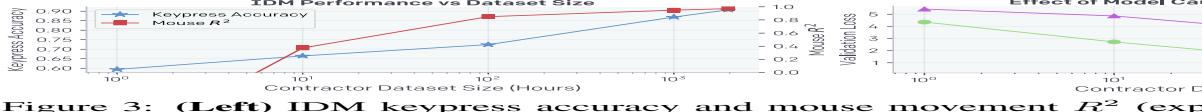


Figure 3: **(Left)** IDM keypress accuracy and mouse movement R^2 (explained variance) as a function of dataset size. **(Right)** IDM vs. behavioral cloning data efficiency.

IDM trains on only 1962 hours of data (compared to the ~ 70 k hours of clean data from the internet) and achieves 90.6% keypress accuracy and a 0.97 R^2 for mouse movement on a held-out validation set of contractor-labeled data (Figure 3 left).

Figure 3 (right) validates our hypothesis that IDMs are far more data efficient than behavioral cloning because inverting environment mechanics is far easier than modeling the entire environment. The IDM is two orders of magnitude more data efficient than a Behavior Cloning model trained on the same data and improves more quickly with more data. This evidence supports our hypothesis that it is more effective to use contractor data within the VPT pipeline by training an IDM foundation model from contractor data directly (Sections 4.5 and 4.6 provide empirical support for this claim).

4.2 VPT Foundation Model Training and Zero-Shot Performance

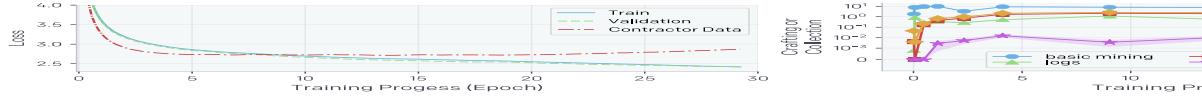


Figure 4: **(Left)** Training and validation loss on the `web_clean` internet dataset, and loss on the main IDM contractor dataset, which has ground-truth labels, and loss on the main IDM contractor dataset, which has ground-truth labels, and loss on the main IDM contractor dataset, which has ground-truth labels. **(Right)** Amount a given item was collected per episode over time. We show the mean and standard error of the mean for various game states and collection rates (Fig. 4, right). The VPT foundation model quickly learns to collect logs, a task we found near impossible for an RL agent to achieve in zero-shot (Sec. 4.4). It also learns to craft those logs into wooden planks and stone.

We now explore the emergent behavior learned by a behavioral cloning policy. We collected a large, but noisy, internet dataset labeled with our IDM. To collect the unlabeled data, we searched for publicly available videos of Minecraft play with search terms such as “survival for beginners.” These searches resulted in ~ 270 k hours of video, which we used to extract “clean” video segments yielding an *unlabeled* dataset of ~ 70 k hours, which we used to train the IDM (Appendix A has further details on data scraping and filtering). We then generated contractor labels for `web_clean` with our best IDM (Section 3) and then trained the VPT foundation model using behavioral cloning. Preliminary experiments suggested that our model could learn to stay in the survival mode regime⁶³ for that training duration (Appendix H), which took ~ 9 days on 7700 parameter model. We evaluate our models by measuring validation loss (Fig. 4, left) and zero-shot performance in the Minecraft environment. Unless otherwise noted, in all environment evaluations, we run the model in standard survival mode game where they play for 60 minutes, i.e. 72000 commands. We plot the mean and shade the standard error of the mean for various game states and collection rates (Fig. 4, right). The VPT foundation model quickly learns to collect logs, a task we found near impossible for an RL agent to achieve in zero-shot (Sec. 4.4). It also learns to craft those logs into wooden planks and stone.

to craft a crafting table, which are required to unlock most other technologies. A human proficient in Minecraft approximately 50 seconds (970 consecutive actions) to craft these behaviors are fairly complex in the native human action space, the VPT model collects these items 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 a non-negligible amount of wooden sticks, which are required to make wooden tools, and collects flowers and crafts dyes from them; kills zombies that appear during the night; collects various berries and mushrooms and eats them; and finds game-generated items such as pillars to collect various rare items from chests. The model also learned to navigate through the game and pillar jump, which involves the agent repeatedly jumping and quickly placing wooden sticks such that it climbs upward by making a pillar.^(iv)

While training and validation loss decrease healthily over training (Fig. 4, left), the validation loss for the contractor_house dataset (which the VPT model does not train on) begins increasing after 7 epochs. This could be out-of-distribution because our contractors may have a different visual domain than the contractor_house dataset. Because there is some impactful visual domain shift compared to videos from the web_clean dataset, we could have expected this would be predictive of declining evaluation performance. We provide notable game statistics from the VPT foundation model rollouts (Figure 5) and in the next section we show that BC fine-tuning performance compares favorably to VPT foundation model trains. We provide more insight into this curious phenomenon in the next section.

4.3 Fine-Tuning with Behavioral Cloning

Foundation models are designed to have a broad behavior profile and be generalizable to a wide variety of tasks. To incorporate new knowledge or allow them to specialize in a specific task or distribution, it is common practice to fine-tune these models to smaller, more specialized datasets. In this work, we fine-tuned the VPT foundation model trained on the broad web_clean dataset had nontrivial difficulty learning to craft a crafting table yet unable to go past this in the technical report. To study into BC fine-tuning, we attempt to improve the VPT foundation model’s performance and craft these “early game” items by fine-tuning to two narrower datasets. The contractor_house dataset contains contractor-related behaviors within the first few minutes of players starting in a fresh world. Contractors can build a basic house from logs, planks, and dirt. Collecting contractor data can be difficult and expensive, so we collected the contractor_house dataset earlygame_keyword by searching for videos online with descriptive terms such as “new world”, “let’s play episode 1”, etc.; this is a subset of web_clean. See Appendix B.4 and A.3 for full descriptions of both datasets.



Figure 5: (Left) Collection and crafting rates for three policies: the Foundation model, and the VPT foundation model BC fine-tuned to the earlygame_keyword and contractor_house datasets. BC fine-tuning to either dataset improves performance compared to the contractor_house dataset) yielding wooden and stone tools. Proficiency in the contractor_house dataset (1390 actions) to construct wooden tools and 2.3 minutes to construct stone tools. (Right) Collection and crafting rates for VPT foundation model rollouts throughout training after they are BC fine-tuned to the contractor_house dataset. Collection and crafting-related behaviors increase throughout foundation model training. The contractor_house dataset contains contractor-related behaviors within the first few minutes of players starting in a fresh world. Contractors can build a basic house from logs, planks, and dirt.

^(iv) Sample videos: https://www.youtube.com/playlist?list=PLNAO1b_agjf3U3rSv

Fine-tuning to `earlygame_keyword` results in a large boost compared to the model: 2.5x more crafting tables, 6.1x more planks, 4.3x more logs, and 5.1x more cobblestones (Fig. 5). However, when fine-tuning to this dataset we did not see any significant improvement over the contractor_house dataset: 213x more crafting tables, 59x more wooden pickaxes, and 59x more crafting over all. In addition, we saw the emergence of crafting sequences that require placing a crafting table on the ground, opening it to reveal a new crafting table, and using it to craft wooden tools. This entire sequence takes a proficient human player a median of 2.3 minutes (1390 consecutive actions) to accomplish. The model goes further and crafts stone tools, requiring a wooden pickaxe to mine, and crafts stone tools, requiring a stone pickaxe to mine. This takes a proficient human player a median of 2.3 minutes (2790 consecutive actions). We also saw this model more frequently raiding villages that randomly spawn animals for food, in addition to many behaviors we saw performed by the foundation model.

Despite the foundation model's zero-shot rollout performance plateauing (Fig. 5, right), fine-tuning performance *does* continue to increase throughout fine-tuning (Fig. 5, left). Additionally, there is a stark difference in performance when fine-tuning from the VPT foundation model (Fig. 5 right, comparing the leftmost bar to the rightmost bar).

4.4 Fine-Tuning with Reinforcement Learning



Figure 6: Typical sequence of items for obtaining a diamond pickaxe. Below each item is the number of actions required to obtain it, the median time to obtain it, and the percentage of contractors that got the item within 10 minutes. The median time to obtain a diamond pickaxe is 2.3 minutes (2790 actions), while the median time to obtain an iron ingot is > 20 minutes because contractors obtained this item in less than 50% of cases.

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 requires a sequence of difficult-to-obtain items that require complex skills like mining at depths, crafting with and without a crafting table, tool use, operating a furnace, and smelting items. These items can be easily lost by dropping items, destroying items, or dying. Obtaining a diamond pickaxe often takes a proficient human over 20 minutes (24,000 actions).

Agents are rewarded for each item obtained in the sequence, with lower rewards for items that are easier to obtain and higher rewards for items near the end of the sequence. We used the phasic policy gradient⁶⁴ RL algorithm for ~1.3 million episodes (round-trip) with the phasic policy gradient⁶⁴ RL algorithm for ~1.3 million episodes (round-trip). Episodes last for 10 minutes. See Appendix G.1 for reward function and RL details. Due to computational constraints, RL experiments use a ~248 million parameter VPT model.

A major problem when fine-tuning with RL is catastrophic forgetting^{65,66} because skills can be lost before their value is realized. For instance, while our VPT model exhibits the entire sequence of behaviors required to smelt iron zero-shot, it forgets how to smelt iron after some training. Players smelting with furnaces. It therefore may have some latent ability to smelt iron, but the prerequisites to do so have been forgotten. To combat the catastrophic forgetting, we use a Kullback-Leibler (KL) divergence loss between the RL model and the frozen VPT model.

Training from a randomly initialized policy fails to achieve almost *any* reward. The hard exploration challenge the diamond pickaxe task is for RL in the natural language setting (Fig. 7a). The model never learns to reliably collect logs, typically the first item required to obtain a diamond pickaxe (Fig. 7b). RL fine-tuning from the VPT foundation model achieves significantly better (Fig. 7a), learning everything up to mining iron ore and crafting furnaces. However, this agent fails at smelting an iron ingot, the next item required to get further.

⁶⁵Sample Videos: https://www.youtube.com/playlist?list=PLNAO1b_agjf2yDSSs4

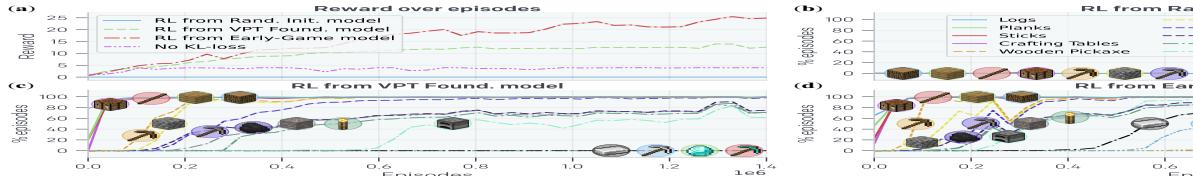


Figure 7: RL Fine-tuning results. (a) RL from a randomly initialized model fails to learn any reward, RL fine-tuning from the VPT foundation model performs suboptimally, RL fine-tuning from the early-game model performs surprisingly well. When training the early-game model without a KL loss to the original policy, it stalls after 100,000 episodes, suggesting that the skills necessary to make further progress have catastrophically forgotten. (b) RL from a randomly initialized model occasionally collects breaking leaves (an easy but inefficient method of getting sticks that does not require mining) 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 smelt iron ingots. (d) RL fine-tuning from the early-game model learns to obtain every item in the sequence towards a diamond pickaxe and crafts a diamond pickaxe.

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

Results further improve by first BC fine-tuning the VPT Foundation Model on the `earlygame_keyword` dataset (the *early-game model*, Sec. 4.3) and then RL fine-tuning on the contractor house dataset (Fig. 7a), which in preliminary experiments we found to perform better than the contractor house dataset (pretraining, BC fine-tuning, and then RL fine-tuning) succeeds in learning to craft diamonds. It achieves over 80% reliability on iron pickaxes, almost 20% reliability on wooden pickaxes, and 2.5% reliability on obtaining a diamond pickaxe (Fig. 7d). For comparison, humans achieve the objective of obtaining a diamond pickaxe by collecting these items in 57%, 15%, and 1% respectively, meaning our model is human-level for crafting iron pickaxes and wooden pickaxes. Others have managed to obtain diamonds with $\sim 0.1\%$ reliability in 15 minutes using a simplified action space designed to ease exploration. To the best of our knowledge, no one has reported non-zero success rates on crafting a diamond pickaxe. Qualitatively, our model can learn useful skills for diamond mining, such as efficient mining patterns, cave exploration, previously placed objects like crafting tables, and advanced techniques like using sticks as fuel when moving on to iron tools.^(vi)

Finally, we validated the importance of the KL loss to the pretrained model. The treatment without a KL loss obtains only items early in the sequence (e.g., logs, planks, sticks, crafting tables) limiting its reward (Fig. 7a). This failure to progress further is likely because, while the initial skills of chopping logs and crafting planks are learned, subsequent skills like crafting a wooden pickaxe 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 easier to scale a foundation model from a small dataset of labeled contractor data to train an IDM within the VPT method than it is to scale a foundation model from that same small contractor dataset. If we could collect a contractor dataset of a similar order of magnitude as `web_clean`, then this would be a significant win; however, collecting that scale of data would have cost millions of dollars. Foundation models trained on increasing orders of magnitude of data from 1 to 100k hours of `web_clean` dataset. Foundation models trained up to and including 1k hours of

^(vi) Videos found at https://www.youtube.com/playlist?list=PLNAO1b_agjf3e_UKv



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

contractor data, and those trained on 5k hours and above are trained on subsets of the entire `web_clean` dataset. The `contractor_house` dataset does not contain any IDM contractor data. Scaling training data increases log crafting capabilities. The zero-shot model only begins to start crafting after ~100 hours of training data. When fine-tuning each foundation model to contractor data, the zero-shot crafting rates for crafting tables and wooden tools increase by orders of magnitude. The zero-shot crafting rates for the entire $\sim 70\text{k}$ hour `web_clean` dataset. We furthermore only see the emergence of the contractor data 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, after which there are few to no gains. Similarly, crafting tables are only crafted after 50 or more hours of data, and again gains plateau after 100 hours. While in all previous experiments we used 1962 hours of data, these results suggest we could reduce that number to

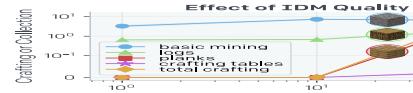


Figure 9: Zero-shot performance of models trained from scratch on the `earlygame_keyword` dataset labeled with IDMs increasing amounts of contractor data.

5 Discussion and Conclusion

The results presented in this paper help pave the path to utilizing the wealth of the web for sequential decision domains. Compared to generative video modeling that would only yield representational priors, VPT offers the exciting possibility to act during pretraining and using these learned behavioral priors as extremely strong priors for RL. VPT could even be a better general representation learning method than GPT-3. The downstream task is not learning to act in that domain—for example, fine-tuning a model to play a game—but rather learning to predict what is happening in a video—because arguably the most important information in a video is what is present in features trained to correctly predict the distribution over future frames. This is an intriguing direction for future work.

Future work could improve results with more data (we estimate we could train models on larger, better-tuned datasets). Furthermore, all the models in this work condition on closed captions only; we cannot ask the model to perform specific tasks. Appendix I presents an experiment where we condition our models on closed captions (text transcripts of speech in a video).

become weakly steerable; we believe this a rich direction for future research. VPT also consistently correlates with downstream evaluation metrics (Sec. 4.2), which is slow and hard-won. Another fruitful future direction would be to investigate various training metrics and downstream evaluations. Finally, while we do not mitigate negative societal impacts from the models trained in this work, as VPT improves domains it will be important to assess and mitigate harms that emerge with or without internet datasets, such as emulating inappropriate behavior.⁶⁷

In conclusion, VPT extends the paradigm of training large and general purpose models on freely available internet-scale data to sequential decision domains. Our model achieves zero-shot behavior and, when fine-tuned with RL, achieves an unprecedented performance on the diamond pickaxe in Minecraft (all the more difficult given the human interface). We show that contractor data is far better used within the VPT pipeline than to train a few models, and that only a small amount of contractor data (about \$2000 USD) was required to learn large amounts of unlabeled online data for use in BC. Finally, learning with the human interface is highly general and allows losslessly modeling the entire distribution of behavior. While we only experiment in Minecraft, we believe that VPT provides a general framework for learning behavioral priors in hard, yet generic, action spaces in any domain that has access to large amounts of available unlabeled data, such as computer usage.

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*, 2021.
- [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] Yinhua Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, 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, 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 Ramanathan, Kaiming He, Mahadevan Paluri, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of pre-training. 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 Dubey, Ishay Amitai, Scott Wu, Andrej Karpathy, Stefan Alama, Zhaoqi Dong, Julian Bern�ić, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning long-range language models from natural language supervision. In *International Conference on Learning Representations (ICLR)*, pages 8748–8763. PMLR, 2021.
- [7] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simon Sissons, Daniel Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Bowring, and Samy Bengio. Opportunities and risks of foundation models. *arXiv preprint arXiv:2106.04560*, 2021.
- [8] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. S-Former: A multi-scale backbone for image classification. *CoRR*, abs/2106.04560, 2021. URL <https://arxiv.org/abs/2106.04560>.
- [9] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and Thore Graepel. Mastering the game of go with deep neural networks and tree search. *Nature*, 2016.
- [10] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Daan Wierwille, John Schulman, Tianqi Chen, Ilya Sutskever, and Pieter Abbeel. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 2019.

- [11] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [12] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Dennison, David Farhi, Quirin Fischer, Shariq Hashmi, Chris Hesse, et al. Scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [13] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Gordon, and Igor Mordatch. Emergent tool use from multi-agent autocuriosity. *arXiv:1909.07528*, 2019.
- [14] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guillaume Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avradip Majumdar, et al. Human-level performance in 3d multiplayer games with population-based training. *Science*, 364(6443):859–865, 2019.
- [15] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sosik, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming humans on a suite of tasks. In *International Conference on Machine Learning*, pages 411–420, 2017.
- [16] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Munos, et al. Unifying count-based exploration and intrinsic motivation in reinforcement learning. *Information processing systems*, 29, 2016.
- [17] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- [18] Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Hod Lipson. Nature, 590(7847):580–586, 2021.
- [19] Peter C Humphreys, David Raposo, Toby Pohlen, Gregory Thornton, Rajat Sen, Michael Lai, David Muldal, Josh Abramson, Petko Georgiev, Alex Goldin, Adam Santoro, et al. Deep reinforcement learning for learning to control computers. *arXiv preprint arXiv:2202.00001*, 2022.
- [20] Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Li Fei-Fei. Bits: An open-domain platform for web-based agents. In Doina Precup, Yair Wolfson, and Yitong Zhou, editors, *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 3135–3144. PMLR, 2017. URL <https://proceedings.mlr.press/v70/shi17a.html>.
- [21] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *Advances in neural information processing systems*, volume 1, page 2, 2000.
- [22] Faraz Torabi, Garrett Warnell, and Peter Stone. Recent advances in behavioral cloning from observation. *arXiv preprint arXiv:1905.13566*, 2019.
- [23] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Information processing systems*, 29, 2016.
- [24] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from raw video. *arXiv preprint arXiv:1805.01954*, 2018.
- [25] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation learning to imitate behaviors from raw video via context translation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 11118–11125. IEEE, 2018.
- [26] Twinfinite Staff. Most played games in 2021, ranked by peak concurrent players. *Twinfinite*. URL <https://twinfinite.net/2021/01/most-played-games-in-2020-ranked-by-peak-concurrent-players/>.
- [27] William H Guss, Brandon Houghton, Nicholas Topin, Phillip Wang, Oriol Vinyals, Pedro Domingos, 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 Yonatan Zohary. A hierarchical approach to lifelong learning in minecraft. In *Proceedings on Artificial Intelligence*, volume 31, 2017.
- [29] Christian Scheller, Yanick Schraner, and Manfred Vogel. Sample efficient learning through learning from demonstrations in minecraft. In *NeurIPS Demonstration Track*, pages 67–76. PMLR, 2020.
- [30] Ingmar Kanitscheider, Joost Huizinga, David Farhi, William Hebgen Gao, Raul Sampedro, Peter Zhokhov, Bowen Baker, Adrien Ecoffet, Jie Li, et al. Curriculum learning in a complex, visual, hard-exploration domain: Minecraft. *arXiv:2106.14876*, 2021.
- [31] Junhyuk Oh, Valliappa Chockalingam, Honglak Lee, et al. Control of movement and action in minecraft. In *International Conference on Machine Learning*. PMLR, 2016.
- [32] Vihang P Patil, Markus Hofmarcher, Marius-Constantin Dinu, Matthias Johannes Brandstetter, Jose A Arjona-Medina, and Sepp Hochreiter. Learning from few demonstrations by reward redistribution. *arXiv preprint arXiv:2106.09939*, 2021.
- [33] Alexey Skrynnik, Aleksey Staroverov, Ermek Aitygulov, Kirill Aksentev, and Aleksandr I Panov. Forgetful experience replay in hierarchical reinforcement learning from demonstrations. *arXiv preprint arXiv:2006.09939*, 2020.
- [34] Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, and Ming Tang. Playing minecraft with sample-efficient hierarchical reinforcement learning. *arXiv:2112.04907*, 2021.
- [35] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Journal of 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] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browne. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):46–58, 2009.
- [38] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jones. A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(3):1–35, 2018.
- [39] Claude Sammut, Scott Hurst, Dana Kedzier, and Donald Michie. *Learning Proceedings 1992*, pages 385–393. Elsevier, 1992.
- [40] Alessandro Giusti, Jérôme Guzzi, Dan C Cireşan, Fang-Lin He, Julian Fontana, Matthias Faessler, Christian Forster, Jürgen Schmidhuber, and Carsten Stachniss. A machine learning approach to visual perception of forest trails for mobile robots. *IEEE Transactions on Robotics and Automation Letters*, 1(2):661–667, 2015.
- [41] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Massimiliano Frasconi,天鹅湖 Lawrence D Jackel, Matheus Monfort, Urs Muller, Jiakai Zhang, et al. End-to-end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [42] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Daniel Ciresan. End-to-end driving via conditional imitation learning. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 4693–4700. IEEE, 2018.
- [43] Rémi Coulom. Computing “elo ratings” of move patterns in the game of go. *Computers and games*, 4(4):198–208, 2007.
- [44] Todd Hester, Matej Večerík, Olivier Pietquin, Marc Lanctot, Tom Schaul, John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning for markov decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.

- [45] Ashley Edwards, Himanshu Sahni, Yannick Schroecker, and Charles Pmlr, 2019.
- [46] Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, a Reinforcement learning of physical skills from videos. *ACM Transactions on Graphics*, 37(6):1–14, 2018.
- [47] Feryal Behbahani, Kyriacos Shiarlis, Xi Chen, Vitaly Kurin, Sudhanshu Joao Gomes, Supratik Paul, Frans A Oliehoek, Joao Messias, et al. Learning in the wild. In *2019 International Conference on Robotics and Automation*, pages 775–781. IEEE, 2019.
- [48] Yusuf Aytar, Tobias Pfaff, David Budden, Thomas Paine, Ziyu Wang, and Peter Welinder. Playing hard exploration games by watching youtube. *Advances in neural information processing systems*, 31, 2018.
- [49] Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Ivan Shelhamer, Jitendra Malik, Alexei A. Efros, and Trevor Darrell. Zero-shot learning via visual analogy. In *ICLR*, 2018.
- [50] Ashvin Nair, Dian Chen, Pulkit Agrawal, Phillip Isola, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Combining self-supervised learning and imitation for vision-based manipulation. In *ICRA*, pages 2146–2153, 05 2017. doi: 10.1109/ICRA.2017.7989247.
- [51] Duy Nguyen-Tuong, Jan Peters, Matthias Seeger, and Bernhard Schölkopf. Comparison of policy gradient methods for reinforcement learning with dynamics: a comparison. In *European symposium on artificial neural networks*, pages 1–6, 2008.
- [52] David Abel, Alekh Agarwal, Fernando Diaz, Akshay Krishnamurthy, and John Langford. Exploratory gradient boosting for reinforcement learning in complex environments. *arXiv:1603.04119*, 2016.
- [53] Dilip Arumugam, Jun Ki Lee, Sophie Saskin, and Michael L Littman. Policy learning from policy-dependent human feedback. *arXiv preprint arXiv:1706.05001*, 2017.
- [54] Alexander Trott, Stephan Zheng, Caiming Xiong, and Richard Socher. Solving sparse reward tasks using self-balancing shaped rewards. In *Information Processing Systems*, 32, 2019.
- [55] Stephan Alaniz. Deep reinforcement learning with model learning and transfer in minecraft. *arXiv preprint arXiv:1803.08456*, 2018.
- [56] Hiroto Udagawa, Tarun Narasimhan, and Shim-Young Lee. Fighting zombies with deep reinforcement learning. Technical report, Technical report, Technical report, Seoul National University, 2016.
- [57] Tianmin Shu, Caiming Xiong, and Richard Socher. Hierarchical and iterative multi-task reinforcement learning. *arXiv preprint arXiv:1712.07370*, 2017.
- [58] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot learning with multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pages 2661–2670. PMLR, 2017.
- [59] Zhengxiang Shi, Yue Feng, and Aldo Lipani. Learning to execute or as a function. *arXiv preprint arXiv:2204.08373*, 2022.
- [60] Tambet Mattiisen, Avital Oliver, Taco Cohen, and John Schulman. Team learning. *IEEE transactions on neural networks and learning systems*, 2022.
- [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. Physics-informed neural networks. In Nando de Freitas, Marina Meila, and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 1227–1237. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/cobbe21a.html>.
- [65] Dhireesha Kudithipudi, Mario Aguilar-Simon, Jonathan Babb, Max Blackiston, Josh Bongard, Andrew P Brna, Suraj Chakravarthi Raja, Nishant Chaturvedi, et al. Biological underpinnings for lifelong learning machines. *Nature*, 3(3):196–210, 2022.
- [66] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Cenzer, 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 Shmuel Rubinstein. The dangers of stochastic parrots: Can language models be too big??. In *2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 1–12, 2021.
- [68] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12(Oct):2825–2830, 2011.
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Lee, 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:1607.06450*, 2016.
- [71] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. In *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:1412.6980*, 2014.
- [73] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Alban Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: A high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://paperswithcode.com/sota/pytorch-an-imperative-style-high-performance-deep-learning-library>.
- [74] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Yoshua Bengio. Identifying and attacking the saddle point problem in convex optimization. *Advances in neural information processing systems*, 29, 2016.
- [75] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How do deep neural networks? *Advances in neural information processing systems*, 29, 2016.

- [76] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv:1901.02860*, 2019.
- [77] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and John Schulman. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [78] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv:1506.02438*, 2015.
- [79] Ronald J Williams. Simple statistical gradient-following algorithms for reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- [80] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1921–1929, 2016.
- [81] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachael Houthooft, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Scalable impala: Bootstrapping from soft actor-critic. *Advances in neural information processing systems*, 30, 2017.
- [82] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal function approximators. In *International conference on machine learning*, pages 1025–1033, 2015.
- [83] Open Ended Learning Team, Adam Stooke, Anuj Mahajan, Catarina Gomes, Jakob Bauer, Jakub Sygnowski, Maja Trebaacz, Max Jaderberg, Michael Wellman, and Michael Wellman. How does open ended learning leads to generally capable agents. *arXiv preprint arXiv:1909.00001*, 2019.
- [84] Jelena Luketina, Nantas Nardelli, Gregory Farquhar, Jakob Foerster, Michael Grefenstette, Shimon Whiteson, and Tim Rocktäschel. A survey of neural architectures informed by natural language. In *Proceedings of the Twenty-Eighth International Conference on Artificial Intelligence, IJCAI-19*, pages 6309–6317. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.2496/00000000000000000000000000000000. URL <https://doi.org/10.2496/00000000000000000000000000000000>.
- [85] DeepMind Interactive Agents Team, Josh Abramson, Arun Ahuja, Alán Carnevale, Mary Cassin, Felix Fischer, Petko Georgiev, Alex Goldin, Timo Janzen, and Michael Littman. Multiagent reinforcement learning for multimodal interactive agents with imitation and self-supervised learning. *arXiv:2112.03763*, 2021.
- [86] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Michael Neumann. Text-to-image generation with clip latents. *arXiv preprint arXiv:2204.00682*, 2022.
- [87] Dong Yu and Li Deng. *Automatic speech recognition*, volume 1. Springer, 2014.
- [88] Daulet Nurmanbetov. rpnunet, May 25 2021. URL <https://github.com/dnurman/rpnunet>. Accessed 2022-04-22.
- [89] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michalek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text-to-image synthesis 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 Tenenbaum, Jonathan Gordon, and the OpenAI supercomputing team, especially Christopher Berner.

Supplementary Information

A Collecting Internet Data

A.1 Initial Unclean Dataset Curation

Our goal was to curate a video dataset of Minecraft gameplay from the survival mode, we prefer the data come from game modes as close as possible to our meaning, preferably coming from Minecraft version 1.16, being on a computer and keyboard vs. video game controllers with keypads and other buttons), being a player, and having the default look of the game (vs. modifications that alter it to make it look realistic). To try to accomplish these goals, we collect a dataset of searches of publicly available videos on the internet. A list of search queries is shown in Table 1.

minecraft survival longplay
minecraft gameplay no webcam
minecraft gameplay survival mode
minecraft survival tutorial
minecraft survival guide
minecraft survival let's play
minecraft survival for beginners
minecraft beginners guide
ultimate minecraft starter guide
minecraft survival guide 1.16
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 survival
how to play minecraft
minecraft survival basic
minecraft survival for noobs
minecraft survival for dummies
how to play minecraft for beginners
minecraft survival tutorial series
minecraft survival new world
minecraft survival a new beginning
minecraft survival episodio 1
minecraft survival episode 1
minecraft survival 1. bölüm
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 filtering to eliminate videos that do not fit our target distribution. In this step, we locate keywords in the video title and description and reject videos that contain them. The keywords we use are: {ps3, ps4, ps5, xbox 360, playstation, timelapse, mod, pocket edition, skyblock, realistic minecraft, how to install, how to download}. This process yielded us ~ 270 k hours of unlabeled data, which we filter down as described in the next section.

A.2 Training a Model to Filter out Unclean Video Segments

We restrict the scope of this work to the Minecraft Survival game mode, and our training dataset to clips that are obtained from this mode that are relatively free of clutter.

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

A.2.1 Label Collection

We asked 5 workers on Amazon Mechanical Turk (mTurk) that we selected with a task to label random screen capture images to be used in training the classification 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 (with visual examples of each class):

1. **Minecraft Survival Mode – No Artifacts:** Video frames (images) from the Minecraft Survival game mode that do not contain any non-game elements (e.g., subscribe buttons, channel logos, advertisements, picture-in-picture overlays).
2. **Minecraft Survival Mode – with Artifacts:** Video frames from the Minecraft Survival game mode 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 provided visual examples from each category in the worker instructions, similar to the sample shown in Figure 11):

Please help us identify screenshots that belong only to the survival mode in Minecraft. Screenshots that belong to other modes (Minecraft creative mode, other games, music videos, etc.) should be labeled 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

Valid survival mode videos have health/hunger bars and an item hotbar at the bottom of the screen.

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 are from the Minecraft survival mode gameplay without any noticeable artifacts.
- **Minecraft Survival Mode – with Artifacts:** These images are from the Minecraft survival mode screenshots, but with some added artifacts. Typical artifacts include picture-in-picture overlays (a logo/brand), text annotations, a picture-in-picture of the player's face, or other UI elements.
- **None of the Above:** Use this category when the image is not a Minecraft screenshot. It may be a non-Minecraft frame or from a different game mode such as the creative mode, the health/hunger bars may not be present, or the item hotbar may or may not be still present.

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

Since we perform rigorous keyword and metadata based filtering of videos (and images) on which we served sample images to be labeled, serving offensive content to workers is not a concern.

low risk and no such images were detected during our manual checks. We also note that no personally identifiable information (PII) was collected.



Figure 10: Amazon Mechanical Turk worker interface showing an example of a Minecraft Survival Mode screenshot.



Figure 11: (Left) Sample image for Class 1: Minecraft Survival Mode – with Health Bar. (Middle) Sample image for Class 2: Minecraft Survival Mode – with Artifacts and None of the Above – contains annotations and picture-in-picture of the narrator. (Right) Sample image for Class 3: Minecraft Survival Mode – with None of the Above – Image is missing the hotbar as well as health and armor bars.

A.2.2 SVM Training

With the image labels collected as described in the previous section, we train a Support Vector Machine (SVM) classifier to distinguish between video segments that consist of frames from the Minecraft Survival Mode category. Given a set of labeled images, we obtain embeddings for each image using a pre-trained ResNet CLIP Model.⁶ This is a ResNet-based CLIP model that is scaled up by a factor of 64x the compute of a ResNet-50. We then train a Support Vector Machine (SVM) with a Gaussian kernel to obtain a frame classifier. We use the Scikit-learn⁶⁸ SVM implementation with the configuration given in Table 2.

Finally, we apply the classifier to frames of raw video sequences at a rate of 3 frames per second. We consider two classes: 'Survival Mode' and 'None of the Above'. For videos that consist of at least 80% "clean" frames at this stage (Classes 1 and 2), we apply a median filter (with a kernel size of 7) to the labels and segment them into "clean" segments that are at least 5s in duration. The result of this is our final dataset.

A.3 early_game Dataset

The `early_game` dataset is a ~3000 hour subset of `web_clean` targeted at new players. It consists of instances where players start in a fresh world with no items. The dataset includes a short text that accompanies the videos in `web_clean` and determine whether any of the following expressions match:

CLIP Model Specification	RN50x64 (see table)						
CLIP Input Image Resolution	448x448x3						
CLIP Embedding Feature Length	1024						
SVM Parameters	<table> <tr> <td>Kernel</td><td>rbf</td></tr> <tr> <td>C</td><td>20</td></tr> <tr> <td>Gamma</td><td>scale</td></tr> </table>	Kernel	rbf	C	20	Gamma	scale
Kernel	rbf						
C	20						
Gamma	scale						
Sample Size	<table> <tr> <td>Class 1</td><td>2200</td></tr> <tr> <td>Class 2</td><td>2200</td></tr> <tr> <td>Class 3</td><td>4400</td></tr> </table>	Class 1	2200	Class 2	2200	Class 3	4400
Class 1	2200						
Class 2	2200						
Class 3	4400						

Table 2: Feature Extraction Details and SVM Configuration. The parameter implementation in Scikit-learn⁶⁸.

- `(episodename|day|session|series|chapter|chap .|seriespart|partelpt|roundel|daylt|täpl|bölüm|episodio|episode|эпизод)*(\.\d{1}|#\d{1}|.\d{1}|\#01|\d{1}|one[^\-0-9]|$)`
 - **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.

B Contractor Data

B.1 Recording Contractor Play

Our contractors use a custom Minecraft recorder that we built that records their feeds as they play. The recorder is implemented using the MCP-Reborn (github Reborn) modding package. To ensure that the recorder environment is the same as the Minecraft environment used for RL rollouts and evaluations (Appendix), we run the underlying game engine for both. The recorder is a Java app that runs in constant resolution of 1280x760. Brightness is set to 0 (the "gloomy" setting) as it is the default setting. Other graphics settings (field of view, GUI scale) are set in the Minecraft environment (C.1); we explicitly prevented users from changing them. Unlike the environment, the recorder allows all keyboard key presses and continuous binned mouse actions. On every game step (or "tick") the frame buffer under the window is downsized to 640x360 and written into a video file. In-game actions are stored in a separate JSONL file (a text file where each line is a JSON-formatted string) and chunked into 5 minute clips; after each 5 minute segment of contractor actions, the recorder automatically uploads the video file, the JSONL file with actions, as well as the contractor's name to our cloud storage. To ensure that contractors cannot corrupt each other's data, we provided each contractor with an individual cloud bucket, as well as with credentials giving write access only to their bucket. We also included adjective-adjective-noun names (e.g. grumpy-amethyst-chipmunk) and a namegenerator python package to ensure contractor anonymity when we upload their data.

B.2 Contractor Contract

We recruited contractors by posting the following offer on the UpWork free

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

We had the applications open for a day, and then randomly selected 10 applicants. Later in the project, as we needed more data and as some contractors left their contracts, we added more applicants from the original pool as well as recent working contractors. The contractors were paid \$20 per hour (minus Upwork applicable taxes). All of the results presented in this paper are based on about 400 hours of contractor data (including data recorded to gather statistics of human play that was not used by us around \$90,000. Over the course of the project, we collected some data for bugs in the recorder and for some ideas we ultimately did not pursue. In total, we could likely obtain most of our results with an IDM trained using only \$2000 worth of foundation VPT model, BC fine-tuning to the `earlygame_keyword` dataset, and contractor results. Collecting the `contractor_house` dataset cost about \$8000. Because we trained on about 2000 hours of contractor data, the actual cost of contractor data is around \$40,000.

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

- Collect as many units of wood as possible, using only wooden or stone tools
- Start a new world every 30 minutes of game play
- Build a basic house in 10 minutes using only dirt, wood, sand, and stone tools (`contractor_house`, more details below in Appendix B.4).
- Starting from a new world and an empty inventory, find resources to obtain a diamond pickaxe in 20 minutes (`obtain_diamond_pickaxe`). This data (and the statistics for how long it takes humans on average 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 information that contractors could theoretically use Minecraft to generate personally identifiable information and/or offensive content (e.g., blocks to write their name or offensive messages, then finding a spot from where they can be visible). In practice, we have not seen any attempts to do so in the contractor videos we watched. Of course, we train our BC models on videos from the internet of people, and if such behavior is in those videos our model could also potentially learn it. However, such behavior is rare enough that our model would not be likely to reproduce it.

B.3 Data for the Inverse Dynamics Model.

Since the IDM’s task is to infer actions given the video, any labelled data is used for training. In practice, we included general gameplay as well as the `earlygame_keyword` data in the previous section, which amounted to a total of 1962 hours. Due to the fact that contractors only used the `contractor_house` task at late stages of the project, they were not included in the final dataset.

B.4 `contractor_house`.

The `contractor_house` contains about 420 hours of data. We asked contractors to build a basic house in 10 minutes, using only basic dirt, wood, and sand, blocks. Each trial starts in a generated world and a timer forcibly ends a trajectory after a 20 minute time limit. Contractors chose to begin their trajectories by crafting basic tools and building a foundation. It was common for the first 2 minutes to be spent crafting a wooden pickaxe and an assortment of stone tools before gathering more building blocks and starting to build a structure.

C Minecraft environment details

Our Minecraft training environment is a hybrid between MineRL²⁷ (github.com/Hexception/MCP-Reborn) Minecraft modding package. Unlike MineRL, our environment is not a single world, but rather a collection of worlds, each with its own set of rules and objectives. This allows us to experiment with different environments without having to restart the entire system.



Figure 12: (**Left**) Sample of a Minecraft frame in the original resolution (640x360). The mouse cursor can be seen in the center of the image. This provides the player’s inventory and can be used to craft very basic items. (**Middle**) We downsample the frame to 128x128 for computational reasons. Shown is a downsampled observation without the GUI. (**Right**) A 128x128 observation with the GUI open. The health, hunger, hotbar overlays, and agent’s position are shown in the lower part of the image.

game, in which the server (or the “world”) always runs at 20Hz and the client can complete (typically at 60-100Hz), in our version the client and server run at the same frequency. This allows us to run the environment slower or faster than the server, avoiding artifacts like missing chunks of the world. The action and observation spaces are identical to those of MineRL environments and are described in more detail in the MineRL paper. The environment also returns diagnostic information, such as in-game stats, inventory, whether any in-game GUI is open, etc., which we use for tracking and evaluating the model’s performance. The agent can “die” in a number of ways, such as staying in water and drowning, being killed by hostile mobs, or falling from a tall structure. The episode ends when the episode on agent “death”. Instead, just as for humans in the regular Minecraft game, the agent loses all its items when it dies and respawns at a random spot close to the initial spawn point in the Minecraft world. The policy state is not masked on death, so the model can still act accordingly.

C.1 Observation space

The environment observations are simply the raw pixels from the Minecraft world as a human would see. Unlike MineRL, we do not remove overlays like the hotbar, health bar, or the animation of a moving hand shown in response to the attack or “use” action. The field of view is set to 70 degrees, which corresponds to the Minecraft default. GUI scale (a parameter of the in-game GUI) is set to 2, and brightness is set to 2 (which is not a Minecraft setting, but is very frequently used in online videos). The rendering resolution is 640x360, which is then scaled up to 128x128 before being input to the models. We empirically found 128x128 to be the best resolution for which in-game GUI elements are still discernible, and then we can compute costs. Whenever an in-game GUI is open, we additionally render the mouse cursor at the appropriate mouse position to match what a human player’s observation looks like (Figure 12).

C.2 Action space

Our action space includes almost all actions directly available to human players, such as mouse movements, and clicks. The specific binary actions we include are shown in Table C.1. One difference between the human action space and our agent’s is that we do not allow the agent to type letters, which is only useful for entering text into the search bar of the crafting interface. A human can either do that or browse the recipe book with the mouse, the latter of which is something the agent cannot do. However, because we do allow the agent to press letters that are also shown in the search bar outside of the GUI, the “W” key triggers the forward action (agents are able to search the recipe book within the GUI (W, A, S, D, E, Q) that produce letters if the recipe book search bar is empty). Interestingly, we have not seen agents attempt to search the recipe book with these letters. Instead, they tend to search 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
sneak	shift key	Move carefully in current direction. In GUI it acts as a modifier key; when it moves item from/to the inventory bar, and when used with craft it allows crafting more than the standard number of items possible instead of one.
sprint	ctrl key	Move fast in the current direction.
attack	left mouse button	Attack; In GUI, pick up the stack of items in a GUI cell; when click attack - no attack - attack items of the same kind present in stack.
use	right mouse button	Place the item currently held or user is looking at. In GUI, pick up the a single item from a stack held by another player.
drop	Q key	Drop a single item from the stack the agent is currently holding. If the player drops the entire stack. In the happens except to the item the mouse is hovering over.
hotbar . [1-9]	keys 1 – 9	Switch active item to the one in a hotbar slot.

Table 3: Binary actions included in the action space. <https://minecraft.gamepedia.com/Controls> has more detailed descriptions of each action.

In addition to the binary (on/off) keypress actions, our action space also includes continuous actions such as with human gameplay, when in-game GUIs are not open, mouse movements, agent’s yaw and pitch, respectively. When a GUI is open, camera actions are disabled. Mouse movements are relative (i.e. they move the mouse or camera relative to the previous position).

Inventory interaction in Minecraft requires fine-grained mouse movements for crafting and smelting, while mining and navigating the world can be achieved via keyboard actions. To be able to achieve both with the same action space, we implemented them as a set of discrete actions with foveated binning along each axis (Fig. 13). In our 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 is a sequence of consecutive image frames (128 frames of video), each of which has dimensions $16 \times 16 \times 3$. The IDM is tasked with predicting the action at each frame. All image pixels are scaled by 255.0 such that they lie within the range $[0, 1]$. The first layer of the IDM consists of 128 learnable filters with a temporal kernel width of 5 and spatial kernel size 3. The convolution is non-causal, 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 important for learning the dynamics of the game.



Figure 13: Relative camera angle or mouse movement in pixels vs. action bin used for both X and Y coordinates. The binning is foveated, meaning that bins for smaller movements are more coarse-grained than for larger movements. The axis (X and Y). The center of each bin (indicated with green circles) is used for movements (that is, when converting from an action expressed as a bin to a movement).

as it incorporates neighboring temporal information immediately, and we compare IDM performance with and without this layer in Figure 14. This comparison is based on the (1962-hour) IDM dataset.

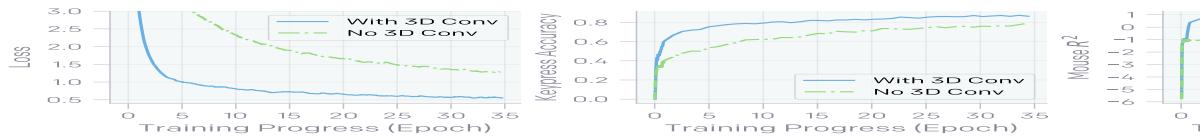


Figure 14: Effect of 3-D Convolution in the IDM Architecture

This initial temporal convolutional layer is followed by a ResNet⁶² image processing part of the model, no extra temporal information is shared between neighbors since each frame was first processed with the temporal convolution, some residual connection is present at this stage. The ResNet image processing network is comprised of three stacks with widths $W = \{64, 128, 128\}$. Each stack is comprised of, in order, (1) a 3x3 convolutional layer with 1-pixel zero padding at the embedding boundary (embedding dimensions are the same as the incoming embedding dimension), (2) a 3x3 max pooling with stride 2 and padding 1 such that the embedding dimension is halved, and (3) two classic ResNet blocks as defined in He et al.⁶² with each block having 64 output channels.

The output of the ResNet stack is flattened into a 1-dimensional vector of 128 vectors (one vector for each frame in the video) such that at this stage there are 128 vectors. Each vector is independently processed with two frame-wise dense layers with 256 and then 4096 output activations, respectively. The result is then fed through four (unmasked) residual transformer⁶⁹ blocks. Each block first has an unmasked attention layer that may attend to future frames, with 32 attention heads of dimension 128 each and a residual connection that skips this layer. The embedding is then passed through another two residual blocks with output dimension 16384 and another with output dimension returning to 128. The final residual connection skips past this pair of frame-wise dense layers (not skipping past the first two residual blocks). All dense layers have their weights tied through time. Finally, independent dense layer heads for each action are pulled from the last two layers of the class on/off categorical parameterized with a softmax for each available action.

categorical for both the discretized horizontal and vertical mouse movement details on the action space).

Each dense layer or convolutional layer in the network is preceded by a layer a ReLU non-linearity. Weights are initialized with Fan-In initialization⁷¹ and zero.

D.2 IDM Training

The total loss for the network is the sum of each independent action prediction key and one for both mouse directions). Each independent loss is the negative of the correct action. We use the ADAM⁷² optimizer with a linear learning rate of 0.003, a learning rate of 0.003, a batch size of 128 (where each item in the batch is a sequence of frames), 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 approximately 4 hours.

We add data augmentation to each video segment; augmentations are randomly applied to each segment such they are temporally consistent. Using the Pytorch⁷³ transform, we randomly scale the image between 0.8 and 1.2, randomly rotate it between -15 and 15 degrees, scale it by a random factor between 0.98 and 1.02, shear it between -10 and 10 pixels, and 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 in this section, we are from one run of training (for IDM, BC, and RL training); this non-ideal because deep learning training tends to be low variance^{74,75} and because we sweep over many hyperparameters (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. The IDM is trained to simultaneously predict the next frame of each video sequence, so the IDM will effectively be causal for frames at time t because future frames are not included in the sequence. For this reason, we generate pseudo labels for the IDM by extracting frames from the middle of a video using a sliding window with stride 64 frames and only use the pseudo labels for frames 32 to 96 (the center 64 frames). By doing this, the IDM prediction for a 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 in Section D.1 except that we modify the architecture so that it is causal (i.e. cannot see future frames). This means the BC architecture does not have the initial non-causal layer that the IDM has (this layer is omitted completely). Furthermore, the residual transformation is causally masked (as is standard in language modeling) and we do 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). This accounts for 35% of all actions they take. Among other reasons, a player may take a long pause, wait for something in the game to finish, to pause between actions, or to take a break (e.g. drink a glass of water). Early on in the project we found that the BC model would take a majority of null actions, often upwards of 95%. In order to prevent this behavior, we remove all null actions from the dataset. We compare a few different treatments: we can filter out null actions after every 1, 3, or 21 frames of consecutive null actions, and include a treatment where we do not filter any null filtering. Null action filtering generally helps, increasing all metrics.

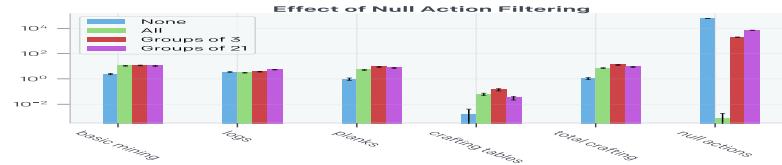


Figure 15: (Left) Effect of Null Action Filtering during training. We compare the number of sampled null actions (y-axis, log scale) and number of sampled null action during rollouts (rightmost group of columns) for four treatments: no null action filtering (blue), filtering all null actions (green), 3 or more null actions (red), and filtering only groups of 21 or more null actions. **(Right)** Joint Hierarchical versus Factorized Action Spaces.

Filtering only groups of 3 performed slightly better than filtering all null actions. In our experiments indicated that filtering all null actions was better; however, after we had already trained our largest models, we found that filtering only 3 or more null actions performed best. Due to compute constraints we were not able to do this 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 on or off, and this choice was independent of whether the mouse was being moved or not. This choice led to issues for modeling the human behavior distribution exactly. Say for a given item, there is a 50% probability (a) move forward and attack or with 50% probability (b) move forward and pick up the item. The best a factored distribution can do is to assign 50% probability to both of these actions because it chooses to press each button simultaneously and independently on the entire action space.

For this reason, we implemented a joint distribution over actions; however, there are over 20 binary buttons and two mouse movement dimensions discretized into 11 bins each, which result in $2^{20} \times 11^2 \approx 1.2 \times 10^8$ possible combinations. This is far too large for the final layer from the transformer stack with a dimension of 4096, so we used a hierarchical joint action space resulting in $4096 \times 1.2 \times 10^8 \approx 5.2 \times 10^{12}$ parameters for the final layer. In order to reduce this we noted that many buttons in Minecraft have no effect if they are pressed simultaneously. For example, if a player tries to move forward and backward at the same time, nothing happens. Below we list the sets of mutually exclusive actions. Furthermore, each action is only ever present in one of the three sets of mutually exclusive actions.

Mutually Exclusive Actions
forward, back
left, right
sprint, sneak
hotbar, [1-9]

Even reducing the joint action space to reflect these mutually exclusive combinations still results in a huge action space when combined with the discretized mouse movements, i.e., $5.2 \times 10^{12} \times 121^2 \approx 5.2 \times 10^{15}$. This calculation results from 3^3 for the 3 sets of 2 mutually exclusive mouse movements, 121^2 for the 4 mouse movement keys, and 5.2×10^{12} for the remaining 4 keys: use, drop, attack, and jump, plus $\times 11^2$ for mouse movement. The inventory button which is mutually exclusive with all other actions is quite large so we chose to implement a hierarchical binary action for camera control. If 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 there is no movement; the secondary mouse movement action is masked during training, and the factored action space need not be sampled during evaluations. While this no longer models the full joint action space, it is quite a bit better than the factored action space since dependencies between the primary action (whether or not to move the mouse (although not which mouse movement) and the secondary action (camera movement)) have been removed. The resulting action space has dimension $3^3 \times 10 \times 2^4 \times 2 + 1 = 8461$ (the 11-dimensional action space for camera movement has been replaced by a multiplier of 2 here, corresponding to the factored action space for whether or not to move the mouse) with an additional 121-dimension hidden state dimension for mouse movements. In the future it would be interesting to implement sequential control over the mouse to more completely model the joint distribution.

In Figure 15 (right) we compare environment rollout performance between the hierarchical joint action space and with the factored action space. Environment rollout performance is comparable; however, we see that the factored action space model samples many more null actions. This is an important example of the factored action space failing to correctly sample the dataset because, due to null action filtering, there are 0 null actions in the dataset that the model can train on. Despite this, the factored model samples many null actions because the key press 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 that it uses 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 that we use a subset of all the videos (early-game dataset, described in A.3) with pseudo labels (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 10 epochs of fine-tuning on contractor-house dataset, and 16 A100 GPUs for about 10 epochs of 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) algorithm⁷⁷, which is based on the proximal policy optimization (PPO) algorithm⁷⁷ that increases the performance of the reinforcement learning agent by performing additional passes over the collected data to optimize the value function.

auxiliary value function. These algorithms have been described extensively so here we describe them only briefly. A major inefficiency when training that, to remain on-policy, one can only take a single gradient step before to be gathered to continue optimization. To alleviate the potentially destructive multiple optimization steps in a single iteration, PPO prevents the policy from taking multiple optimization steps in a single iteration by clipping the loss before the update becomes too large.⁷⁷ We also use generalized advantage estimation to speed-up credit assignment by looking more than 1 step into the future. The advantage of an action, with the look-ahead being determined by hyperparameter γ .

PPG improves the sample efficiency of PPO when the policy and value function are shared representation. PPG splits optimization in two phases: a wake phase and a sleep phase. In the wake phase, the policy and value function are optimized as in normal PPO, except that every sample is used at most once, which prevents the policy from being trained multiple times on the same sample. In the sleep phase PPG optimizes the value function and an auxiliary value function (which is optimized with the exact same loss as the regular value function, but during training), while keeping a Kullback-Leibler (KL) divergence loss to the policy fixed. The purpose of the sleep phase to ensure that the policy does not change. Because the policy is frozen during the sleep phase, PPG does allow samples to be reused multiple times in this phase. The reason behind optimizing the value function during the sleep phase is that value function is less sensitive to being trained multiple times on the same sample. Optimizing the auxiliary value function does not directly affect either the value function or the policy, but it does affect the representation of both functions (the assumption being that predicting the auxiliary value function requires encoding all features that are important for distinguishing states). The combined losses (value function loss, auxiliary value function loss, and KL loss) are used to train the policy. In experiments a single iteration consists of two sleep cycles and one wake cycle. Because the value and auxiliary value functions are not optimized during the sleep phase, they are initialized at the start of RL fine-tuning. Each value function is implemented as a fully connected layer on top of the last residual transformer block of the pretrained network. The weights of the auxiliary value function are randomly initialized while the weights of the value function are initialized with zero weights, which appeared to prevent overfitting in training that could happen with a randomly initialized value function. To prevent the auxiliary value function from having gradients that depend greatly on the magnitude of the reward, the auxiliary value function target is updated by subtracting the mean and dividing by the standard deviation estimated through an exponentially weighted moving average.

To prevent catastrophically forgetting the skills of the pretrained network we apply an auxiliary KL divergence loss between the RL model and the frozen pretrained policy. This loss is defined as:

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

Where π_θ is the the policy being trained, π_{pt} is the frozen pretrained policy, and ρ is the coefficient to weight this loss relative to other losses.

In the fine-tuning experiments, this KL divergence loss replaces the common cross-entropy loss, which is often added to RL experiments to encourage exploration.^{79,80} The goal of entropy maximization is that, when all actions appear to have equal value, such as when the agent has learned about the next reward, it should maximize its entropy to increase the uncertainty of the next reward. Blindly exploring by maximizing entropy is effective when state and action spaces are sufficiently small or the reward is sufficiently dense, but becomes ineffective when state and action spaces are large and rewards are sparse, which is the case for our task. Instead of blindly exploring through uniform-random actions, we assign a policy to explore interestingly new states, and thus, in states where the agent assigns equal probability to all actions, it should mimic the action-distribution of the pretrained policy instead of uniform random action distribution. In experiments with a randomly initialized policy we found that using a maximization loss with a coefficient of 0.01, which has been an effective setting for PPO work.³⁰ Empirically, we found that a high coefficient ρ for this KL divergence loss prevents the agent from properly optimizing the reward function while a low coefficient ρ prevents the agent from properly optimizing the policy function.

Hyperparameter	Value
Learning rate:	2×10^{-5}
Weight decay:	0.04
Batch size:	40
Batches per iteration:	48
Context length:	128
Discount factor (γ):	0.999
GAE λ :	0.95
PPO clip:	0.2
Max Grad norm:	5
Max Staleness:	2
PPG sleep cycles:	2
PPG sleep value-function coefficient:	0.5
PPG sleep auxiliary value-function coefficient:	0.5
PPG sleep KL coefficient:	1.0
PPG sleep max Sample Reuse:	6
KL divergence coefficient ρ :	0.2
Coefficient ρ decay:	0.999

Table 6: Hyperparameters for RL experiments. These are the hyperparameters for two exceptions. First, when fine-tuning from the early-game model without setting out of a sweep over 5 different learning rates), as we found that performance was lower with the standard learning rate of 2×10^{-5} and the agent did not even learn. We suspect that the reason that the learning rate needed to be lowered when the KL loss is that the KL loss prevents making optimization steps that change the single step, especially in early iterations when the value function has not learned. The KL loss thus makes it possible to optimize with a higher learning rate. In the RL from a randomly initialized policy there is no KL divergence loss or KL coefficient, instead we use an entropy bonus of 0.01, which reportedly worked well in practice.

protecting the learned skills of the pretrained policy and preventing catastrophic forgetting. We start with a relatively high coefficient ρ and decay it by a fixed factor after each epoch. This method protects policy skills in early iterations while guaranteeing that the agent maximizes the reward function, regardless of how different its behavior has become from the pretrained policy.

For the reward function we estimated the rough quantities of each item that the agent gathers when trying to craft a diamond pickaxe, and we reward the model a quantity for each item. We started these estimates by iterating over the technology tree, adding a diamond pickaxe and adding the requirements for each item to the reward function for crafting a diamond pickaxe, then we added the 3 diamonds required for crafting a diamond pickaxe, then we added the 1 iron pickaxe requirement (and so on). Then we added coal and torches to the reward function, with a base reward when smelting iron and for crafting torches while the torches themselves prevent enemies from spawning. Finally, we reward the model for bringing additional logs required to craft all items in the reward function, but we reward up to 8 logs if the agent has crafted into a crafting table or sticks if the agent runs out. In practice, the agent places the torches, or uses coal as fuel when smelting, and the reward was based on human expectations on what would be useful to execute this task. Around how an RL model behaves after training. Finally, to encourage the agent to gather diamonds and crafting diamond pickaxes after it has crafted its first diamond pickaxe, we set a limit on the number of diamonds or diamond pickaxes that would be rewarded.

The rewards for the different items are separated into 4 tiers, roughly depending on how many items the agent would usually get the relevant item. The first tier consists of all wooden items with a base reward of 1, the second tier consists of all items requiring coal with a base reward of 2, the third tier consists of all items requiring iron with a base reward of 4, and the final tier consists of all items requiring diamonds with a base reward of 8. Thus items later in the sequence of items towards a diamond pickaxe will have a higher reward than earlier items.

Item	Quantity rewarded	Reward per item
Log	8	0.000000
Planks	20	0.000000
Stick	16	0.000000
Crafting table	1	0.000000
Wooden pickaxe	1	0.000000
Cobblestone	11	0.000000
Stone pickaxe	1	0.000000
Furnace	1	0.000000
Coal	5	0.000000
Torch	16	0.000000
Iron ore	3	0.000000
Iron ingot	3	0.000000
Iron pickaxe	1	0.000000
Diamond	inf	0.000000
Diamond pickaxe	inf	0.000000

Table 7: Reward per item and total quantity rewarded

give a higher reward. To make sure that the agent does not over-value items which can cause the agent to focus on planks at the expense of creating a crafting table, we reward the agent for each item by the total quantity that the agent gets rewarded for determining the reward, the total quantity for diamonds is 3 and the total quantity for wooden pickaxes is 1, even though we did not put a limit on the number of these items. For example, the agent is rewarded for 3 iron ore, which has a base reward of 4 and up to 3 blocks of iron ore are rewarded, thus the reward per block of iron ore and reward for each item are listed in Table 7.

While every item in the sequence towards a diamond pickaxe is rewarded, the sequence is sparse and, in some cases, even deceptive. The sparsity comes from the fact that the agent needs to take many actions to find the next reward, even after the agent has acquired all the items in the sequence. For example, the agent needs to mine 10,000 blocks of cobblestone (e.g. human players often take more than 10,000 actions to find a diamond pickaxe). The reward function can be deceptive when the most efficient method to get the next item is to mine down immediately after crafting a wooden pickaxe, while leaving up to 9 blocks of cobblestone unmined. For example, the agent can make it far more difficult to get the next item. For example, a good strategy to craft a stone pickaxe quickly is to mine (i.e. spend a few seconds to pick up 10,000 blocks of cobblestone) and then immediately start crafting a wooden pickaxe, such that the agent has immediate access to a crafting table. This makes it easier for the agent to craft a stone pickaxe quickly. However, the fastest way to get a reward for a stone pickaxe is to mine down immediately after crafting a wooden pickaxe, while leaving up to 9 blocks of cobblestone unmined. Thus following the optimal strategy for gathering cobblestone makes it more difficult for the agent to craft a stone pickaxe.

Experiments ran for approximately 6 days (144 hours) on 80 GPUs (for pre-training) and 56,719 CPUs (mostly for collecting rollouts from Minecraft). In this time the agent collected roughly 4,000 optimization iterations and collected roughly 1.4 million Minecart frames, with an average of 12,000 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 experiments are presented in this section. First, we show the items-over-training figure without the early-game model without a KL loss (Fig. 16). When training without the early-game model, the agent only learns to obtain the four items that the early-game model is capable of obtaining: logs, planks, sticks, and crafting tables. Second, we present preliminary experiments where we directly compare RL fine-tuning from the house-building model and RL fine-tuning from 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 was 0.001, and the reward for each item was $1/\text{quantity}$ for all items (i.e. items closer to the end of the sequence did not have an increased reward). While RL fine-tuning from the house-building model

worked better than RL fine-tuning from the early-game model, fine-tuning from the early-game model worked better after 800,000 episodes and showed signs of smelting iron ingots. The early-game model was chosen for the main experiments.

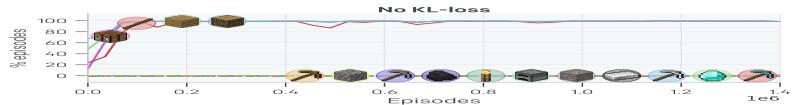


Figure 16: Items obtained when RL fine-tuning from the early-game model. The early-game model learns to obtain all items that the early-game model can craft zero-shot: a wooden stick, a wooden torch, a crafting table, and an iron ingot. In contrast to the treatment with a KL-penalty, it takes the model beyond these initial four, likely because skills that are not performed zero-shot, which the model thus does not initially see any reward, are catastrophically forgotten when the rewards are learned.

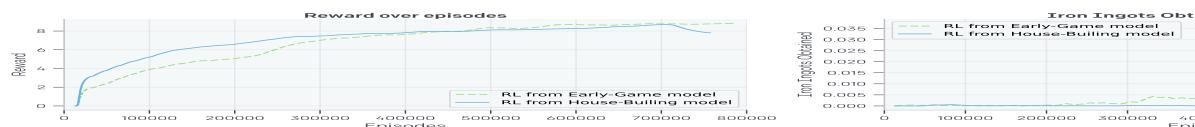


Figure 17: Preliminary experiments when RL fine-tuning from the early-game model. (Left) While reward initially increases faster when fine-tuning from the house-building model, fine-tuning from the early-game model ends up with a slightly higher reward. (Right) RL fine-tuning from the early-game model is much better at learning how to smelt an iron-ingot, which is why the early-game model was chosen for the main experiments.

H Foundation Model Scaling

In early experiments we found that increasing model size led to models learning regime longer into training.⁶³ Here we compare the 0.5B model described above with both a 248M and 71M parameter model. Both of these models are trained for the same number of epochs, up to the 30 epochs the 0.5B model trained for. These models have the same number of layers as the 0.5B model but each layer in the 248M parameter model has 1/2 the width and the 71M parameter model 1/3 the width. The 71M model was trained with an initial learning rate of 0.0001831, batch size of 480, and weight decay of 0.044506. The 248M model had a learning rate of 0.0001831, batch size of 640, and weight decay of 0.051376.

In Figure 18 we show validation loss on web_clean with IDM pseudo-labels. We also show the validation environment performance for the 71M, 248M, and 0.5B models. While the 71M model has the highest validation loss on web_clean, these results do not tell the clear story that one model is better than its smaller counterparts. The 71M model has the lowest contractor data loss and the best validation loss, while the 248M model has the highest web_clean loss, and it also has the best zero-shot environment performance. It is interesting to note that the 71M model even had non-zero wooden tool crafting (Fig. 18 bottom row), which also appears to be better at crafting than the 0.5B, and also has lower contractor data loss.

While the zero-shot results suggest smaller models are better, fine-tuning to contractor_house, model size rank ordering reverses and the 0.5B model performs best both in validation loss (Fig. 19 left) and in environment performance (Fig. 19 right).

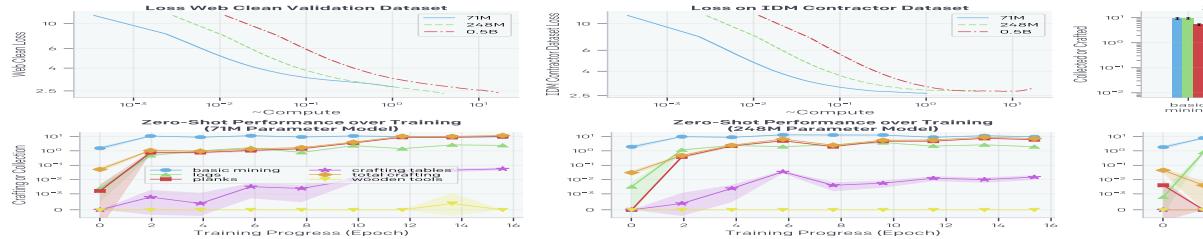


Figure 18: Training and Zero-Shot Performance versus Model Scale. In the is compute normalized to that used by the 71M parameter model, such that after the 71M model has used 1 “compute”. The 248M parameter model and the 71M model have used the same amount of data (15 epochs), and the 0.5B parameter model is trained on the same amount of data (1 epoch). **(Top Left)** Loss on the web_clean validation dataset; note that these models were trained only on web_clean and not contractor_house. **(Top Middle)** Loss on the contractor_house validation dataset; note that these models were trained only on web_clean and not contractor_house. **(Top Right)** Zero-shot environment rollout performance at the end of training for the 71M model (bottom left), 248M model (bottom middle), and 0.5B model (bottom right).

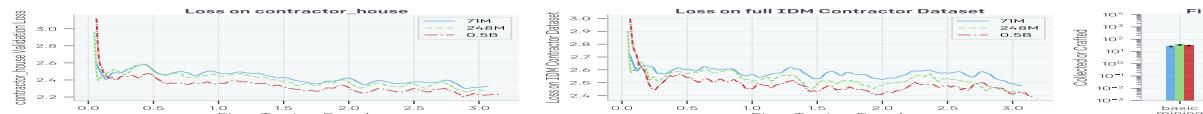


Figure 19: contractor_house fine-tuning performance versus model size. We fine-tune the 71M model on the contractor_house holdout validation set. **(Left)** Loss on the contractor_house holdout validation set. **(Middle)** Loss on the full contractor_house dataset. **(Right)** Environment rollout performance at the end of fine-tuning.

followed by the 248M model and then the 71M model. Environment models trained using the same game engine that we use to collect contractor data, which comes from videos taken from the web. It is plausible that the larger models capture peculiarities in web data during pretraining since they have worse contractor_house performance (middle), and this causes them to perform more poorly in the environment. We hypothesize that because the contractor_house dataset we fine-tune to is generated by the same game engine, the larger models that are a better overall Minecraft prior (as indicated by their validation loss in Fig. 18 top left) can quickly shift their low level features to coming from our game engine, resulting in better environment rollout performance. This is further supported by Fig. 19 (middle) showing loss on the contractor_house dataset after training, which has no overlap with contractor_house. After just a few epochs on contractor_house, all models quickly improve in loss on the full IDM, with the larger models now performing best. While not conclusive, we believe this provides some intuition for future studies of model scaling for sequential decision making.

I Text Conditioning

Goal-conditioned policies^{81,82} make it possible for a single agent to perform goals in a single environment, which is particularly relevant in open-ended Minecraft. In recent work, goal specification has increasingly taken the form of languages⁸³, or even natural language^{84,85}. The benefits of language-conditioned agents are tremendous, especially natural-language-conditioned agents, as their goal variety of potentially very complex tasks. Text conditional models have shown to perform tasks zero-shot (or learn them few-shot) including generalizing in compositional and combinatorial possibilities allowed by natural language (e.g. 2⁸⁶). We hypothesize that we should expect similar capabilities to emerge in conditioned virtual agents, if they are similarly trained on enormous amounts of natural language description to a sequence of actions that completes the section we take preliminary steps toward that future. Our preliminary experiments show that it is possible to pretrain a natural-language-conditioned model for Minecraft using the approach presented in this paper (VPT) plus conditioning on the speech videos.

In online videos, the human actor sometimes indicates their intent in the (e.g. “Let’s go chop some trees to make a wooden axe” or “now let’s learn Photoshop”). Conditioning on this closed caption data could produce a *steerable* agent, i.e., it may later be possible to condition the model with text such as “I am going to pickaxe,” or “I am going to build a house,” and have the agent perform those actions than simply follow typical human behavior (as was investigated in the rest of the paper). The way to produce a steerable agent is via RL fine-tuning, which we could have done by adding a bit indicating the task to be completed, as has been done in [1]. Conditioning on natural language offers many benefits over that approach. First, it is powerful, being able to express any task. Second, one does not need to pre-specify what needs to be completed ahead of time. This would allow for general, capable, zero-shot agents extending those capabilities to embodied tasks such as completing tasks on computers in 3D worlds. Third, text conditioning can be used even when tasks are difficult to define functions (e.g. “Let’s build a house” or—if the agent is capable of doing it—“I will now build a castle surrounded by a moat”). In the limit, VPT+text could produce powerful, capable, natural-language-conditional agents with the powers of GPT-3, able to understand instructions, and complete tasks zero or few shot, but in the form of agents that can move in worlds, complete tasks on computers, and in other similar embodied sequences. We do not reach those lofty goals in this work, but we began a first step toward that direction.

Many Minecraft videos feature audio commentary from the player. This could be present in the form of closed captions for the videos, or could be extracted via speech recognition (ASR).⁸⁷ Our dataset features about 17k hours of commentary captions.

We fine-tuned the 220 million parameter VPT foundation model used in experiments (chosen vs. 0.5B for the same reason: to reduce compute cost). Given text-conditioning input on the subset of our data for which closed captions are available, we first split videos into 30 second chunks. The same text-conditioning input, we first split videos into 30 second chunks. The same text frame in a given chunk, and is made up of all the closed captions occurring within the chunk. The text embedding vector is obtained by concatenating the text as the line of text preceding and following the chunk (if any). Because the vast majority (95%) of our closed caption data lacks capitalization and punctuation, it is purified using a library⁸⁸. We then obtain a text embedding vector of length 4,096 from the OpenAI GPT-2 model, which is processed by a randomly initialized multi-layer perceptron (MLP) with hidden size 2,048. The resulting activations are added for each frame to the pretransformerActivations before the transformer layers (pretransformerActivations += mlp(textEmbedding)). The model is fine-tuned for four epochs.

Our model shows evidence of steerability. When conditioned on sentences explore (such as "I'm going to explore" and "I'm going to find water") the agent farther from its spawn point (Figure 20a). Additionally, we can steer the agent

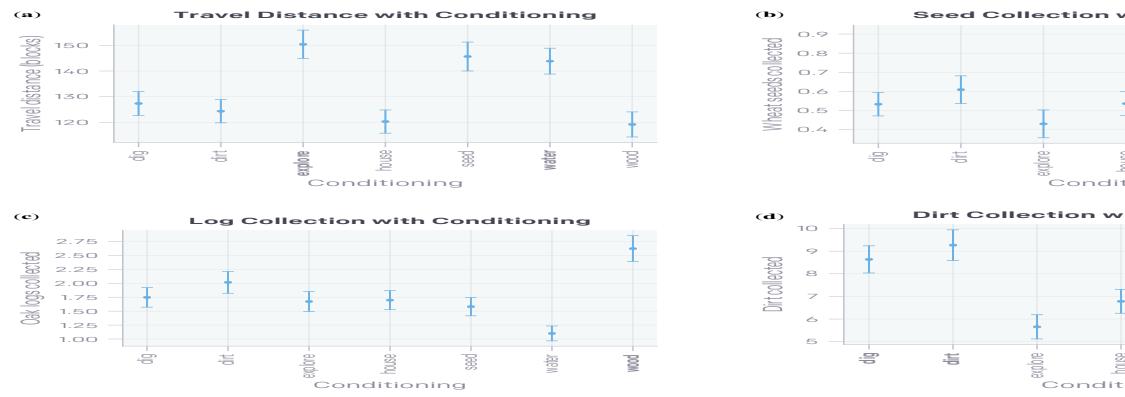


Figure 20: Evidence for conditioning. In each plot, the variants expected to be the most common are shown in bold. The strings corresponding to each variant are shown in Table 8. Statistics are shown for 100 episodes per minute. (a) Distance traveled by the agent. Both “explore” and “water” variants encourage a steerable agent to move more than when doing other tasks, which is needed to get seeds. (b) Collection of wheat seeds. This produces more travel (as the agent sometimes needs to move to a biome) and distance is the Euclidean distance from the spawn point to the farthest point in the episode on the horizontal (x-z) plane. (c) Collection of logs. This is substantially more than other variants, as expected of a steerable agent. (d) Collection of wood. This is the most common type of wood (logs). The “wood” variant collects significantly more wood than other variants, as expected of a steerable agent (we speculate that the “water” variant collects no trees in water). (e) Collection of dirt. The “dirt” and “dig” variants collect dirt, while the variants that are (indirectly in the case of “dig”) conditioned to collect dirt aim at the ground rather than at grass or trees when collecting seeds or wood. The slightly higher amount of dirt collected by these variants. In all cases, 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.

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 with text such as “gather seeds/chop wood/collect dirt” (Figure 20b,c,d).

While our results show some level of steerability, more work is required to improve it. We were not able to successfully steer agents to gather flowers or to hunt, both of which are common in the early game, but less common (and, in the case of hunting animals, more difficult) than gathering dirt, wood, or seeds. Likewise, an experiment in which the agent was conditioned to craft a given item (e.g., “craft a wooden axe”) failed to show that the conditioning had a significant effect on what was crafted. Instead, it seemed the agent was more influenced by the prior, unconditioned behavior of what human players would craft next given the resources available, which makes sense since in Minecraft, especially in the early game, there is a relatively consistent pattern where resources in a specific order go produce more powerful tools (Fig. 6). For example, we asked the agent to craft a stone pickaxe and it instead made a wooden pickaxe, which it often would make the stone pickaxe anyway. Finally, looking at videos of a human player can convince us that the “house” conditioning causes the agents to take more steps to build a house than other variants.

Thus, our results show that it is possible to train a somewhat steerable natural language processing agent. However, its steerability is still too weak to be practically useful, and much more research could be accomplished with more research, data, and training. Our primary research direction is to have the model predict future text as well as just the current text.