

# MOSCAR: A LARGE-SCALE MULTILINGUAL AND MULTIMODAL DOCUMENT-LEVEL CORPUS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Multimodal Large Language Models (mLLMs) are trained on a large amount of text-image data. While most mLLMs are trained on caption-like data only, Alayrac et al. (2022) showed that additionally training them on interleaved sequences of text and images can lead to the emergence of in-context learning capabilities. However, the dataset they used, M3W, is not public and is only in English. There have been attempts to reproduce their results but the released datasets are English-only. In contrast, current multilingual and multimodal datasets are either composed of caption-like only or medium-scale or fully private data. This limits mLLM research for the 7,000 other languages spoken in the world. We therefore introduce mOSCAR, to the best of our knowledge the first large-scale multilingual and multimodal document corpus crawled from the web. It covers 163 languages, 303M documents, 200B tokens and 1.15B images. We carefully conduct a set of filtering and evaluation steps to make sure mOSCAR is sufficiently safe, diverse and of good quality. We additionally train two types of multilingual model to prove the benefits of mOSCAR: (1) a model trained on a subset of mOSCAR and captioning data and (2) a model trained on captioning data only. The model additionally trained on mOSCAR shows a strong boost in few-shot learning performance across various multilingual image-text tasks and benchmarks, confirming previous findings for English-only mLLMs. The dataset will be made publicly accessible.

## 1 INTRODUCTION

Multimodal large language models (mLLMs) are trained on large amounts of text-image data (Radford et al., 2021; Yu et al., 2022; Li et al., 2023; Wang et al., 2023; OpenAI, 2023; Gemini Team et al., 2023; Chameleon Team, 2024). The main paradigm until recently was to train a model from a large collection of web-crawled images and their captions (Li et al., 2021; Wang et al., 2022; Chen et al., 2023b). Models such as Flamingo (Alayrac et al., 2022) challenged this paradigm by being additionally trained on interleaved sequences of text and images from web documents, showing state-of-the-art results on various tasks and in-context learning capabilities that are not present in models trained on caption-like data only. Additionally, McKinzie et al. (2024) recently proved that including interleaved text-image data during training was necessary to get good few-shot learning performance. However, the datasets used to train mLLMs are either private (Alayrac et al., 2022), monolingual or multilingual but only medium-scale (Srinivasan et al., 2021). Some attempts have been made to reproduce these datasets (Zhu et al., 2023; Laurençon et al., 2023) but the resulting datasets are only available in English.

Few image-text datasets are multilingual and most of them are obtained by translating English caption-like datasets, such as multilingual Conceptual Captions (Sharma et al., 2018), into multiple languages using neural machine translation (NMT) systems (Surís et al., 2022; Maaz et al., 2024). This presents some drawbacks such as some languages still being poorly translated by current state-of-the-art NMT models (Liu et al., 2020; Costa-jussà et al., 2022) and some cultural subtleties inherent in each language not being fully conveyed. Some efforts have been conducted to collect large-scale multilingual image captioning datasets, such as LAION-5B (Schuhmann et al., 2022), but they are limited to caption data too, are relatively noisy and more importantly contain a non-negligible share of “not safe for work” (NSFW) content such as pædopornographic images (Schuhmann et al., 2022).

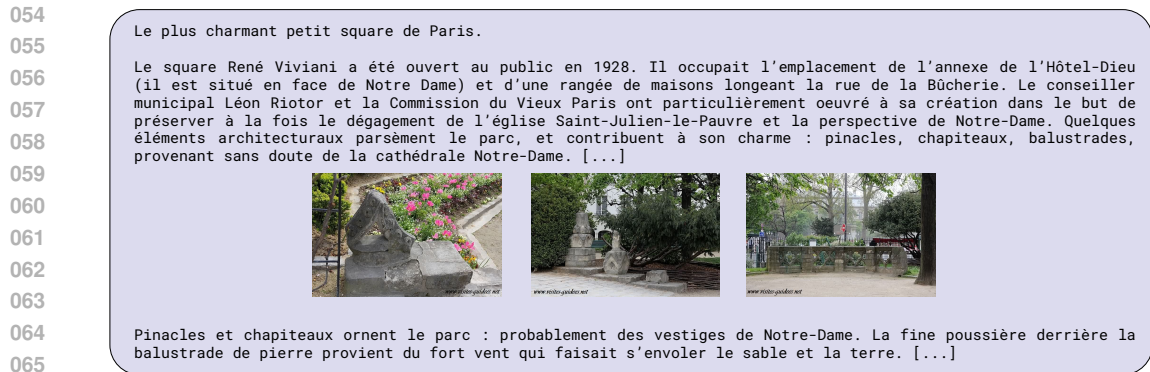


Figure 1: Example of a French document from mOSCAR.

066  
 067  
 068  
 069 This motivated us to collect and release the first large-scale multilingual and multimodal document  
 070 dataset derived from Common Crawl.<sup>1</sup> Our dataset, multimodal OSCAR (mOSCAR), follows the  
 071 OSCAR initiative (Ortiz Suárez et al., 2019; Abadji et al., 2021; 2022) and covers 303M documents  
 072 in 163 languages, 200B tokens and 1.15B images. Figure 1 shows an example of a document, more  
 073 can be found in the Supplementary Material. We carry out extensive filtering to increase its safety  
 074 and quality. To prove mOSCAR’s utility, we train a multilingual OpenFlamingo (Awadalla et al.,  
 075 2023) from a Gemma-2B language model (Gemma Team et al., 2024) on a subset of mOSCAR and  
 076 captioning data from LAION-400M (Schuhmann et al., 2021), recaptioned with BLIP (Li et al.,  
 077 2022), filtered with CLIP (Radford et al., 2021) and translated with NLLB (Costa-jussà et al., 2022).  
 078 We compare against a similar model trained on captioning data only and show we obtain a strong  
 079 boost in few-shot learning, confirming previous findings for English (Alayrac et al., 2022; McKinzie  
 080 et al., 2024; Laurençon et al., 2024). The dataset and models will be made publicly available.

## 082 2 RELATED WORK

083  
 084 **Large-scale web-based datasets** Numerous datasets have been created by filtering web-crawled  
 085 data. These include large-scale text-only datasets (Ortiz Suárez et al., 2019; Raffel et al., 2020;  
 086 Wenzek et al., 2020; Gao et al., 2020; Abadji et al., 2021; Xue et al., 2021; Laurençon et al., 2022;  
 087 Abadji et al., 2022; Penedo et al., 2023) and multimodal ones (Sharma et al., 2018; Changpinyo et al.,  
 088 2021; Jia et al., 2021; Schuhmann et al., 2021; 2022; Byeon et al., 2022; Laurençon et al., 2023;  
 089 Zhu et al., 2023; Gadre et al., 2024). Even if these datasets are not as high quality as smaller and/or  
 090 hand-crafted ones, they are now the standard to pretrain foundation models, as it has been shown that  
 091 training bigger models on more data leads to better downstream performances (Brown et al., 2020;  
 092 Hoffmann et al., 2022; Touvron et al., 2023a;b).

093  
 094 **English image-text datasets** The first open-source image-text datasets were manually created,  
 095 small-scale and English-only (Ordonez et al., 2011; Lin et al., 2014; Plummer et al., 2015; Krishna  
 096 et al., 2017). Scaling up these datasets was an appealing solution to overcome limitations of previous  
 097 image-text models; a few works (Sharma et al., 2018; Changpinyo et al., 2021) proposed to collect  
 098 millions of image-text pairs from the web before filtering them with well-designed steps. Relaxing the  
 099 filtering steps enabled the collection of more data and led to large-scale datasets to train image-text  
 100 foundation models (Radford et al., 2021; Li et al., 2021; Schuhmann et al., 2021; 2022; Byeon et al.,  
 101 2022). However, these datasets generally contain caption-like image-text pairs only, and it is therefore  
 102 difficult to observe in-context learning abilities similarly to text-only language models trained on raw  
 103 documents (Raffel et al., 2020). Alayrac et al. (2022) overcome this issue by training their model  
 104 directly on documents with interleaved image-text data. While their results are promising, their M3W  
 105 dataset is English-only and private. Recently, open-source efforts (Zhu et al., 2023; Laurençon et al.,  
 106 2023) have been made to release a similar dataset but they are still monolingual.

107 <sup>1</sup><https://commoncrawl.org/>. The Common Crawl Foundation is a non-profit organization that  
 crawls the web on a monthly basis.

**Multilingual image-text datasets** Only a few image-text datasets are available in multiple languages. One of the first focused on collecting Google images from short queries based on word frequencies from Wikipedia pages in 98 languages (Hewitt et al., 2018). Later, Srinivasan et al. (2021) proposed the WIT dataset, an image-text dataset composed of Wikipedia pages. Although of high quality, it is only medium-scale even for high-resource languages and there are fewer than 50k unique images for most languages. Another approach lies in bootstrapping multilingual and multimodal data from a model trained with English-only data (Mohammed et al., 2023). While effective for captioning, it is computationally expensive to implement in practice. Other multilingual image-text datasets exist but focus on captions only and are highly domain-specific (Kosar et al., 2022; Leong et al., 2022).

### 3 DATASET CREATION PIPELINE

#### 3.1 DATA COLLECTION

We collect mOSCAR from the Web ARchive Content (WARC) files of three 2023 Common Crawl dumps, processing them using the FastWARC library (Bevendorff et al., 2021). We remove documents smaller than 500 bytes (50% of the documents), as we find they are usually too small to be considered documents and tend to contain noisy text. We then navigate through the entire Document Object Model (DOM) tree with a depth first search algorithm and ChatNoir library (Bevendorff et al., 2018) to extract nodes of interests corresponding to specific HTML tags.

Following previous work, we extract text from the tags that usually contain the main content of web pages (we refer to them as DOM text nodes), i.e. `<p>`, `<h*>`, `<title>`, `<description>`, `<ul>`, `<ol>`, `<aside>`, `<dl>`, `<dd>`, `<dt>`. Similarly to (Laurençon et al., 2023), we choose to remove `<table>` content as most often it is irrelevant and difficult to render. We extract all `<img>` tags (we refer to them as DOM image nodes). We then remove documents with fewer than 3 text nodes (as they do not contain enough text) and more than 30 image nodes (as we found them to be too noisy).

#### 3.2 LANGUAGE IDENTIFICATION

We identify the language of each document using the state-of-the-art open-LID language detector (Burchell et al., 2023), covering 201 languages. We apply open-LID to each DOM text node and keep the three most probable languages with their respective probabilities. The language of the document is then determined by summing over the probabilities of each language detected for each text segment, weighted by the number of characters in the segment<sup>2</sup> and taking the language with the highest score.

#### 3.3 TEXT-ONLY FILTERING

We apply a series of filtering steps to the text content of each document independently of the images, with the aim of discarding poor quality documents and cleaning text as best as possible. We first filter at the text-node level and then at the whole document level, before running near-deduplication to keep unique text nodes within a document and unique documents in the dataset.

**Text node filtering** We use a set of heuristics (see Supplementary Material) to extract as much human-generated content as possible while discarding noisy text related to ads and website functions (e.g. “Instagram”, “Facebook”). We then keep DOM text nodes with content over 10 bytes. This step, designed to improve the quality of extracted text, removes on average 55% of text nodes.

**Document filtering** We mostly filter “not safe for work” (NSFW) content at the document level. We use an English regular expression to detect adult content, similar to the one used by the Université Toulouse 1 Capitole<sup>3</sup> and remove the entire document if there is a match with any of the DOM text nodes’ contents, removing on average 0.5% of documents (mostly English ones). We acknowledge that there is a high probability that this also discards safe content, e.g. we could remove content from

<sup>2</sup>This is to avoid mis-assigning the language due to the presence of many short, non-informative DOM text nodes in the same language (e.g. “Cookies”, “Subscribe”, “Newsletter” etc.) and because language identification is generally less reliable for short segments.

<sup>3</sup>[https://dsi.ut-capitole.fr/blacklists/index\\_en.php](https://dsi.ut-capitole.fr/blacklists/index_en.php)

certain communities who use some explicit words in a non-sexual way (Sap et al., 2019). However, we explicitly favour recall over precision to minimise the risk of unsafe content. We additionally remove documents containing fewer than five DOM text nodes and fewer than 300 characters after the previous filtering steps, removing 70.6% of documents.

**Deduplication** We conduct several types of per-language deduplication at different levels, as this has been shown to improve training efficiency (Abbas et al., 2023). First, we keep unique documents only by removing exact duplicates at the document level. We also remove exact duplicates of text nodes within the same document (4% of text nodes) and near-duplicate text nodes (1% of text nodes) by computing the Levenshtein ratio (Levenshtein, 1966) between all text nodes within the same document and applying a threshold of 0.95. If near-duplicates are found, we keep the first one in the document. Finally, we conduct per language near-deduplication at the document level with MinHashLSH (Broder, 1997; Gionis et al., 1999) following Smith et al. (2022), removing on average 19% of documents:<sup>4</sup> we turn documents into hashing vectors, compute min hashes from these vectors and perform Locality Sensitive Hashing to remove duplicates<sup>5</sup> (see Supplementary Material for more details).

**Toxicity filtering** Toxic content targeting individuals or groups of people is widespread on the internet and can therefore be found in large-scale web-crawled datasets like mOSCAR without appropriate filtering steps. To alleviate this issue, we apply the same method used by Costa-jussà et al. (2022) and remove documents from mOSCAR based on the presence of a list of "toxic" words for each language<sup>6</sup>. As some words in the list can also be used in a non-toxic way based on the context (e.g.: 'breast' in English), we tag the document as toxic and remove it from mOSCAR if it contains at least two distinct words in the list. This filtering step removes 0.95% of the documents for very high-resource languages (>5M documents), 2.13% for high-resource languages (<5M, >500K), 0.47% for mid-resource languages (<500K, >50K) and 0.64% for low-resource languages (< 50K). When manually analysing 100 random documents removed by this filtering step in each of the 2 (high-resource) languages we are native speakers of (English and French), we found 58 documents with toxic content.

**Personal Identifiable Information** Personal Identifiable Information (PII) can be found in large-scale web-crawled datasets, we therefore conducted an additional filtering step to replace all detected PII by place holder strings using regular expressions (see Supplementary Material for the list of regular expressions we used). Concretely, we replaced all detected email addresses, phone numbers, credit card numbers, IP addresses and passport numbers.

### 3.4 IMAGE-ONLY FILTERING

We downloaded images from the URLs in DOM image nodes using a modified version of the `img2dataset` toolkit (Beaumont, 2021) that includes an antivirus scan and follows `robots.txt` instructions to respect the Robots Exclusion Protocol. We then apply a series of filtering steps, first removing images based on heuristics, and then applying multiple NSFW detection models to remove undesirable content. Finally, we conduct a set of deduplication steps.

**Rule-based filters** Similarly to previous works (Schuhmann et al., 2021) and to avoid extracting low-resolution images and favicons, we keep images with a minimum height and width of 150 pixels. We restrict the aspect ratio to be between 3 and 1/3 (to remove banners), we remove images if their URLs contain the words "logo", "banner", "button", "widget", "icon" or "plugin" or if the image name from the URL matches "twitter", "facebook" or "rss" (to remove logos). This step removes 13.6% of the URLs. At this stage, we downloaded 2.5B images with an average success rate of 55%.

**NSFW detection** We use multiple NSFW automatic models to remove as much unsafe content as possible. We first combine two NSFW detectors: `nsfw-detector` (Laborde), a 5-class classifier with a

<sup>4</sup>With some disparity among languages as we found more duplicates for low- than high-resource languages.

<sup>5</sup>We performed this using the `datasketch` python library.

<sup>6</sup>The list of these words for each language can be found here: <https://github.com/facebookresearch/flores/tree/main/toxicity>

MobileNet (Howard et al., 2017) backbone fine-tuned on 60GB of annotated data and NudeNet,<sup>7</sup> an object detector trained to detect different types of nudity in images. We combined the two models as we found the first to be gender-biased while the second gives a large number of false positives for non-human images. Concretely, we consider an image an NSFW candidate if the sum of the probabilities for the classes ‘porn’ and ‘hentai’ is superior to 0.8 using nsfw-detector. We then tag the image as NSFW if one of the sensitive ‘exposed’ classes of NudeNet gets a probability superior to 0.5. If a document contains an image with an NSFW tag, we remove the entire document from the dataset, which removes 0.5% of images. We manually inspecting 1,000 images of the remaining data and found no NSFW content. We manually inspected 1,000 images of the removed content and found 63.4% of NSFW images.

**CSAM content** Child Sexual Abuse Material (CSAM) is widespread on the internet and is therefore likely to be found in such a large-scale dataset crawled from the web. Removing CSAM is challenging as there is no training data nor open-source detection models available as these could be used in a harmful way. We rely on Thorn’s CSAM classifier<sup>8</sup>, a proprietary classifier trained to detect CSAM content in images. We tag the image as CSAM if the probability of the class CSAM is superior to 0.4 to favour recall over precision. As mentioned above, if a document contains an image with a CSAM tag, we remove it from the dataset. This step removes 0.07% of the images.

**Deduplication** To avoid memorisation issues often seen in models trained on datasets with many duplicated images (Somepalli et al., 2023; Carlini et al., 2023; Webster et al., 2023; Somepalli et al., 2024), we perform deduplication at the image level. We first remove duplicate images within the same document by URL matching (removing 8.7% of URLs). We then compute a perceptual hash (pHash) for each image using the imagehash library<sup>9</sup> and remove images with the same pHash within the same document, keeping only the first occurrence. We also limit the number of times an image can appear in the dataset per-language to 10 using both URL matching and perceptual hashing (this removes 2.5% of images). We do this per-language and not across languages as having the same images in documents from different languages could encourage cross-lingual transfer.

**Personal Identification Information** To protect PII in images, we use a lightweight face detector<sup>10</sup> and apply a threshold of 0.99 to detect faces in the images. We apply such a high threshold as we found the model to be biased towards detecting faces with high probability in images without any human. For each image in mOSCAR, we distribute the bounding boxes of the detected faces so that users can blur them when downloading the images. More details are provided in the Supplementary Material.

### 3.5 DATA DECONTAMINATION

LLMs and mLLMs are trained on web-crawled data that can contain the benchmarks they are tested on (Dodge et al., 2021). As they are good at memorizing training data (Carlini et al., 2023), this data contamination is problematic. We therefore discard all images with the same perceptual hash as any of the images from the evaluation benchmarks (and their training sets) we use (see Section 5.1). This step removes on average 126,016 images for high-resource languages (up to 300K images for English), 6,862 images for mid-resource languages and 45 images for low-resource languages.

### 3.6 TEXT-IMAGE JOINT FILTERING

Our aim is to obtain truly multimodal documents where all images are related to at least one of the text nodes in some way<sup>11</sup> and vice versa. We choose to apply joint text-image filtering to discard images and/or text nodes that are irrelevant to the rest of the document (e.g. the case of ads and website functionalities). To do this, we use NLLB-SIGLIP<sup>12</sup> (Visheratin, 2023), a multilingual version

<sup>7</sup><https://github.com/vladmandic/nudenet>

<sup>8</sup><https://safer.io/>

<sup>9</sup><https://github.com/JohannesBuchner/imagehash>

<sup>10</sup><https://github.com/Linzaer/Ultra-Light-Fast-Generic-Face-Detector-1MB>

<sup>11</sup>We do not limit ourselves to caption-like relation and instead allow all types of text-image relation.

<sup>12</sup>siglip-base-patch16-224 as vision encoder and nllb-distilled-600M as text encoder.

of SIGLIP (Zhai et al., 2023) trained with the encoder of NLLB (Costa-jussà et al., 2022), which covers all mOSCAR languages.<sup>13</sup> We compute cosine similarity scores between all images and all paragraphs<sup>14</sup> within a same document. To remove irrelevant text nodes or images in a document, we mimic a text-image retrieval task, which means we avoid using arbitrary cosine similarity thresholds for each language and can reduce length biases and those in favour of caption-like paragraphs. For each candidate pair we randomly sample 63 negative images and 63 negative similar-length paragraphs from the same language but other documents. We tag the text node (resp. image) as valid if the cosine similarity of the pair is among the top 8 of the text-to-image (resp. image-to-text) similarity scores computed with the candidate text node (resp. image) and all the negative images (resp. text nodes). This means that we tag the text node (resp. image) as valid if it has a significantly higher score than a score computed with a random image (resp. text) for at least one of the images (resp. text node) in the document. We then discard text nodes and images not tagged as valid (on average 35% of the DOM text nodes and 10% of the images within a document).

After this filtering step, we apply additional text-only filters to keep documents superior to 100 bytes. We also reapply the open-lid language detector (Burchell et al., 2023) as described in Section 3.2 as we found the last filtering step to change the major language of some documents.

#### 4 MULTIMODAL OPEN SUPER-LARGE CRAWLED AGGREGATED CORPUS (MOSCAR)

mOSCAR is extracted from three Common Crawl dumps from 2023. Due to computational constraints and in order to extract a maximum number of documents for low-resource languages, we extracted all languages from the first dump only. We removed the 6 most high-resource languages from the second dump and only extracted the languages with fewer than 1M documents for the last dump. Table 1 shows a distribution of the total number of languages and their number of documents.

To avoid data poisoning (Carlini et al., 2024), we release a hash (sha512) with each mOSCAR image.

#documents	10M	5M	1M	500K	200K	50K	10K	5K	1K
#languages	10	15	36	46	56	75	119	133	154

Table 1: Number of languages with at least  $N$  documents

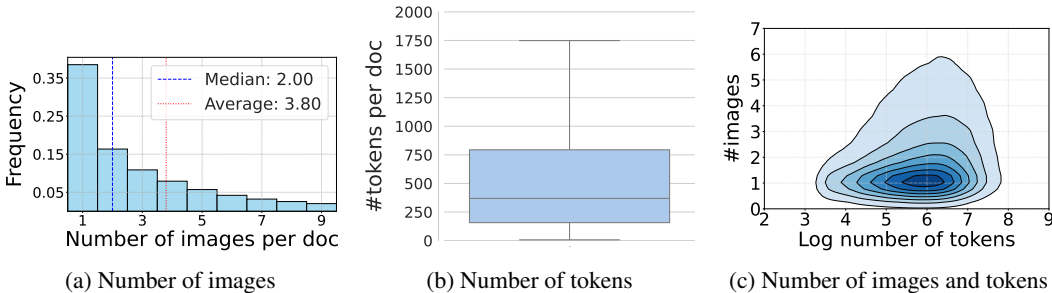


Figure 2: Distributions of numbers of tokens and images per document

mOSCAR is composed of 303M documents (200B tokens, 1.15B images) from 163 languages. Figure 2 shows the distribution of images and tokens per document and their joint distribution. As shown in Figure 2a, the mean and median number of images per document is 2 and 3.80.

##### 4.1 QUALITY VS DIVERSITY

While improving overall data quality, the filtering steps we applied (see Section 3) necessarily have a negative impact on diversity. We therefore study the trade-off between quality and diversity and compare against previously published, well-used datasets.

<sup>13</sup>We use the open-clip (Ilharco et al., 2021) model version and the transformers (Wolf et al., 2020) library.

<sup>14</sup>We refer to paragraph as the text content in a DOM text node.

#### 4.1.1 TEXT CONTENT

**Diversity** By construction, mOSCAR is diverse in terms of number of languages, so we focus on the diversity of mOSCAR’s English documents and compare against mmc4 (Zhu et al., 2023), OBELICS (Laurençon et al., 2023) and the English subset of WIT (Srinivasan et al., 2021). We compute the Vendi score (Friedman and Dieng, 2023) on a set of SimCSE embeddings (Gao et al., 2021) with a RoBERTa encoder (Liu et al., 2019) to evaluate the content diversity. Since embedding-based diversity metrics target content diversity well but are less relevant for lexical diversity (Tevet and Berant, 2021), we measure lexical diversity via the distinct  $n$ -gram ratio (Li et al., 2016).

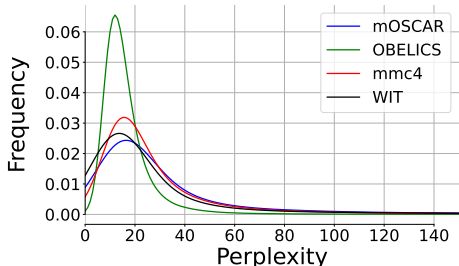


Figure 3: Perplexity of 100K random documents from different datasets.

	Vendi score	Dist. $n$ -gram ratio
mOSCAR	69.05 ( $\pm$ 0.14)	0.472 ( $\pm$ 0.002)
mmc4	67.93 ( $\pm$ 0.12)	0.494 ( $\pm$ 0.002)
OBELICS	58.49 ( $\pm$ 0.09)	0.488 ( $\pm$ 0.001)
WIT	<b>73.30</b> ( $\pm$ 0.09)	<b>0.530</b> ( $\pm$ 0.001)

Table 2: Average text diversity scores ( $\pm$  standard error) of text documents.

**Comparison with other datasets** For content diversity, we randomly sample 30M documents for mOSCAR, mmc4 and OBELICS and 3M documents for WIT and represent the documents by their SimCSE embedding. We compute the Vendi Score with cosine similarity on a randomly sampled subset of 65,536 documents. Table 2 shows that mOSCAR English content is more diverse than mmc4 and OBELICS but less diverse than WIT. For lexical diversity, we randomly sample 3M documents for mOSCAR, mmc4, OBELICS and WIT and compute the distinct  $n$ -gram ratio on a subset of 8,192 documents for  $n$  from 1 to 4. Table 2 shows that mOSCAR is slightly less lexically diverse than OBELICS and mmc4, while WIT is by far the most diverse.

**Quality** To evaluate document quality, we focus on English documents and compute their perplexity using Gemma-2B (Gemma Team et al., 2024). Figure 3 shows the kernel density estimation of the distribution of the perplexity of 100K randomly sampled documents from different datasets: mOSCAR is comparable to mmc4 and WIT, while OBELICS appears to be the of the highest quality. mOSCAR is therefore comparable to other interleaved image-text dataset in terms of quality and diversity of its English subset. It is however more diverse than English-only datasets by its multilingual construction and more than 10 times larger than existing multilingual interleaved image-text datasets such as WIT.

#### 4.1.2 IMAGE DIVERSITY

mOSCAR	LAION-400M	WIT	English	All
55.74 ( $\pm$ 0.16)	<b>67.59</b> ( $\pm$ 0.16)	36.14 ( $\pm$ 0.08)	52.36 ( $\pm$ 0.18)	<b>54.78</b> ( $\pm$ 2.29)

(a) Comparison of different datasets.

(b) mOSCAR (English vs. any language).

Table 3: Average Vendi score ( $\pm$  standard error) of images.

**Comparison with other datasets** We compute the Vendi Score on random samples of images for different datasets, comparing the images from English mOSCAR documents with those from Conceptual Captions (Changpinyo et al., 2021), LAION-400M (Schuhmann et al., 2021) and WIT (Srinivasan et al., 2021). We represent each image by its SigLIP<sup>15</sup> (Zhai et al., 2023) embedding and compute the Vendi score on batches of size 65,536 and a total of 1M images for each dataset. In

<sup>15</sup>We use siglip-base-patch16-224.

Table 3a, we notice that the set of images in mOSCAR documents are more diverse than images from WIT documents but less diverse than LAION-400M.

**Multilingual diversity** We also compare the diversity of images from English documents and of images sampled from documents of any language (English included). We use multilingual SigLIP (Chen et al., 2023a) trained on WebLI (Chen et al., 2023b) to compute image embeddings used to get the Vendi score. We again use a batch of size 65,536 and a total of 3M images, and we do not sample multiple images from a same document. For the multilingual setting, we randomly sample 50 languages and an equal number of images for each language to build the batch. As we did not do any image deduplication across languages, we could expect to have less diversity in the multilingual setting. However, Table 3b shows that the set of images is on average more diverse when sampled from all documents than from English-only documents. This means that the distribution of images is not exactly the same across languages, potentially due to cultural differences.

## 5 TRAINING A MULTILINGUAL MULTIMODAL LANGUAGE MODEL

We train a multilingual Flamingo-like model on mOSCAR that we call multilingual Open Flamingo. As adding captioning data to training data has been shown to improve zero-shot performance (McKinzie et al., 2024), we additionally train on LAION-400M, which we re-captioned using BLIP (Li et al., 2022), filtered with CLIP score (Radford et al., 2021) and translated using distilled NLLB-600M (Visheratin, 2023) following the proportion of languages found in mOSCAR. We use Gemma-2B (Gemma Team et al., 2024) as the underlying language model and we train the model on 50M mOSCAR documents and 100M randomly sampled image-text pairs. We also train a model on 300M image-text pairs, a model trained on 35M WIT (Srinivasan et al., 2021) documents and 70M text-image pairs and a model trained on 50M mOSCAR documents from the English subset and 100M English image-text pairs as comparison baselines. We additionally compare with OpenFlamingo-3B-MPT (Awadalla et al., 2023) as the *translate-test* baseline. The full list of languages for training and the implementation details can be found in the Supplementary Material.

### 5.1 EVALUATION SETUP

We evaluate the models using a broad set of image-text multilingual tasks and benchmarks. We use the IGLUE benchmark (Bugliarello et al., 2022) composed of XVNLI, MaRVL (Liu et al., 2021) to test reasoning, xGQA (Pfeiffer et al., 2022) to test visual question answering capabilities and xFlickr&CO (Young et al., 2014; Karpathy and Fei-Fei, 2015; Yoshikawa et al., 2017) for captioning. We also include Crossmodal-3600 (XM3600) (Thapliyal et al., 2022) and MaXM (Changpinoy et al., 2022) as they cover a broader range of languages. To test to what extent models trained on mOSCAR can perform zero-shot multimodal machine translation (MMT), we also test on Multi30K (Elliott et al., 2016; 2017; Barrault et al., 2018) and CoMMuTE (Futeral et al., 2023). For captioning we compute the CideR (Vedantam et al., 2015) score and we tokenize references and model outputs with the Stanford Core NLP tokenizer for English and Stanza (Qi et al., 2020) tokenizers for other languages. To evaluate Multi30k, we compute BLEU (Papineni et al., 2002) score from Sacrebleu (Post, 2018) with *l3a* tokenization and default parameters. We use accuracy for CoMMuTE. More details can be found in the Supplementary Material.

### 5.2 RESULTS

Tables 4 and 6 show the average results across all languages. Full results are available in the Supplementary Material. We notice that the multilingual OpenFlamingo trained additionally on mOSCAR gets better results than the model trained on captioning data only while having seen fewer image-text pairs during training. More importantly, when increasing the number of few-shot examples from 0 to 16, it sees gains of on average +6.71 points on VQA benchmarks and +19.39 CideR points on captioning benchmarks. In contrast, the model trained on text-image pairs only sees gains of +2.82 and +9.08 points respectively. In cross-modal machine translation, the model additionally trained on interleaved data is again much better than the one trained on just captioning data, which is not able to translate the Multi30k benchmark at all.<sup>16</sup> Moreover, mOSCAR helps the model to learn to zero-shot

<sup>16</sup>Most of the time, the model is not able to follow the prompt and only outputs the end of sequence token.



disambiguate translations as shown by the improved average score on CoMMuTE (63.54) compared to the model trained on captions only (61.36).

Multilingual Open Flamingo trained on mOSCAR & text-image pairs is also better than Open-Flamingo 3B MPT evaluated on translate test benchmarks<sup>17</sup>. However, we obtain the best results (except for MaXM) by evaluating our multilingual Open Flamingo on the translate-test benchmarks since the underlying language model (Gemma-2B) is far better in English than other languages. We also notice that all models struggle with reasoning classification tasks (MaRVL, XVNLI) where they obtain scores close to random guessing.

Table 5 additionally shows that Multilingual Open Flamingo trained on mOSCAR obtains much better results than the same model trained on WIT for equivalent training data seen during training<sup>18</sup> (except for Multi30K benchmark) which means mOSCAR is better suited than WIT for training multilingual mLLMs. Eventually, Table 7 shows that we don't face a drop in performances in English performances when training the model on 43 languages (multilingual Open Flamingo) in comparison to training it on the English subset of mOSCAR and English text-image pairs.

Additional comparison results with InternVL2 (Chen et al., 2024), Llava-NeXT (Li et al., 2024), PaliGemma (Beyer\* et al., 2024) and Idefics2 (Laurençon et al., 2024) can be found in the Supplementary Material.

	#shots	xFlickR&CO	XM3600	xGQA	MaXM	MaRVL	XVNLI	Multi30K	CoMMuTE
	0	16.91	7.45	26.95	22.33	49.56	33.88	22.91	63.34
Multi. OF	4	34.80	22.18	32.33	26.33	49.64	34.07	23.27	63.22
<i>mOSCAR + cap.</i>	8	36.90	23.48	34.24	27.08	<b>51.48</b>	<b>36.60</b>	23.59	<b>63.54</b>
	16	<b>39.46</b>	<b>23.67</b>	<b>35.23</b>	<b>27.47</b>	49.84	34.85	<b>23.85</b>	62.78
	0	9.57	4.21	8.62	4.01	<b>49.88</b>	33.76	0.00	61.36
Multi. OF	4	13.20	9.26	13.45	4.15	49.54	32.04	0.00	61.13
<i>cap. only</i>	8	18.00	10.35	12.82	4.88	49.65	33.71	0.01	60.90
	16	19.87	12.07	13.37	4.89	49.79	32.70	0.74	60.25

Table 4: Results averaged over all languages. Multi. OF refers to multilingual Open Flamingo, *mOSCAR + cap.* refers to the model trained on text-image pairs and mOSCAR while *cap. only* refers to the model trained only on text-image pairs. **Bold** is best result.

	#shots	xFlickR&CO	XM3600	xGQA	MaXM	MaRVL	XVNLI	Multi30K	CoMMuTE
	0	19.07	8.73	25.08	19.64	<b>49.77</b>	33.01	22.70	<b>63.75</b>
Multi. OF (35M)	4	34.32	20.59	31.90	23.90	49.67	36.07	22.79	63.65
<i>mOSCAR + cap.</i>	8	36.77	22.15	33.9	24.41	49.72	<b>37.16</b>	23.21	63.00
	16	<b>37.63</b>	<b>22.24</b>	<b>35.71</b>	<b>25.38</b>	49.73	35.36	23.48	62.77
	0	9.39	4.67	19.81	14.63	49.71	32.78	<b>26.99</b>	56.75
Multi. OF (35M)	4	7.68	2.99	25.68	16.12	49.72	33.51	<b>26.99</b>	53.27
<i>WIT + cap.</i>	8	8.91	3.63	27.06	16.81	49.74	32.77	<b>26.99</b>	55.33
	16	9.74	4.14	28.14	16.34	49.74	33.63	<b>26.99</b>	54.04

Table 5: Results averaged over all languages and comparison between a model trained on WIT and a checkpoint of multilingual Open Flamingo trained on 35M mOSCAR documents (full model was trained on 50M mOSCAR documents). Both models were trained on 35M documents from their respective training datasets and 70M text-image pairs for fair comparison. Multi. OF (35M) refers to multilingual Open Flamingo trained on 35M documents. **Bold** is best result.

We additionally compare results at different training steps, defined by the number of images seen during training. Figure 4 shows the difference of averaged scores between the model trained on all data and the model trained only on text-images pairs. We notice that the gap first decreases until 20M images seen and keep increasing over training at all training steps after that. Particularly, the gap is wider for few-shot learning.

<sup>17</sup>This means benchmarks were translated from local languages to English, using Google Translate API

<sup>18</sup>We select the checkpoint of multilingual Open Flamingo trained on 35M documents and 70M captions to have fair comparison.

	#shots	xGQA	MaXM	MaRVL	XVNLI
OF-3B MPT	0	18.34	7.68	49.75	32.73
	4	22.97	7.82	49.70	35.82
	8	28.57	8.32	49.71	31.29
	16	31.82	9.04	49.72	33.29
Multi. OF <i>mOSCAR + cap.</i>	0	30.16	10.06	49.93	34.66
	4	35.55	9.89	48.99	36.10
	8	36.78	10.12	<b>50.54</b>	<b>39.69</b>
	16	<b>37.75</b>	<b>11.49</b>	49.57	37.97

Table 6: *Translate-test* results averaged over languages where all benchmarks were translated from local languages into English using Google Translate API. Multi. OF *mOSCAR + cap.* refers to Multilingual Open Flamingo trained on mOSCAR and text-image pairs while OF-3B MPT refers to Open Flamingo (Awadalla et al., 2023) based on MPT (Team, 2023) and trained on mmc4 (Zhu et al., 2023) and English text-image pairs.

	#shots	xFlickr&CO	XM3600	xGQA	MaXM	XVNLI
Multilingual OF <i>mOSCAR + cap.</i>	0	29.64	42.57	34.24	36.58	34.62
	4	51.47	77.98	37.91	38.13	33.59
	8	56.75	77.64	39.44	<b>38.52</b>	<b>38.75</b>
	16	<b>59.89</b>	78.18	40.09	35.80	36.60
English OF <i>English mOSCAR + English cap.</i>	0	32.70	43.75	34.71	36.19	35.82
	4	51.39	75.33	37.48	37.96	34.88
	8	51.44	77.73	39.64	38.35	36.86
	16	59.24	<b>78.38</b>	<b>40.36</b>	37.35	37.11

Table 7: Results on the English subsets of the test sets and comparison between multilingual Open Flamingo and an Open Flamingo trained on the English subset of mOSCAR and English text-image pairs (English OF). Both models were trained on 50M documents from their respective training datasets and 100M text-image pairs for fair comparison. **Bold** is best result.

## 6 CONCLUSION, LIMITATIONS AND SOCIETAL IMPACTS

We introduce mOSCAR, a large-scale multilingual and multimodal dataset covering 163 languages and composed of 303M documents, 200B tokens and 1.15B images. We show that mOSCAR is of good quality, diverse and can be used to train a multilingual and multimodal LLM. We ensure that mOSCAR is as safe as possible by applying a series of filtering steps to remove NSFW and toxic content. We however did not conduct any analysis of its biases as this is challenging in a multilingual setting. As it is crawled from the internet, it is indeed possible that mOSCAR reflects biases widespread on it. Training a model on mOSCAR must therefore be combined with additional alignment training steps to mitigate potential biases towards groups of people. Nevertheless, by its multilingual nature, mOSCAR is a step towards the inclusion of more languages, cultures and people in accessing mLLMs.

## REFERENCES

- Julien Abadji, Pedro Javier Ortiz Suárez, Laurent Romary, and Benoît Sagot. Ungoliant: An optimized pipeline for the generation of a very large-scale multilingual web corpus. In Harald Lungen, Marc Kupietz, Piotr Bański, Adrien Barbaresi, Simon Clematide, and Ines Pisetta, editors, *Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-9)*, pages 1–9, Limerick, 2021. Leibniz-Institut für Deutsche Sprache. doi: 10.14618/fids-pub-10468. URL <https://nbn-resolving.org/urn:nbn:de:bsz:mh39-104688>.
- Julien Abadji, Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. Towards a cleaner document-oriented multilingual crawled corpus. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache,

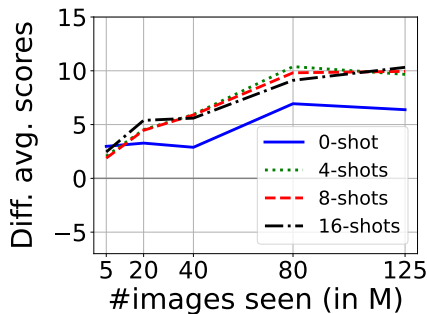


Figure 4: Score differences averaged over benchmarks and languages between the model trained on mOSCAR + text-image pairs and the model trained only on text-image pairs. **Bold** is best result.

- 540 Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard,  
541 Joseph Mariani, H el ene Mazo, Jan Odijk, and Stelios Piperidis, editors. *Proceedings of the*  
542 *Thirteenth Language Resources and Evaluation Conference*, pages 4344–4355, Marseille, France,  
543 June 2022. European Language Resources Association. URL [https://aclanthology.org/](https://aclanthology.org/2022.lrec-1.463)  
544 [2022.lrec-1.463](https://aclanthology.org/2022.lrec-1.463).
- 545 Amro Abbas, Kushal Tirumala, D aniel Simig, Surya Ganguli, and Ari S Morcos. Semdedup: Data-  
546 efficient learning at web-scale through semantic deduplication. *arXiv preprint arXiv:2303.09540*,  
547 2023.
- 548 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel  
549 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language  
550 model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–  
551 23736, 2022.
- 552 Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,  
553 Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel  
554 Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for  
555 training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- 556 Lo c Barrault, Fethi Bougares, Lucia Specia, Chiraag Lala, Desmond Elliott, and Stella Frank.  
557 Findings of the third shared task on multimodal machine translation. In *Proceedings of the Third*  
558 *Conference on Machine Translation: Shared Task Papers*, pages 304–323, 2018.
- 559 Romain Beaumont. img2dataset: Easily turn large sets of image urls to an image dataset. [https:](https://github.com/rom1504/img2dataset)  
560 [://github.com/rom1504/img2dataset](https://github.com/rom1504/img2dataset), 2021.
- 561 Janek Bevendorff, Benno Stein, Matthias Hagen, and Martin Potthast. Elastic ChatNoir: Search  
562 Engine for the ClueWeb and the Common Crawl. In Leif Azzopardi, Allan Hanbury, Gabriella Pasi,  
563 and Benjamin Piwowarski, editors, *Advances in Information Retrieval. 40th European Conference*  
564 *on IR Research (ECIR 2018)*, Lecture Notes in Computer Science, Berlin Heidelberg New York,  
565 2018. Springer.
- 566 Janek Bevendorff, Martin Potthast, and Benno Stein. Fastwarc: optimizing large-scale web archive  
567 analytics. *arXiv preprint arXiv:2112.03103*, 2021.
- 568 Lucas Beyer\*, Andreas Steiner\*, Andr e Susano Pinto\*, Alexander Kolesnikov\*, Xiao Wang\*, Daniel  
569 Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello,  
570 Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Grit-  
571 senko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias  
572 Bauer, Matko Bo snjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Bal-  
573 azevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah  
574 Harmsen, and Xiaohua Zhai\*. PaliGemma: A versatile 3B VLM for transfer. *arXiv preprint*  
575 *arXiv:2407.07726*, 2024.
- 576 A.Z. Broder. On the resemblance and containment of documents. In *Proceedings of the Compression*  
577 *and Complexity of Sequences 1997*, pages 21–29. IEEE Computer Society, 1997. doi: 10.1109/  
578 SEQUEN.1997.666900.
- 579 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
580 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
581 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 582 Emanuele Bugliarello, Fangyu Liu, Jonas Pfeiffer, Siva Reddy, Desmond Elliott, Edoardo Maria  
583 Ponti, and Ivan Vuli c. Iglue: A benchmark for transfer learning across modalities, tasks, and  
584 languages. In *International Conference on Machine Learning*, pages 2370–2392. PMLR, 2022.
- 585 Laurie Burchell, Alexandra Birch, Nikolay Bogoychev, and Kenneth Heafield. An open dataset  
586 and model for language identification. In Anna Rogers, Jordan Boyd-Graber, and Naoaki  
587 Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computa-*  
588 *tional Linguistics (Volume 2: Short Papers)*, pages 865–879, Toronto, Canada, July 2023. As-  
589 sociation for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.75. URL [https:](https://aclanthology.org/2023.acl-short.75)  
590 [://aclanthology.org/2023.acl-short.75](https://aclanthology.org/2023.acl-short.75).

- 594 Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Sae-  
595 hoon Kim. Coyo-700m: Image-text pair dataset. [https://github.com/kakaobrain/  
596 coyo-dataset](https://github.com/kakaobrain/coyo-dataset), 2022.  
597
- 598 Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr,  
599 Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models.  
600 In *Proceedings of the 32nd USENIX Conference on Security Symposium, SEC '23, USA, 2023*.  
601 USENIX Association. ISBN 978-1-939133-37-3.  
602
- 603 Nicolas Carlini, Matthew Jagielski, Christopher Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum  
604 Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. Poisoning web-scale training  
605 datasets is practical. In *Proceedings of the 2024 IEEE Symposium on Security and Privacy (SP)*,  
606 pages 179–179, Los Alamitos, CA, USA, 2024. IEEE Computer Society. doi: 10.1109/SP54263.  
607 2024.00179. URL [https://doi.ieeecomputersociety.org/10.1109/SP54263.  
608 2024.00179](https://doi.ieeecomputersociety.org/10.1109/SP54263.2024.00179).  
609
- 609 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing  
610 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the  
611 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, Nashville,  
612 TN, USA, 2021.  
613
- 613 Soravit Changpinyo, Linting Xue, Idan Szepes, Ashish V Thapliyal, Julien Amelot, Michal Yarom,  
614 Xi Chen, and Radu Soricut. Maxm: Towards multilingual visual question answering. *arXiv  
615 preprint arXiv:2209.05401*, 2022.  
616
- 617 Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, Jialin Wu, Paul Voigtlaender, Basil  
618 Mustafa, Sebastian Goodman, Ibrahim Alabdulmohsin, Piotr Padlewski, et al. Pali-3 vision  
619 language models: Smaller, faster, stronger. *arXiv preprint arXiv:2310.09199*, 2023a.  
620
- 621 Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian  
622 Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual  
623 language-image model. In *Proceedings of the International Conference on Learning Representa-  
624 tions*, Kigali, Rwanda, 2023b.  
625
- 625 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi  
626 Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial  
627 multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.  
628
- 629 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan,  
630 Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling  
631 human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.  
632
- 632 Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld,  
633 Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the  
634 colossal clean crawled corpus. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and  
635 Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural  
636 Language Processing*, pages 1286–1305, Online and Punta Cana, Dominican Republic, November  
637 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.98. URL  
638 <https://aclanthology.org/2021.emnlp-main.98>.  
639
- 639 Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30k: Multilingual english-  
640 german image descriptions. In *Proceedings of the 5th Workshop on Vision and Language*, pages  
641 70–74. Association for Computational Linguistics, 2016. doi: 10.18653/v1/W16-3210. URL  
642 <http://www.aclweb.org/anthology/W16-3210>.  
643
- 644 Desmond Elliott, Stella Frank, Loïc Barrault, Fethi Bougares, and Lucia Specia. Findings of the  
645 second shared task on multimodal machine translation and multilingual image description. In *Pro-  
646 ceedings of the Second Conference on Machine Translation, Volume 2: Shared Task Papers*, pages  
647 215–233, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.  
URL <http://www.aclweb.org/anthology/W17-4718>.

- 648 Dan Friedman and Adji Bousso Dieng. The vendi score: A diversity evaluation metric for machine  
649 learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=g970HbQyk1>.  
650
- 651 Matthieu Futral, Cordelia Schmid, Ivan Laptev, Benoît Sagot, and Rachel Bawden. Tackling  
652 ambiguity with images: Improved multimodal machine translation and contrastive evaluation. In  
653 Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual  
654 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5394–  
655 5413, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/  
656 2023.acl-long.295. URL <https://aclanthology.org/2023.acl-long.295>.  
657
- 658 Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen,  
659 Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the  
660 next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36,  
661 2024.
- 662 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,  
663 Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for  
664 language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 665 Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence  
666 embeddings. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau  
667 Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language  
668 Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic, November 2021.  
669 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL <https://aclanthology.org/2021.emnlp-main.552>.  
670
- 671 Aristides Gionis, Piotr Indyk, Rajeev Motwani, et al. Similarity search in high dimensions via hashing.  
672 In *Proceedings of the 25th VLDB Conference*, volume 99, pages 518–529, Edinburgh, Scotland,  
673 UK, 1999.  
674
- 675 John Hewitt, Daphne Ippolito, Brendan Callahan, Reno Kriz, Derry Tanti Wijaya, and Chris Callison-  
676 Burch. Learning translations via images with a massively multilingual image dataset. In *Proceed-  
677 ings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
678 Papers)*, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- 679 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
680 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas  
681 Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Au-  
682 relia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and  
683 Laurent Sifre. Training compute-optimal large language models. In S. Koyejo, S. Mo-  
684 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural In-  
685 formation Processing Systems*, volume 35, pages 30016–30030. Curran Associates, Inc.,  
686 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
687 file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf).
- 688 Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,  
689 Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for  
690 mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- 691 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,  
692 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali  
693 Farhadi, and Ludwig Schmidt. OpenCLIP, July 2021. URL [https://doi.org/10.5281/  
694 zenodo.5143773](https://doi.org/10.5281/zenodo.5143773).
- 695 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,  
696 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with  
697 noisy text supervision. In *Proceedings of the Thirty-Eighth International Conference on Machine  
698 Learning*, pages 4904–4916, online, 2021. PMLR.  
699
- 700 Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions.  
701 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–  
3137, 2015.

- 702 Vaclav Kosar, Antonín Hoskovec, Milan Šulc, and Radek Bartyzal. GLAMI-1M: A Multilingual  
703 Image-Text Fashion Dataset. In *Proceedings of the 33rd British Machine Vision Conference*,  
704 London, UK, 2022. BMVA Press. URL <https://bmvc2022.mpi-inf.mpg.de/0607.pdf>.  
705 pdf.  
706
- 707 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie  
708 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language  
709 and vision using crowdsourced dense image annotations. *International journal of computer vision*,  
710 123:32–73, 2017.
- 711 Gant Laborde. Deep nn for nsfw detection. URL [https://github.com/GantMan/nsfw\\_model](https://github.com/GantMan/nsfw_model).  
712 model.  
713
- 714 Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral,  
715 Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen,  
716 Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella  
717 Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen,  
718 Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan  
719 Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Muñoz, Jian Zhu, Daniel  
720 Van Strien, Zaid Alyafeai, Khalid Almubarak, Minh Chien Vu, Itziar Gonzalez-Dios, Aitor Soroa,  
721 Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Adelani, Long  
722 Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, Margaret Mitchell,  
723 Sasha Alexandra Luccioni, and Yacine Jernite. The bigscience roots corpus: A 1.6tb composite  
724 multilingual dataset. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh,  
725 editors, *Advances in Neural Information Processing Systems*, volume 35, pages 31809–31826. Curran  
726 Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/  
727 paper/2022/file/ce9e92e3de2372a4b93353eb7f3dc0bd-Paper-Datasets\\_  
728 and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/ce9e92e3de2372a4b93353eb7f3dc0bd-Paper-Datasets_and_Benchmarks.pdf).
- 729 Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov,  
730 Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, Matthieu Cord, and  
731 Victor Sanh. OBELICS: An Open Web-Scale Filtered Dataset of Interleaved Image-Text  
732 Documents. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors,  
733 *Advances in Neural Information Processing Systems*, volume 36, pages 71683–71702. Curran  
734 Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/  
735 paper/2023/file/e2cfb719f58585f779d0a4f9f07bd618-Paper-Datasets\\_  
736 and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/e2cfb719f58585f779d0a4f9f07bd618-Paper-Datasets_and_Benchmarks.pdf).
- 736 Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building  
737 vision-language models? *arXiv preprint arXiv:2405.02246*, 2024.  
738
- 739 Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel  
740 Whitenack. Bloom library: Multimodal datasets in 300+ languages for a variety of down-  
741 stream tasks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of  
742 the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8608–8621,  
743 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.  
744 doi: 10.18653/v1/2022.emnlp-main.590. URL [https://aclanthology.org/2022.  
745 emnlp-main.590](https://aclanthology.org/2022.emnlp-main.590).
- 746 Vladimir Iosifovich Levenshtein. Binary codes capable of correcting deletions, insertions and  
747 reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966. Doklady Akademii Nauk SSSR, V163  
748 No4 845-848 1965.
- 749 Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang,  
750 Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal cap-  
751 abilities in the wild, May 2024. URL [https://llava-vl.github.io/blog/  
752 2024-05-10-llava-next-stronger-llms/](https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms/).  
753
- 754 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting  
755 objective function for neural conversation models. In Kevin Knight, Ani Nenkova, and Owen  
Rambow, editors, *Proceedings of the 2016 Conference of the North American Chapter of the*

- 756 *Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San  
757 Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/  
758 N16-1014. URL <https://aclanthology.org/N16-1014>.
- 759  
760 Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven  
761 Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum  
762 distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- 763 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image  
764 pre-training for unified vision-language understanding and generation. In *Proceedings of the 39th*  
765 *International Conference on Machine Learning*, pages 12888–12900, Baltimore, Maryland, USA,  
766 2022. PMLR.
- 767 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image  
768 pre-training with frozen image encoders and large language models. In *Proceedings of the 40th*  
769 *International Conference on Machine Learning*, pages 19730—19742, Honolulu Hawaii USA,  
770 2023.
- 771 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
772 Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Proceedings of*  
773 *the 13th European Conference on Computer Vision*, pages 740–755, Zurich, Switzerland, 2014.  
774 Springer.
- 775  
776 Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond  
777 Elliott. Visually grounded reasoning across languages and cultures. In Marie-Francine Moens,  
778 Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Con-*  
779 *ference on Empirical Methods in Natural Language Processing*, pages 10467–10485, Online  
780 and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguis-  
781 tics. doi: 10.18653/v1/2021.emnlp-main.818. URL [https://aclanthology.org/2021.](https://aclanthology.org/2021.emnlp-main.818)  
782 [emnlp-main.818](https://aclanthology.org/2021.emnlp-main.818).
- 783 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
784 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
785 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 786 Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike  
787 Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation.  
788 *Transactions of the Association for Computational Linguistics*, 8:726–742, 2020. doi: 10.1162/  
789 [tacl\\_a\\_00343](https://aclanthology.org/2020.tacl-1.47). URL <https://aclanthology.org/2020.tacl-1.47>.
- 790 Muhammad Maaz, Hanoona Rasheed, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M  
791 Anwer, Tim Baldwin, Michael Felsberg, and Fahad S Khan. PALO: A Polyglot Large Multimodal  
792 Model for 5B People. *arXiv preprint arXiv:2402.14818*, 2024.
- 793 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint*  
794 *arXiv:2405.09818*, 2024.
- 795  
796 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,  
797 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly  
798 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 799 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,  
800 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models  
801 based on Gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 802  
803 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufer,  
804 Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. MM1: Methods, Analysis & Insights  
805 from Multimodal LLM Pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- 806 Owais Khan Mohammed, Kriti Aggarwal, Qiang Liu, Saksham Singhal, Johan Bjorck, and Subhojit  
807 Som. Bootstrapping a high quality multilingual multimodal dataset for Bletchley. In Emtiyaz Khan  
808 and Mehmet Gonen, editors, *Proceedings of The 14th Asian Conference on Machine Learning*,  
809 volume 189 of *Proceedings of Machine Learning Research*, pages 738–753. PMLR, 12–14 Dec  
2023. URL <https://proceedings.mlr.press/v189/mohammed23a.html>.

- 810 OpenAI. GPT-4 Technical Report. *ArXiv*, abs/2303.08774, 2023. URL <https://arxiv.org/abs/2303.08774>.  
811  
812
- 813 Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2Text: Describing Images Using 1 Mil-  
814 lion Captioned Photographs. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q.  
815 Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Asso-  
816 ciates, Inc., 2011. URL [https://proceedings.neurips.cc/paper\\_files/paper/2011/file/5dd9db5e033da9c6fb5ba83c7a7e9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/5dd9db5e033da9c6fb5ba83c7a7e9-Paper.pdf).  
817
- 818 Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. Asynchronous pipelines for processing  
819 huge corpora on medium to low resource infrastructures. In Piotr Bański, Adrien Barbaresi,  
820 Hanno Biber, Evelyn Breiteneder, Simon Clematide, Marc Kupietz, Harald Lungen, and Caroline  
821 Iliadi, editors, *Proceedings of the Workshop on Challenges in the Management of Large Corpora*  
822 (*CMLC-7*), pages 9 – 16, Cardiff, UK, 2019. Leibniz-Institut für Deutsche Sprache. doi: 10.14618/  
823 ids-pub-9021. URL <http://nbn-resolving.de/urn:nbn:de:bsz:mh39-90215>.  
824
- 825 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
826 evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin, editors,  
827 *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages  
828 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.  
829 doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 830 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Hamza Alobeidli,  
831 Alessandro Cappelli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refined-  
832 web dataset for falcon llm: Outperforming curated corpora with web data only. In A. Oh,  
833 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural*  
834 *Information Processing Systems*, volume 36, pages 79155–79172. Curran Associates, Inc., 2023.  
835 URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/fa3ed726cc5073b9c31e3e49a807789c-Paper-Datasets_and_Benchmarks.pdf)  
836 [fa3ed726cc5073b9c31e3e49a807789c-Paper-Datasets\\_and\\_Benchmarks.](https://proceedings.neurips.cc/paper_files/paper/2023/file/fa3ed726cc5073b9c31e3e49a807789c-Paper-Datasets_and_Benchmarks.pdf)  
837 pdf.
- 838 Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin Steitz, Stefan Roth, Ivan Vulić,  
839 and Iryna Gurevych. xGQA: Cross-lingual visual question answering. In Smaranda Mure-  
840 san, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association for Com-*  
841 *putational Linguistics: ACL 2022*, pages 2497–2511, Dublin, Ireland, May 2022. Associa-  
842 tion for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.196. URL <https://aclanthology.org/2022.findings-acl.196>.  
843  
844
- 845 Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and  
846 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer  
847 image-to-sentence models. In *Proceedings of the IEEE international conference on computer*  
848 *vision*, pages 2641–2649, Santiago, Chile, 2015.
- 849 Matt Post. A call for clarity in reporting BLEU scores. In Ondřej Bojar, Rajen Chatterjee, Chris-  
850 tian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno  
851 Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéal, Mariana Neves, Matt  
852 Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Con-*  
853 *ference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, Oc-  
854 tober 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6319. URL  
855 <https://aclanthology.org/W18-6319>.
- 856 Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza: A Python  
857 natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual*  
858 *Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020. URL  
859 <https://nlp.stanford.edu/pubs/qi2020stanza.pdf>.  
860
- 861 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
862 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
863 models from natural language supervision. In *Proceedings of the 38th International Conference*  
*on Machine Learning*, pages 8748–8763, online, 2021. PMLR.



- 864 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
865 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
866 transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.  
867
- 868 Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. The risk of racial bias in  
869 hate speech detection. In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings*  
870 *of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678,  
871 Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1163.  
872 URL <https://aclanthology.org/P19-1163>.
- 873 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,  
874 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: Open Dataset  
875 of CLIP-Filtered 400 Million Image-Text Pairs. In *Proceedings of the Data Centric AI NeurIPS*  
876 *Workshop 2021*, online, 2021.
- 877 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
878 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,  
879 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia  
880 Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. In  
881 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural*  
882 *Information Processing Systems*, volume 35, pages 25278–25294. Curran Associates, Inc., 2022.  
883 URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/a1859debf3b59d094f3504d5ebb6c25-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/a1859debf3b59d094f3504d5ebb6c25-Paper-Datasets_and_Benchmarks.pdf).  
884  
885
- 886 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A  
887 cleaned, hypernymed, image alt-text dataset for automatic image captioning. In Iryna Gurevych  
888 and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for*  
889 *Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia,  
890 July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1238. URL  
891 <https://aclanthology.org/P18-1238>.
- 892
- 893 Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared  
894 Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, et al. Using Deep-  
895 Speed and Megatron to Train Megatron-Turing NLG 530B, A Large-Scale Generative Language  
896 Model. *arXiv preprint arXiv:2201.11990*, 2022.
- 897
- 898 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion  
899 art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the*  
*IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- 900
- 901 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-  
902 standing and mitigating copying in diffusion models. *Advances in Neural Information Processing*  
*Systems*, 36, 2024.
- 903
- 904 Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. WIT:  
905 Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceedings*  
906 *of the 44th International ACM SIGIR Conference on Research and Development in Information*  
907 *Retrieval*, pages 2443–2449, online, 2021.
- 908
- 909 Dídac Surís, Dave Epstein, and Carl Vondrick. Globetrotter: Connecting languages by connecting  
910 images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
911 pages 16474–16484, Canada, 2022.
- 912
- 913 MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable  
914 llms, 2023. URL [www.mosaicml.com/blog/mpt-7b](http://www.mosaicml.com/blog/mpt-7b). Accessed: 2023-05-05.
- 915
- 916 Guy Tevet and Jonathan Berant. Evaluating the evaluation of diversity in natural language generation.  
917 In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty, editors, *Proceedings of the 16th Conference*  
*of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages  
326–346, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.  
eacl-main.25. URL <https://aclanthology.org/2021.eacl-main.25>.

- 918 Ashish V. Thapliyal, Jordi Pont Tuset, Xi Chen, and Radu Soricut. Crossmodal-3600: A mas-  
919 sively multilingual multimodal evaluation dataset. In Yoav Goldberg, Zornitsa Kozareva, and  
920 Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural*  
921 *Language Processing*, pages 715–729, Abu Dhabi, United Arab Emirates, December 2022.  
922 Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.45. URL  
923 <https://aclanthology.org/2022.emnlp-main.45>.
- 924 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
925 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and  
926 Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 927 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
928 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
929 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 931 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image  
932 description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern*  
933 *recognition*, pages 4566–4575, 2015.
- 934 Alexander Visheratin. NLLB-CLIP – train performant multilingual image retrieval model on a budget.  
935 *arXiv preprint arXiv:2309.01859*, 2023.
- 937 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,  
938 Lei Zhao, Xixuan Song, et al. CogVLM: Visual Expert for Pretrained Language Models. *arXiv*  
939 *preprint arXiv:2311.03079*, 2023.
- 940 Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal,  
941 Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language:  
942 Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022.
- 943 Ryan Webster, Julien Rabin, Loic Simon, and Frederic Jurie. On the de-duplication of laion-2b. *arXiv*  
944 *preprint arXiv:2303.12733*, 2023.
- 945 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán,  
946 Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets  
947 from web crawl data. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri,  
948 Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani,  
949 Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the*  
950 *Twelfth Language Resources and Evaluation Conference*, pages 4003–4012, Marseille, France,  
951 May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.494>.
- 952 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
953 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick  
954 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,  
955 Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural  
956 language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*  
957 *Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Associa-  
958 tion for Computational Linguistics. URL [https://www.aclweb.org/anthology/2020.](https://www.aclweb.org/anthology/2020.emnlp-demos.6)  
959 [emnlp-demos.6](https://www.aclweb.org/anthology/2020.emnlp-demos.6).
- 960 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya  
961 Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In  
962 Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven  
963 Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021*  
964 *Conference of the North American Chapter of the Association for Computational Linguistics:*  
965 *Human Language Technologies*, pages 483–498, Online, June 2021. Association for Computational  
966 Linguistics. doi: 10.18653/v1/2021.naacl-main.41. URL [https://aclanthology.org/](https://aclanthology.org/2021.naacl-main.41)  
967 [2021.naacl-main.41](https://aclanthology.org/2021.naacl-main.41).
- 968 Yuya Yoshikawa, Yutaro Shigeto, and Akikazu Takeuchi. STAIR captions: Constructing a large-scale  
969 Japanese image caption dataset. In Regina Barzilay and Min-Yen Kan, editors, *Proceedings of the*  
970  
971

972 *55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*,  
973 pages 417–421, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi:  
974 10.18653/v1/P17-2066. URL <https://aclanthology.org/P17-2066>.  
975

976 Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual  
977 denotations: New similarity metrics for semantic inference over event descriptions. *Transactions*  
978 *of the Association for Computational Linguistics*, 2:67–78, 2014. doi: 10.1162/tac1\_a\_00166. URL  
979 <https://aclanthology.org/Q14-1006>.

980 Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu.  
981 Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*,  
982 2022.

983 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language  
984 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*  
985 *(ICCV)*, pages 11975–11986, October 2023.  
986

987 Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae  
988 Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale  
989 corpus of images interleaved with text. *arXiv preprint arXiv:2304.06939*, 2023.  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025