
ReMI: A Dataset for Reasoning with Multiple Images

— Supplementary Material

Mehran Kazemi¹, Nishanth Dikkala², Ankit Anand¹, Petar Devic³,
Ishita Dasgupta¹, Fangyu Liu¹, Bahare Fatemi², Pranjal Awasthi²,
Dee Guo², Sreenivas Gollapudi², Ahmed Qureshi³
¹Google DeepMind, ²Google Research, ³Google

1 Dataset documentation: datasheet

In this section, we follow the recommendations in Gebru et al. [1] to provide comprehensive documentation for our dataset, available at <https://huggingface.co/datasets/mehrankazemi/ReMI>. The dataset croissant can be found at <https://huggingface.co/api/datasets/mehrankazemi/ReMI/croissant>.

1.1 Motivation

- **For what purpose was the dataset created?** Reasoning with multiple images (possible interleaved with text) is an important emerging capability of the LLMs; we created a dedicated dataset to enable tracking progress in this area.
- **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?** The dataset was created by a team of researchers at Google all listed as authors of the paper.
- **Who funded the creation of the dataset?** The dataset was created by Google employees as part of their work at the company, so Google funded its development.

1.2 Composition

- **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?** Each instance within the dataset represents a textual question as well as multiple accompanying photos. The photos have a wide range: matplotlib charts and function graphs, NetworkX graphs, latex rendered tables and shapes, screenshots from Google Maps, images from COCO [2], etc.
- **How many instances are there in total (of each type, if appropriate)?** ReMI contains 13 tasks, each having 200 examples for evaluation. All datasets contain photos.
- **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?** Parts of the dataset have been created programatically. So it is possible to create larger sets. But the released dataset corresponds to the entire data that was used for the paper.
- **What data does each instance consist of?** Each example in the dataset contains a textual reasoning question with image markers representing where the images must be inserted, a list of accompanying images, and the label.
- **Is there a label or target associated with each instance?** Each data point has a label corresponding to the ground truth response to be used for evaluation.
- **Is any information missing from individual instances?** No. All instances have both their text, images and label available.

- 34 • **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?** There is no relationship between different instances.
- 35
- 36 • **Are there recommended data splits (e.g., training, development/validation, testing)?**
- 37 This dataset serves as an evaluation benchmark. We have released a test set and a fewshot
- 38 set.
- 39 • **Are there any errors, sources of noise, or redundancies in the dataset?** The dataset has
- 40 been quality checked over multiple rounds by multiple authors until no errors were found.
- 41 • **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
- 42 **(e.g., websites, tweets, other datasets)?** The dataset is self-contained and does not link to
- 43 external resources.
- 44 • **Does the dataset contain data that might be considered confidential (e.g., data that**
- 45 **is protected by legal privilege or by doctor-patient confidentiality, data that includes**
- 46 **the content of individuals' non-public communications)?** The dataset does not contain
- 47 any confidential or sensitive information.
- 48 • **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
- 49 **threatening, or might otherwise cause anxiety?** The dataset does not contain any data
- 50 that could be considered offensive, insulting, threatening, or anxiety-inducing.
- 51 • **Does the dataset identify any subpopulations (e.g., by age, gender)?** The dataset does
- 52 not identify or contain any information that would allow for the identification of subpopula-
- 53 tions based on attributes.
- 54 • **Is it possible to identify individuals (i.e., one or more natural persons), either directly**
- 55 **or indirectly (i.e., in combination with other data) from the dataset?** The dataset does
- 56 not contain any information that could be used to directly or indirectly identify individuals,
- 57 either on its own or in combination with other data.
- 58 • **Does the dataset contain data that might be considered sensitive in any way (e.g., data**
- 59 **that reveals race or ethnic origins, sexual orientations, religious beliefs, political opin-**
- 60 **ions or union memberships, or locations; financial or health data; biometric or genetic**
- 61 **data; forms of government identification, such as social security numbers; criminal**
- 62 **history)?** The dataset does not contain any sensitive data that could reveal attributes.

63 1.3 Collection

- 64 • **How was the data associated with each instance acquired?** The text of the questions
- 65 have been created automatically using code, or, in the case of one task, by prompting a
- 66 language model. The images are acquired from three sources: 1- Automatically generated
- 67 using visualization libraries, 2- screenshots from Google Maps, and 3- from the COCO
- 68 dataset.
- 69 • **What mechanisms or procedures were used to collect the data (e.g., hardware appa-**
- 70 **ratutes or sensors, manual human curation, software programs, software APIs)?** The
- 71 data collection process primarily involved a combination of automated procedures and
- 72 manual human input. The majority of the data collection was performed using software
- 73 programs and scripts that were developed and executed by the authors. These programs
- 74 included algorithms and techniques designed to generate and curate the specific types of
- 75 data required for the dataset.
- 76 • **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,**
- 77 **deterministic, probabilistic with specific sampling probabilities)?** NA
- 78 • **Who was involved in the data collection process (e.g., students, crowdworkers, con-**
- 79 **tractors) and how were they compensated (e.g., how much were crowdworkers paid)?**
- 80 The data collection process was conducted solely by the authors of the paper.
- 81 • **Over what timeframe was the data collected?** The data collection process took place
- 82 over a four-month period.
- 83 • **Were any ethical review processes conducted (e.g., by an institutional review board)?**
- 84 No.
- 85 • **Did you collect the data from the individuals in question directly, or obtain it via third**
- 86 **parties or other sources (e.g., websites)?** We collected data from no individuals.

- 87 • **Were the individuals in question notified about the data collection?** NA.
- 88 • **Did the individuals in question consent to the collection and use of their data?** NA.
- 89 • **If consent was obtained, were the consenting individuals provided with a mechanism**
- 90 **to revoke their consent in the future or for certain uses?** NA.
- 91 • **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
- 92 **a data protection impact analysis) been conducted?** NA.

93 1.4 Uses

- 94 • **Has the dataset been used for any tasks already?** The dataset is used for benchmarking
- 95 model performances on multi-image reasoning, in this paper.
- 96 • **Is there a repository that links to any or all papers or systems that use the dataset?**
- 97 The dataset is generated and used in this paper only.
- 98 • **What (other) tasks could the dataset be used for?** To measure model performance on
- 99 multi-image reasoning.
- 100 • **Is there anything about the composition of the dataset or the way it was collected and**
- 101 **preprocessed/cleaned/labeled that might impact future uses?** No.
- 102 • **Are there tasks for which the dataset should not be used?** The main purpose of the
- 103 dataset is to be used for evaluation and it should not be used for training.

104 1.5 Distribution

- 105 • **Will the dataset be distributed to third parties outside of the entity (e.g., company,**
- 106 **institution, organization) on behalf of which the dataset was created?** Yes, the dataset
- 107 is available publicly in Huggingface.
- 108 • **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** The
- 109 dataset is distributed through Huggingface.
- 110 • **When will the dataset be distributed?** The dataset is already available.
- 111 • **Will the dataset be distributed under a copyright or other intellectual property (IP)**
- 112 **license, and/or under applicable terms of use (ToU)?** In the spirit of open science and
- 113 collaboration, we have released the datasets under a Creative Commons Attribution 4.0
- 114 International (CC BY 4.0) license. For comprehensive details about the terms of the CC
- 115 BY 4.0 license, please visit the Creative Commons website: <https://creativecommons.org/licenses/by/4.0/>.
- 116
- 117 • **Have any third parties imposed IP-based or other restrictions on the data associated**
- 118 **with the instances?** Please refer to the copyright.
- 119 • **Do any export controls or other regulatory restrictions apply to the dataset or to indi-**
- 120 **vidual instances?** Please refer to the copyright.

121 1.6 Maintenance

- 122 • **Who will be supporting/hosting/maintaining the dataset?** The authors will support, host,
- 123 and maintain the dataset.
- 124 • **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
- 125 The owner (Mehran Kazemi) can be contacted through mehrankazemi@google.com.
- 126 • **Is there an erratum?** No. If errors are found in the future, we will release errata using the
- 127 same link.
- 128 • **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
- 129 **instances)?** If any issues are found in the dataset, it will be updated to ensure correctness.
- 130 • **If the dataset relates to people, are there applicable limits on the retention of the data**
- 131 **associated with the instances (e.g., were the individuals in question told that their data**
- 132 **would be retained for a fixed period of time and then deleted)?** NA.
- 133 • **Will older versions of the dataset continue to be supported/hosted/maintained?** Yes,
- 134 older versions of the dataset will continue to be maintained and hosted.

135 • **If others want to extend/augment/build on/contribute to the dataset, is there a mecha-**
136 **nism for them to do so?** Yes, others are welcome to extend, augment, build on, or contribute
137 to the dataset. They are encouraged to download the dataset, create their own modified
138 versions, and publish their work on their preferred platform.

139 **2 Accessibility**

140 The datasets created in this research are available for download at [https://huggingface.co/](https://huggingface.co/datasets/mehrankazemi/ReMI)
141 [datasets/mehrankazemi/ReMI](https://huggingface.co/datasets/mehrankazemi/ReMI). In the spirit of open science and collaboration, we have released
142 the datasets under a Creative Commons Attribution 4.0 International (CC BY 4.0) license. For
143 comprehensive details about the terms of the CC BY 4.0 license, please visit the Creative Commons
144 website: <https://creativecommons.org/licenses/by/4.0/>.

145 **3 Reproducibility**

146 In our pursuit of reproducibility, we conducted experiments with five prominent models: Gemini
147 Ultra, Gemini 1.5 Pro, Gemini Flash, Claude Sonnet, and GPT4 Turbo. All of these variants are
148 publicly accessible through various platforms, ensuring transparency and enabling further research.
149 For the Gemini and Claude models, we utilized the Google Cloud Platform (Vertex AI) as our
150 computational infrastructure to execute the experiments. For GPT4, we used the OpenAI API.

151 The author bear all responsibility in case of violation of rights.

152 **References**

- 153 [1] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,
154 Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64
155 (12):86–92, 2021.
- 156 [2] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
157 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
158 *Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014,*
159 *Proceedings, Part V 13*, pages 740–755. Springer, 2014.