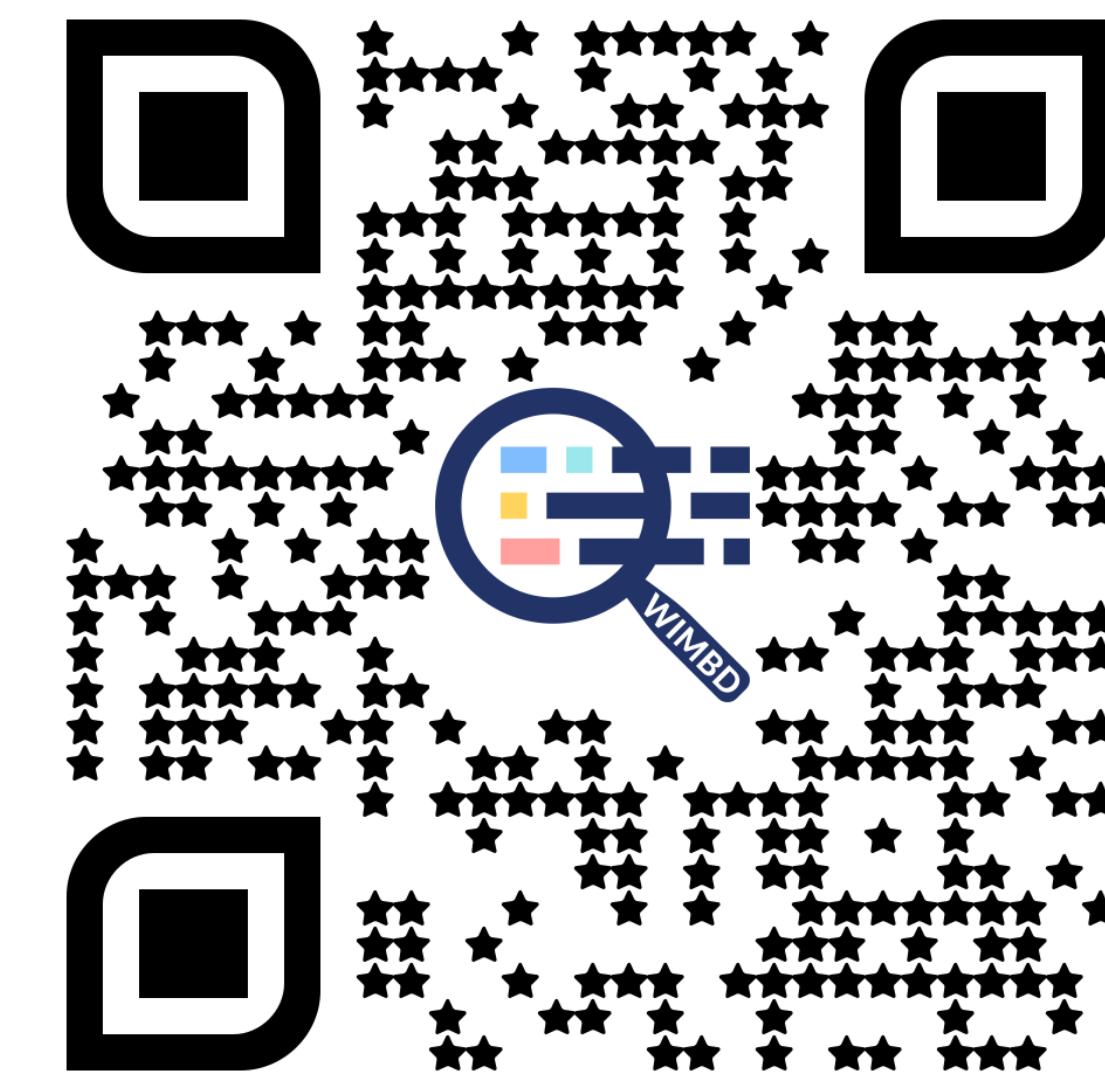
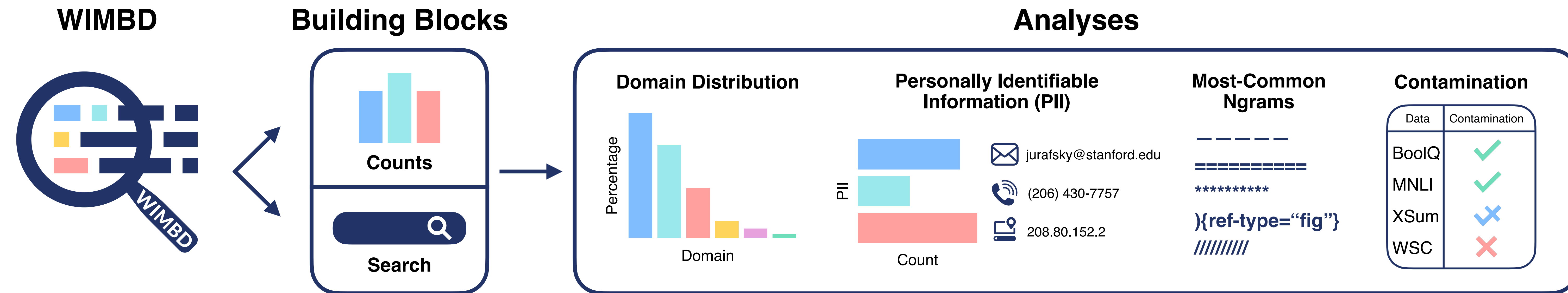


What's In My Big Data?


Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhinav Ravichander, Dustin Schwenk, Alane Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hanna Hajishirzi, Noah A. Smith, Jesse Dodge



(1) Motivation

- Datasets are the foundation of ML models
- To understand model behavior, we must understand their underlying data
- How do we analyze the contents of terabytes of unstructured text data?!

(2) The Platform

- Search
 - Counts
 - Tokens
 - Domains
 - ...
- 
- ```
from wimbd.es import count_documents_containing_phrases

count_documents_containing_phrases("c4", "artificial intelligence")
6,065,714
```

## (3) Datasets & Analyses

| Corpus      | Model              | Size (GB) | # Documents   |
|-------------|--------------------|-----------|---------------|
| OpenWebText | GPT-2*             | 41.2      | 8,005,939     |
| C4          | T5                 | 838.7     | 364,868,892   |
| mC4-en      | umT5               | 14,694.0  | 3,928,733,374 |
| OSCAR       | BLOOM*             | 3,327.3   | 431,584,362   |
| The Pile    | GPT-J/Neo & pythia | 1,369.0   | 210,607,728   |
| RedPajama   | LLaMA*             | 5,602.0   | 930,453,833   |
| S2ORC       | SciBERT*           | 692.7     | 11,241,499    |
| peS2o       | -                  | 504.3     | 8,242,162     |
| LAION-2B-en | Stable Diffusion*  | 570.2     | 2,319,907,827 |
| The Stack   | StarCoder*         | 7,830.8   | 544,750,672   |

### 1.Data Statistics

- High-level statistics
- Internet domains distribution
- Dates distribution

### 2.Data Quality

- Common n-grams
- Duplicates

### 3.Community/Society Measurements

- Contamination
- PII

### 4.Cross-data Analysis

- Distributional similarity
- Overlapping documents

*More analyses in the paper!*

## (4) Results

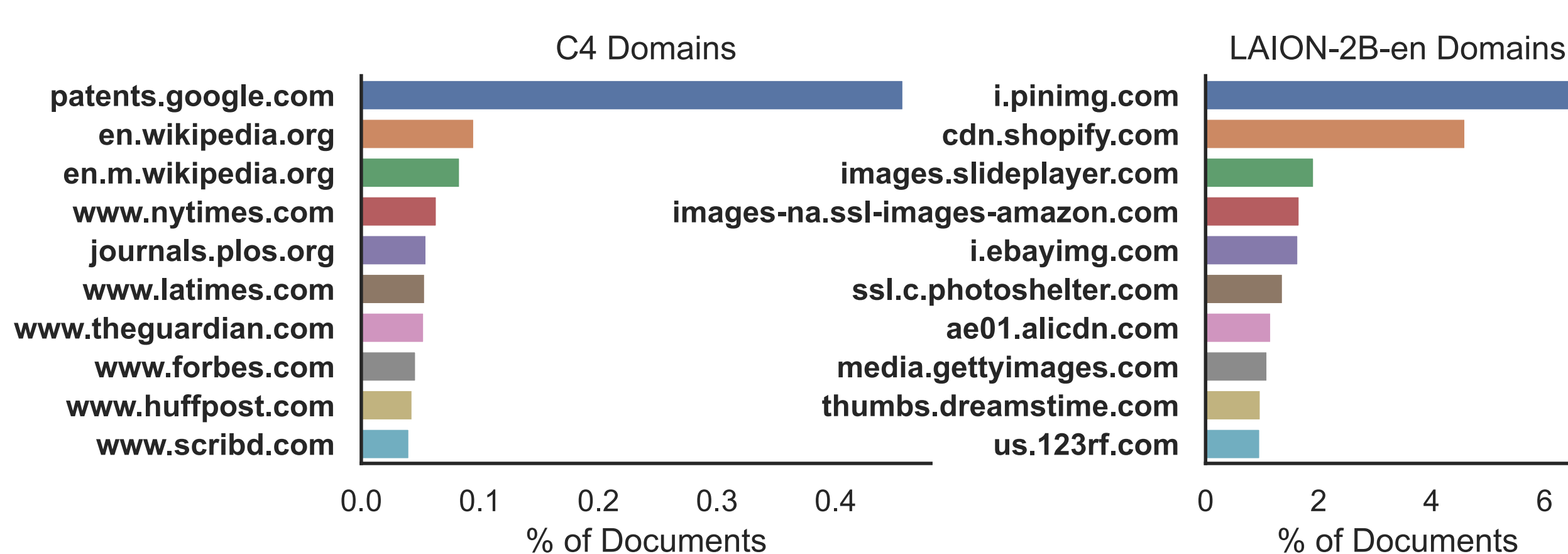
### (i) Most common n-grams

| n-gram                        | OpenWebText | Count | n-gram                        | C4    | Count |
|-------------------------------|-------------|-------|-------------------------------|-------|-------|
| ????????                      | 3.4M        | 9M    | ????????                      | 7.27M |       |
| =====                         | 1.05M       |       | =====                         | 4.41M |       |
| *****                         | 830K        | 3.87M | *****                         | 3.87M |       |
| *****                         | 595K        | 1.91M | *****                         | 1.91M |       |
| #####                         | 302K        | 784K  | #####                         | 784K  |       |
| amp ; amp ; amp ; amp ; amp ; | 278K        | 753K  | amp ; amp ; amp ; amp ; amp ; | 753K  |       |
| amp ; amp ; amp ; amp ; amp ; | 265K        | 752K  | amp ; amp ; amp ; amp ; amp ; | 752K  |       |
| amp ; amp ; amp ; amp ; amp ; | 249K        | 752K  | amp ; amp ; amp ; amp ; amp ; | 752K  |       |
| amp ; amp ; amp ; amp ; amp ; | 88.1K       | 752K  | amp ; amp ; amp ; amp ; amp ; | 752K  |       |
| amp ; amp ; amp ; amp ; amp ; | 83.3K       | 748K  | amp ; amp ; amp ; amp ; amp ; | 748K  |       |

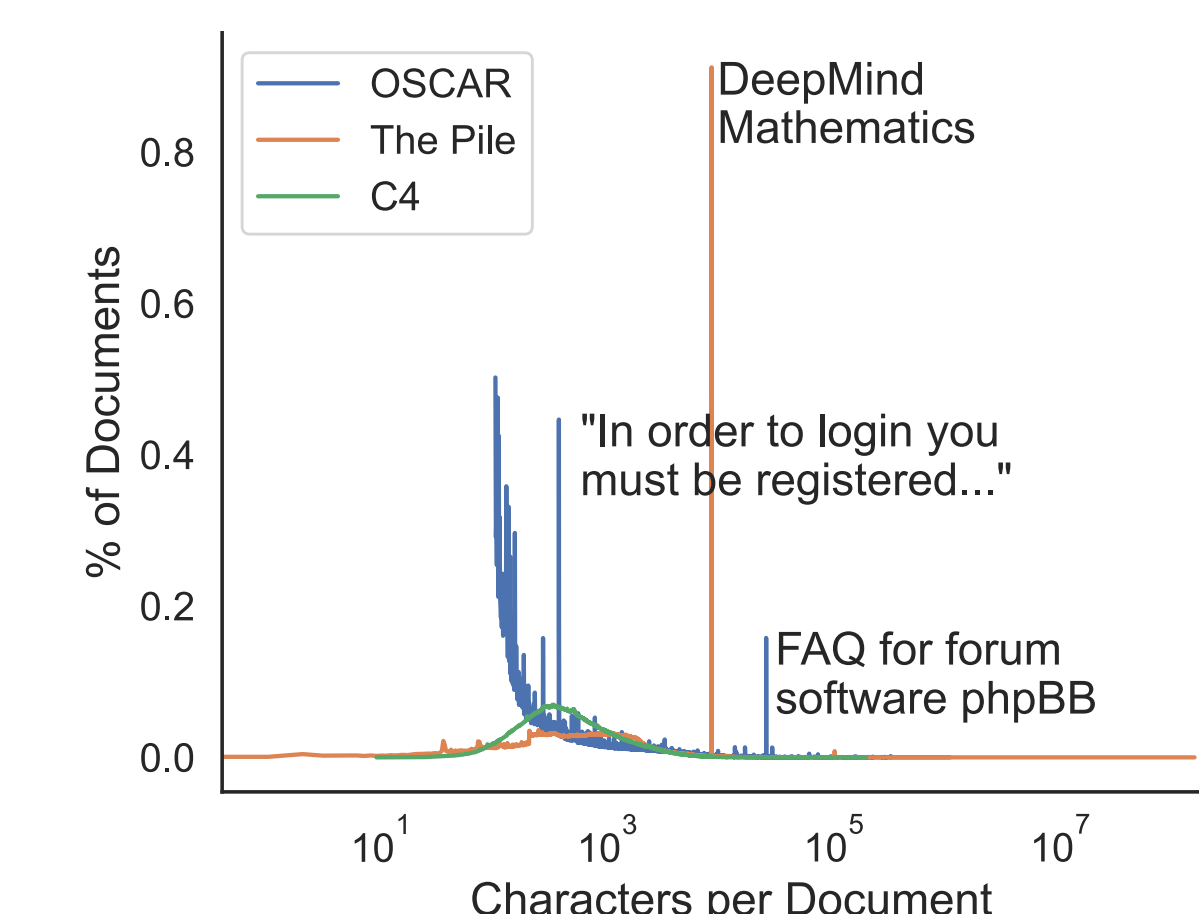
| n-gram                          | LAION-2B-en | Count | n-gram                          | The Stack | Count |
|---------------------------------|-------------|-------|---------------------------------|-----------|-------|
| #####                           | 1.65M       | 4.29B | #####                           | 3.87B     |       |
| #####                           | 1.43M       | 2.75B | #####                           | 2.75B     |       |
| #####                           | 1.15M       | 2.62B | #####                           | 2.62B     |       |
| #####                           | 809K        | 1.46B | #####                           | 1.46B     |       |
| < br /> < br /> < br /> < br /> | 797K        | 1.46B | < br /> < br /> < br /> < br /> | 1.46B     |       |
| < br /> < br /> < br /> < br /> | 796K        | 1.42B | < br /> < br /> < br /> < br /> | 1.42B     |       |
| < br /> < br /> < br /> < br /> | 796K        | 1.42B | < br /> < br /> < br /> < br /> | 1.42B     |       |
| < br /> < br /> < br /> < br /> | 576K        | 1B    | < br /> < br /> < br /> < br /> | 1B        |       |
| < br /> < br /> < br /> < br /> | 437K        | 938M  | < br /> < br /> < br /> < br /> | 938M      |       |
| < br /> < br /> < br /> < br /> | 437K        |       | < br /> < br /> < br /> < br /> |           |       |

### (ii) Internet-domain distribution

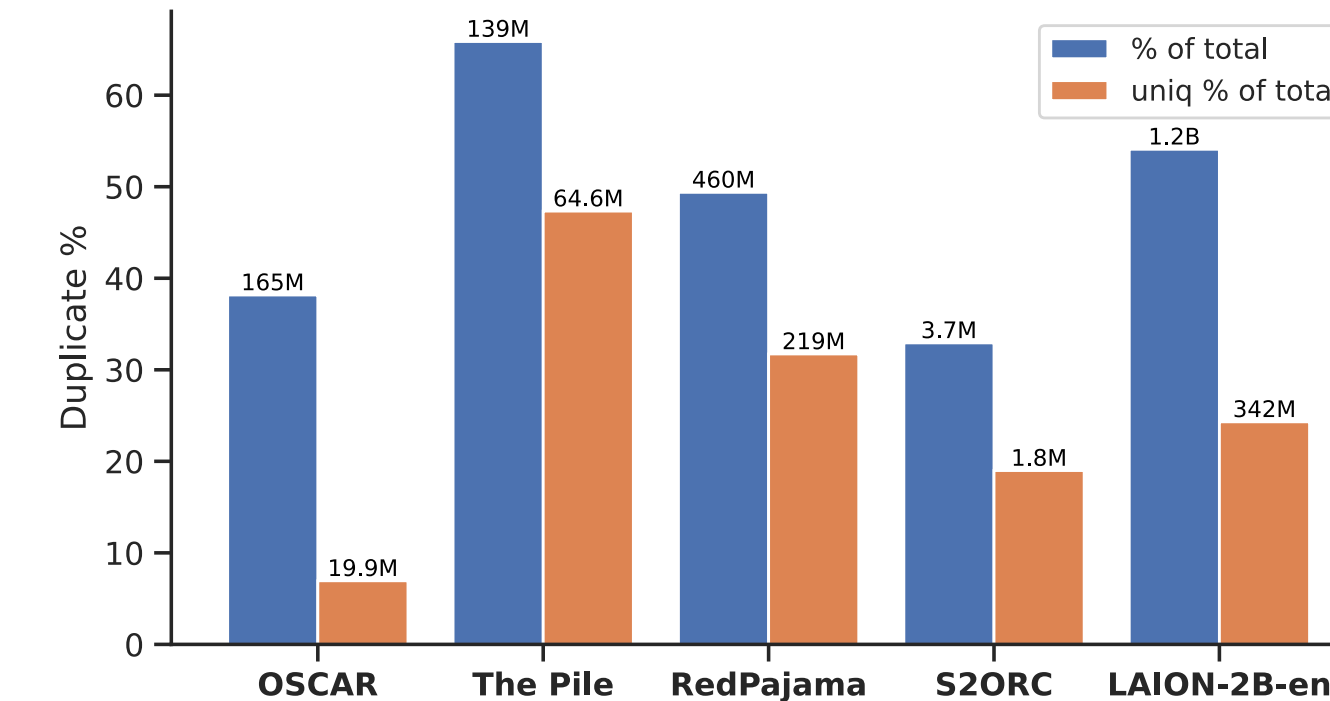
| Domain                  | Corpus    | Rank    | Tokens  | % of All Tokens |
|-------------------------|-----------|---------|---------|-----------------|
| www.geteasysolution.com | C4        | 473082  | 49,859  | 0.000032%       |
| www.geteasysolution.com | RedPajama | 472159  | 49,859  | 0.000023%       |
| www.geteasysolution.com | mC4-en    | 1658921 | 156,174 | 0.0000056%      |



### (iii) Length distribution



### (iv) Duplicate documents



### (v) Personally Identifiable Information (PII)

| Corpus      | Email Addresses | Phone Numbers | IP Addresses |
|-------------|-----------------|---------------|--------------|
|             | Count           | Count         | Count        |
| OpenWebText | 364K            | 533K          | 70K          |
| OSCAR       | 62.8M           | 107M          | 3.2M         |
| C4          | 7.6M            | 19.7M         | 796K         |
| mC4-en      | 201M            | 4B            | 97.8M        |
| The Pile    | 19.8M           | 38M           | 4M           |
| RedPajama   | 35.2M           | 70.2M         | 1.1M         |
| S2ORC       | 630K            | 1.4M          | 0K           |
| peS2o       | 418K            | 227K          | 0K           |
| LAION-2B-en | 636K            | 1M            | 0K           |
| The Stack   | 4.3M            | 45.4M         | 4.4M         |

### (vi) Benchmark contamination

