

unique
and other
in the labels
internet-scale
semi-supervised
semi unlabeled videos.
unlabeled data we can train an
unlabeled source of online
unlabeled – from which we can then
unlabeled, with both imitation learning and
behavioral prior has nontrivial zero-
shaped tasks that are impossible to learn from
to report computer agents that can craft
sufficient humans upwards of 20 minutes (24,000
play to accomplish.

For sequentially act within an environment (e.g. robotics, game playing, and computer
in e.g. natural language). In a few rare settings, such as Chess, Go, and StarCraft, there
is much less in the form of *unlabeled* video (i.e. without the actions taken
in these domains taken a wealth of data also exists on the
without a general foundation

His was a large effort by a dedicated team. Each author made huge contributions on many fronts over long
periods. All members were full time on the project for over six months. BB, IA, PZ, and JC were on the
original VPT project team and were thus involved for even longer (over a year). Aside from those original team
members, author order is random. It was also randomized between IA and PZ.

[†]OpenAI

[‡]University of British Columbia

...
s
ings,
within
causal BC
distribution, which does
ation bottlenecks
am tasks with either

ts
and other us-
(c) it
a research
st the native hu-
e domain shift between
use and keyboard
dragging items to specific slots or navigating human contractors to
attacking macros,^{30,32-34} using the native human interface that uses a lower frame rate and
exploration task of gathering a single wooden log while already facing a tree takes 60
ack actions with the human interface, meaning the chance for a naive random policy to
simple task of gathering a single wooden log while already facing a tree with RL from
making most simple tasks near impossible to learn with RL alone, such as crafting planks and crafting tables, accom-
section 4 we show that the VPT foundation model has nontrivial zero-shot performance, accom-
ishing tasks impossible to learn in Minecraft a median of 50 seconds or ~970 consecutive actions (tasks
requiring a human proficient in Minecraft to behavioral cloning to smaller datasets that target more specific actions).
Through fine-tuning with distributions, our agent is able to push even further into the technology tree, crafting stone tools



Figure 1: Example Minecraft crafting GUI. Agents use the mouse and keyboard to navigate the menus and drag and drop items. This choice means that our models play at 20 frames per second, making most simple tasks near impossible to learn with RL from scratch only, the VPT method is general and applied to any domain.

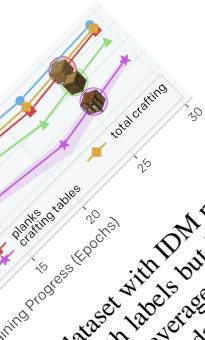
of roll
variations
learning is
trajectories
video games.^{10,44}
not work; however,
actions.²² For instance,
ency to exhibit behaviors,
pose to first learn a latent
e motion-capture methods to track
waypoints. Similarly, Behbahani
the current state towards expert-provided goal
(IDM),⁵¹ which aims to uncover the underlying
past and future timesteps, e.g. $p_{\text{pw}}(a_t | o_t, o_{t+1})$, and
ectors of observations labeled with the IDM. Data to
any point in training if there are sequences in the dataset that both
ures that the BC model in the target environment such that both
we first train an IDM on a small number of labeled trajectories collected for a
correctly label them. Therefore, if the BC model does not explore
they play the game as would normally as we record their keypresses and
throughout BC training.
because human contractors reach most relevant parts of the state space, we can
ost previous work in semi-supervised imitation learning.^{1,27,28,30–34,52–60}
x and open-ended environment of Minecraft. Minecraft is a voxel-based 3D much
due its popularity and wide variety of mechanics, has attracted a vast amount of video
tasks such as navigation,^{53,60} block placing,^{54,55} instruction following,^{58,59} combat,⁵⁶ worlds
ers.^{28,31,57} Work operating,³² automated curriculum learning³⁰ and, most closely related to the RL
experiments presented in Sec. 4.4, diamond mining.^{27,32–34} However, to the best of our knowledge,
there is no published work that operates in the full, unmodified human action space, which includes,
drag-and-drop inventory management and item crafting.

Indeed,
This IDM
harder task of
details.

ing the web for related
ts, such as a video feed
from platforms other than
modes, e.g. in Minecraft we
ust gather or craft all their items.
s from survival mode, and call all
el, and enough training compute, a BC
y still perform well in a clean Minecraft
n clean and unclean segments). We do this by
unclean (Appendix A.2).

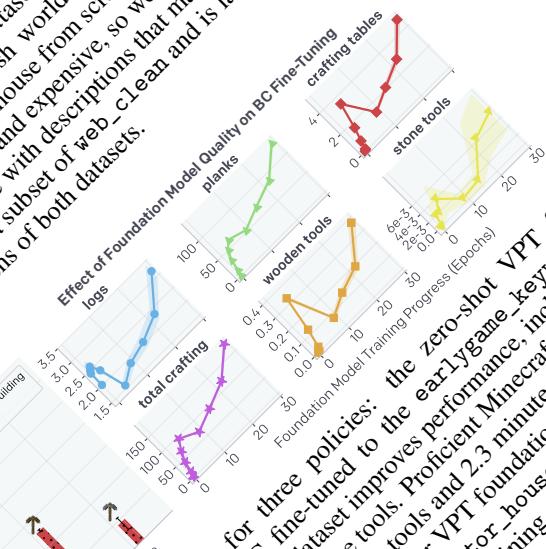
ndation model with standard behavioral cloning, i.e. mini-
actions predicted by the IDM on clean data. For a particular
 $a_t | o_1, \dots, o_t$, where $a_t \sim p_{\text{IDM}}(a_t | o_1, \dots, o_t, \dots, o_T)$ (1)
lowing sections, this model exhibits nontrivial zero-shot behavior and can be
imitation learning and RL to perform even more complex skills.

Performance of the Inverse Dynamics Model
The IDM architecture is comprised primarily of a temporal convolution layer, a ResNet⁶² image
processing stack, and residual unmasked attention layers, from which the IDM simultaneously
predicts keypresses and mouse movements (see Appendix D for IDM architecture and training
details). A key hypothesis behind our work is that IDMs can be trained with a relatively small amount
of labeled data. While more data improves both mouse movement and keypress predictions, our best



or learned by a behavioral cloning policy trained on an extremely labeled videos of our IDM. To collect the unlabeled internet dataset, we search results in ~270k hours of video, which we refer to as “minecraft web-clean” (Section 3) and then trained the VPT foundation model with a 0.5 billion parameter model (Appendix H), which took ~9 days on 720 V100 GPUs. Preliminary experiments suggested that our model could benefit from 30 epochs of training duration (Fig. 4, left) and rolling them out in the standard survival mode game where they play for 60 minutes, i.e. 72000 consecutive actions, and we plot the mean and shade game collection rates (Fig. 4, right). The VPT foundation model quickly learns to chop down trees to collect logs, a task we found near impossible for an RL agent to achieve with the native human interface (Sec. 4.4). It also learns to craft those logs into wooden planks and then use those planks

task terms (logs, planks, crafting tables, crafting tables, crafting tables) across a range of narrower datasets. The model's ability to collect a wide variety of items from scratch using primarily house descriptions, so we also construct a subset of web-clean and is labeled with descriptions that match keywords from both datasets.



(Left) Collection and crafting rates for three policies: the zero-shot VPT found on the VPT foundation model BC fine-tuned to the earlygame keyword or either house dataset. (Right) Collection and crafting rates for VPT foundation model BC fine-tuned to either dataset improves performance, including for task terms (logs, planks, crafting tables, and total crafting). Proficient Minecraft players take 1.2 minutes (1390 actions) to construct wooden and stone tools. Proficient Minecraft players take 2.3 minutes (2790 actions) throughout training after they are BC fine-tuned to the contractor_house dataset. In general, crafting-related behaviors increase throughout foundation model training. Fig. 4 defines the other task terms (logs, planks, crafting tables, and total crafting).

(iv) Sample videos: https://www.youtube.com/playlist?list=PLNAOlb_agf3U3rSvG_BCWqJ869Ndhcp

Diamond Pickaxe
24000+ act.
200+ minutes
12% 10 min
1 act.
1 minute
1% 10 min

Item is the median
percentage of contractors
axe is unknown (except
of 20-minute episodes).

ing goal of obtaining a diamond pickaxe is world. Doing so involves acquiring items like mining, inventory management, or dying. Obtaining a diamond pickaxe more (Fig. 6). Adding to the difficulty, progress is measured by the number of actions taken in the sequence, with lower rewards for items that have to be obtained in fewer actions. Agents are optimized for reward function and RL training details. Due to the value is realized. It therefore may have some behaviors required to smelt iron zero-shot, it did train on examples of previously learned behaviors. If the RL divergence loss between the RL model and the frozen pretrained policy fails to achieve almost any reward, underscoring how better (Fig. 7a), learning everything up to mining an iron ingot, the next item required to get further into the tech tree, likely this agent fails at smelting an iron ingot, the next item required to get further into the tech tree, likely

(a). The model never learns to diamond pickaxe (Fig. 7b). RL fine-tuning from the VPT foundation model does substantially better (Fig. 7c). However, this agent fails at smelting an iron ingot, the next item required to get further into the tech tree, likely

Sample Videos: https://www.youtube.com/playlist?list=PLNAOib_agf2yDSs4AqcoyPv4z_eWUiKn

is too low, even
when fine-tuning to the
RL Model to the
three-phase training to
learning extremely difficult tasks:
items in 57%, 15%, and 12% of episodes,
rafting iron pickaxes and mining diamonds, and
diamond **pickaxe**. To the best of our knowledge, we are **the first to**
efficient mining patterns, cave exploration, returning to
rafting a wooden pickaxe. Qualitatively, the model developed
initial skills of chopping logs and crafting planks are lost due to catastrophic forgetting.

Properties of the Foundation Model

In we validate a core hypothesis behind this work: that it is far more effective to use contractor data to train an IDM within the VPT method than it is to directly train a BC contractor dataset from that same small contractor dataset. If we could cheaply collect a labeled foundation model from a similar order of magnitude as web-clean, then this would not 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 data from 1 hour up to the full $\sim 70k$ web-clean dataset. Foundation models trained up to and including 1k hours are trained on the IDM found at https://www.youtube.com/playlist?list=PLNAOlb_agf3e_UKweM5pQUStw8r-Wfc

Inclusion

This paper help pave the path to utilizing the wealth of unlabeled data on the task is not learning and using these learned behavioral priors, VPT offers the exciting possibility of directly learning in a video—because arguably the most important information in any given scene would be in features trained to correctly predict the distribution over future human actions. We leave Future work could improve results with more data (we estimate we could collect $>1M$ hours) and larger, better-tuned models. Furthermore, all the models in this work condition on past observations only; we cannot ask the model to perform specific tasks. Appendix I presents preliminary experiments on conditioning our models on closed captions (text transcripts of speech in videos), showing they

After which there are few to no gains and differences are seen in all previous experiments we use our best IDM trained on increasing amounts of contractor data.

IDM Training Data (Hours)	Training (Zero-Shot) (IDM)	Training (Zero-Shot) (VPT)	Training (Zero-Shot) (VPT + BC)	Training (Zero-Shot) (VPT + BC + Keyword)
10^0	10^0	10^0	10^0	10^0
10^1	10^1	10^1	10^1	10^1
10^2	10^2	$10^{2.5}$	$10^{2.5}$	10^2
10^3	10^3	10^3	10^3	10^3

After which there are few to no gains and differences are seen in all previous experiments we use our best IDM trained on increasing amounts of contractor data.

- [1] Naman Dua, Arunveer Singh, and Sargur Nallapati. Multi-task learning for multi-domain text classification. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1901–1909, 2017.
- [2] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [4] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [5] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [6] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [7] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [8] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [9] Ming Tang, Yuxin Chen, Omer Levy, Mike Lewis, and Llion Jones. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [10] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michael Mathieu, Andrew Dudzik, Jun Young Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.

- al
y random
June. First return,
achita Chhaharia, Alistair
Antoro, et al. A data-driven
and Perez-Landez, and Percy Liang. World of
In Doina Precup and Yee Whye Teh,
ience on Machine Learning, volume 70 of
3135–3144. PMLR, 06–11 Aug 2017. URL
117a. html.
gorithms for inverse reinforcement learning. In *Icml*,
Peter Stone. Recent advances in imitation learning from
y:1905.13566. 2019.
mon. Generative adversarial imitation learning. Advances in neural
systems, 29, 2016.
ut Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv*
05.01954, 2018.
to imitate behaviors from raw video via context translation. *arXiv*
ence on Robotics and Automation (ICRA), pages 1118–1125. IEEE, 2018.
Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation:
winfinite Staff.
most-played-games-in-2020-ranked-by-peak-concurrent-players/.
rent players.
Most played games in 2021, ranked by peak concurrent
Twirfinite. URL <https://twirfinite.net/2021/12/>.
William H Guss, Brandon Houghton, Nicholas Topin, Phillip Wang, Cayden Codei, Manuela
Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations.
[27] arXiv preprint arXiv:1907.13440, 2019.

- Juewu-mc:
arXiv preprint
network. Advances in
Trends in cognitive sciences,
and Brett Browning. A survey of robot
ous systems, 57(5):469–483, 2009.
Elyan, and Christina Jayne. Imitation learning:
ng Surveys (CSUR), 50(2):1–35, 2017.
Zier, and Donald Michie. Learning to fly. In Machine
o visual perception of forest trails for mobile robots. Gianni Di Caro, et al. A
Christian Forster, Jürgen Schmidhuber, Gianni Di Caro, et al. A
D Jackel, Mathew Monfort, Urs Müller, Daniel Dvorakowski, Bernhard Fimner, Beat Flepp, Prasoon
evelilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy.
-end driving via conditional imitation learning. In 2018 IEEE international conference
(4):198–208, 2007.
Rémi Coulom. Computing “elo ratings” of move patterns in the game of go. ICGA journal, 30
Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations. In
Proceedings of the AAAI Conference on Artificial Intelligence, volume 32, 2018.

- [1] Michael L Littman, Deep reinforcement learning with model learning and monte carlo tree search. arXiv preprint arXiv:1902.04257, 2019.
- [2] Mingzhang Xiong, and Richard Socher. Fighting zombies in minecraft with hierarchical and interpretable skill acquisition. arXiv preprint arXiv:1712.07294, 2017.
- [3] Mingxian Shi, Yue Feng, and Aldo Lipani. Learning to execute or ask clarification questions. In International Conference on Machine Learning, pages 261–2670. PMLR, 2017.
- [4] Tambe Matissen, Avital Oliver, Taco Cohen, and John Schulman. Teacher–student curriculum learning. IEEE transactions on neural networks and learning systems, 31(9):3732–3740, 2019.
- [5] Robert George Douglas Steel, James Hiram Torrie, et al. Principles and procedures of statistics. Principles and procedures of statistics., 1960.
- [6] , and Robert E Schapire. Learning inverse problems in complex domains. arXiv preprint arXiv:1808.08456, 2018.
- [7] Yannan, and Shim-Young Lee. Keeping your discounted rewards. Advances in Neural Information Processing Systems, 31:2661–2670. PMLR, 2018.
- [8] Mingzhang Xiong, and Richard Socher. Hierarchical and interpretable skill acquisition in deep reinforcement learning. Technical report, Technical report, Stanford University, 2017.
- [9] Mingzhang Xiong, and Richard Socher. Fighting zombies in minecraft with hierarchical and interpretable skill acquisition. Technical report, Technical report, Stanford University, 2017.
- [10] Mingzhang Xiong, and Richard Socher. Hierarchical and interpretable skill acquisition with model learning. arXiv preprint arXiv:1902.04257, 2019.

- [74] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, 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 32, 2019. URL <http://papers.neurips.cc/paper/748035CurranAssociates,Inc.,2019>.
- [75] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? Advances in neural information processing systems, 27, 2014.

We thank the following
Joel Lehman, III
Jonathan Goode,
and Christopher

Acknowledgements

We thank the following people for helpful discussions and support: Bob McGrew, Ken Stanley, Joel Lehman, Ilya Sutskever, Wojciech Zaremba, Ingmar Kanitscheider, David Fathi, Glenn Powell, Jonathan Gordon, and the OpenAI supercomputing team, especially Christian Gibson, Ben Chess, and Christopher Bemner.

on ap-
. PMLR,
-107.12808, 2021.
doi: 10.24963/ijcai.2019/880. URL
<https://github.com/Felflflare/rpunct>.

erson, Arun Ahuja, Arthur Brussee, Federico Georgiev, Alex Goldin, Tim Harley, et al. Creating a neural with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

punct, May 25 2021. URL <https://github.com/Felflflare/rpunct>.

antau, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings pre-training. *arXiv preprint arXiv:2201.10005*, 2022.

on and self-supervised learning. *arXiv preprint arXiv:2204.06125*, 2022.

er, Jacob Andreas, Edward H. Durbin, and Michael L. Littman. A survey of reinforcement learning. In *Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI) 2019*, pages 6317–6317. International Joint Conferences on Artificial Intelligence (IJCAI), 2019. URL <https://doi.org/10.24963/ijcai.2019/880>. URL <https://doi.org/10.24963/ijcai.2019/880>.

ies
beginners
orial series
al new world
a new beginning
vival episodio 1
t survival episodio 1
e a new minecraft bolüm
arch terms used for generating the initial web dataset.

that do not fit our target distribution. In this step, we look for a list of filtering terms that contain these terms. The blacklist use are: {ps3, ps4, ps5, xbox 360, playstation, timelapse, multiplayer, minecraft pe, skyblock, realistic minecraft, how to install, how to download, realmcraft, animation}. described in the next section.

We restrict the scope of this work to the Minecraft Survival mode and therefore limit our training dataset to clips that are obtained from this mode that are relatively free from visual artifacts.

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

16

the image provided in
Minecraft. Everything
marked as None of the
above will be classified as None of the screen.

- No Artifacts: These images will be clean screenshots
- Mode without any noticeable artifacts.
- Mode with some added artifacts: These images will be valid survival
but with some added artifacts. Typical artifacts may include image
brand), text annotations, a picture-in-picture of the player, etc.
- Above: Use this category when the image is not a valid Minecraft survival
mode such as the creative mode or from a different game mode. In non-survival
modes such as the creative mode, the health/hunger bars will be missing from the
item hotbar may or may not be still present.

We spent \$319.96 on human labeling experiments on mTurk, of which \$159.98 was directly
spent by workers. The remaining amount was spent towards Amazon platform fees. The workers
received \$0.01 per labeled image, at an hourly compensation of \$7.20 (based on an estimated labeling
time of 5 seconds/image – in our internal sample run of the same task, we found the average labeling
time to be <3 seconds).
Since we perform rigorous keyword and metadata based filtering of videos (as described in A.1) from
which we served sample images to be labeled, serving offensive content to workers was extremely



The early-game Dataset
The early-game dataset is a ~3000 hour
text that accompanies the videos in w
behavior, i.e. instances where players s
expressions match.

The early-game dataset is a ~3000 hour subset of `web-clean` targeted at “early game” Minecraft. We obtain the metadata from the previous section, i.e. instances where players start in a fresh world with no items. We obtain the following regular expression matches:

elements that are at least 5s in duration. The result of this is our final `web-clean`.
elements that are at least 80% “clean” frames at this stage (Classes Minecra
rtifacts and None of the Above are both considered not clean).
classifier to frames of raw video sequences at a rate of 3 fram
es, we obtain embeddings from the Minecraft
0. We then train a Support Vector Ma
Scikit-learn⁶⁸ SVM implem
ne 2.

Contractor Contract

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

or stone
is used to obtain
coal.
K (and the subtasks
clude personally identifiable
content (e.g. by using Minecraft
so in the contractor videos that we
also potentially learn it, although we expect
it to be likely to reproduce it.

Even the video, any labelled data is appropriate for IDM
gameplay as well as the treechop task described
stages of the project, they were not included in IDM training.

se contains about 420 hours of data. We asked contractors to build a basic
using only basic dirt, wood, and sand, blocks. Each trajectory starts in a newly
ose to begin their trajectories by crafting after a 20 minute time limit. For this task, many
portion of stone tools before to be spent crafting a wooden pickaxe and then mining stone
ortment of stone tools before gathering more building blocks and beginning to create their

Our Minecraft training environment is a hybrid between MineRL²⁷ and the MCP-Reborn
(github.com/Hexception/MCP-Reborn) Minecraft modding package. Unlike the regular Minecraft

C Minecraft environment details

includes almost all actions directly available to human players, such as keypresses, events, and clicks. The specific binary actions we include are shown in Table 3.

which is only useful for entering text into the search bar of the crafting recipe book. Humans either do that or browse the recipe book with the mouse, the latter of which our agent can still do. However, because we do the forward action (e.g. "W" key triggers the agent to press letters that are also shortcuts for actions (e.g. "A", "S", "D", "E", "Q") that produce the forward action) agents are able to press a few keys outside of the GUI (W, A, S, D, E, Q) that produce the forward action). Agents that have not seen agents attempt to search the recipe book with the mouse or craft by dragging items around the crafting window. Instead, our agents navigate the recipe book if the recipe book search bar is selected. We

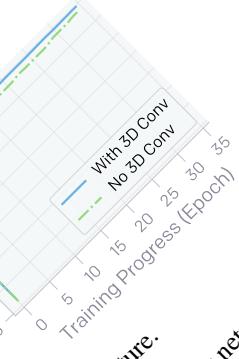
includes almost all actions directly available to human players, such as keypresses, events, and clicks. The specific binary actions we include are shown in Table 3.

which is only useful for entering text into the search bar of the crafting recipe book. Humans either do that or browse the recipe book with the mouse, the latter of which our agent can still do. However, because we do the forward action (e.g. "W" key triggers the agent to press letters that are also shortcuts for actions (e.g. "A", "S", "D", "E", "Q") that produce the forward action) agents are able to press a few keys outside of the GUI (W, A, S, D, E, Q) that produce the forward action). Agents that have not seen agents attempt to search the recipe book with the mouse or craft by dragging items around the crafting window. Instead, our agents navigate the recipe book if the recipe book search bar is selected. We

A collage of screenshots from the game Minecraft. The top left shows a player's hand holding a sword. The top right shows a player standing next to a horse. The bottom left shows a close-up of a sword blade. The bottom right shows a player's hand holding a sword, with a small text overlay reading 'Crafting'.

Machine Learning Dynamics Model Training Details

22



the IDM Architecture.

a ResNet 62 image processing network. In this processing network some temporal information is shared between neighboring frames; however, each stack is comprised of three subsequent convolutional layers in order, (1) an initial 3×3 padding 1 such that the embedding boundary (such that the outgoing dimension) with W output channels, net blocks as defined in He et al. 62 with each layer also having W width and height are flattened into a 1-dimensional vector of size $2^{17} = 131072$ (one activation) such that at this stage there are 128 vectors of size $2^{17} = 131072$ (one future frames, with 32 attention heads of dimension 16384 and another with output dimension 16384 and another with output dimension 128 each and a surrounding non-causal connection skips past this pair). All dense layers have their weights tied through time, so each frame in the video is processed with the same weights.

Finally, independent dense layer heads for each action are pulled from the final embedding – a 2 class on/off categorical parameterized with a softmax for each available key as well as a 11-way

Null Action Filtering

Training

results
mitigated
e data points

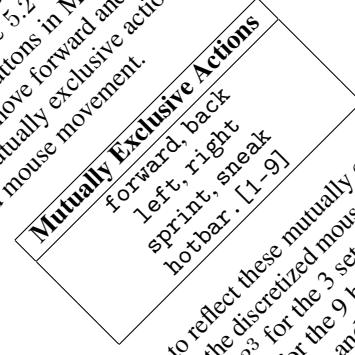
than behavioral cloning
ly predict all 128 actions for
imes at the end of the video clip
s reason, we apply the IDM over a
y use the pseudo-label prediction for
e IDM prediction at the boundary of the
es of a full video.

ture is the same as the IDM architecture described in Appendix
ecture so that it is causal (i.e. cannot see the future when making
d completely). Furthermore, the initial non-causal convolution
ard in language modeling. Additionally, we do Transformer layers are now
to keys and values from past batches and we do Transformer-XL-style⁷⁶ training
relative attention position embedding.

for common action humans take is the null action (no keypresses or mouse movements), which
or something in the game to finish. Among other reasons, a player may take the null action to
water. Early on in the project we found that the BC model would take a break to grab a glass
35% of null actions, often upwards of 95%. In order to prevent this behavior we removed frames than
null actions from the dataset. We compare a few different treatments: we filter nulls if there have
been 1, 3, or 21 frames of consecutive null actions, and include a treatment that does not perform
any null filtering. Null action filtering generally helps, increasing all crafting rates (Figure 15 left).

ng
more
ruments

be independently on
en state. This could cause
(b) move left and drop their
d independently of each of the 4 constituent
actions; however, the full joint distribution
dimensions. This is far too large for many reasons, e.g.
 5.2×10^{11} actions. Furthermore, the inventory button is
move forward and backward have no effect when simultaneously
mouse movement.



the joint action space to reflect these mutually exclusive combinations still results in a
neither in the set is an option, $\times 10$ for the 3 sets of 2 mutually exclusive movements, i.e. $3^3 \times 10 \times 2^4 \times 11^2 + 1 \approx$
remaining binary 4 keys: use, drop, attack, and jump, $\times 11^2$ or no hotbar keypress, $\times 2^4$ for the
quite large so we chose to implement a secondary hierarchical binary action for camera being moved or not. If
this action is on, then there is a secondary discrete action head with 121 classes (the joint distribution
of mouse movements because each discretized mouse direction has 11 bins) that determines where

ing
asted in

training

in model training, except we either use a focused
with ground-truth labels, or contractor data
Table 5. We used 16 A100 GPUs. The hyperparameters used
et, and 16 A100 GPUs for about 2 days when fine-tuning

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

Table 5: Hyperparameters for behavior cloning fine-tuning

1 Reinforcement Learning Fine-Tuning

RL experiments were performed with the phasic policy gradient (PPG) algorithm,⁶⁴ an RL algorithm
based on the proximal policy optimization (PPO) algorithm⁷⁷ that increases sample efficiency by
performing additional passes over the collected data to optimize the value function as well as an

BC pre-training, they are trained as a single, fully connected model (Appendix D.1). While the weights of the regular function. To prevent the value-function divergence by the standard deviation, we normalize the standard deviation, which are

the RL model and the frozen pretrained policy.¹⁰ This loss is defined, π_{pt} is the frozen pretrained policy, π_θ is the policy being trained and the pretrained policy, $KL(\pi_{pt}, \pi_\theta)$ is the KL divergence loss replaces the common entropy maximization loss, this KL divergence loss is the frozen pretrained policy, and ρ is a coefficient to encourage exploration.⁷⁹⁸⁰ The idea behind entropy maximization is to have equal entropy to increase the chance that the agent has not explored all actions appear to have equal exploration.⁷⁹⁸⁰ The entropy maximization loss is sufficiently small or the reward is sufficiently large and rewards are sparse, which is the case in the diamond-pickaxe action distribution. In experiments with a randomly initialized policy instead of a uniform-random action distribution, it should mimic the action-distribution in states where the agent assigns equal value to each of its actions, it is much more likely to take sequences of actions that lead to other losses.

Empirically, we found that a coefficient of 0.01, which has been an effective setting in other Minecraft work.³⁰ Interestingly, we found that a high coefficient ρ for this KL divergence loss while a low coefficient ρ was ineffective at encouraging the agent from properly optimizing the reward function while a low coefficient ρ was ineffective at

ing catastrophic forgetting. As such, guaranteeing that the policy can eventually fix its behavior has to be do so relative to current quantities of each item that a human player might have, and we reward the model for gathering up to that item function, then we added the 1 iron pickaxe required for mining diamonds, and torches to the reward function, with coal being useful as fuel, and torches to the reward function, with 3 diamonds and 2 sticks required for crafting the table or sticks if the agent runs out. In practice the agent rarely collects 8 logs, places the torches themselves improve visibility and an RL model expectations on what would be useful to execute this task, rather than collects 5 logs and an RL model behaves after training. Finally, to encourage the agent to keep mining diamonds or diamond pickaxes after it has crafted its first diamond the agent to keep mining base reward of 1, the second tier consists of all items requiring coal with a base reward of 4, and the final tier is diamond pickaxe generally on the number of diamonds or diamond pickaxes after training. Thus items later in the sequence of items towards a diamond pickaxe generally would usually get the different items are separated into 4 tiers, roughly depending on how late a player would reward for the relevant item. The first tier consists of all items requiring coal with a base reward of 4, and the final tier is diamond pickaxe, we did not put base reward of 8. Thus items later in the sequence of items towards a diamond pickaxe generally

the fact that it can take thousands of seconds to craft all the necessary prerequisites and a diamond after crafting an iron bar. For example, a good strategy for getting one item immediate access to a crafting table as soon as it is gathered cobblestone, while leaving the crafting table behind. After 6 days (144 hours) on 80 GPUs (for policy optimization) and 16 days (384 hours) on 80 GPUs (for fine-tuning), the agent has learned rollouts from Minecraft. In this time the algorithm performed 1.4 million Minecraft episodes consisting of 16.8 billion frames.

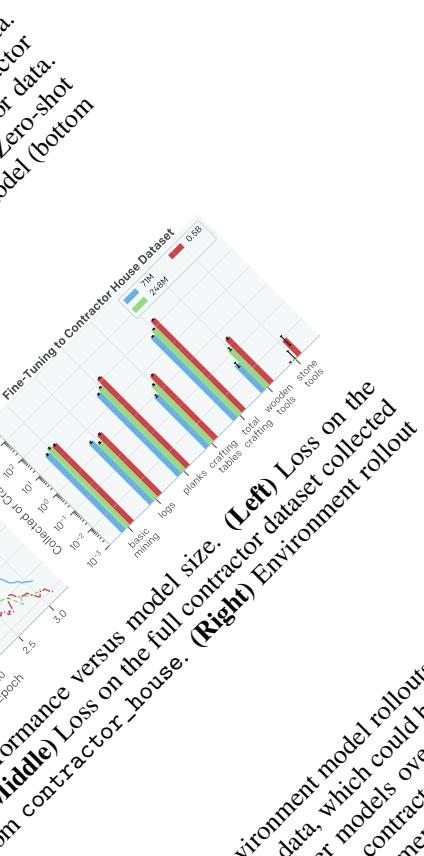
Learning Fine-Tuning Additional Data

In this section, we show the main results of the RL fine-tuning experiments that are helpful for understanding the main items-over-training figure (Fig. 16). When training figure without a KL loss, the model is capable of getting zero-shot, which directly compare RL fine-tuning from the early-game model with the early-game model (Fig. 17). These experiments differ from the house-building model and RL fine-tuning from the early-game model that the early-game model is trained without a KL loss (Fig. 16). Second, we present preliminary experiments in which treatments shown here, the KL loss coefficient was set to 0.4, the learning rate was set to 6×10^{-5} , and the reward for each item was 1/quantity for all items (i.e. items closer to the diamond pickaxe did not have an increased reward). While RL fine-tuning from the house-building model initially

y-game model compared to initially increases faster when game model eventually obtains a higher likelihood was chosen for future RL fine-tuning

el. Here we compare the 0.5B model staying in the efficient M parameter model has 1/2 the width and each layer in the 0.5B validation loss on web-clean with ground truth labels collected during contractor play, and zero-shot performance for the 71M, 248M, and 0.5B models. While the 71M model even had non-zero wooden tool crafting (Fig. 18 bottom left), and also has the best zero-shot environment dataset loss while having the 0.5B, and also has the best zero-shot contractor dataset loss. In fact, we see fine-tuning to contractor-house, model size rank ordering reverses and now the 0.5B model performs best both in validation loss (Fig. 19 left) and in environment performance (Fig. 19 right)





that because the models that are a better overall Minecraf t have worse contractor data, which could be visually distinct from our game engine, resulting in better environment loss in Fig. 18 top left) can quickly shift their low level features to perform better on the rollout features. After just a few steps of fine-tuning, which has no overlap with contractor-house, all models quickly improve in loss on the contractor dataset collected from the web. It is plausible that the larger models perform more poorly in the environment due to zero-shot. However, we further supported by Fig. 19 (middle) showing loss on the contractor-house, all models now performing best. While not conclusive, we believe this investigation provides some intuition for future studies of model scaling for sequential decision making problems.

in
model:
ically wooden
e in Section 4.4
work³⁰. An alternate
reconceive of the task and
s on computers or in simulated
VPT+text could conceivably produce
the powers of GPT to meta-learn, follow
in the form of agents that can act in virtual
similar embodied agents that can act in virtual
at we began a first step towards decision domains.
mentary from the player. This commentary is sometimes
the videos, or could be extracted post-hoc using automated
set features about 17k hours of content with associated closed
parameter VPT foundation model used in the RL-fine-tuning ex-
for the same reason: to reduce compute costs) with an additional ex-
first split videos into 30 second chunks. The same text is associated with every
closed caption and following the closed captions. The same text is available. To obtain the
preceding and is made up of all the closed chunks (if any). Because the vast majority (around
the subset of our data for which closed captions are available. To obtain the
lacked capitalization and punctuation (if any).
e then obtain a text embedding vector of length 4,096 from the OpenAI embedding API⁸⁹
processed by a randomly initialized multi-layer perceptron (MLP) with two hidden layers of
2,048. The resulting activations are added for each frame to the pretrained using the punct
model is fine-tuned for four epochs.
ore the transformer layers (pretrained for four epochs). When conditioned on sentences that incite the agent to
explore (such as ‘I’m going to explore’ and ‘I’m going to find water’ the agent travels significantly
farther from its spawn point (Figure 20a). Additionally, we can steer the agent to preferentially collect

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

Table 8: Strings corresponding to each conditioning variant.

