
Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos

Sebastian¹ Byrd² Peter Doherty¹ Joel Briley¹
bensch@csail.mit.edu jbyrd@mit.edu pdoherty@mit.edu joel@mit.edu

Jelley¹ Abishek Balaji¹ Brantley Boggs¹ Reid Campbell¹
jelley@mit.edu abishek@mit.edu brantley@mit.edu reidc@mit.edu

MitSci¹
mitsci@mit.edu

Abstract

Recent work in reinforcement learning has heavily relied on a technique for training agents with basal general capabilities in text, images, and other modalities.¹⁻³ However, for many important decision domains such as robotics, video game, and computer vision, publicly available datasets often contain feedback required to train behavioral goals in the same way. We extend the internet-scale pretraining paradigm to important decision domains through semi-supervised imitation learning where agents learn to act by watching unlabeled videos. Specifically, we show that with a small amount of labeled data we can train an inverse dynamics model accurate enough to let a large unlabeled source of video data—here, millions of people playing Minecraft—train which we can then train a general behavioral prior. Exploiting its native hierarchical structure (macc and leighfield et al. 2018), we show that this behavioral prior has substantial general capabilities and that it can be fine-tuned with both imitation learning and reinforcement learning to hard-expert tasks that are impossible to learn from scratch via reinforcement learning. For many tasks we match human-level performance, and we are the first to report a complete agent that can solve difficult tasks, which can take up to hours (speeds of 20 minutes (2400 environment actions) of people to accomplish).

1 Introduction

Recent years have demonstrated the efficacy of pre-training large and general “internet-scale” models on noisy internet-scale datasets for use in downstream tasks in natural language⁴⁻⁶ and computer vision.⁷⁻⁹ For important decision domains (e.g., robotics, game playing, and computer vision) where agents must repeatedly act within an environment, a wealth of data like exists on the web; however, most of this data is in the form of unlabeled videos (i.e., without the actions taken at each frame); making it much less straightforward to train a behavioral prior in these domains than it is in e.g., natural language. In a few rare settings, such as Chess, Go, and StarCraft, there

¹This work was partly funded by grants from the National Science Foundation (NS0904710, NS0904711, and NS0904712) and the Defense Advanced Research Projects Agency (DARPA) under contracts HR0017-09-2-0001 and HR0017-09-2-0002.

²Present

³University of Illinois Urbana-Champaign

already exist large datasets with action labels from various video platforms that researchers have used for imitation learning.¹³ When large labeled datasets do not exist, the canonical strategy for training agents against reinforcement learning (RL),¹⁴ which can be simple infections and complex behavioral-exploitation problems,¹⁵⁻¹⁷ may entail tasks, e.g., navigating vehicles using Photoshop, hacking flights, etc., can be very hard to train with RL and is at the edge of currently available sources of labeled data.^{18,19} In this paper, we seek to extend the paradigm of training large, general-purpose simulation tasks to expand their usefulness by utilizing freely available internet-wide multi-modal video datasets with a simple cross-exploited imitation learning method. We call this method Video Following (VF) and demonstrate its efficacy in the domain of Minecraft.

Existing cross-exploited imitation learning methods are to train with a no-explicit action label; however, they generally rely on the policy's ability to explore the environment throughout training, making them susceptible to explosive bottlenecks.²⁰⁻²² Furthermore, most prior cross-exploited imitation learning work has been in the relatively low-dimensional space; because no experiments with larger data (~700 hours of available video), we hypothesize that we can achieve good performance with a much simpler method, a result that has proven true for performing in other modalities such as text.²³ In particular, given a large but unlabeled dataset, we propose incrementally pre-training by gathering a small amount of labeled data to train an inverse dynamics model (IDM) that predicts the action taken at each timestep in our video. Behavioral cloning (BC) can require a large amount of data because the model must learn to take intent and the distribution over future behaviors from only past observations. In contrast, the inverse dynamics modeling task is simpler because it is one-class, meaning it can look at both past and future frames to take actions. In most settings, environment mechanics are so simple that the health of human behavior that can take place within the environment, supporting that non-exploited IDM could explicitly be able to take the case BC would. Using pre-trained generative IDM, we can train a model to simulate distribution of behavior in the previously unlabeled dataset with standard behavioral cloning methods, which does not require any model robustness and thus does not suffer from any potential explosive bottlenecks in the environment. Finally, we show we can fine-tune the model to interact more fully with other behavioral cloning procedures and learning.

We chose to test our method in Minecraft because: (i) it is one of the most actively played games in the world²⁴ and has a wealth of currently available video data online; (ii) it is a fairly open-ended sandbox game within extremely wide variety of potential things to be built, collected, mining, or crafts more applicable to real-world applications such as computer setup, which also tends to be visual and open-ended; and (iii) it has already pre-trained by the RL community on a specific domain due to its complexity and correspondingly difficult exploitation challenges.²⁵⁻²⁷ In this work, we use the native human interface for Minecraft so that we can (i) most accurately model the human behavior distribution and reduce domain shift between video data and the environment; (2) make this collection easier by allowing the human controller to play the game without notification, and (3) eliminate the need to hand-explicit session interface for models to interact with the environment. This choice means that our model (a IDM) can pre-train and just use a mouse and keyboard interface to interact with human (UI) for collecting, crafting, tooling, trading, etc., including dropping items in specific slots or swapping the recipe book with the mouse cursor (Fig. 1). Compared to prior work in Minecraft that uses a lower dimensional and abstracted policy and training metric,²⁸⁻³⁰ our approach leverages the human interface directly across the environment-exploitation difficulty, making most simple tasks more compatible with RL from scratch. Even the simple task of gathering single wooden log while already being a tree takes 60 consecutive stuck actions with the human interface; meaning the chance for a naive random policy to succeed is ~4%. While this paper shows results in Minecraft only, the VF method is general and could be applicable any domain.

In Section 4 we show that the VF learning method has material two-fold performance: accomplishing tasks impossible to learn with RL alone, such as crafting plants and collecting files (just requiring a human patient in Minecraft a median of 9 iterations ~150 consecutive actions). Through this, we try to build and refine the technology to target more specific behavioral distributions, our agent is able to push even further into the technology tree, solving more tasks



Figure 1: Example Minecraft crafting UI. Agents use the mouse and keyboard to script mouse and keyboard inputs.

(target human actions of 2 minutes or <200 actions). Finally, the *tiny* and *IL* policies for most human improvements are specific to well-known tasks, as expected and as in Microsoft such can pose challenges by using human-like metrics. This requires a policy human action sprawl of 20 minutes or <2000 actions. The two contributions of this work are (1) we are the first to have provided results applying semi-supervised learning to extremely long, noisy, and freely available video data for expert human actions, (2) we show that our proposed tiny *far-tiny* model is quite able to do the task that were otherwise impossible to learn, (3) we show that labeled contact data is the more efficiently used within the *TFI* method than it would be by directly training a baseline solution and (4) we open source our contact data, trained model weights, and Microsoft datasets for future research on learning to act via semi-supervised imitation learning at scale.

1 Preliminaries and Related Work

Imitation learning methods¹⁻⁴ seek to construct a policy that correctly models the distribution of behavior in some dataset $D = \{s_t, a_t\}$, $t \in [1, T]\}$ of state-action pairs. In order to roll out these policies in an environment, they must be causal, meaning they condition on observations from the current timestep t and past timesteps only, i.e., $\pi = \pi(s_t | s_{t-1})$. Imitation learning is simpler than imitation and rollout with corresponding actions. Imitating labeled trajectories have been shown in social robots,⁵⁻⁸ self-driving cars,⁹⁻¹² board games,¹³ and video games.¹⁴

When labeled demonstrations are unavailable, standard behavioral cloning will not work; however, there is a large body of work in imitating behavior from unlabeled demonstrations.¹⁵ For instance, GAIL,¹⁶ constructs an adversarial objective by training the learned policy to exhibit behavior indistinguishable from those in the target dataset. Ghoshal et al.¹⁷ propose to first learn a latent policy using unlabeled demonstrations and then map the learned latent action to real actions with a nonlinear discriminative function. PolicyCF¹⁸ first performs expert methods to find appropriate policies in videos and then train IL agents to match these experts. Similarly, Doshi-Velez et al.¹⁹ and Agapiou et al.²⁰ both self-train to match experts, however they control experts but are restricted to unsupervised learning methods. Pohlen et al.²¹ and Narasimhan et al.²² use policy gradient methods to indicate that above the current state reward expert-provided policy states approach best theoretical local experts. However, it is overruled. Trisch et al.²³ simultaneously train (1) a inverse dynamics model (IDM),²⁴ which aims to recover the underlying action between timestep t and observation of $t+1$ and later timesteps, e.g. $p_{\text{idm}}(a_t | s_t, s_{t+1}, \dots)$ and (2) stochastic control (SC) model to trajectory of observations labeled with the IDM. Data from the IDM is collected by rolling out the SC model in the target environment such that both could improve tandem. However, it may waste training if there are experts in the dataset but the IDM performs poorly on it, it implies that the SC model performs those in similar experts in order to let IDM to improve and correctly label them. Therefore, if the SC model does not update efficiently, it will severely slow down learning. In order to avoid this potential issue we propose a simple two-stage approach to first train a IDM on a small subset of labeled trajectories collected from human controllers (they play the game as well normally with small fine-tunings and noise movements). Because human controllers teach and collect parts of the state space, we can build the IDM fast through the SC training.

Comparing and prior work in semi-supervised imitation learning, we experiment in the much more complex and unstructured environment of Microsoft. Microsoft's crowd-sourced 3D video game dataset, due to its popularity and wide variety of activities, has started a renaissance of IL research.²⁵⁻²⁸ A large body of work focuses on real, robot-scale Microsoft webcams with tasks such as navigation,²⁹ hand placing,³⁰ interaction following,³¹ cooking,³² and others.³³⁻³⁵ Hand grasping with robotic, publicly generated environments of Microsoft itself has included full-grasping,³⁶ annotated collision learning³⁷ and more closely related to the IL experiments presented in this paper, fine-grasping.³⁸⁻⁴⁰ However, to the best of our knowledge, there is no published work that operates in the full, unstructured human action space, which includes day-to-day interactive navigation and item collecting.

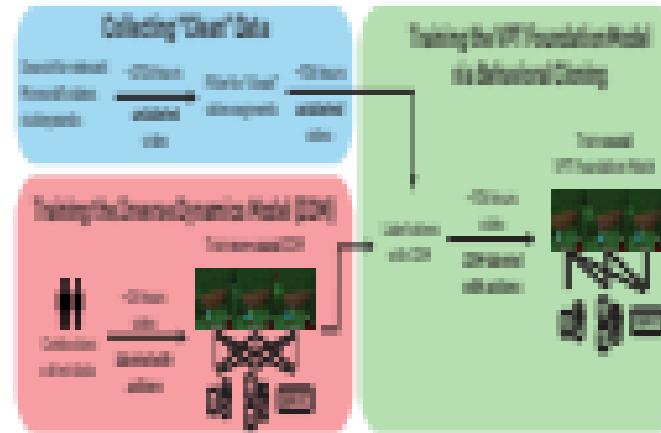


Figure 1: Video Processing/WI Model Overview.

3 Methods

User Dynamics Model (UDM): WI, illustrated in Figure 1, requires us first collect small amounts of labeled metadata data with which to train an inverse dynamics model $\pi_{\text{UDM}}(a|v_1, \dots, v_t)$, which enables us to make the negative log-likelihood of an action at time step t given a history of T observations $v_j = 1 \in [1, T]$. In contrast to a imitation learning policy, the DNN can be untrained, meaning its predictions are a combination of both past and future actions, i.e., a_{t+1} . Compared to the behavioral cloning objective of matching the distribution of human data given past frame only, we hypothesize that encoding environment dynamics is easier and more data efficient to learn. Indeed, Sec. 4.1 will show that the DNN objectives need fewer frames, and furthermore Sec. 4.1 will show that with very little labeled data (as few as 100 frames) we can train a fairly accurate DNN. The DNN can be used to generate values, providing the large amount of data required for the training of behavioral cloning. See Appendices D and E for DNN training and data collection details.

Data Filtering: We gather a large dataset of Minecraft videos by searching the web for related keywords (Appendix A). Online video sites (1) include several actions, such as a writer tool of the player's face, chess, chess tips, mathematics, etc., (2) are collected from platforms where the user interacts with different people, or (3) are from different game modes, e.g., in Minecraft we only want “survival” when players start from scratch and must gather resources to build their items. We call this “bias.” It does not contain visual artifacts and is a low-level task, and could also have “noise.” Through this, sleep, song, and music training songs, a MC mode trained on both modes and other video would likely still perform well in a clean Minecraft environment. Hence, in explicitly accounting negative diversity, we choose to filter out random aspects of video (note that a video may contain both clean and random aspects). We do this by training a model to filter out random aspects using a small dataset (10000 images sampled from online video library) containing no clean or random (Appendix A.2).

WI Foundation Model: We train a baseline model with standard behavioral cloning, i.e., minimizing the negative log-likelihood function predicted by the DNN on clean data. For a particular trajectory of length T we minimize:

$$\min_a \sum_{t=1}^T -\log \pi_{\text{UDM}}(a|v_1, \dots, v_t) \text{ s.t. } a \sim \pi_{\text{UDM}}(a|v_1, \dots, v_t) \quad (1)$$

As we will see in the following section, this model exhibits material two-shot behavior and can be fine-tuned with both imitation learning and RL to perform even more complex tasks.

4 Results

4.1 Performance of the User Dynamics Model

The DNN architecture is composed primarily of a temporal convolutional layer, a ResNet² image processing stack, and residual masked attention layers, from which the DNN simultaneously predicts keypresses and mouse movements (see Appendix D for DNN architecture and training details). Also, hypothesis selection and initial DNN can be trained with a relatively small amount of labeled data. While more data improves both accuracy and keypress prediction, we find

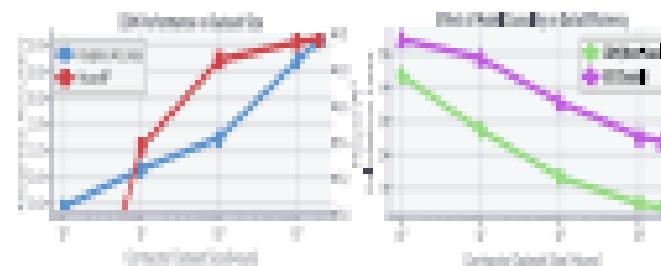


Figure 3: (Left) DNN layers accuracy and cross-entropy loss (negative log-likelihood) as a function of dataset size. (Right) DNN cross-entropy loss vs. efficiency.

DNNs on our 100k base-of-line (equivalent to the ~2M base-of-line data we collected from the internet) achieve 90% layers accuracy and a 0.15 L^2 loss metric evaluated on held-out validation set of another 100k data (Figure 3 left).

Figure 3 right validates our hypothesis that DNNs are less accurate than K-means, likely because clustering assignment makes it easier for training to make full utilization of these features. The DNNs (two orders of magnitude more efficient than a K-means trained on the same data and improve exponentially with more data). This further supports the hypothesis that it is more effective to use customer data within the NFT pipeline by training an DNN than it is to train a baseline model from customer data directly (Sections 4.3 and 4.6 provide additional evidence).

4.2 NFT Foundation Model Training and Cross-Net Performance



Figure 4: (Left) Training and validation loss on the web_base dataset with DNN pre-trained, and loss on the raw DNN customer dataset, which has pre-trained DNN but no self-attention (see text). (Right) Across a given loss we collect the epoch corresponding 200 thousand reward epochs as a function of training epoch, starting with the initial run of the new basic training code (baseline of first green), and adding 100 reward that can be gained when each task. Loss is obtained by equally taking two in the reward, a difficult loss for an RL agent (which we show in Sec. 4.4). Baseline can be called base loss, and taking two called base plus. Collecting reward step is generating QBs, and policy loss takes a factor of 5 seconds (Q-learning action) to take a scaling table.

We now explore the expected behavior based by a theoretical strategy that has extremely large, but noisy, inter-task losses (labeled with no DNN). To collect the related inter-task losses, we recorded the publicly available videos of Minecraft play with conditions such as “mineshaft survival beginner.” These results resulted in a ~2M hours of video, which we filtered down to “bad” video epochs yielding an additional dataset of ~2M hours, which we take to be a web_base dataset (Appendix A has better details on data capping and filtering). We then generated pre-trained DNN web_base with our test DNN (Section 3) and then trained the NFT foundation model with behavioral cloning. Preliminary experiments suggested that we could only benefit from 20 epochs of training and that a 10 billion parameter model was required to stay in the efficient learning regime²⁴ so that training finished (Appendix B), which took ~4 days on 72 V100 GPUs.

We evaluate our model by measuring validation loss (Fig. 4, left) and policy loss on the Minecraft environment. Data is pre-averaged, and we present evaluations to open agents in a standard survival mode game where they play for 10 minutes, i.e. 1200 consecutive actions, and we plot the mean and std of the standard error of the mean for various game statistics such as collecting and collecting items (Fig. 4, right). The NFT foundation model quickly learns to play more items in collecting, which we found too impossible for all agents to achieve with the naive baseline model (Sec. 4.4). It also learns to collect less bags from broken plates and less no-things plates

to collect calling tiles, which are required to collect most often technology in the game and take a human participant in Microsoft’s approach (Sheeran & Chater, 2008) to collect. While these behaviors are likely complex in the entire human action space, the VPT baseline model collects these items at a rate that is over 100% of our baseline collection rate, e.g., on average an contractor will call 5.4 calling tiles in 90 minutes of play versus 4.8 from the baseline model. The model also collects a non-negligible amount of worker tasks, which are required to make worker tools, collect various flowers and materials from them, kill zombies that appear during the night, build settlements, collect various items and materials and eat them, and build pre-generated villages from which to collect various rare items from them. The model also learned to complete several terrain, river, and pillar jumps which enables the agent repeatedly jumping and quickly placing a hill block itself such that it quickly spreads by making spikes.¹⁰

While training and validation loss decrease (though very slowly) (Fig. 4 left), loss on our contractor dataset (which the VPT model does not train on) begins to exceed the VPT model. One reason this could be out-of-distribution because our contractor may have a different distribution of play or because there is some sequential visual bias in difficult compared to tiles from the test set. While we could have expected the available predictor of collecting evaluation performance, we do not see notable gains statistics from the VPT baseline model without (Figure 4, right) decrease over training, and in the test setting we show that BC fine-tuning performance continually improves on the VPT baseline model train. We provide more insight on this analysis process in Appendix B.

4.3 Fine-Tuning with Behavioral Cloning

Baseline models are designed to have global behavior goals and to generally capture across a wide variety of tasks. To incorporate new knowledge or allow them to specialize on a narrow task distribution, it is common practice to fine-tune their model to make more specific choices. The VPT baseline model trained on technological, class historical materializes their performance; it was able to collect a calling tile yet unable to go past this in the technology tree. As a case study into BC fine-tuning, we attempt to improve the VPT baseline model’s ability to collect and call how “body part” items by fine-tuning on two separate datasets targeted at Microsoft behavior within the first few minutes of players starting in a local world. In the first dataset, contractor, Jesus, contractors have 10 minutes to build their base from scratch using mostly wood, sand, and dirt. Gathering contractor data can be difficult and expensive, so we use the contractor dataset we’ve previously reported by casting for video clips with the query for search keywords such as “new world”, “the begin episode”, etc.; this is a subset of the ~1.2m multi-modal data in the IBM. See Appendix B.4 and A.3 for full descriptions of both datasets.

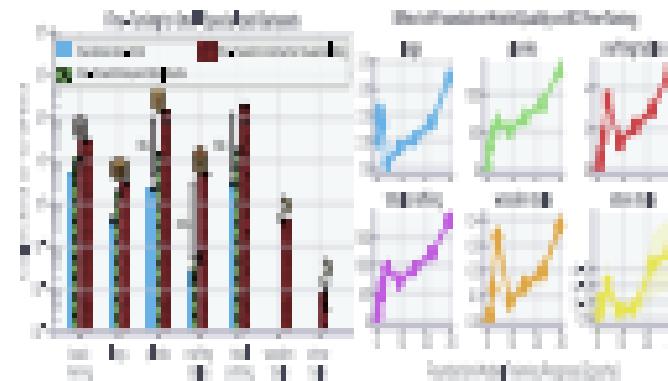


Figure 5: (Left) Collection and calling rates for three policies: the zero-shot VPT baseline model, and the VPT baseline model BC fine-tuned to the contractor Jesus dataset BC fine-tuned to the contractor Jesus dataset BC fine-tuning to other data improves performance, including for the contractor Jesus dataset yielding worker and stone tools. Indeed Microsoft players take a median of 1.2 minutes (20% faster) to construct worker tools and 1.3 minutes (20% faster) to construct stone tools. (Right) Collection and calling rates for VPT baseline model complete throughout training after they are BC fine-tuned to the contractor Jesus dataset. In general, calling/collect behavior increase throughout baseline model training. By 400 steps the other tool items (log, plants, calling tiles, and metal calling).

¹⁰https://www.vision.csail.mit.edu/pubs/ICML2018_gpt2.pdf

The training to set types (jagged tools) in day four compared to the zero-shot baseline model (2% success rate), 13% success rate, 43% success, and 5% success rate overall (Fig. 5). However, when the agent is in the dataset it did not see any saw blades except early in the sequence of cutting skills. We see an even larger improvement when the training is in the contractor dataset (base dataset: 2% success rate), 5%, 36% success rate, 76% success, and 98% success overall. In addition, we see the concept of cutting wooden tools, which requires placing a cutting table on the ground, opening it to creatures cutting interface, and then using it to cut wooden tools. This entire sequence takes approximately human player median of 12 minutes (200 consecutive actions) to accomplish. The model performs substantially faster, which requires a wooden puzzle to open, and cut wood tool, requiring it to open the cutting table (the table requires human player a median of 23 minutes (300 consecutive actions)). We also see the model more frequently cutting villages that suddenly appear in the game; hunting animals for food, in addition to many behaviors were performed by the baseline model.¹⁷

Despite the baseline model's zero-shot offline performance (planning 10 min training (Fig. 4, right), fine-tuning performance does continue to increase throughout baseline model training (Fig. 5 right). Additionally, there is a small difference in performance when training from scratch vs. fine-tuning from the VPT baseline model (Fig. 5 right comparing the left and rightmost points).

4.3 Fine-Tuning with Reinforcement Learning

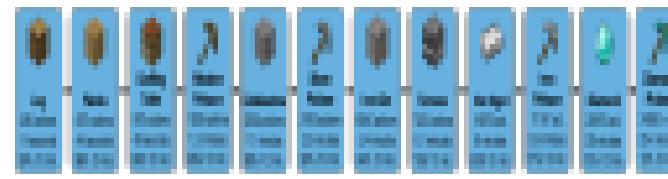


Figure 6: Typical sequence of items for training a diamond puzzle. Below each item is the median time and number of actions required to obtain the item and percentage of contracts that contain the item within 10 minutes. The median time to obtain a diamond puzzle is unknown (except that it is > 10 min); because contracts obtained this item is less than 10% of 10-min epochs.

To demonstrate the efficacy of RL fine-tuning, we chose to challenging and obtaining a diamond puzzle within 10-minutes starting from a zero-shot baseline model. Doing so while incurring a expense of difficult-to-obtain items, but requiring simple skills like mining inventory management, cutting wood and without a cutting table, and so, opening a hammer, and mining at the lowest depth, where many hard-to-reach items will be seen (Fig. 6). Adding to the difficulty, progress can be easily lost by dropping items, destroying items, or dying. Obtaining a diamond puzzle more often than not takes approximately human over 20 minutes (3000 actions).

Ages are provided for each item obtained in the sequence, with lower values in items that are to be collected first and higher values for items near the end of the sequence. Ages are updated with the policy update.¹⁸ RL experiments (~1.1 million epochs, roughly 1 to 10³ times). Episodes last for 10-minutes. See Appendix G.1 for visualization and RL training details. Due to computational constraints, RL experiments use ~10 million parameters VPT model (Appendix B).

An interpretation about the fine-tuning with RL is catastrophic forgetting,¹⁹ because previously learned skills can be forgotten when trained. For instance, while our VPT baseline achieves 100% of the entire sequence of items required to reach the zero-shot, it obtains no example of placing cutting interfaces. It therefore may have some latent ability to recall how to use the sawing projectiles, but has not been performed. To combat the catastrophic forgetting of items with which they can continually improve exploration through RL fine-tuning we allow an auxiliary Label-Shift-Loss (LL) discrepancy loss between the RL model and the human ground truth policy.²⁰

Training from a randomly initialized policy fails to refine about any reward, understanding how far an exploration challenge the diamond puzzle task is from RL to realize human action space (Fig. 7a). The model never learns to reliably collecting, typically the best it may do is obtaining a diamond puzzle (Fig. 7b). RL fine-tuning from the VPT baseline model does significantly better (Fig. 7c), learning everything up to mining iron ore and cutting houses (Fig. 7d). However, this approach is much slower in speed, the next item requires a puzzle to be solved too. Likely

¹⁷<https://www.youtube.com/watch?v=OJfC2X40Bsg&t=25s>

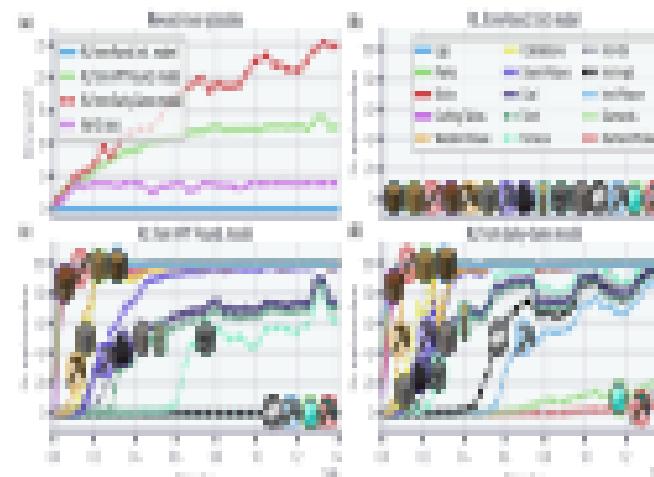


Figure 7: RL fine-tuning results. (a) RL fine-tuning from a randomly initialized model fails to get above any reward. (b) RL fine-tuning from the VPT Foundation model performs substantially better with a reward near 10, and RL fine-tuning from the early-gate model performs best with a reward of 25. When training the early-gate model without RL, loss in the original policy (No RL) is approximately 10000 episodes, suggesting that it will require many more episodes to reach convergence. (c) RL fine-tuning from a randomly initialized model occasionally selects sticks by breaking them (as they are inefficient method of getting sticks that does not require drops or planks) and never loses visibility collecting. (d) RL fine-tuning from the VPT Foundation model learns everything in the environment to collect sticks and broken frames, but fails to learn to use the frame to reach tree logs. (e) RL fine-tuning from the early-gate model learns to obtain 10 broken-level items in the sequence towards a diamond picture and ends a diamond picture in 125 of episodes.

hence the zero-shot probability for the VPT Foundation model ranks as first (logit is 10.0), even when given the pre-synthesized material.

Results further improve by first BC fine-tuning the VPT Foundation Model to the early-gate dataset (the early-gate model, Sec. 4.2) and then fine-tuning with RL (Fig. 7d), which is preliminary evidence we found to perform better than first the tuning to contractor Jason followed by fine-tuning with RL (Appendix E2). The three-phase tuning (pretraining, BC fine-tuning, and then RL fine-tuning) results in learning extremely efficient tasks: it achieves over 80% visibility on sticks, about 20% visibility on collecting diamonds, and 23% visibility on obtaining a diamond picture (Fig. 7d). Furthermore, human players give the objective of obtaining a diamond picture value these items in 70, 100, and 125 of episodes, respectively, meaning our model is human-level for collecting raw materials and mining diamonds. Others have managed to train humans with ~ 4.7 million hours¹⁰ training on Foundation¹¹, but do so with a simplified action space (dropping raw materials). To the best of our knowledge, we are the first to report zero-shot success rates on collecting a diamond picture. Qualitatively, the model developed useful skills for diamond mining, such as efficient mining patterns, core exploitation, attention to previously placed objects like mining tables, and other techniques. Mining woods pictures is hard when carrying an item tool.¹²

Finally, we validate the importance of the RL loss in the pre-trained model during RL fine-tuning. The test set without any RL loss obtains only items early in the sequence (logs, planks, sticks, and mining tables) limiting its reward (Fig. 7e). The hidden hyperparameters limit the sequence to likely frames, while the initial skills of dropping logs and collecting planks are being learned with RL, otherwise with fine-tuning a working pipeline we lost these strategic insights.

4G Data Scaling Properties of the Foundation Model

In this section we validate a core hypothesis behind this work: that it is far more effective to use labeled contractor data to train an RLMM within the VPT method than it is to directly train a BC Foundation model from that same real contractor dataset. If we could cheaply collect a labeled contractor dataset that is similar with the target task as a soft class, then this would be important; however, collecting that scale of data would have cost millions of dollars. Figure 8 compares Foundation models trained on increasing orders of magnitude of datasets. Training with 100-1000 soft-class dataset Foundation models trained up to and including 10 hours are trained with 10M

¹⁰https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2348149/ (last checked 2023-01-10).

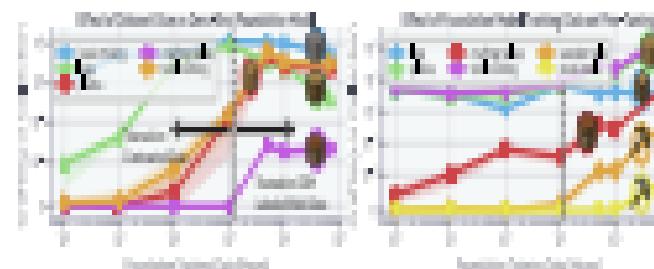


Figure 3: (Left) Cross-validation performance of baseline models trained on varying amounts of data. Models to the left of the dashed line (e.g., points < 50k lines) were trained on contractor data (green) and client data (purple), and models to the right were trained on 500k point-labeled sets of sub-clients. Due to sample limitation, this analysis was performed with smaller (1 million parameter) models except for the final point, which is the 1.5 billion parameter NFT baseline model. (Right) The corresponding performance of each model after BC fine-tuning each model on the contractor Jones dataset.

contractor data, and those trained on 50k lines and above are trained on subsets of sub-client, which does not contain any DNN contractor data. Seeing this implies access to collection, mining, and calling capabilities. The zero-shot model only begins to start calling calling data after 500k lines of training data. When the fine-tuning each baseline model to contractor Jones, we see that calligraphic fine-tuning takes and model has access by rules of nuptials: the average rate: 2.0 lines sub-client lines. We believe only in the concept of calling the tool to be kept data safe.

4.6 Effect of Baseline Model Quality on Basical Calling

This section investigate how downstream BC performance is affected by DNN quality. We train DNNs on increasingly large datasets and see each independently fine-tuned (green, Jayant dataset) to smaller dataset and then to the full (purple) (Fig. 4). We find that all BC could have written out dataset and repeat process starting from scratch, a reflection of DNN contractor dataset (Fig. 4).

DNNs trained on at least 10 lines of data are required for any calling and the calling rate increases quickly up until 10 lines of data, after which there are less to no gain and differences are likely due to noise. Similarly, calling data are only called after 100 millions of DNN lines, and again gains plateau the 100 lines. Which all points experiments we see varied DNN trained on 100k lines of data, since results suggest to call in the first number is about 100 lines.



Figure 4: Zero-shot performance of BC models trained on the sub-client (Jayant) dataset trained with DNNs that were trained on increasing amounts of contractor data.

5 Discussion and Conclusion

The results presented in this paper clearly provide path to utilizing the wealth of unlabeled data in the real-world expert decision-making. Compared to generic rule-making or extensive methods that would only yield representative rules, NFT allows the mining possibility of directly learning to automatically generate and store these learned learned prior, a extremely effective acquisition prior for RL. TPI could even be a better general representation learning mechanism when the downstream task is not learning to write function rules—for example, the ability to explain what is happening our rule—because equally the most important information in any given scene would present in features trained to correctly predict the function over later human action. We leave this intriguing direction to future work.

Future work: and could improve results with more data (we estimate we could collect 10M lines) and larger, better-trained models. Furthermore, all the models in this and previous experiments are only relevant to one sub-client and is specific to their specific tasks. Appendix I presents preliminary experiments investigating our models on closed captions (text transcripts of specific videos), showing they

more readily available; we believe this will facilitate further research. Also, how can we most effectively combine will investment evaluation metrics (Sec. 4.2) with other model properties outlined here? Another fruitful future direction would be to investigate the correlation between return trading metrics and investment evaluation. Finally, while we do not anticipate any direct negative social impacts from the methodology in this work, as TPI improves and expands to other domains it will be important to assess and mitigate harm that may result after the loss of privacy or personal datasets, such as enabling inappropriate behavior.²⁷

In conclusion, TPI extends the paradigm of learning key and general property behaviors from freely available internet-scale data to a general finance domain. Our model exhibited impressive two-class behavior and, when fine-tuned with EC, achieved an unprecedented result of calling a financial pipeline in Microsoft (all the time Microsoft gives the human intuition). We further showed that outcome data is far better used within the TPI pipeline than to train a classifier model directly and that only a small amount of outcome data (about 20M ES) was required while maintaining an array of additional valuable data (see Sec. 4.2). Finally, learning with the human hybrid and neural interface is highly generalized across industry settings for more sophisticated human behaviors. While we only experiment in Microsoft, we believe that TPI provides a general recipe for training behavioral agents in hard, yet generic, action spaces in many domains that have a large amount of freely available unlabeled data, such as computing usage.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models as few-shot learners. *Advances in neural information processing systems*, 35:367–380, 2022.
- [2] Junhong Chen, Ming-Wei Chang, Xunyu Li, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Yuhao Fan, Myrto Maniatis-Tay, Nandkishore Deep Das, Jian Li, Michael Lewis, Luke Vilnis, and Noam Slonim. Robert: A robustly optimized fast pretraining approach. *arXiv preprint arXiv:1907.11601*, 2019.
- [4] Christopher J. O’Hearn, Sami Ben Romdhane, Katherine Lee, Sharan Narayan, Michael Nitkin, Yang Zhou, Wei Li, and Pratej Li. Explaining the limits of transfer learning with masked language models. *arXiv preprint arXiv:1910.14632*, 2019.
- [5] Dario Kraljević, Bojan Šešić, Nenad Krstović, Luka Čirjak, Mirko Palić, Vojin Lić, Adrijan Blagojević, and Lazar Vučić-Milutinović. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pages 11–26, 2021.
- [6] Alex Radford, Jong-Wook Kim, Cenili Bailey, Ariana Barone, Gilad Gol, Sathya Jagadish, Girish Sastry, Amanda Askell, Peter Molino, Jack Clark, et al. Llama: too much fine-tuning hurts performance. *International Conference on Machine Learning*, pages 1034–1051, PMLR, 2021.
- [7] Ross Bougares, Drew A. Hodson, Ethan Hirsh, Ross Allen, Simon Amos, Sydney Wu An, Michael F. Bernstein, Jameson Bell, Austin Bowditch, Emma Brunskill, et al. On the opportunities and risks of learning models. *arXiv preprint arXiv:1906.07313*, 2019.
- [8] Yoshua Bengio, Dumitru Erhan, Yoshua Bengio, and Karim Fréchette. Scaling feature extraction. *CiCP*, abs/2016/0261/321, 2021. <https://arxiv.org/abs/2016/0261>.
- [9] David Silver, Alp Tüzün, Chris Daupe, John Cox, Laurent Sifre, George Van Den Driessche, Marc Pilatowski, Iman Alaviyagh, Yash Pannirselvam, Marc Lanctot, and Timothy Lillicrap. OpenSpiel: a python library for reinforcement learning. *Nature*, 529(760):484–489, 2016.
- [10] David Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Ó Domhnaill, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Nature*, 529(760):59–65, 2016.

- [1] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An introduction*. MIT press, 2018.
- [2] Christopher Beal, Guy Lachman, Baekil Cho, Yuki Ono, Paragmoy Dabir, Craig Denison, David Fohi, Qiang Fu, Ming He, Ming Hu, Eric Hsu, and Ben Shafrazi. Multi-player reinforcement learning. *arXiv preprint arXiv:1912.06881*, 2019.
- [3] Bruce Bonet, James Lederer, Tom Moller, T. N. Quan, Paul Riedl, Ben Shafrazi, and Ilya Shvarts. Emergent role for multi-agent interaction. *arXiv preprint arXiv:1909.07221*, 2019.
- [4] Max Jaderberg, Wojciech M. Czarnecki, Ludwig Dalle Molle, Guy Lever, Timo Mann, George van den Oord, Charles Blundell, Neil C. Heess, Tom E. Erez, Avradip Saha, et al. Human-level performance in 11 atari games with population-based reinforcement learning. *Nature*, 524(7567):529–535, 2015.
- [5] Alexey Radulovich Balu, Rui Hu, Steven Karpov, Kishore Namachchivaya, Alex Vostrikov, Zhezong David Guo, and Charles Blundell. Agent2: Cooperating for continuous learning. In *International Conference on Machine Learning*, pages 57–65. PMLR, 2019.
- [6] Max Jaderberg, Simon Osband, George Tsewang, Tim Salim, David Silver, and Karen Simola. Unifying model-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 31, 2018.
- [7] Yen Chen, Harrison Shum, Anna Shrestha, and Oleg三人. Exploration by random action. *arXiv preprint arXiv:2001.1394*, 2020.
- [8] Alireza Fard, Ismail Hacioglu, Ismail Ustunay, Kenneth O’Hearn, and Jeff Donahue. Iterative refinement. *NIPS*, 2018, 158–166, 2018.
- [9] Paul Chouleye, Ismail Hacioglu, Illy Hahn, George Tsewang, Lubin Chouy, Koen Meers, Jonathon Hare, Eddie Geng, Alex Gold, Alex Pentland, et al. A baseline approach for learning social games. *arXiv preprint arXiv:2004.05732*, 2020.
- [10] Taehan Kim, Andrey Krapush, Linna Fa, Justin Brondum, and Perry Liang. World of Warcraft: An open-domain platform for web-based agents. In Doina Precup and Yee-Wing Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1135–1144. PMLR, 16–21 Aug 2017. URL <https://proceedings.mlr.press/v70/kim17a.html>.
- [11] Andrew TNg, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *IJCAI*, volume 1, page 1, 2000.
- [12] Finn Turpin, Gauthier Gidel, and Peter Stone. Deep abstractions in infinite horizon learning from observation. *arXiv preprint arXiv:1805.13966*, 2018.
- [13] Jonathan Shulman and Stefanos Zafeiriou. Generative adversarial imitation learning: Abstraction and information processing systems, 31, 2018.
- [14] Finn Turpin, Gauthier Gidel, and Peter Stone. Reinforced learning from observation. *arXiv preprint arXiv:1811.00544*, 2018.
- [15] William Lin, Shubham Gupta, Peter Abbeel, and Sergey Levine. Unsupervised domain adaptation: Learning cross-domain distributions via adversarial registration. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.
- [16] Dominic Bell. Most played game in 2021, voted by paid community players. URL <https://fivethirtyeight.com/2022/01/most-played-game-in-2021-voted-by-paid-community-players/>.
- [17] William D. Gees, Brandon Bogatin, Nicholas Topp, Billy Rose, Caylor Gold, Michael Wilson, and Nathan Saldanha. Most A long calculation of microtubule rotation. *arXiv preprint arXiv:1907.11440*, 2019.

- [2] Omer Tsoori, Shlomo Groves, Tomi Zilberman, David Mandelzon, and Shai Ben-Zvi. A deep hierarchical approach to lifelong learning in neural net. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [3] Christian Schaffter, David Scherer, and Michael Veerle. Simple efficient reinforcement learning through learning loss decomposition is robust. In *NIPS 2016 Computer and Decision Science Track*, pages 67–76. PMLR, 2016.
- [4] Hugo Larochelle, Joelle Pineau, David Precup, William Vaughan, Balaji Lakshminarayanan, Raftis Sutskever, Peter Abbeel, Doina Precup, Alex Szepesvári, Fei Tang, et al. Multi-task episodic learning in complex, real-life applications from scratch. *arXiv preprint arXiv:2001.14076*, 2020.
- [5] António M. Vilas-Boas, Domingos, Braga, et al. Context-awareness, active perception, and action in mobile. In *International Conference on Mobile Learning*, pages 279–279. PMLR, 2016.
- [6] Ying Peng, Maksim Shatskikh, Maria-Cristina Diaz, Tatiana Dvurechenskaya, David M. Rhee, Marcos Brondum, Jocelyn Noyola-Molina, and Scott Beibman. Map-aware learning from live annotations by visual collaboration. *arXiv preprint arXiv:2001.04008*, 2020.
- [7] Alenay Szymek, Abby Stavrou, Brad Stiggle, Emil Kostov, Yuliya Savenkov, Yuliya Savenkov, and Arashashir Pasztor. Exploit-exploit: reply-a hierarchical reinforcement learning loss decomposition. *arXiv preprint arXiv:2006.09951*, 2020.
- [8] Ethan Xu, Jiaoyang Li, Junyu Yu, Jiaoyang Yu, Qiang Yu, and Wei Tang. Iterative: Playing chess with single-client federated reinforcement learning. *arXiv preprint arXiv:2102.14077*, 2021.
- [9] David Pritchard. When Anesthetics last which is a good effect. *Anesthesia and information processing system*, 1, 1991.
- [10] Steve Small. Is iterative learning the route to learned what? *Task in cognition*, 34(2):25–32, 1999.
- [11] Bruno D’Alessandro, Luca D’Amato, Massimo Viggiani, and Renzo Ruggiero. Learning of robot learning loss decomposition. *Robotics and computer systems*, 57:46–61, 2014.
- [12] Alessandro Giordani, Michael Nesterov, Eyal Geva, and Omri Ziv. Iterative learning: A survey of learning methods. *ACM Computing Survey (CSUR)*, 50(3):1–35, 2017.
- [13] Clark Glymour, Scott Flanagin, Dan Klahr, and Donald Nisbett. Learning to be like. In *Machine Learning Proceedings 1992*, pages 35–36. Elsevier, 1992.
- [14] Alireza Gholami, Hamed Ghasemi, Farid Gheysari, Farid Gheysari, Parham Fattah, Marzieh Farokhi, Christian Freude, Hugo Schmidhuber, Gauri D’Cunha, et al. A machine learning approach to visual perception of the visual scene outside. *IEEE Robotics and Automation Letters*, 10(2):63–67, 2021.
- [15] Maciej Chojnicki, Bartłomiej Laski, Dawid Dziedzic, Bartłomiej Laski, Paweł Czajka, Paweł Czajka, Łukasz Włodarczyk, Małgorzata Laski, Małgorzata Laski, Józef Tęczak, Józef Tęczak, and Bartłomiej Laski. Iterated learning in self-driving cars. *arXiv preprint arXiv:1904.07314*, 2019.
- [16] Felipe Galván, Matías Milic, Antonio López, Valerio Valera, and Alvaro Barrientos. End-to-end learning via combined imitation learning. In *2018 IEEE International conference on robotics and automation (ICRA)*, pages 485–490. IEEE, 2018.
- [17] Niels Drotz. Computing “the entropy” of state patterns in the ground truth. *ICDAR journal*, 30(1):39–50, 2017.
- [18] Tomi Zilberman, Omer Tsoori, Shlomo Groves, Michael Mandelzon, Tomi Zilberman, Omer Tsoori, Shlomo Groves, Michael Mandelzon, Tomi Zilberman, Shlomo Groves, Michael Mandelzon, and Shai Ben-Zvi. Deep learning loss decomposition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

- [45] Aditya Elgondi, Renuka Iyer, Daniel Schuster, and Charles Weil. Unifying latent policy functionals. In *International conference on machine learning*, pages 175–182. PMLR, 2011.
- [46] Yuxin Fang, Jiyang Liawala, Shreyas Malhotra, Peter Abbeel, and Sergey Levine. Self-reinforced learning of physical skills from videos. *ICML workshop on Robotics*, 10(2), 2018.
- [47] Przemysław Błaszczyk, Grzegorz Stach, Xi Chen, Naiqi Ren, Sotiris Loukas, Giacomo Cipolla, Dan Gutfreund, Sagnik Paul, Paul A. Hertz, Ivan Viola, et al. Learning handimanagement in the wild. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 75–80. IEEE, 2019.
- [48] Yann LeCun, Yoshua Bengio, David Hinton, Thien-Hau Pham, Zhi Wang, and Yoshua Bengio. Playing backgammon games by watching experts. *Advances in neural information processing systems*, 31, 2018.
- [49] Dejal Patel, Paweł Włodarczyk, Guangyu Lin, Pali Agrawal, Xi Chen, Yit-Sian Tan, Paul Schulman, Shreyas Malhotra, Alex A. Aja, and Peter Abbeel. Toward visual imitation. In *ICLR*, 2018.
- [50] Alvin Sun, Dan Chen, Pali Agrawal, Pali Agrawal, Peter Abbeel, Shreyas Malhotra, and Sergey Levine. Distilling self-supervised learning of motion for visual handimanagement. pages 346–355. ICML, 2019. id: 111992A201907002.
- [51] Duy Nguyen-Huu, Jia Pan, Mathias Sojka, and Bernt Schiele. Learning more dynamic responses. *Indepth report on visual visual search*, number 0196, 2018.
- [52] David Held, Arsh Agrawal, François Fleuret, Riley Endresen, and Robert E. Mahon. Epilepsy patient learning to ride a car and learning to speak German. *arXiv preprint arXiv:1803.04010*, 2018.
- [53] Duy Nguyen-Huu, Jia Li, Sojka Sehn, and Michael Ullman. Deep reinforcement learning from policy-dependent human feedback. *arXiv preprint arXiv:1803.04377*, 2018.
- [54] Alexander Tsvetkov, Stephen Zhang, Clinton Wang, and Richard Socher. Keeping you honest: Solving sparse reward tasks using self-balancing depth rewards. *Advances in Neural Information Processing Systems*, 32, 2019.
- [55] Stephen Zhang. Deep reinforcement learning with soft boundary constraints to avoid overconfid. *arXiv preprint arXiv:1802.08456*, 2018.
- [56] Ernesto Iglesias, Luis Rosales, and Shin-Yung Lee. Fighting policies in neural reinforcement learning. Technical report, Technical report, Technical report, Stanford University, 2018.
- [57] Samira Das, Clinton Wang, and Richard Socher. Hierarchical step policy: Efficient learning in multi-task reinforcement learning. *arXiv preprint arXiv:1712.07944*, 2017.
- [58] Duy Nguyen-Huu, Sankar Subramanian, and Peter Abbeel. Iterated task generalization with multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pages 363–371. PMLR, 2017.
- [59] Duy Nguyen-Huu, Shreyas Malhotra, and Peter Abbeel. Learning to create an admissible problem. *arXiv preprint arXiv:204.05731*, 2020.
- [60] Timothy Nicaise, Arsalane Chou, Tom Cramer, and Olli Simula. Task-adaptive curriculum learning. *EETI transactions on neural networks and learning systems*, 31(9):752–764, 2019.
- [61] Robert Gray, Douglas Red, James Ross Bell, and Priscilla and practices of statistics. *Principles and practices of statistics*, 1996.

- [6] Liangyu Si, Dongya Zhang, Shengyu Liu, and Jun Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 73–79, 2015.
- [7] David Lopis, Sean McDaniel, Tom Rognes, Ben J. Brown, Benjamin Chen, Steven D. Bell, Scott Gray, Michael E. McCoy, Wu, and David Andric. Scaling law for machine learning models. *arXiv preprint arXiv:2001.09613*, 2020.
- [8] Karl N. Leibniz, Jacob Elieser, Peter Glavin, and Mike Schlosser. Basic policy gradient. In Maria Nahm and Ying Tang, editors, *Proceeding of the 20th International Conference on Machine Learning*, volume 129 of *Proceeding of Machine Learning Research*, pages 203–212 PMLR, 11–14 Jul 2013. URL <https://proceedings.mlr.press/v129/leibniz20a.html>.
- [9] Dimna Ladidipati, María-Ayuso-Sanz, Isabela Ribeiro, Mário Sáto, Douglas Radcliffe, José Miguel, Pedro Mira, Suryakumar Raju, Niall Corry, Jeff Cen, et al. Biological entropy in lifelong learning machines. *Neuro Machine Intelligence*, 4(3):26–34, 2021.
- [10] Jason Karpowicz, Karina Pocino, Scott Rothkopf, Joel Yau, Galina Bogatin, Antti Virolainen, Daniela Diaz, Wei Qian, Troy Daniels, Agustín Gómez-Benito, et al. Decreasing synaptic latency in neural networks. *Proceedings of national academy of sciences*, 114(2):523–528, 2017.
- [11] Emily M. Bender, Tomi Kain, Anya W. Miller-Rizzo, and Shreyas Chaitanya. On the dangers of stochastic points: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 624–635, 2021.
- [12] Fabio Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Rémi Weiss, Trevor Oliphant, et al. Scikit-learn: Machine learning in python. *The Journal of machine Learning research*, 12:2825–2830, 2011.
- [13] Aditi Parikh, Naman Dua, Shikhar Maji, Nitish Khurana, Umar Khan, Akash Srivastava, Nitish Patel, and Milind Pai. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [14] Jimmy Le Viêt, Jacob Baudéon, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [15] Yann A. LeCun, Leon Bottou, Genevieve Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 1–48. Springer, 2012.
- [16] Dilek Hacıoglu and Liqiang Yu. Hierarchical stochastic optimization. *arXiv preprint arXiv:1412.6901*, 2014.
- [17] Alan Puth, Ian Cox, Francisco Flores, Alan Lee, James Radford, Gregory Charvat, Trevor Kilian, Zengin Yuksel, Nadeem Ghaffari, Luis Antón, Alan Donkin, John Kop, Shuai Tang, Zicheng De-Pin, Martin Raison, Alysha Trott, Senthil Chandrasekaran, Rohit Srivastava, Li-Feng Jing, and Senthil Chandrasekaran. PyTorch: An imperative style, high-performance deep learning library. In H. Wallach, R. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8023–8033. Curran Associates, Inc., 2019. URL <https://paperswithcode.com/paper/8023-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [18] Yann N. Dauphin, Karissa Petersen, Capil Correia, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in neural information processing systems*, 27, 2014.
- [19] Ivan Trifunovic, Jeff Cen, Tolka Rengaraju, and Balaji Lakshminarayanan. The vanishing gradient in deep neural networks! *Advances in neural information processing systems*, 32, 2019.

- [16] Zenglin Zhu Tang, Tong Tang, Jian Cheng, Qun Wu, and Huaishan Zhang. Reinforcement learning based policy search based on local graph search. *arXiv preprint arXiv:1901.02901*, 2019.
- [17] Michael A. Bowling, Eric Blasius, Pauline Durance, Alex Koller, and Guy Van den Broek. Neural policy optimization algorithm. *arXiv preprint arXiv:1903.04731*, 2019.
- [18] Alan Shrivastava, Policy Matrix, Sajay Lotia, Michael Mathis, and Peter Abbeel. High-dimensional actions control using predicted strategy selection. *arXiv preprint arXiv:1904.02011*, 2019.
- [19] Ronald Williams. Simplified statistical policies: theory for stochastic reinforcement learning. *Machine learning*, 10(3):229–256, 1992.
- [20] Volodymyr Mnih, Adriaan Pritzel, Mirza Micić, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1123–1131. PMLR, 2016.
- [21] Maria Andrychowicz, Ilja Nemenman, Krzysztof Jozefowicz, Radford H. M. Paszke, Bartłomiej Wąs, Mikołaj Bober, Mikołaj Wierwucik, Maciej Mieczkowski, Timothy Lillicrap, David Silver, and Krzysztof Kowalewski. Asynchronous methods for deep reinforcement learning. In *Advances in neural information processing systems*, 30, 2017.
- [22] Tom Schaul, David Silver, Karl Gregor, and David Silver. Universal value function approximators. In *International conference on machine learning*, pages 1811–1819. PMLR, 2015.
- [23] Open Esports League Team, Adam Stach, Amy Kheiraj, Christian Gross, Charlie Beck, David Bass, Matt Segalow, Rajiv Tereja, Mac McElroy, Michael Mathis, and Spencer Ladd. Reinforcement learning leads to precisely adaptive agents. *arXiv preprint arXiv:2007.03988*, 2020.
- [24] Max Laskin, Xiangyu Zhou, Gregory Frosini, Michael Freitas, Jacob Vossler, Shuang Goh, Scott Wilson, and Tim Rockford. A survey of reinforcement learning informed by robotics research. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 639–647. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/80. URL: <https://doi.org/10.24963/ijcai.2019/80>.
- [25] DeepMind Research Agents Team, Isidore Bonnefond, Ann Haug, Arthur Bruce, Fabien Geurts, May Cen, John Frandsen, Peter Granger, Michael He, Tim Harley, and Daniel Whiteside. Reinforcement learning agents with imitation and self-supervised learning. *arXiv preprint arXiv:2012.07021*, 2020.
- [26] Alpár Renzeli, Pauline Durance, Alex Koller, Guy Van den Broek, and Mati Den. Hierarchical reinforcement learning: precision with big batch. *arXiv preprint arXiv:2004.01211*, 2020.
- [27] Dong Yu and Li Deng. Automatic speech recognition, volume 1. Springer, 2016.
- [28] DataBiosphere. *spack*. May 2020. URL: <https://github.com/llnl/spack>. Accessed 2020-04-22.
- [29] Aravind Srinivasan, Bo Li, Rui Pan, Alex Koller, Jose Mikel Ria, Jerry Liang, Qingqiang Yan, Michael Tzou, Deyi Wang, Xian Bai, Zhen Bai, and Yuxin Chen. Learning to control by contrastive pre-training. *arXiv preprint arXiv:2001.06615*, 2020.

Acknowledgements

We thank the following people for helpful discussions and support: Bob Freitas, Ken Stach, Isidore Bonnefond, Rajiv Tereja, Wojciech Laski, Ignacio Santesteban, David Freitas, Sean Freitas, Jonathan Gordon, and the OpenAI supercomputing team, especially Christian Olsson, Ben Oord, and Christopher Brown.

Supplementary Information

A. Collecting Internet Data

A.1 Initial Video Dataset Creation

Our pipeline creates a video dataset of Microsoft employee live tv recordings from 2014 until 2016. Additionally, we prefer the data over from past years as close as possible to our evaluation environment, meaning publicly available live Microsoft videos (i.e. being on a computer) (which uses a more natural interface, with open windows with live feeds and other buttons, being single (i.e. multi-) player, and having the default look of the past year), modifications that alter interface, such as to make it less realistic. To try to avoid bias, first, we collect a randomly performing keyword search of publicly available videos on the internet. A list of search queries tested are given in Table 1.

| |
|-------------------------------------|
| msft annual budget |
| msft yearly revenue |
| msft yearly revenue total |
| msft annual total |
| msft annual profit |
| msft annual last play |
| msft annual live segments |
| msft live segments |
| msft annual state profit |
| msft annual profit 2016 |
| msft live in that year annual total |
| msft annual total last |
| msft annual last play episode 1 |
| last play annual episode 1 |
| msft annual last |
| msft annual having to play |
| how to play msft annual |
| how to play msft annual |
| msft annual last |
| msft annual for sales |
| msft annual for losses |
| how to play msft live segments |
| msft annual total wins |
| msft annual new total |
| msft annual new beginning |
| msft annual episode 1 |
| msft annual 2014 1 |
| msft annual 2016 |
| last play msft annual total |

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

For videos that have segments available, we perform an additional step of automatically filtering technician videos that do not have segment information. In this step, we look for a list of Microsoft keywords in the video title and description and ignore videos that contain those terms. The Microsoft keywords we use are: (pol, pol, pol, the, Bl, playdate, finale, multiple, annual, pc, point editor, d3d12, enable msft, how to install, how to download, released, msft),. This process yields a ~70% based validation loss, which we filter down to only a “clean” subset as described in the next section.

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

We restrict the scope of this work to the Microsoft Surface year books and therefore limit our training dataset to clips that are obtained from the books that are relatively free from visual artifacts.

Therefore, we asked annotators to label a set of videos with faces (image) from Microsoft videos (MSVTT). These images are from a random subset of the videos we collected toward the beginning of the project (Section A).

A.2.1 Label Collection

Microsoft workers on Amazon Mechanical Turk (mTurk) were selected with a simple qualification task to find certain scene capture images in front of existing tv shows. A sample worker interface that the workers are provided is given in Figure H.

Revised workers to label videos as being on screen following three categories (see Figure 11 for visual examples of each class):

1. **Research: Neutral Face - No attributes:** Video frame (image) that correspond to the Microsoft neutral pose task and contains no pose visualization (e.g., subtitles, character, character, advertisement, picture-in-picture of the scene, etc.).
2. **Research: Neutral Face - with attributes:** Video frame (image) of the Microsoft neutral pose task that includes visual entities.
3. **Face of the Show:** Video frame (image) that are not from the Microsoft neutral pose task, including those from the Microsoft pose tasks and a creative task or consider yourself neutral.

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

Please help us identify screenshots that belong only to the neutral task in Microsoft. Everything else (Microsoft neutral task, other poses, music video, etc.) should not be a focus of the show. Neutral task is identified by the lack of the faces after screen:

- always face (no face)
- always face (one or more faces visible)
- showing item held

Neutral task:

Only neutral task video have buildings has either one below a fraction of the screen.
Creative task:

Creative task only has a item below and doesn't include a focus of the show.

Label Descriptions:

- **Research: Neutral Face - No attributes:** The image will be due screenshots from the Microsoft neutral pose people without any specific action.
- **Research: Neutral Face - with attributes:** The image will be addressed with screens, but will not include others. Typical entities may include image, video, play button, text annotation, picture as depicted by player, etc.
- **Face of the Show:** The first category when the image is not a valid Microsoft neutral screenshot. It may have no faces and focus on a different pose task. In some cases pose tasks such as the creative task, the buildings task will be missing from the image; the last letter say or may not be present.

In total, we paid \$29.95 on human labeling experiments as a total, of which \$29.95 was directly paid to workers. The remaining amount was spent towards Amazon platform fees. The workers earned \$0.01 per labeled image, at an hourly compensation of \$0.24 based on a estimated labeling time of 1.5 seconds. As an internal sample test of the annotation, we found the average labeling time to be < 1 second.

Since no platforms provide automatic functionality of this, we located it in AI form, where we need sample images to be labeled, owing which annotators take care of this.

low risk and no risk images were detected during our initial trials. We only collected high-risk images during our experiment, and for which user fully assigned it to failed plates; further no generally identifiable information (PI) was added.



Figure 11: Amazon Mechanical Turk worker interface showing a sample image.

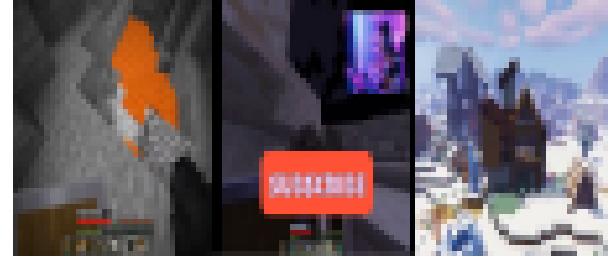


Figure 12: (Left) Sample image for Class 1: Minecraft Survival Mode - No Kitchen. (Middle) Sample image for Class 2: Minecraft Survival Mode - with Kitchen - image contains scratches and noise as poster of the website. (Right) Sample image for Class 3: One of the items - image showing kitchen as well as health and armor bars, indicating that item was captured during survival mode gameplay.

A.2.2 SVM Training

While images which collected as described in the previous section are not included in the final video dataset, but exist in those from the Minecraft: Survival Mode - No Kitchen category, given a set of labeled images, we train a model for each image using the LIBSVM `LibLinearCLPModel`.¹ This is a linear-based CLP model that is scaled up to have approximately 10x the capacity of a `Logistic`. We then train a Support Vector Machine (SVM) using the LIBSVM kernel to obtain a linear classifier. We use the `Scikit-Learn`² SVM implementation with the parameter configuration given in Table 2.

Finally, we apply the classifier to frames of an video sequence to create the final dataset. We filter the frames that consist of at least 80% 'kitchen' frames at this stage (Class 1 consists of Survival Mode - with Kitchen and 80% of the items are both considered clean). From the rest, we apply another filter (with a kernel size of 7) to the labels and expand videos by splitting the 'kitchen' frames that are at least 5% in duration. The result of this is our final `clean` dataset.

A.2.3 Only_grey Dataset

The `only_grey` dataset has about 100,000 short video clips from a 'only grey' Minecraft behavior, i.e. instances where players stand in their world with no items. We obtain the metadata that accompanies the video in `only_grey_clean` dataset via the following regular expression match:

| | |
|----------------------------|-------------------|
| CDF Model Specification | 2020-04 (version) |
| CDF Language Definition | 460-461 |
| CDF Encoding Format Length | 104 |
| SMI Parameters | None |
| SMI Parameters | None |
| Sample Size | None |
| Sample Size | None |
| Sample Size | None |

Table 1: Frame Extraction Details and SMI Configuration. The parameters are for the SMI implementation in `MediaBox`¹⁰.

- `(skip_start, skip_end, clip_start, clip_end)`
- `(skip_start, skip_end, clip_start, clip_end, frame_start, frame_end)`
- `clip`
- `beginning`
- `(start_timestamp, end_timestamp)`
- `from scratch`

From this set of videos, we take only the first 1 minute of each video.

II. Contracting Data

II.1 Recording Contractor Pay

Our contractors receive their full reward for what they work their action and provide back to the play. The rewards is implemented using the MC3-Rewards (`rewards.Rewards`) MC3-Rewards recording package. To ensure that the rewards environment is as close as possible to the Microsoft environment used for RL training and evaluation (Appendix C), we use the same underlying game engine for both. The rewards is a language that runs in a windows mode, with constant resolution of 1280x720. Rewards is set to tell the ‘playout’ setting in Microsoft, which is the default setting. Other specific settings (field size, ODE resolution) specific to rewards used in the Microsoft environment (C.1), we explicitly present here from playout specific settings. Unlike the environment, the rewards allows all keyboard key presses and continuous (as opposed to binary) mouse actions. The very first step in “Set” the base buffer used to display the game window is located in the `446380` and written into a video file. In-game actions are recorded in a separate `BSML`. We extract file video and from a `BSML` (contracting). All recordings are divided into 1 minute steps: after each 5 minute segment of contractor play, the rewards automatically uploads the video file, the `BSML` file with actions, as well as a Microsoft state file. To ensure that contractors cannot corrupt each other’s data, we provide every contractor with an individual download link, a unique identification string with access only to themselves. Contractors also included subjective subjective score (e.g., `group-wins-losses-deposits`) generated with the `mc3agent` python package to ease contractor monitoring when we publish the file.

II.2 Contractor Content

We recruit contractors by posting the following file on the UpWork freelancing platform.

“We are collecting data for training AI models in Microsoft. You’ll need to install your, download the modified version of Microsoft (that collects and uploads your play data), and play Microsoft environments! Paid per hour of gameplay. Price depends on Microsoft set currency. We don’t collect any data that is unrelated to Microsoft how you complete.”

We built applications open to public, and then randomly selected 10 applicants for the list of contractors. Later in the project, we've decided more than 1000 user contracts would be terminated, so called new applicants from the original pool and the client has to manually review contractors. The contractors were paid \$21 per hour (minus OpenShift platform fees and application fees). All of the costs presented in this paper are based on about 4.5M hours of time (including time needed to gather statistics of game play that we can used for training), which cost us around \$90,000. Over the course of the project, we collected some data we didn't use for the to-buy-in-the-marketplace and for which it was ultimately did not prove. Instead, we spent about \$500k for contractor compensation over the course of the project. However, as we discuss in Sec. 4.6, we could likely obtain most of our results with a \$100k budget using only 1000 worth of data (i.e., the baseline MDP model, EC file-timing in the each game, Javacard dataset, and the EC file-timing results). Collecting the contractor Javacard dataset cost about \$8000. Because we used the EC dataset in about 200 hours of contractor data, the total cost of contractor data in the research was around \$16000.

In early stages of the project, we were planning to use contractors data only for the purpose of training the DNN. In such approach tasks were given after that "play the survival mode of Minecraft like you normally would". Later in the project, we expected that contractor to perform specific tasks in Minecraft, such as:

- Collect a noisy multi-task puzzle; using only wood or stone tools (troweling)
- Start a new world every 30 minutes of game play
- Build a brick house in 30 minutes using only dirt, wood, sand, and other tools or stone tools (contractor Javacard, more details below in Appendix B.4).
- Starting from a new world and an empty inventory, find resources and craft a diamond pickaxe in 30 minutes (dirt, stone, diamond, pickaxe). This dataset was used to obtain statistics for how long it takes human on average to complete the task (and the individual requires complete it) when obtaining a diamond pickaxe in their pad.

Since we only recorded in-game events and videos, the data does not include personally identifiable information. That being said, contractors could potentially use Minecraft's open-world property to generate personally identifiable information under different contexts (e.g., by using Minecraft blocks to write their names or offensive messages, then taking a screenshot while the message will be visible). In practice, we have never seen any attempts to do so in the contractor videos that we watched. Of course, we trained EC models on videos from thousands of people playing Minecraft, and if such behavior is in those videos we could still be potentially forced to, although we expect such behavior is rare enough that it would not result in a likely to happen event.

B.3 Data for the Inverse Dynamics Model.

Since the DNN's task is to make actions given the video, any labeled data is appropriate for DNN training. In practice, we included general gameplay as well as the troweling task data described in the previous section, which amounted to about 10K hours. Due to collecting dataset like contractor Javacard only at the stages of the project, they were not included in DNN training.

B.4 contractor Javacard.

The contractor Javacard contains about 23 hours of data. We asked contractors to build a house in 30 minutes, using only basic dirt, wood, and sand blocks. Each trajectory data is manually generated and automatically adds trajectory after 30 minutes time limit. For the test, many contractors chose to keep their trajectories by collecting basic tools and building blocks, specifically it was common for the last 2 minutes to keep on collecting a wooden pickaxe and then trying to use an assortment of stone tools before gathering enough building blocks and trying to construct a structure.

C. Minecraft environment details

Our Minecraft training environment is a hybrid between MineRL⁷ and the MC2 dataset (which uses Experience Replay) Minecraft training package. Unlike the regular Minecraft

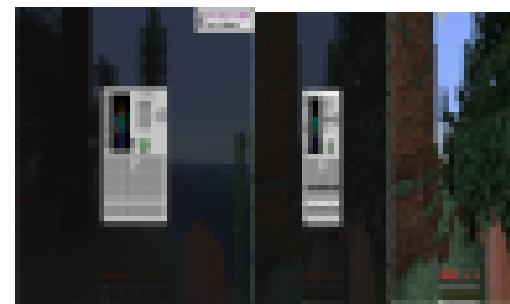


Figure 2: (Left) Sampled observation in the original resolution (4Kx2K) with an n-pixel GII agent. The noise areas can be seen in the centre of the image. This particular GII shows the player's inventory and can be used to eat very basic items. (Middle) We downsample images to 128x128 for computational costs. Shows a low-sampled observation with an n-pixel GII for eating. This is the condition used by our model. (Right) A 128x128 observation generated by our model without n-pixel GII. The health, hunger,饱食度 (satiation), and agent food can be seen in the bottom of the image.

pose, in which the error (in the ‘wall’) always runs at 30% and the distance to get to a building can complete (typically until 100%) in one visit to the closest server run in the same thread at the same frequency. This allows us to run the environment slower than the real time, while maintaining the training quality of the model. The action and observation spaces are similar to those of D4RL environments and are described in more detail in the following subsections. The environment also stores diagnostic information, such as n-pixel size, content of the agent's inventory, whether an n-pixel GII is active, etc., which we use in training and evaluating but not in input to the model. The episode length is 10 minutes for RL experiments and 10 minutes for BC model evaluation. The agent can ‘die’ in a number of ways, such as staying idle for too long and drowning, being killed by hostile mobs, or falling from a tall structure. We do not terminate the episode on agent ‘death’. Instead, just as the human in the sample Minecraft game, the agent keeps alive after it has died and requires a new spawn spot due to the initial spawning spot in the same Minecraft world. The player can be attacked in death, so the model can consider that the agent has died and act accordingly.

C.1 Observation space

The environment observations are simply the raw pixels from the Minecraft game that a human would see. Unlike Model 1, we do not encode velocity, the health/hunger/birth indicators, and the location of a nearby land shore in response to the attack or ‘use’ actions. The field of view is 70 degrees, which corresponds to the Minecraft default. GII scale (parameters controlling the size of the n-pixel GII) is set to 2, and brightness is set to 0.5 (which is a Minecraft default, but is very bright and visible visually). The rendering resolution is 4Kx2K, which is downsampled to 128x128 before being input to the model. We empirically found 128x128 to be the smallest resolution for which n-pixel GII elements are still discernible, and therefore that is minimum compute costs. Whenever an n-pixel GII is active, we additionally render an image of a noise mask at the appropriate noise position to match what a human player’s eye sees (see also Fig. 1).

C.2 Action space

Our action space includes all actions directly available to human players, such as key presses, mouse movements, and clicks. The specific binary actions we include are shown in Table 3.

The difference between the human action space and an agent’s is that we include typographical letters, which is only useful for interacting with the world by breaking traps/blocks. These can either be in the form of the escape tool with lowercase, the letters of which our agent can kill in. However, because we shall use the agent’s type letters that are also short actions (e.g. outside of the GII), the ‘W’ key (agent’s leftward action) agents are able to press the key within the GII (R, L, S,D,E, Q) (the police letter). If the escape tool used by a robot. We have not seen agents attempt to break the trap tool with these other letters. Instead, our agents script the trap tool with the mouse or switch by dragging items around the crafting window.

| Action | Description | Details |
|---------------|--------------------|---|
| forward | Walk | Move forward. |
| back | Stop | Move backward. |
| left | Turn | Turn left. |
| right | Turn | Turn right. |
| jump | Jump | Jump. |
| inventory | Equip | Open a slot inventory and the 3D crafting grid. |
| use | Use key | More variability in consumption of actions. In the 3DCrafts controller key, when used with attack it moves item from the inventory to the hotbar, and when used with craft it adds the selected item to the player's inventory. |
| upward | Attack | Move forward in current direction of motion. |
| attack | Left mouse button | Attack. In 3D, pick up the stack of items or place the stack of items in a 3D cell; when used as a double click (attack-on-attack), collects all items of the same kind present in inventory as a single stack. |
| use | Right mouse button | Place the item currently held or use the block-like place-in-holding slot. In 3D, picking up the stack of items or placing items from a stack. Left by mouse. |
| drop | Q key | Drops a single item from the stack of items the player is currently holding. If the player presses shift then it drops the entire stack. In the 3D, the user drops items except to the item he wants to keep. |
| hotbar: [0-9] | Keys 0 - 9 | Switch active item to the one in a given hotbar cell. |

Table 1: Binary actions included for the action space: <https://microsoft.github.io/rlkit/>. Controller key: more detailed description of selection.

In addition to the binary (switch) actions, our action space also includes movement actions. As with most people, when a user 3D's are not open, move L and R actions change the angle you are picking (e.g. jump). When 3D's are open, these actions move the same way. These movements are relative (i.e. they move the user in a vector relative to the current position, and thus their effect depends on the current position).

Inventory interaction is different requires fine-grained movement controls to achieve tasks such as crafting and reading, while moving and navigating the world can be achieved with coarse movement actions. Difficulties in interacting with the game state space, well-implemented movement controls and efficient actions will facilitate learning along and aids (Fig. 1), which is primary experiments we found to improve learning performance.

D Inverse Dynamics Model Training Details

D.1 DDM Architecture

The DDM model has approximately 10 million trainable weights. The input to the DDM is 128 consecutive images from 12 frames of video, each of which has dimensions $12 \times 12 \times 3$. The DDM is tasked with predicting the action at each frame. All image pixel values are then divided by 255.0 such that they lie within the range [0, 1]. The first layer of the DDM is a 2D convolution with 128 kernel size with a stride kernel width of 3 and stride kernel height of 1. This convolution is cross-correlation for controlling at the index 1 for batches of pixel values at times $t = 2, 3, \dots, 12, t + 1, \dots, T$. We had to crop to fit correctly important in DDM training.

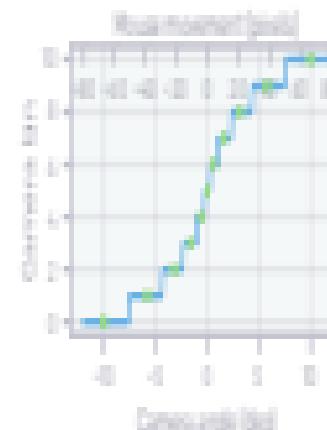


Figure 11. Relative cross-entropy losses measured by pixels vs. activation. The same binary is used but Loss Function. The binary is forced, meaning the binary is not bin-pool for smaller neurons while more bin-pool for large neurons. Thus each bin has breadth (X and Y). The state should be balanced with pixels (bin) is used when re-balancing neurons (that is, when creating bins in size especially when two neurons are measured).

as it incorporates applying temporal information immediately and to allow multi-step query DNN performance without a bottleneck layer as in Figure 14. This comparison was made on the Inception v3 DNN dataset.

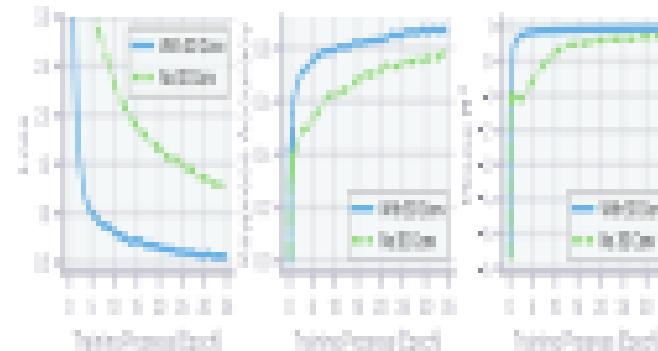


Figure 12. Effect of 50% Correlation in the DNN Architecture.

The initial temporal convolutional layer is followed by a stack of $^{(2)}\text{conv}$ layers processed sequentially. In this period the model, no static temporal information is shared between neighboring frames; however, since each frame was first processed with the temporal convolution, static temporal information is present at this stage. The flexible image processing network is composed of three subsequent stacks with widths $W = [34, 19, 19]$. Each stack is composed of (i) one 1D 1x1 convolutional layer with 16 output channels overlapping at the embedding boundary (such that the overlapping embedding dimensions will have a 16-neuron embedding dimension with 16 output channels), (ii) a 3D max pooling with stride 2 overlapping 1 such that the embedding width and height are halved, and (iii) two classic ReLU blocks with kernel size $3 \times 3 \times 3$ with each layer also having 16 output channels.

The output of the ReLU block is flattened into a 1-dimensional vector of size $2^3 \times 16 \times 16$ (one vector for each frame in the video) such that this step takes an 18 vector of size 1344. Each vector is independently processed with two frame-wise dense layers with 256 input activations and then 496 output activations, respectively. The result is then followed by 4 subsequent non-local (masked) residual connections¹⁰ blocks. Each block consists of masked attention layers (i.e., frame-aggregated latent frames, with 12 attention heads of dimension 12) followed sequentially by residual connection blocks (skip blocks). The embeddings are then passed through a frame-wise dense layer with output dimension 1024 and another with output dimension extending to 4096, a single residual connection block per the pair of frame-wise dense layers (but skipping past each layer separately), two skip blocks per pair. All dense layers have their weights tied throughout time, so each frame in the video is processed with the same weights.

Finally, independent dense layers lead to evaluation as predicted for final embedding - a 1-class multi categorical predicted with a softmax head available as well as a 11-way

categorical for both the behavioral behavioral and verbal assessments (see Appendix C2 for details on fraction space).

Each individual's convolutional layer is fine-tuned independently, along with 3 additional fully connected layers initialized with the initializers¹ (softmax initialized to zero).

B2 DNN Training

The total loss for the network is the sum of softmax prediction loss (over broad key and each behavior function). Each independent loss is computed by flattened off the convolution. We use the Adam² optimizer with a base learning rate of 0.001. We use an initial learning rate of 0.001, a batch size of 12 (over each item in the data) in a mini-batch of 128 frames, and an epoch length of 100. Hyperparameters were tuned in preliminary experiments. The DNN is trained on an extracted dataset (30 epochs). The total Adam loss is 22.4097%.

We add the augmentation mask when applying expectation over each weight matrix per request activity as temporally consistent. Using the PyTorch³ function `linear`, we apply linear layer weights between -0.1 and 0.1, activation between 0.0 and 1.2, brightness between 0.0 and 1.2, and contrast between 0.0 and 1.2. We also randomly rotate the image between -2 and 2 degrees, randomly translate later between 0.0 and 1.0, scale between 0.5 and 1.0 degrees, and translate it between -2 and 2 pixels in both the x and y directions.

Due to the large computational cost of running all of the experiments in this paper, training results are shown for training for DNN, BC, and RL training. The standard situation is assigned because deep learning training results in the lowest DNN inference latency relative to other approaches (e.g., no latent for the sequential task).

B3 Generating Reward Labels with the DNN

Section 4.1 shows that inverse dynamics modeling is a much easier tool than behavioral cloning because it can be non-linear. The DNN is trained simultaneously predict all 13 actions for each video sequence, so the DNN will effectively be used for generating behavioral cloning because later frames are not included in the sequence. For this reason, we apply the DNN over a video using a sliding window with stride 10 frames and only use the pseudo-label prediction for frames 10 to 16 (the center 14 frames). By doing this, the DNN predictions at the boundary of the video clip is associated to only the last and last frame of a 160-video.

C Foundation Model Behavioral Cloning

C.1 Foundation Model architecture

The behavioral cloning model architecture is the same as the DNN architecture described in Appendix B1 except for two modify the architecture so that it is causal (i.e., cannot predict future video frame predictions). This means the BC architecture does not have the initial non-causal convolution the DNN has (this layer is omitted completely). Furthermore, the residual connection layers are now causally masked (this is standard in language modeling) and we do Transformer-XL-style⁴ training where losses can distinguish and video frame predictions within the same video. We also use Transformer-XL-style relative position encoding.

C.2 Full Action Filtering

The most common action has over 100 million (in logspace) unique occurrences, which accounts for 99% of all actions by size. In my experiments, a player may take the final action in their last recording into their first, in pure betweenness, or in idle sleep to play a game of chess. Early on in the project we found that the BC model would take a much longer time than 30% of all actions, often upwards of 100%. In order to prevent this behavior we removed from our full action filter dataset. We compare a few different test statistic we filters out if they have less than 1, 10, or 20 hours of unique occurrences, and which is best for this application may not always. Full action filtering greatly helps increasing the inference rate (Figure 17(B)).

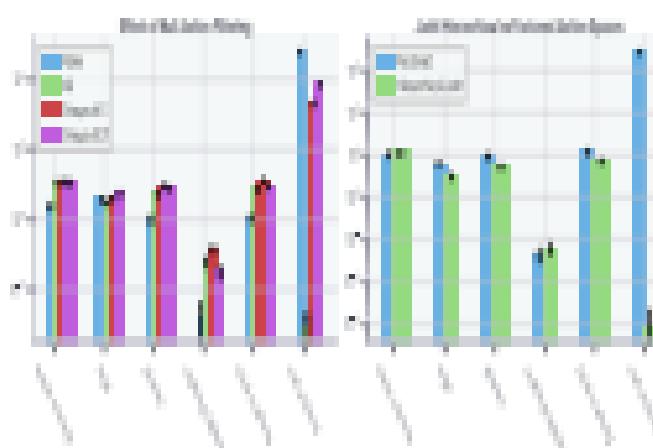


Figure 15: (Left) Effect of Full Action Filtering (by group). We compare estimated actions and counts of sampled full actions using either (optimal pre-filtering) for the filtering treatment: no full action filtering (blue), filtering direct actions (red), Direct+ (purple) group of 3 or more collisions (red), and filtering only groups of 3 or more collisions (purple). (Right) Standard versus Filtered Action Space.

Filtering only group of 3 performed slightly better than filtering direct actions in group of 3. Total experiments indicated that filtering direct actions was better; however, after filtering nothing and filtering had already trained no legal actions, we found that filtering only group of 3 was still actions performed best. Due to compute constraints we cannot do individual experiments with the setting, but this is as well as a reasonable choice for my interests.

13 Joint Biunivocal Action Space

We originally worked with a latent action space: where each type could be independently on or off, and the choice was independent of whether the same was being used. This will come down to making the latent belief distribution much like this prior state. However, also with 99% probability it's more logical and standard with 99% probability it's more likely to pick this item. The latent belief distribution can do this by giving 99% probability to each of the 4 constituent actions because it chooses to press each button simultaneously and independently. See Appendix C.2 for details on the static action space.

For this reason, we implemented a joint distribution over actions; however, the full joint distribution uses 20 binary buttons, and two more movement directions described later. It has total recall rate $2^4 \times 2^2 \times 12 \times 10^3$ possible configurations. This is far too large for many reasons, e.g. the final layer of the transformer model with dimension of 4096 would need to be mapped to each combination resulting in $4096 \times 12 \times 10^3 \times 12 \times 10^3$ parameters for this final layer alone. In order to reduce this to something more manageable in memory we can do other clever simultaneously good (for example, if a player tries to move left and right at the same time, they cannot do place). Below we list the sets of mutually exclusive actions. Furthermore, the inventory button is exclusive with all other buttons and vice versa.

| Action Exclusion Actions |
|--------------------------|
| Inventory back |
| left, right |
| up, down |
| lateral - [D, I] |

Even though the joint action space is effective, mostly exclusive configurations still result in a large action space, which is difficult to be learned from experience (i.e., $2^{16} \times 2 \times 12^2 \times 10^3 \times 12 \times 10^3$). The calculation results from Jester 3 set of 20 mutually exclusive buttons show that: inventory is forced to be on, while left, right buttons are mutually exclusive, up, down, attack, and jump are mutually exclusive movements, and finally -[D, I] for the inventory button which is mutually exclusive with all other actions. $\sim 2^{16} \times 10^3$ is still quite large, so we chose to implement a hierarchical memory, which becomes heavily nested at first. If the action is on, then there is a secondary classification with 12 classes (the joint distribution of those movements because each selected move distribution has 12 bins) for learning when

In more technical terms, if the learned action is all 0s, then there is no cause movement, however the secondary action component is encoded during training and the secondary action field needs to be sampled during evaluation. While this is a large model field for distribution, it is quite a bit better than the learned action space since it depends on between keypresses as well as releases with more fine tuning (although not much has been done) (see [this link](#)). The resulting action space has dimension $1^2 \times 10 \times 2 \times 2 \times 1 = 160$ (the 1^2 dimensional multiplier for causes over each character replaced by a multiple of 2 here, corresponding to primary actions for whether or not to move the mouse). When additional 1D-dimensions lead to further causes increases. In the future it would be interesting to implement sequential conditional action spaces to more completely model the joint distributions.

In Figure 15 (right) we compare environment editor performance between BC models with the learned joint action space and with the learned action space. Environment statistics are fairly comparable; however, we see that the learned action space model samples far more valid actions. This is a important example of the learned action space failing to correctly model the distribution in the dataset because the learned action filtering becomes finalization in the dataset these models fail on. Despite this, the learned model samples many valid actions because its prediction breadth key is not conditioned on other keypresses.

1.4 Foundation Model Training

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

| Hypoparameter | Value |
|---------------|--------|
| Learning rate | 0.0001 |
| Batch size | 100 |
| Epochs | 30 |
| Height/depth | 1000 |

Table 4: Hyperparameters for foundation model training

2 Behavior Cloning Fine-Tuning

Behavior cloning fine-tuning is similar to the DRL training, except we either use a shared robot validate robot (only_gym dataset, described in A.3) with pre-trained, or continue the contractor dataset (same dataset, described in A.4) with ground-truth block. The hyperparameters used for behavior cloning fine-tuning are listed in Table 5. We used 16 ALIVECPUs for about two hours for tuning the only_gym contractor (contractor, same dataset), and 16 ALIVECPUs for about 2 days when fine-tuning on only_gym dataset.

| Hypoparameter | Value |
|---------------|--------|
| Learning rate | 0.0001 |
| Batch size | 10000 |
| Epochs | 1 |
| Height/depth | 10 |

Table 5: Hyperparameters for behavior cloning fine-tuning

3 Reinforcement Learning Fine-Tuning

3.1 Reinforcement Learning Fine-Tuning Training Details

ILC agents are trained with policy gradient (PG) algorithm,¹⁰ and algorithms based on the proximal policy optimization (PPO) algorithm.¹¹ Both are sample efficient by performing additional passes over the collected data to optimize the value function as well as an

policy value function. These algorithms have been described elsewhere previously,¹⁰ where we describe them only briefly. A major inefficiency when training policy algorithms is that, to remain on-policy, one can only take a single gradient step before any update has tends to be pulled to another optimisation. To alleviate the potentially disastrous effects of taking multiple optimisation steps in a single iteration, PG prevents the policy from changing too much in a single step by clipping the loss when the difference between the current policy and the policy before the update becomes too large.¹¹ We also use generalised advantage estimation (GAE), which can speed up credit assignment by looking further than 1 step into the future when determining the advantage of an action, with the look-ahead being determined by hyperparameter γ .¹²

PG improves the sample efficiency of PPO when the policy and value function that it is activated by following different optimisation processes for the policy, the value function, and their shared representation. PPO splits optimisation in two phases: a value phase and a step phase. In the value phase, the policy and value function are optimised in a normal PPO setting, with the only exception being that every sample is used at most once, which prevents the policy from switching on these samples. In the step phase, PPO optimises the value function and an auxiliary value function (which is optimised with the exact same loss as the regular value function, but its output is never used during training) while keeping all Kullback–Leibler (KL) divergence loss in the policy before the start of the step phase to ensure that the policy does not change. Because the policy is not optimised in this step, PPO does allow samples to be reused multiple times in this phase. The assumption behind optimising the value function during the step phase is that value function optimisation is less sensitive to being trained multiple times on the same sample. Optimising the auxiliary value function does not directly affect either the value function or the policy, but it can improve the shared representation of both functions (the assumption being that predicting the value function requires knowing all features that are important in distinguishing states). The coefficients for the three losses (value function loss, auxiliary value function loss, and KL divergence loss) in Table 6.1 in our experiments a single iteration consists of two step cycles and one value cycle.

Because the value and auxiliary value functions are not optimised during BC pre-training, they are initialised at the start of KL fine-tuning. Each value function is implemented as a single fully connected layer on top of the last encoded features (last) of the pre-trained model (Appendix D). The weights of the auxiliary value function are randomly initialised while the weights of the regular value function are initialised with zero weights, which apparently prevent disastrous updates early in training that could happen with a randomly initialised value function. To prevent the value function loss from having gradients that depend greatly on the magnitude of the reward, we normalise the value function target by subtracting the mean and dividing by the standard deviation, which are calculated through an exponentially weighted moving average.

To prevent catastrophically forgetting the skills of the pre-trained network when KL fine-tuning, we apply an auxiliary KL divergence loss between the KL model and the former pre-trained policy.¹³ This loss is defined as:

$$L_{\text{aux}} = \beta L(\pi_t, \pi_0) \quad (6)$$

Where π_t is the policy being trained, π_0 is the former pre-trained policy, $L(\pi_t, \pi_0)$ is the Kullback–Leibler divergence between the policy being trained and the pre-trained policy, and β is a coefficient to weight this loss relative to other losses.

In the bc-tiny experiments, the KL divergence loss replaces the entropy maximisation loss, which is often added to KL experiments to encourage exploration.¹⁴ The described entropy maximisation is that, when all actions appear to have equal value, such as when the agent has not learned about the environment, it should encourage entropy to increase the chance that it chooses the next action. Randomly exploring by maintaining entropy is effective when the state and action spaces are sufficiently small as the reward is sufficiently dense, but becomes ineffective when the state and action spaces are large and reward is sparse, which is the case in the Diamond-poker task. Instead of randomly exploring through uniform-random actions, we assume that the pre-trained policy has an action distribution that is much more likely to take sequences of actions that lead to interesting new states, and thus, in states where the agent assigns equal value to each of its actions, it should move the action distribution of the pre-trained policy instead of uniform-random action distribution. In experiments with a randomly initialised policy we do include the entropy maximisation loss with a coefficient of 0.01, which has been an effective setting in other Diamond tasks.¹⁵ Empirically, we found that a high coefficient for the KL divergence loss would prevent the agent from properly optimising the visual function while also sufficient, was ineffective at

| Hypothese | Wk |
|--------------------------------------|-------|
| Lone play: | 15.07 |
| High play: | 0.00 |
| Baseline: | 0 |
| Double play: | 0 |
| Low play: | 12 |
| Discounted (-): | 0.00 |
| LL-1: | 0.00 |
| EW-1: | 0.00 |
| Max bonus: | 1 |
| Min bonus: | 1 |
| EW play cycle: | 1 |
| EW play rule-fair coefficient: | 0.0 |
| EW sequential rule-fair coefficient: | 0.0 |
| EW play LL coefficient: | 1.0 |
| EW play max/best rule: | 1 |
| LL discount coefficient: | 0.0 |
| LL discount play: | 0.00 |

Table 6 Hypotheses in LL experiments. These are the hypothesized for all treatments with two exceptions. First, when the training rule is only one rule without LL, the hypothesis, in addition to the LL-discounted hypothesis, the learning rate was set to 1×10^{-2} (the best setting and a step over 1 different learning rate), we found that performance was obviously lower with the standard learning rate of 2×10^{-2} and the potential causes from two other hypotheses. It is known that the training rule works better for lower values for training without a LL, because that the LL has greater scaling optimization steps that change the policy too much in a single step, especially in early iteration, when the value function has not been optimized yet, and the LL has this makes it possible to optimize with adaptive learning rate. Second, when using LL, it is a randomly initialized policy here is no LL-discount but in LL-discount play, but instead we use an entropy bonus of 0.01, which repeatedly switched its parameter.

protecting the learned skill of the persisted policy and preventing catastrophic forgetting. In such, we start with relatively high coefficient γ and decay it by a schedule after each iteration (Table 5). This method protects policy drift in early iteration while guaranteeing that the policy can eventually dominate the recall function, regardless of how different it behavior has to be in comparison to the persisted policy.

For increased function reactivation through qualities of each item that a human player might prefer when trying to craft a learned policy, and we used the model for gathering up to that quantity for each items. We start the reactivation by testing every technology technique whether a learned policy will satisfy the requirement for each items to be reactivated (e.g., have added a learned policy to the reactivation, then we add the Learned Materials and Learned Input required for crafting a learned policy, then we added the Learned Input required for mining diamond, and so on). Then we added each materials in the reactivation, with each being added a task when reading one and in crafting tasks while the tasks themselves improve visibility and prevent damage from spawning. Finally, we tested the model for crafting additional log (Crafting requires to craft all items in the recall function, but no crafting specific log), which can be used as fuel or crafted into a crafting table or tools if the agent runs out. In particular, upon each selection for additional log, places the tasks, or uses coal and other when crafting, but the reactivation can base on human expectation on what would be useful to execute his task, rather than designed around how an RL model behaves after training. Finally, to encourage the agent to keep mining diamond and crafting learned policies, when he crafted a learned policy, we allocated a limit on the number of times each learned policy that will be reacted.

The rewards for the different items are: expectation items, simply depending when the player would usually get the relevant items. The last term consists of all reward and star items and has a base reward of 1, whereas the reward of all items requiring coal will also be reward of 1, so that the reward of all items requiring coal will also reward of 1, and the total for a diamond with a base reward of 1. The item base in the sequence of items towards a learned policy generally

| Item | Quality reward | Reward per item |
|----------------|----------------|-----------------|
| Spade | 0 | 10 |
| Club | 20 | 120 |
| Diamond | 30 | 180 |
| Cards | 1 | 1 |
| Woden pole | 1 | 1 |
| Chopping stick | 1 | 1 |
| Stone | 1 | 1 |
| Flint | 1 | 1 |
| Bone | 1 | 1 |
| Leather | 1 | 1 |
| Leather pole | 1 | 1 |
| Bone | 1 | 1 |
| Gold | 1 | 1 |
| Flint | 20 | 120 |
| Bone | 20 | 120 |
| Leather | 20 | 120 |
| Leather pole | 20 | 120 |
| Bone | 20 | 120 |
| Diamond | 30 | 180 |
| Diamond pole | 30 | 180 |

Table 2: Reward per item and total quality reward.

give a higher reward. To make sure that the agent learns one-value items that are needed in the gathering task (e.g. the agent is rewarded for up to 20 plastic but only up to 1 chopping stick, which can cause the agent to focus on plastic at the expense of creating a cutting table), we divide the reward of each item by the total quality that the agent gets rewarded for (in the process of determining the reward, the total quality for diamonds is 3 and the total quality for diamond poles is 1, even though it did not put a limit on the number of these items being rewarded). For example, if the agent is rewarded for 3 diamonds, which has a reward of 60 for being in the task and up to 1 black diamond are awarded, the reward per item of diamond is 40. The quality and reward for each item are listed in Table 2.

While every item in the sequence towards a diamond pole is rewarded, the revitalization is still sparse and occurs once every despite. The quality comes from the fact that it is only the end of actions to find the next reward, even after the agent has explored all the necessary pre-requisites (e.g. bone plays often the role that 100 actions to find a diamond the cutting or bone pole). The revitalization can be deceptive when the most efficient method for gathering items can make it far more difficult to get the next item. For example, a good strategy for the agent to call a stone pickaxe quickly is to make it a special item reward to pick up its cutting table after calling a wooden pickaxe, such that the agent has immediate access to a cutting table a reward for collecting enough�石质的. However, the fastest way to profit a wooden pickaxe is to move directly to the cutting a wooden pole; while leaving the cutting table behind. This following the optimal strategy for gathering sufficient makes it more difficult to locate a stone pickaxe.

Experiments run for approximately 6 days (14 hours) on 100 CPU (for policy optimisation) and 3,200 CPU (simply for collecting actions from Microsoft). In this time the algorithm performed roughly 400 optimization iterations and collected roughly 14 million Microsoft episodes consisting of 2,000 frames each (for a total of 16 Million frames).

C.2 Reinforcement Learning Fine-Tuning: Additional Data

Additional figures that help better understand the main results of Sec. 3.2. The tuning experiments are presented in this section. First, we show the item-over-tuning figure when RL fine-tuning from the early-game model without a RL loss (Fig. 16). When training without a RL loss, the model only learns to obtain the five items for the early-game model (i.e. update objecting item-set), which are log, plastic, sticks, and cutting tables. Second, we present preliminary experiments in which we directly compare RL fine-tuning from the base-building model and RL fine-tuning from the early-game model (Fig. 17). These experiments differ from the main experiments in that, for both treatments shown here, the RL loss coefficient was set to 1.0 (learning rate was set to 4×10^{-5}), and increased the end item reward quantity for all items (i.e. items closer to the diamond pole did not have an increased reward). While RL fine-tuning from the base-building model initially

and better than IL learning from the only-gate model. IL learning from the only-gate model tended better than 10000 epochs of direct signal matching from inputs, which is why the only-gate model was chosen in the main experiments.

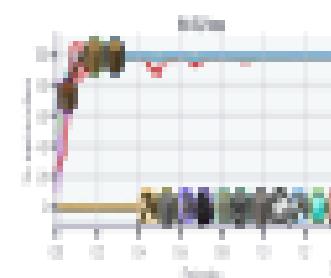


Figure 16: Loss vs. Epochs when IL learning from the only-gate model achieves 0.15 loss. The model learns to obtain all losses for the only-gate model can not be beaten, whether logic, gate, or direct signal matching. In contrast to the test cases with a IL learning, it does not learn any items beyond the initial loss, likely because skills that are not performed are not used, and therefore the model that does not initially see any reward, are statistically bypassed while the first few items are learned.

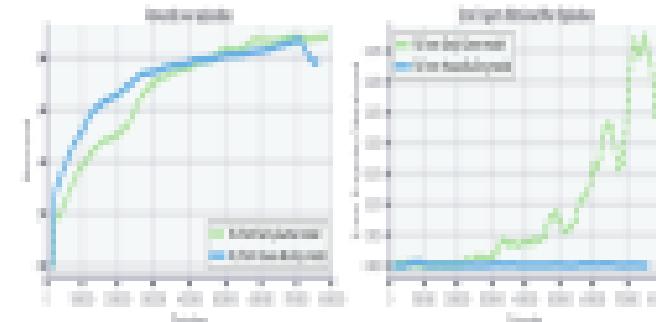


Figure 17: Preliminary experiments when IL learning from the only-gate model compared to IL learning from the base building model. (Left) With regard initially increase faster when IL learning from the base building model, IL learning from the only-gate model eventually shows a slightly higher trend. (Right) IL learning from the only-gate model has a higher likelihood of selecting a correct input, which is why the only-gate model has discrete lines in the loss experiment.

III. Evaluation Model Scaling

In early experiments we found that increasing model size led to much staying in local optima learning regions during training.¹² Here we compare the ILG model described in Section 4.2 with 1000 and 2000 parameter model. Both of these models are trained for 10 epochs as compared to the 30 epochs the ILG model trained for. These models have the same architecture as the ILG model but each layer in the 2000 parameter model has 1024 units and each layer in the 1000 parameter model has 512 units. The 2000 model was trained with an initial learning rate of 0.00056, batch size of 30, and weight decay of 0.049%. The 1000 model had an initial learning rate of 0.00033, batch size of 30, and weight decay of 0.053%.

In Figure 18 we show validation loss on only_class with ILG model-based loss on bottleneck dataset and with the ILG with gated truth labels collected during counterfactual and zero-shot environment performance for the 1000, 2000, and ILG models. While larger models have better validation loss on only_class, these models tend to do worse than the ILG model in terms of validation loss on only_class, due to the fact that they have worse environment performance. In fact, we see that the 2000 model even had lower validation loss than training (Fig. 18(a), left). The 2000 model also appears to be better at scaling than the ILG, as it has lower validation loss than the ILG.

While the overall results suggest smaller models are better, IL learning still scales linear. When the training is counterfactual, models with small validation losses and zero-shot ILG model perform best both in validation loss (Fig. 18(a)) and environment performance (Fig. 18(b)).

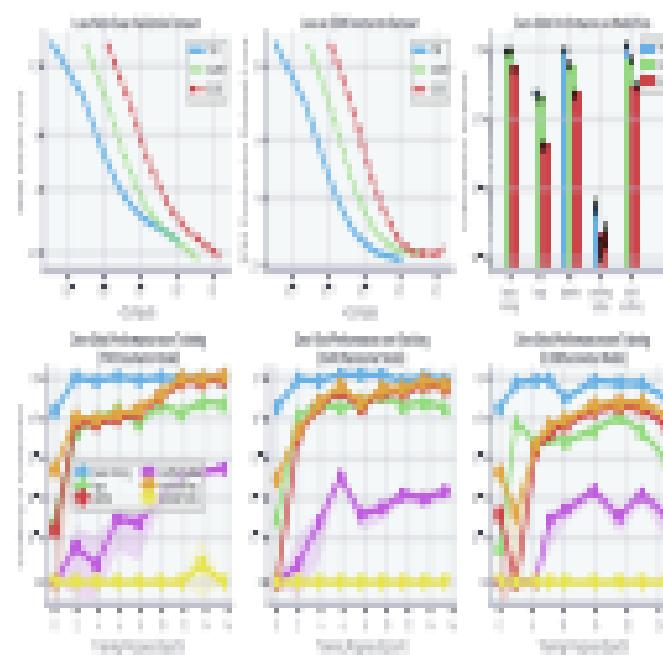


Figure 11: Training and Cross-Val Performance across Model Size. In the first two plots the ratio is compute normalized is learned by the 700 parameter model, and therefore 10 epochs training for the 700 model has used 1 ‘compute’. The 1000 parameter model and the 7000 model are trained on the same amount of data (5 epochs), while the 10000 model is trained on 30 epochs of data. (Top Left) Loss on the val. class validation dataset. (Top Middle) Loss on the EBM validation dataset; note that these models were trained only on val. class validation by contractor Joss. (Top Right) Env. Val. env. validation val. performance after end of training. (Bottom) Env. Val. env. validation val. performance over training for the 700 model (bottom left), 1000 model (bottom middle), and 10000 model (bottom right).

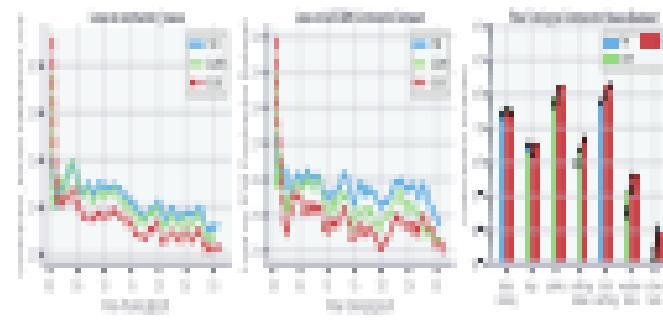


Figure 12: Contractor Joss Env. Val. Performance across Model Size. (Left) Loss on the contractor Joss val. validation. (Middle) Loss on the full contractor Joss dataset collected to train the EBM; this dataset is disjoint from contractor Joss. (Right) Env. Validation val. performance at the end of fine-tuning.

followed by the 2000 model and then the 700 model. Environment model val. loss are performed using the same gate weights that we use to collect contractor data, which could be visually distinct from videos taken from the web. It is plausible that the large models overfit to the visual predictions involved during pre-training; they have more contractor data loss (Fig. 11 top middle), and this causes them to perform worse poorly in the environment val. loss. However, we hypothesize that because the contractor Joss dataset video frames collected from the web require the large models that are later used. These inputs are likely to be mostly low val. class validation loss in Fig. 11 top left, compared with the contractor Joss loss (indicates no performance on this setting from our perspective, resulting in low environment val. performance). This hypothesis is further supported by Fig. 11 middle showing loss on the contractor dataset collected for EBM training, which has no overlap with contractor Joss. After just 10 steps of fine-tuning in contractor Joss, all models quickly improve its loss on the full EBM contractor dataset, with larger models now performing less. With no contractor, we believe this investigation provides some intuition for future studies of switching the input modalities in fine-tuning problems.

I Test Conditioning

Goal-conditioned policies^{10,11} make it possible for a single agent to perform a wide variety of goals in a single environment, which is particularly relevant in open-world environments such as Minecraft. In recent work, goal specification has increasingly taken the form of domain-specific languages¹², or even natural language^{13,14}. The benefits of language-conditioned agents can be tremendous, especially natural-language-conditioned agents, as their goal space contains a wide variety of potentially very complex tasks. Test condition matching has shown its amazing ability to perform tasks two-fold (or more than two-fold) in half the time by specifying a sequence of steps via the compositional and conditional possibilities offered by natural language (e.g. GPT¹⁵ and DALL-E¹⁶). We hypothesize that we should expect similar capabilities to emerge with natural-language-conditioned visual agents. They are visually trained on various sequences of data that pair from a natural language description to a sequence of actions that complete the specified goals. In this section we take preliminary steps toward that line. Our preliminary experiments provide evidence that it is possible to produce a natural-language-conditioned model for Minecraft using the general approach presented in this paper (NFT) plus conditioning on the speech that often accompanies videos.

In other videos, the human actor sometimes indicates their intent in their verbal commentary (e.g. “Let’s go play some tennis to make a tennis rac” or “now let’s have her to copy photos in Photoshop”). Conditioning on the closed captions associated with a visually pre-trained model (i.e., it may already possess the capability to condition the model with test text as “I am going to call someone police” or “I am going to build a house,” and have the agent perform those tasks specifically rather than simply follow typical human behavior (as was investigated in the rest of this paper)). A human may be pushed to produce a describable agent in the NFT fine-tuning, which we could have done in Section 4.1 by allowing it to indicate the task to be completed, as has been done in prior work.¹⁷ However, conditioning on natural language often may benefit over that approach. First, it is flexible and powerful, being able to express any task. Second, one does not need to pre-specify the nature of the task in the completed form of text. This will allow for greater cognitive flexibility of agents like GPT but extending their capabilities to extended tasks such as completing goals as sequences in a situated Minecraft. Third, test conditioning conditioned over other tasks are difficult to specify in natural language (e.g. “Let’s build a house” vs “The agent is capable of doing a non-copying task like ‘I will now build a castle surrounded by a road’”). In fact, NFT-conditioned models may produce powerful, cognitive natural-language-conditioned agents with the power of GPT to transduce, follow instructions, and complete tasks over or after the data, but in the form of agents that can act in visual worlds, complete tasks as sequences, and in other similar extended repeatable decision domains. We do not reach this lofty goal in this work, but we hope our first step towards exploring what direction.

Many Minecraft video datasets also commentary from the players. This commentary is sometimes present in the closed descriptions for the video, or would be extracted post-hoc using automated speech recognition (ASR).¹⁸ One class of features about the human commentaries with associated closed captions:

We fine-tuned the 2B million parameter TPT baseline model used in the NFT fine-tuning experiments (shown in Table 1) for the same reason: to gather enough data with an additional test-conditioning input on the subset of video files for which closed captions are available. To obtain the conditioning input, we took open-videos into 30-second chunks. The user text is associated with every frame in a given chunk, and is made up of all the closed captions occurring within that chunk, as well as the line of text preceding and following the chunk, if any. Because the vast majority (around 95%) of the closed captions that lack capitalization and punctuation, this is consistent with the speech theory¹⁹. We then obtain user embedding vectors of length 4,096 from the OpenAI embeddings API²⁰, which is produced by a randomly initialized multi-layer perceptron (MLP) with ten hidden layers of size 3,072. The resulting activations are added for each frame to be passed through multi-activations before the final linear layer (pre-softmax activation = $\text{clip}(x, \text{clip}(x, -10, 10), 0)$). The model is fine-tuned for language tasks.

On multi-class multilabelled testability. When conditioned on sentences that isolate the agent to specific goals like “I am going to copy” and “I am going to build a” the agent tends significantly further from its query point (Figure 2b). Additionally, we can see the agent is preferentially selected

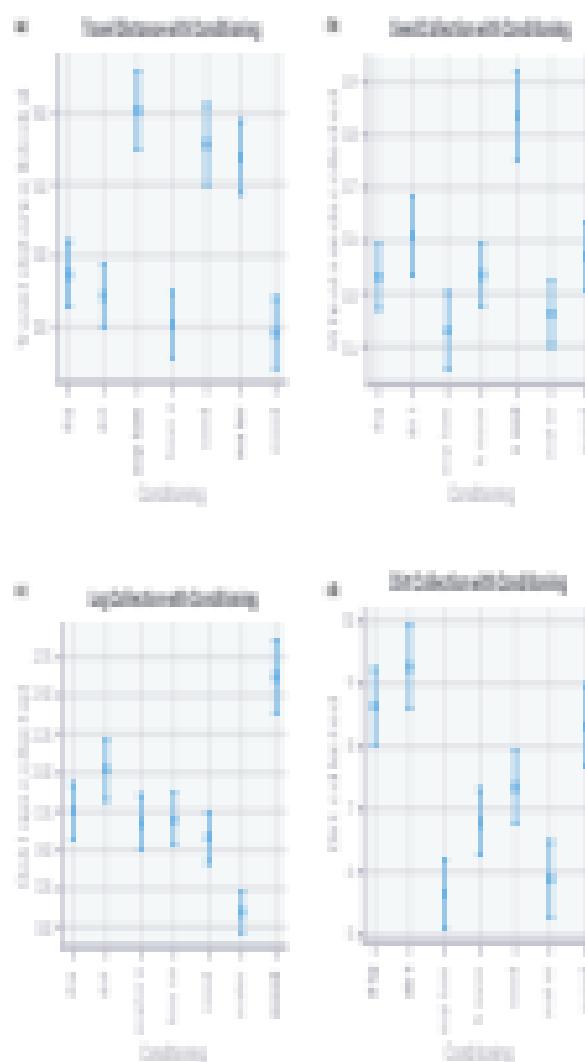


Figure 20: Evidence for conditioning. In each plot, the various expected treatment are shown in bold. The steps corresponding to each outcome shown in Table 1. Statistics are measured over 10000 episodes. (a) Distance travelled by the agent. Both “regular” and “water” test steps should encourage a sterile agent to move more than when doing other tasks, which is what occurs. Gaseous water is needed to get enhanced movement in all fixtures, which is likely why the “leaf” condition produces more travel (as the agent sometimes tends to move to a fixture with grass). The travel distance is the Euclidean distance from the open point to the latest point the agent reached during the episode on the horizontal (x-y) plane. (b) Collection of stones task. The “soil” variant collects substantially more than other variants, as expected of a sterile agent. (c) Collection of leafs (the most common type of “leaf” leaf). The “leaf” variant collects significantly more leafs, which is to be expected of a sterile agent (as speculated for “water” variant which has lower disease infection rates). (d) Collection of list. The “list” and “leaf” variants collect a large amount, while the variants that are uniformly in the case of “leaf” conditionally collect list. This can be attributed again to the ground rather than air gas or tree when collecting seeds or leaf, which likely explains the slightly higher amount of list collected by tree variants. In all cases, the error bars are 95% confidence intervals of three to six, over 1000 episodes per conditioning variant. Treatments for which the box is not bold for plot do not overlap are statistically significantly different at $p < 0.05$ level.

| Variant name | Step |
|--------------|--------------------------------|
| bg | Engaging in step with no gas |
| bs | Engaging to collect soil |
| sgs | Engaging to explore |
| sls | Engaging in a microbial source |
| sd | Engaging to collect seeds |
| sw | Engaging to find water |
| swd | Engaging to deposit seed |

Table 4: Steps corresponding to each conditioning variant.

can give less and a risk, and similarly condition will tell such as “To play is not reducing multi-player” (Page 38, ch.).

With an agent does some level of credibility, one need to expect to increase it. For example, we can not do it necessarily the agent to give fines or to fine, both of which are possible in theory price, but has certain limit, in the case of having enough, machine credibility has giving fine, word, or cash. Likewise, an experiment in which the agent is presented with a collection values and various resources, and additional incentives given him (e.g. “To play is not reducing multi-player”) which shows that the conditioning had a significant effect on which less get called. Instead, it would like agent was more influenced by the price, unconditional probability of what human players will call not give the resources available, which is not the surprising since in Microsoft, especially in the early price, there is a relatively consistent public policing resources in a specific role in police more powerful tools (Fig. 6). For example, if the agent had the resources to make a stop police, and wouldn't intend to make a ticket from a police, it often would make the stops police saying. Finally, looking at values of agent behaviors related to continue a little “base” conditionality cause the agent to take more steps towards holding a base for other values.

Thus, one needs also that is possible to make a somewhat specific multi-player-conditioned agent. However, its credibility will be much more practically useful, and it's for the relative belief could be accompanied with more research, data, and training sample. Another exciting research direction is to have the multi-player function as well as joint distribution.