# FIT400M Data Card

## FIT400M

**This doc:** GitHub Repo link
**Dataset:** GitHub Repo
**Data Card Authors:**
Burcu Karagol Ayan ,
Emily Denton

### DATASET SUMMARY

This data card describes the FIT400M (Filtered Image-Text 400M) dataset, a cleaned-up version of the AltText dataset.

FIT400M is an internal Google dataset, intended to be used for research purposes by Google teams as a high-quality, large-scale resource for training language and image models, including text-and-image dual encoders and text-to-image generation models.

## Dataset Owners

| TEAM(S) | CONTACT DETAIL(S) | AUTHOR(S) |
|---|---|---|
| FIT400M Project | **Dataset Owner(s):** FIT400M Project<br>**Affiliation:** Google Research<br>**Contact:** fit400m@google.com | <ul><li>Burcu Karagol Ayan, Software Engineer, Google Research</li><li>Yinfei Yang, Research Scientist, (during his time at Google)</li><li>Jason Baldridge, Research Scientist, Google Research</li></ul> |

## Dataset Overview

| DATA SUBJECT(S) | DATASET SNAPSHOT | | CONTENT DESCRIPTION |
|---|---|---|---|
| Non-Sensitive Data about people<br><br>Data about places and objects | FIT400M is a static snapshot. Both images and text are static. | | An image and associated text description coming from the alt-text field.<br><br>FIT400M dataset contains 400 million images and associated text, which are noisily labeled with multiple labels that are algorithmically inferred. In addition to |
| | **Size of dataset** | 10.83TB | |
| | **Number of Instances** | ~400,000,0( | |
| | **Number of Fields** | 7 | |
| | **Labeled Classes** | N/A | |

| Number of Labels | Variable[1] | these labels, each data point is associated with an image_id, image height and image width. |
| --- | --- | --- |
| Algorithmic Labels | 5[2] | |
| Human Labels | Unavailable[3] | |

*Above:* Summary of FIT400M dataset.

[1] Some fields (such as image pixels, text, image width and height) have variable values, no classes. Some internal-google-generated labels have tens of thousands of possible values. 'label_score' is a float value.

[2] 'smeared_image_labels' and 'label_score' fields are machine generated.

[3] All labels are either from the original Alt-Text dataset or from other classifiers. There are no human-annotated labels in this dataset.

- It stands for **Filtered Image-Text 400M**.

- The dataset is a cleaned version of the noisy AltText dataset that was used for training ALIGN ([paper](#)) and MURAL ([paper](#)) models.

## Sensitivity Of Data Fields

| SENSITIVITY TYPE(S) | FIELD(S) WITH SENSITIVE DATA | |
| --- | --- | --- |
| None | **Intentionally Collected Sensitive Data**<br>No sensitive data was intentionally collected.<br><br>**Unintentionally Collected Sensitive Data**<br>S/PII, pornographic content, or images depicting violence were not explicitly collected as a part of the dataset creation process and any fields we found may contain such information have been filtered.<br><br>We used algorithmic methods and relied on other classifiers for identifying S/PII information, pornographic and violence depicting images hence it is possible we may have missed some instances in the process. Specifically we filtered (1) any address, email and phone information; (2) images with high porn scores and (3) images labeled as portraying abuse; (4) text identified as having certain adult content references.<br><br>Fields that may contain such sensitive data are:<br>    - image_data (pixels of the images)<br>    - object_groundtruth (associated text) | |

## Version And Maintenance

| MAINTENANCE STATUS | DATASET VERSION | MAINTENANCE PLAN |
| --- | --- | --- |

| Limited Maintenance | Current Version: 1.0 | FIT400M is a static dataset from a specific point in time and maintenance will be limited. |
|---|---|---|
| The data will not be updated, but any technical issues will be addressed. | Last Updated: 08/2021 | |
| | Release Date: N/A | Feedback: For feedback, reach out to fit400m@google.com. |

## Example Of Data Points

| PRIMARY DATA MODALITY | LINK(S) TO DATA POINT(S) | DATA FIELDS |
|---|---|---|
| Multimodal | Below are examples of kind data in the FIT400M dataset.<br><br><br><br>An injured dog with a cone walking outside | \| Field name \| Type \| Description \|<br>\| ----------------- \| ------\|----------- \|<br>\| `image_id` \| Integer \| Unique id for the data point. \|<br>\| `image_data` \| Bytes \| The pixel data for the image. \|<br>\| `image_meta_data/width` \| Integer \| Width of the image. \|<br>\| `image_meta_data/height` \| Integer \| Height of the image. \|<br>\| `raw_text` \| Bytes \| List of text associated with the image. \|<br>\| `smeared_image_labels` \| Bytes \| Machine generated image labels.These were used for further sampling the data. \|<br>\| `label_score` \| Float \| ALIGN model score that shows the semantic similarity of the text and the image. \| |

A man is cutting carrots

Below is a typical data point.

| Field name | Value |
| --- | --- |
| image_id | 0x0001c2dc4950cb12 |
| image_data | "\xFF\xD8\xFF\xE0\x00\x10JFIF\x00\x01\x01\x00\x00\x01\x00\x01\x...\" |
| image_meta_data/width | 678 |
| image_meta_data/height | 452 |
| raw_text | "person washing hands from faucet outdoors" |
| smeared_image_labels | ["Descriptive", "Indoor", "Drink", ...] |
| label_score | 0.24681078 |

The dataset does not contain atypical data points as far as we know.

## Provenance

### Data Collection & Sources

| METHOD(S) USED | METHODOLOGY DETAIL(S) | SOURCE DESCRIPTION(S) |
|---|---|---|
| Scraped or Crawled<br><br>Taken from other existing datasets | **Taken from other existing datasets:**<br><br>**Source :** FIT400M dataset is a cleaned version of the [AltText dataset](). The AltText dataset was created by applying minimal frequency-based filtering to the image-alt text pairs, following a similar procedure to the Conceptual Captions dataset ([Sharma et al., 2018]()).<br><br>**Is this source considered sensitive or high-risk?** [Yes / **No**]<br><br>**Dates of Collection:** [unknown to 2020-11-13]<br><br>**Primary modality of collected data:**<br>• Multimodal (Image and text)<br><br>**Update Frequency for collected data:**<br>• Static | • **AltText dataset:** A large, noisy image-text pair dataset. It contains 1.8B image-text pairs. |

| COLLECTION CADENCE | DATA INTEGRATION | DATA PROCESSING |
|---|---|---|
| Static<br>Data was collected once from single or multiple sources. | AltText dataset<br><br>**Included Fields**<br>All the fields except 'smeared_image_labels' and 'label_score' are coming from the AltText dataset.<br><br>**Excluded Fields**<br>None | All data is coming from AltText with some added machine-generated fields. |

### Collection Criteria

| DATA SELECTION | DATA INCLUSION | DATA EXCLUSION |
|---|---|---|
| Records from the AltText dataset are chosen according to the following criteria:<br><br>• **No sensitive data:** Using several machine generated signals, | Records that are not excluded are in the final dataset. | The following filters were applied:<br><br>• Records identified as containing address, email or phone numbers are excluded. |

- records identified as having sensitive information are excluded.

- **English only text:** Using machine generated language identification signals, records with high confidence non-English text are excluded.

- **No adult content:** Using both image and text based machine generated signals, records identified as having adult content are excluded.

- **Low image-text semantic alignment:** Using the ALIGN model, records that are identified as having low semantic alignment score are excluded.

- **Text consisting of mainly numbers:** Records with text fields that are mostly numbers are excluded.

- Text identified as non-English with a confidence > 0.7 are excluded.

- Images with porn score > 0.7 are excluded.

- Text identified as having certain adult content references are excluded.

- Records with ALIGN score < 0.21 were excluded.

- Texts that are mostly numbers are excluded.

## Data Sampling

| METHOD(S) USED | CHARACTERISTIC(S) | SAMPLING CRITERIA |
|---|---|---|
| Multi-stage Sampling<br>Random Sampling<br>Stratified Sampling | **Stratified Sampling**<br>Upstream Source        AltText dataset cleaned according to the criteria described in Collection Criteria section.<br>Total data sampled        ~834,220,000<br>Sample size                ~825,400,000<br><br>Using the smeared_image_labels, data is sampled to make sure all tail labels are present.<br><br>**Random Sampling**<br>Upstream Source        Dataset version coming from the stratified sampling described above. | First 'smeared_image_labels' are used to sample the data making sure all tail labels are represented in the cleaned up version mentioned above.<br>Then 400M records were randomly selected. |

| | |
|---|---|
| Total data sampled | ~825,400,000 |
| Sample size | ~400,000,000 |

400M records are randomly selected.

## Sociodemographic Information

| SOCIODEMOGRAPHIC CATEGORIES LABELED IN THE DATASET | INTENTIONALITY |
|---|---|
| None | **Intentionally Collected Information**<br>*No sociodemographic information was labeled or collected as a part of the dataset creation process.*<br><br>**Unintentionally Collected Information**<br>*Sociodemographic information was not explicitly collected as a part of the dataset creation process. However, text in the dataset might reference sociodemographic information and it may be possible to infer some sociodemographic information from image pixels or the text field.*<br><br>Fields that may be used to infer sociodemographic information:<br>    - image_data (pixels of the images)<br>    - object_groundtruth (associated text)<br>*Sociodemographic information was not explicitly collected as a part of the dataset creation process. However, text in the dataset might reference sociodemographic information and it may be possible to infer some sociodemographic information from image pixels or the text field.*<br><br>Fields that may be used to infer sociodemographic information:<br>    - image_data (pixels of the images)<br>    - object_groundtruth (associated text) |

## Transformations

| TRANSFORMATION(S) APPLIED | FIELD(S) TRANSFORMED | LIBRARY(IES) AND METHOD(S) USED |
|---|---|---|
| Others (Cleaning text fields) | **Cleaning text fields**<br><br>Some irrelevant text is removed from the 'raw_text' field (and all other fields that are replicas of it). These are very frequently occurring text pieces that do not add value to the text semantics. | **Cleaning text fields**<br><br>**Method:** Using hand written rules, certain prefixes and suffixes are removed from the 'raw_text' field. If the remaining text is too short (less than 3 tokens), the record is not included in the final version of the dataset.<br><br>**Platforms, tools, or libraries:**<br>• Find most common prefixes and suffixes<br>• Remove the irrelevant prefixes and suffixes<br><br>**Transformation Results**<br>Below are some samples of text removed:<br>• free png download<br>• image \\d+ of \\d+<br>• jpeg |

## Breakdown Of Transformations

### Other Transformations (Cleaning Text Field)

| DESCRIPTION | METHOD(S) USED | COMPARATIVE SUMMARY |
|---|---|---|
| 'raw_text' is the field that holds the text associated with the image. Certain text is removed using regular expressions.<br><br>Other fields that are really replicas of the 'raw_text' field are also updated accordingly. | **Platforms, tools, or libraries:**<br>• Find most common prefixes and suffixes<br>• Remove the irrelevant prefixes and suffixes | Below are some samples of text removed:<br>• free png download<br>• image \\d+ of \\d+<br>• jpeg |