
MINEDOJO Datasheet

Linxi Fan¹, Guanzhi Wang^{2*}, Yunfan Jiang^{3*}, Ajay Mandlekar¹, Yuncong Yang⁴,
Haoyi Zhu⁵, Andrew Tang⁴, De-An Huang¹, Yuke Zhu^{1 6†}, Anima Anandkumar^{1 2†}

¹NVIDIA, ²Caltech, ³Stanford, ⁴Columbia, ⁵SJTU, ⁶UT Austin

*Equal contribution †Equal advising

<https://minedojo.org>

1 Motivation

For what purpose was the dataset created? We create this internet-scale multimodal knowledge base to facilitate research towards open-ended, generally capable embodied agents.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)? This knowledge base was created by Linxi Fan (Nvidia), Guanzhi Wang (Caltech), Yunfan Jiang (Stanford), Ajay Mandlekar (Nvidia), Yuncong Yang (Columbia), Haoyi Zhu (SJTU), Andrew Tang (Columbia), De-An Huang (Nvidia), Yuke Zhu (Nvidia and UT Austin), and Anima Anandkumar (Nvidia and Caltech).

Who funded the creation of the dataset? NVIDIA Corporation.

2 Author Statement

We bear all responsibilities for the licensing, distribution, and maintenance of our datasets. This document follows the Datasheet format [1] whenever applicable.

3 Distribution

Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? Yes, the dataset is publicly available on the internet.

How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? All datasets can be downloaded from <https://zenodo.org/>. Please refer to this table of URL, DOI, and licensing:

Database	DOI	License
YouTube	10.5281/zenodo.6641142	Creative Commons Attribution 4.0 International (CC BY 4.0)
Wiki	10.5281/zenodo.6640448	Creative Commons Attribution Non Commercial Share Alike 3.0 Unported
Reddit	10.5281/zenodo.6641114	Creative Commons Attribution 4.0 International (CC BY 4.0)

Have any third parties imposed IP-based or other restrictions on the data associated with the instances? No.

Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? No.

25 **4 Maintenance**

26 **Who will be supporting/hosting/maintaining the dataset?** The authors will be supporting,
27 hosting, and maintaining the dataset.

28 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?** Please
29 contact Linxi Fan (linxif@nvidia.com), Guanzhi Wang (guanzhi@caltech.edu), and Yunfan
30 Jiang (yunfanj@cs.stanford.edu).

31 **Is there an erratum?** No. We will make announcements if there is any.

32 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete
33 instances)?** Yes. New updates will be posted on <https://minedojo.org>.

34 **If the dataset relates to people, are there applicable limits on the retention of the data associated
35 with the instances (e.g., were the individuals in question told that their data would be retained
36 for a fixed period of time and then deleted)?** N/A.

37 **Will older versions of the dataset continue to be supported/hosted/maintained?** Yes, old
38 versions will be permanently accessible on zenodo.org.

39 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for
40 them to do so?** Yes, please refer to <https://minedojo.org>.

41 **5 Composition**

42 **What do the instances that comprise the dataset represent?** For YouTube videos, our data is in
43 JSON format with video URLs and metadata. We do not provide the raw MP4 files for legal concerns.
44 For Wiki, we provide the text, images, tables, and diagrams embedded on the web pages. For Reddit,
45 our data is in JSON format with post IDs and metadata, similar to YouTube. Users can reconstruct
46 the Reddit dataset by running our script after obtaining an official Reddit API license key.

47 **How many instances are there in total (of each type, if appropriate)?** There are more than 740K
48 YouTube videos with 2.2B words of transcripts, 6,735 Wiki pages with 5.8M bounding boxes of
49 visual elements, and more than 350K Reddit posts with 6.8M comments.

50 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of
51 instances from a larger set?** We provide all instances in our Zenodo data repositories.

52 **Is there a label or target associated with each instance?** No.

53 **Is any information missing from individual instances?** No.

54 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social
55 network links)?** We provide metadata for each YouTube video link and Reddit post ID.

56 **Are there recommended data splits (e.g., training, development/validation, testing)?** No. The
57 entire database is intended for pre-training.

58 **Are there any errors, sources of noise, or redundancies in the dataset?** Please refer to Sec.
59 4 in [supplementary.pdf](#)

60 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
61 **websites, tweets, other datasets)?** We follow prior works [2] and only release the video URLs
62 of YouTube videos due to legal concerns. Researchers need to acquire the MP4 and transcript files
63 separately. Similarly, we only release the post IDs for the Minecraft Reddit database, but we also
64 provide a script that can reconstruct the full Reddit dataset given a free official license key.

65 **Does the dataset contain data that might be considered confidential?** No.

66 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,**
67 **or might otherwise cause anxiety?** We have made our best efforts to detoxify the contents via
68 an automated procedure. Please refer to Sec. 4 in `supplementary.pdf`.

69 **6 Collection Process**

70 The collection procedure, preprocessing, and cleaning are explained in details in Sec. 4 of
71 `supplementary.pdf`

72 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and**
73 **how were they compensated (e.g., how much were crowdworkers paid)?** All data collection,
74 curation, and filtering are done by MINEDOJO coauthors.

75 **Over what timeframe was the data collected?** The data was collected between Dec. 2021 to
76 May 2022.

77 **7 Uses**

78 **Has the dataset been used for any tasks already?** Yes, we have used the MINEDOJO YouTube
79 database for agent pre-training. Please refer to Sec. 5 in our main paper.

80 **Is there a repository that links to any or all papers or systems that use the dataset?** N/A.

81 **What (other) tasks could the dataset be used for?** Our knowledge base is primarily intended
82 to facilitate research in open-ended, generally capable embodied agents. However, it can also be
83 broadly applicable to research in video understanding, document understanding, language modeling,
84 multimodal learning, and so on.

85 **Is there anything about the composition of the dataset or the way it was collected and**
86 **preprocessed/cleaned/labeled that might impact future uses?** Nothing we are aware of.

87 **Are there tasks for which the dataset should not be used?** We do not condone any research
88 that intentionally generates harmful or toxic contents using our YouTube, Wiki, and Reddit data.

89 **References**

90 [1] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M.
91 Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Commun. ACM*, 64(12):
92 86–92, 2021. doi: 10.1145/3458723. URL <https://doi.org/10.1145/3458723>.

93 [2] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
94 narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew
95 Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv: Arxiv-1705.06950*,
96 2017.