

663 **A NeurIPS Datasets and Benchmark Checklist**

- 664 1. For all authors...
- 665 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
666 contributions and scope? [Yes] We demonstrate clearly our main claim that models  
667 trained on web data can outperform models trained on curated corpora by showing in  
668 Figure 1 and Figure 3 that models trained on RefinedWeb outperform models trained  
669 on The Pile and match the performance of the GPT-3 series of models.
- 670 (b) Did you describe the limitations of your work? [Yes] See Section 5.
- 671 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We discuss  
672 potential negative societal impacts in the relevant section of our datasheet, see Table 6.
- 673 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
674 them? [Yes] We have added an opt-out mechanism allowing for authors of web pages  
675 included in the public extract to have them removed. Furthermore, the source of our  
676 data (CommonCrawl) honors opt-out requests made based on the robots.txt file.
- 677 2. If you are including theoretical results...
- 678 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 679 (b) Did you include complete proofs of all theoretical results? [N/A]
- 680 3. If you ran experiments (e.g. for benchmarks)...
- 681 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
682 perimental results (either in the supplemental material or as a URL)? [Yes] We have  
683 publicly released: (1) a 600B tokens extract of RefinedWeb based on which the commu-  
684 nity can train models from scratch to verify our claims; (2) the 1B and 7B parameters  
685 models pretrained on RefinedWeb for this paper, allowing the community to reproduce  
686 our evaluations by running the Eleuther AI Eval Harness [49] on these models.
- 687 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
688 were chosen)? [Yes] The model cards for the 1B and 7B models, available on the  
689 HuggingFace Hub where they are publicly released, contains detailed details of the  
690 architecture and pretraining hyperparameters.
- 691 (c) Did you report error bars (e.g., with respect to the random seed after running ex-  
692 periments multiple times)? [N/A] Pretraining large language models is extremely  
693 compute-intensive, and it is not practically feasible to have multiple complete runs  
694 across multiple seeds for the models in this paper.
- 695 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
696 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B.2.
- 697 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 698 (a) If your work uses existing assets, did you cite the creators? [Yes] We appropriately cite  
699 the works which have inspired and which are leveraged for this paper.
- 700 (b) Did you mention the license of the assets? [Yes] We refer to the CommonCrawl Terms  
701 of Use in our datasheet (Table 6), which users of RefinedWeb have to abide to.
- 702 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
703 Yes, we have linked to all assets released in this work (dataset and models).
- 704 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
705 using/curating? [Yes] All data collection is handled by CommonCrawl; they abide to  
706 no-crawl requests in the robots.txt file of the domain, and we have further included  
707 an opt-out mechanism for our dataset.
- 708 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
709 information or offensive content? [Yes] We further discuss this subject in our datasheet  
710 Table 6. We measured the prevalence of toxic content in Figure 5 finding it is in-line  
711 with curated datasets such as The Pile. We note that all information contained in our  
712 dataset is already available online, but that given its sheer scale it is impossible for us  
713 to guarantee the complete absence of PII.

- 714 5. If you used crowdsourcing or conducted research with human subjects...  
715 (a) Did you include the full text of instructions given to participants and screenshots, if  
716 applicable? [N/A]  
717 (b) Did you describe any potential participant risks, with links to Institutional Review  
718 Board (IRB) approvals, if applicable? [N/A]  
719 (c) Did you include the estimated hourly wage paid to participants and the total amount  
720 spent on participant compensation? [N/A]

## 721 B NeurIPS Datasets and Benchmarks track details

### 722 B.1 Accessing and using the data

723 To host artefacts related to this paper, we are leveraging the HuggingFace Hub. This both guarantees  
724 long-term availability, and standardization allowing for interoperability with tools built by the  
725 community. Our 600B tokens public RefinedWeb extract is made available using the datasets  
726 library [83], and our 1B and 7B models are released with the transformers library [84].

- 727 • **Falcon-RefinedWeb**, the 600B tokens extract we make publicly available: <https://huggingface.co/datasets/tiiuae/falcon-refinedweb> licensed under an ODC-By  
728 1.0 license, and users should also abide to the CommonCrawl ToU;
- 729 • **Falcon-RW-1/7B**, the two models we have pretrained on RefinedWeb-only for this paper:  
730 <https://huggingface.co/tiiuae/falcon-rw-1b> and <https://huggingface.co/tiiuae/falcon-rw-7b>, both licensed under an Apache 2.0 license.

733 The complete model and data cards on the hub contain useful code examples for getting started with  
734 our public released assets. For the cards included in this supplementary, we focused on documentation  
735 rather than technical details.

736 The DOI associated with the public RefinedWeb extract is 10.57967/hf/0737.

### 737 B.2 Compute resources

738 We used resources from AWS Sagemaker, training on P4d instances with eight A100 40GB per  
739 node. Nodes were interconnected with 50Gb/s of EFA interconnect. We estimate to have used  
740 55,000A100-hour for this project: 35,000 for the 7B model; 5,000 for the 1B; and 15,000 for small  
741 scale ablations and earlier experiments not reported in this paper. See Appendix I.4 for details on  
742 CPU processing.

### 743 B.3 Statement on resubmission

744 This work was previously submitted at ICML 2023, and was rejected with scores 8/7/6/3. Compared  
745 to this earlier version, we have made the following improvements based on reviewers' feedback:

- 746 • **Added further comparisons with state-of-the-art models**, such as OPT and Pythia.
- 747 • **Analysed the toxicity of RefinedWeb**, as presented in Figure 5, leveraging the Perspective  
748 API to demonstrate the prevalence of toxic content in RefinedWeb to be similar to The Pile.
- 749 • **Added an ablation on the effect of our pipeline on other web datasets**, as presented in  
750 Appendix G.1, demonstrating its wide applicability and generalizing our findings on the  
751 value of filtering and deduplication.
- 752 • **Improved the overall presentation**, by surfacing content from the Appendix on our  
753 evaluation and on the limitations of our work.

### 754 B.4 Statement on license

755 We release the dataset under the ODC-By-1.0 license, further requesting from users that they abide  
756 by the CommonCrawl Terms of Use. This is inspired by other public massive web crawls, such as C4  
757 [9]. This licensing permits sharing, reuse, and adaptation of the dataset, for any purpose, as long as

758 both CommonCrawl and this paper are acknowledged and provided appropriate credit. The purpose  
759 of this statement is to clarify the responsibilities and liabilities associated with the use of this dataset.  
760 While we have made every effort to ensure the accuracy and legality of the data contained within this  
761 dataset, we cannot guarantee its absolute completeness or correctness due to its scale.

762 Therefore, in the event that any rights, legal or otherwise, are violated through the use of this  
763 dataset, including but not limited to copyright infringement, privacy violations, or misuse of sensitive  
764 information, we, the authors, assume no liability for such violations.

765 By utilizing this dataset, you agree that any consequences, legal or otherwise, arising from the use of  
766 this dataset will be the sole responsibility of the user. You acknowledge that you will exercise due  
767 diligence and adhere to all applicable laws, regulations, and ethical guidelines when using the dataset.

768 By accessing, downloading, or using this dataset, you signify your acceptance of this statement and  
769 your commitment to abide by the terms and conditions of the ODC-By-1.0 license and by Terms of  
770 Use of CommonCrawl.

771 If you do not agree with the terms of this statement, the ODC-By-1.0 license, or the CommonCrawl  
772 Terms of Use, you are not authorized to use this dataset.

| MOTIVATION  |  |
|---|--|
| <b>For what purpose was the dataset created?</b>  | RefinedWeb was created to serve as a large-scale dataset for the pretraining of large language models. It may be used on its own, or augmented with curated sources (e.g., Wikipedia, StackOverflow).  |
| <b>Who created the dataset and on behalf of which entity?</b>   | The dataset was created by the Technology Innovation Institute.  |
| <b>Who funded the creation of the dataset?</b>  | The creation of the dataset was privately funded by the Technology Innovation Institute.   |
| <b>Any other comment?</b>   | RefinedWeb is built on-top of CommonCrawl, using the Macrodata Refinement Pipeline, which combines content extraction, filtering heuristics, and deduplication. In designing RefinedWeb, we abided to the following philosophy: (1) <b>Scale first.</b> We intend MDR to produce datasets to be used to train 40-200B parameters models, thus requiring trillions of tokens [4]. For English-only RefinedWeb, we target a size of 3-6 trillion tokens. Specifically, we eschew any labour intensive human curation process, and focus on CommonCrawl instead of disparate single-domain sources. (2) <b>Strict deduplication.</b> Inspired by the work of [30], which demonstrated the value of deduplication for large language models, we implement a rigorous deduplication pipeline. We combine both exact and fuzzy deduplication, and use strict settings leading to removal rates far higher than others have reported. (3) <b>Neutral filtering.</b> To avoid introducing further undesirable biases into the model [31, 44], we avoid using ML-based filtering outside of language identification. We stick to simple rules and heuristics, and use only URL filtering for adult content. |
| COMPOSITION   |  |
| <b>What do the instances that comprise the dataset represent?</b>   | Instances are text-only documents, corresponding to single web pages.  |
| <b>How many instances are there in total?</b>   | RefinedWeb contains ~10 billion documents, or around 5 trillion tokens. The public version is a subset representing a tenth of the full version.   |
| <b>Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?</b> | RefinedWeb is built using all CommonCrawl dumps until the 2023-06 one; it could be updated with additional dumps as they are released. The public release of RefinedWeb is a 600GT random extract of the 5,000GT of the full dataset (limited to 600GT for commercial reasons). For experiments, we randomly sampled from the public extract, or earlier development versions.   |
| <b>What data does each instance consist of?</b>   | Each instance is a text-only document, with metadata about its origin in CommonCrawl and source page URL. We also distribute a multimodal version of RefinedWeb, containing interlaced links to images.  |
| <b>Is there a label or target associated with each instance?</b>  | No.  |
| <b>Is any information missing from individual instances?</b>  | No.  |
| <b>Are relationships between individual instances made explicit?</b>  | No.  |

|  |   |
|--|---|
| <b>Are there recommended data splits?</b>  | No.   |
| <b>Are there any errors, sources of noise, or redundancies in the dataset?</b>   | Despite our best efforts to filter content that does not qualify as natural language, and to deduplicate documents, our pipeline may let through documents that may be considered as errors or redundant.         |
| <b>Is the dataset self-contained, or does it link to or otherwise rely on external resources?</b>  | The base version of the dataset is self-contained, but the multi-modal version is interlaced with links to images—these are not distributed as part of the dataset, and constitute an external source.            |
| <b>Does the dataset contain data that might be considered confidential?</b>  | All documents in RefinedWeb have been publicly available online.  |
| <b>Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?</b> | Yes, as this type of data is prevalent on the internet, it is likely our dataset contains such content. Notably, we estimate the prevalence of toxic content in the dataset to be similar to The Pile (Figure 5). |

#### COLLECTION

|   |   |
|---|---|
| <b>How was the data associated with each instance acquired?</b>                       | We downloaded with warcio publicly available .WET files from the CommonCrawl foundation.  |
| <b>What mechanisms or procedures were used to collect the data?</b>                   | We refer to the CommonCrawl website ( <a href="https://commoncrawl.org">commoncrawl.org</a> ) for details on how they collect data. |
| <b>If the dataset is a sample from a larger set, what was the sampling strategy?</b>  | Whenever we use subsets, we randomly sample from the original data.   |
| <b>Who was involved in the data collection process and how were they compensated?</b> | The original data collection was performed by CommonCrawl; authors from this paper were involved in retrieving it and preparing it. |
| <b>Over what timeframe was the data collected?</b>                                    | We use all CommonCrawl dumps from 2008 to January/February 2023.  |
| <b>Were any ethical review processes conducted?</b>                                   | No.   |

#### PREPROCESSING

|   |  |
|---|--|
| <b>Was any preprocessing/cleaning/labeling of the data done?</b>                      | Yes, we applied extensive preprocessing and cleaning of the data. We first filter URLs to remove adult content using a blacklist and a score system (Appendix 1.1), we then use <code>trafilatura</code> [46] to extract content from pages, and perform language identification with the <code>fastText</code> classifier from CCNet [27]. After this preprocessing stage, we filter data using heuristics from MassiveWeb [12] and our own line-wise corrections (Appendix 1.2). Finally, we run extensive deduplication, removing URLs revisited across dumps (Section 3.3) and performing subsequently fuzzy and exact substring deduplication, with each stage drawing from [30]. See Section 3 for further details and Table 2 for an outline. |
| <b>Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?</b> | During development, we saved intermediary outputs from our pipeline for investigations and for ablations—intermediary outputs exist for about 5% of RefinedWeb. We did not keep intermediary outputs for the final production version of the dataset due to storage and resource constraints.  |
| <b>Is the software that was used to preprocess/clean/label the data available?</b>    | No.  |

#### USES

|   |   |
|---|---|
| <b>Has the dataset been used for any tasks already?</b>   | Yes, this data has been used to develop large language models: both for scientific experiments (e.g., this paper) and production use. See Almazrouei et al. [82] for details.                               |
| <b>Is there a repository that links to any or all papers or systems that use the dataset?</b>   | On a voluntary/self-reporting basis, the HuggingFace Hub where this dataset is hosted will point to models trained using this dataset.  |
| <b>What (other) tasks could the dataset be used for?</b>  | RefinedWeb was built as a large-scale corpora representative of the web, and as such may see many downstream uses which are difficult to predict.   |
| <b>Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?</b> | For the public extract of RefinedWeb, we chose to only draw from the English version of the dataset, preventing multilingual applications.  |
| <b>Are there tasks for which the dataset should not be used?</b>  | Any tasks which may considered irresponsible or harmful.  |
| <b>DISTRIBUTION</b>   |   |
| <b>Will the dataset be distributed to third parties outside of the entity on behalf of which the dataset was created?</b>                                 | Yes, we make a 600B tokens extract publicly available for NLP practitioners. We currently don't plan to share the full version of the dataset.  |
| <b>How will the dataset will be distributed?</b>  | The dataset will be made available through the HuggingFace Hub, in the datasets format [83].  |
| <b>When will the dataset be distributed?</b>  | The dataset is available immediately.   |
| <b>Will the dataset be distributed under a copyright or other intellectual property license, and/or under applicable terms of use?</b>                    | The public extract is made available under an ODC-By 1.0 license; users should also abide to the CommonCrawl ToU: <a href="https://commoncrawl.org/terms-of-use/">https://commoncrawl.org/terms-of-use/</a> |
| <b>Have any third parties imposed IP-based or other restrictions on the data associated with the instances?</b>   | Not to our knowledge.   |
| <b>Do any export controls or other regulatory restrictions apply to the dataset?</b>  | Not to our knowledge.   |
| <b>MAINTENANCE</b>  |   |
| <b>Who will be supporting/hosting/maintaining the dataset?</b>  | The dataset will be hosted on the HuggingFace Hub, and we will release further versions if necessary based on opt-out requests.   |
| <b>How can the owner/curator/manager of the dataset be contacted?</b>   | falconllm@tii.ae  |
| <b>Is there an erratum?</b>   | No.   |
| <b>Will the dataset be updated?</b>   | Yes, for opt-out requests.  |
| <b>If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?</b>                                       | The license allows for the community to fork, build upon, and modify this dataset, as long as proper attribution is given.  |

Table 6: Datasheet for RefinedWeb, following the framework introduced by [85].

774 **D Falcon-RW Model Cards**

| <b>MODEL DETAILS</b>                             |  |
|--|--|
| <b>Organization</b>                              | The models were created by the Technology Innovation Institute.  |
| <b>Model date</b>                                | Falcon-RW models were trained in December 2022/January 2023.   |
| <b>Model type and information about training</b> | Falcon-RW are autoregressive Transformer models trained with a causal language modeling objective. Architecture based on GPT-3 [2], with ALiBi positional encodings [72] and FlashAttention [73]. See Section 4.1 for details.   |
| <b>Licence</b>                                   | Apache 2.0.  |
| <b>Point of contact</b>                          | falconllm@tii.ae   |
| <b>INTENDED USE</b>                              |  |
| <b>Primary intended uses</b>                     | Research on large language models, and the influence of adequately filtered and deduplicated web data on the properties of large language models (fairness, safety, limitations, etc.).  |
| <b>Primary intended users</b>                    | NLP researchers.   |
| <b>Out-of-scope use cases</b>                    | Production use without adequate assessment of risks and mitigation; any use cases which may be considered irresponsible or harmful.  |
| <b>FACTORS</b>                                   |  |
| <b>Relevant factors</b>                          | Falcon-RW models are trained on English data only, and will not generalize appropriately to other languages. Furthermore, as they are trained on a large-scale corpora representative of the web, they will carry the stereotypes and biases commonly encountered online.  |
| <b>Evaluation factors</b>                        | We evaluated the toxicity of the underlying pretraining dataset and found it to be in line with common curated pretraining datasets such as The Pile (see Figure 5). Note that this only accounts for toxicity under the definition of Perspective API: "content that is rude or disrespectful". Notably, this fails to include concerns about social biases or harmfulness. |
| <b>METRICS</b>                                   |  |
| <b>Model performance measures</b>                | We focus our evaluation on measuring the zero-shot generalization capabilities of our models across a wide range of tasks, leveraging the Eleuther AI language model evaluation harness [49].  |
| <b>Variation approaches</b>                      | Due to the costs associated with training Falcon-RW we cannot train the models multiple times and measure variability across training runs.  |
| <b>EVALUATION DATA</b>                           |  |
| <b>Datasets</b>                                  | We evaluate zero-shot accuracy on 18 varied tasks, detailed in Table 3.  |
| <b>Motivation</b>                                | We selected and aggregated tasks to build comparisons with other models in the literature (see Section 4.1; Appendix H.1 for details).   |
| <b>Preprocessing</b>                             | We use the default prompts and setup of [49].  |
| <b>TRAINING DATA</b>                             |  |
| <b>See the dedicated datasheet in Table 6.</b>   |  |

Table 7: **Model card for Falcon-RW**, following the framework introduced by [86].



775 **E Dataset analysis**

776 The large-scale and diverse nature of web corpora make them difficult to document and analyse  
 777 extensively; we provide some key metrics in the section, focusing on document lengths in Figure 4(a),  
 778 a breakdown of the top domain names in Figure 4(b), and the distribution of toxic content in Figure 5

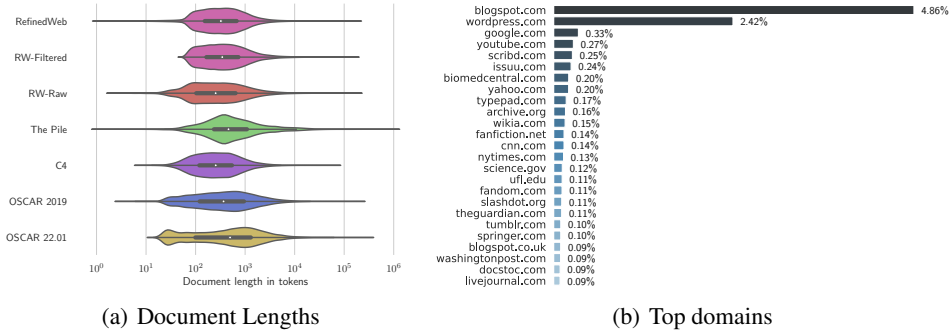


Figure 4: **Make-up of RefinedWeb in document lengths (left) and top domains (right).** (a) We find the OSCAR datasets and RW-Raw to have similar document length distributions; following filtering, most of the short documents are discarded from RW-Filtered. As deduplication removes spans, it shortens documents in RefinedWeb. We note the make-up of C4 and RefinedWeb to be similar, with a longer tail of short documents for RefinedWeb. Finally, The Pile exhibit a unique make-up, with a long tail of both long (books, etc.) and short documents. (b) Top domains in RefinedWeb span from popular content platforms (Blogspot, WordPress, Tumblr, etc.), to news websites (CNN, New York Times, etc.), and include also technical content such as BioMed Central or Springer.

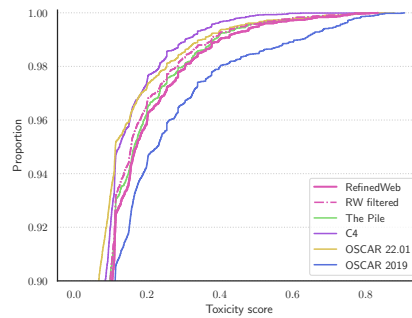


Figure 5: **Toxic content in RefinedWeb is distributed similarly to The Pile.** Cumulative proportion of documents below a given toxicity score, as evaluated by the Perspective API.

779 **F Multilingual RefinedWeb**

780 **Multilingual data.** Using the language identification filter, we classify processed CommonCrawl  
 781 data into 176 languages. Figure 6 shows the top 20 languages present in the data *excluding English*,  
 782 based on their relative contribution in descending order. 58.20% of all documents in CommonCrawl  
 783 were identified as English. We find the distribution of languages in CommonCrawl to only be partially  
 784 aligned with the worldwide distribution of language speakers [87]: Russian is over-represented  
 785 (2nd in CC but only 8th worldwide), Mandarin Chinese is under-represented (6-7th in CC but 2nd  
 786 worldwide), and Hindi does not show-up in the top 20 despite being the 3rd most spoken.

787 **Processing multilingual data.** The MDR pipeline can be used to process all languages: features  
 788 such as text extraction are language-agnostic, whereas specific filters such as line-wise corrections  
 789 need to typically be tuned for each individual language. We also found tuning deduplication param-  
 790 eters for individual languages to be beneficial.



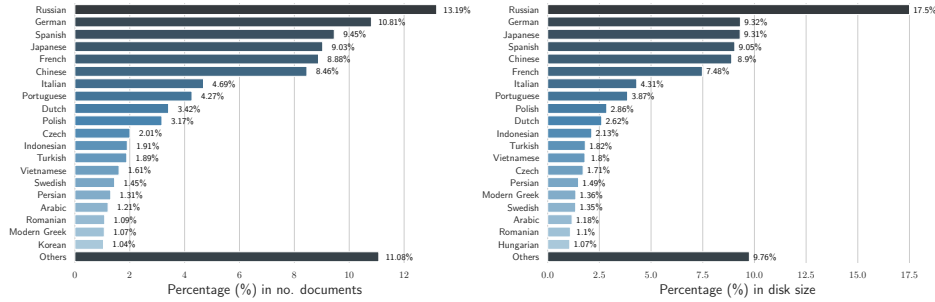


Figure 6: **The representation of languages in CommonCrawl does not align with the worldwide distribution of language speakers.** Top 20 languages (excluding English, which accounts for 58.20%) from processed CommonCrawl based on number of documents and disk size.

## 791 G Additional results

792 In this section, we present additional results obtained during the development of the Macrodata  
 793 Refinement pipeline. For Appendix G.1 and Appendix G.3, these were obtained using earlier  
 794 development versions of the dataset, so results are not directly comparable with the main text. For  
 795 Appendix G.2, this is based on the Falcon-RW models.

### 796 G.1 Small-scale ablations on deduplication approaches

797 We present results in Table 8—setup is similar to our earlier ablations, training 1B models for 30GT:

- 798 • **MinHash alone is insufficient**, it does not match the performance of exact deduplication.  
 799 Conversely, combining it with exact deduplication doesn't improve performance further.
- 800 • **Masking spanned duplicates degrades performance**, systematically underperforming  
 801 other approaches. Dropping and cutting spans perform similarly, although it's likely that  
 802 dropping documents slightly outperforms cutting.

803 We chose to apply MinHash before exact deduplication, as it is easier to scale: approximate dedu-  
 804 plication acts as a pruning phase, enabling us to scale deduplication further. Finally, we choose the  
 805 common option of cutting spans, as dropping resulted in even more stringent rejection rates which  
 806 would have compromised our ability to collect 5 trillion tokens.

Table 8: **MinHash alone is insufficient to match the performance of exact substring deduplication, and combining the two does not significantly improve performance. Of all of the exact substring approaches, masking duplicated spans underperform, but all others exhibit similar performance.**

✓ Minhash + Exact substring-Cut corresponds to our final deduplication setup. Perplexity in bits-per-bytes on The Pile (pile-bpb, lower is better), zero-shot performance aggregated over LAMBADA, PIQA, and HellaSwag (agg-dev). Best results in **bold**, best results with minhash in underline.

| Minhash | Exact substring            | pile-bpb ↓  | agg-dev-1 ↑  |
|---------|----------------------------|-------------|--------------|
|         | <u>RefinedWeb-Filtered</u> | 1.11        | 43.51        |
|         | Mask                       | 1.08        | 45.84        |
| ✓       | Mask                       | 1.07        | 46.28        |
| ✓       |                            | 1.07        | 46.57        |
| ✓       | <u>Cut</u>                 | <b>1.05</b> | 47.11        |
|         | Cut                        | 1.06        | 47.24        |
| ✓       | Drop partial               | <b>1.05</b> | 47.25        |
|         | Drop any                   | 1.07        | 47.77        |
| ✓       | Drop any                   | 1.07        | <u>47.86</u> |
|         | Drop partial               | 1.06        | <b>47.97</b> |
|         | <u>Pile</u>                | 0.88        | 43.70        |

807 **G.2 Language modeling evaluation**

808 Along with our aggregates, we also evaluated perplexity on Wikitext (Table 9). We found that models  
 809 trained on RefinedWeb achieve performance close to that of models trained on The Pile. Importantly,  
 810 we note that RefinedWeb does not contain any content from Wikipedia – it is explicitly filtered out at  
 811 the URL level. We believe this accounts for most of the difference in perplexity, as RW models may  
 812 not be familiar with the idiosyncrasies of Wikitext (e.g., layout of an article, etc.)

Table 9: **Models trained on RefinedWeb achieve performance close to models trained on The Pile on Wikitext, despite not having seen any content from Wikipedia.** Perplexity in bits-per-bytes on Wikitext (wiki-bpb, lower is better.)

| Model size | 1B       | 7B   |      |
|------------|----------|------|------|
| Dataset    | The Pile | RW   | RW   |
| wiki-bpb ↓ | 0.64     | 0.66 | 0.60 |

813 **G.3 Does deduplication help with multiple epochs?**

814 Earlier in this work, we outlined that to scale pretraining data, practitioners had two choices: (1)  
 815 improve data collection, which is the avenue we chose to pursue; (2) train models on multiple epochs  
 816 of the same data. Due to current uncertainties in the ability of larger models to sustain multiple  
 817 epochs without adverse effects [41], we focused on (1). A fairly rational question regarding (2) is  
 818 whether deduplication may improve the situation, and whether deduplicated data may be able to  
 819 sustain more epochs without compromising model quality.

820 We train 1B parameters models on 30GT of RW and RW-Filtered. We keep the number of pretraining  
 821 tokens fixed, but train for 1, 5, 25, and 100 epochs. This is a small-scale, limited set-up, which would  
 822 have to be improved to obtain definitive results. We plot the degradation in performance compared to  
 823 a single epoch in Figure 7(a) and the gap between RW and RW-F in Figure 7(b). We find that the  
 824 absolute degradation is less important for RefinedWeb than for RefinedWeb-Filtered; furthermore,  
 825 the gap widens with increasing number of epochs. However, we observe significant variability across  
 826 tasks.

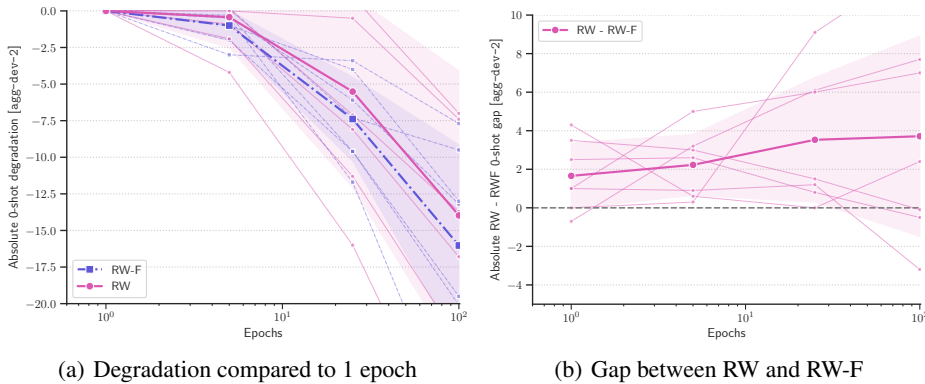


Figure 7: **Deduplication may reduce the degradation in performance incurred by multiple epochs.** However, our experiments were only performed at small-scale (1B models trained on 30GT), and we see high variability in outcomes across tasks. Zero-shot performance measured on the agg-dev-2 aggregate (HellaSwag, PIQA, ARC, BoolQ, COPA, MRPC, SciQ). Individual curves for per-task results and 1-σ standard deviation across all tasks in the aggregate in transparent.

Table 10: **We source evaluation results from a variety of papers across the literature, maximizing task coverage.** Although most results come from the EAI Evaluation Harness [49], results from PaLM and GPT-3 are sourced from their respective papers. Note in Figure 1 that the results from the GPT-3 paper are still ahead of results obtained through the API with the EAI evaluation harness.

| Models             | Aggregates reported | Source of results | EAI eval harness? |
|--------------------|---------------------|-------------------|-------------------|
| Ours               | main, core, ext     | This paper        | ✓                 |
| BS-A&S*            | main, core          | [11]              | ✓                 |
| GPT-Neo*           | main, core          | [11]              | ✓                 |
| PaLM <sup>†</sup>  | main                | [22]              |                   |
| GPT-3 API*         | main, core          | [11]              | ✓                 |
| GPT-3 <sup>†</sup> | main                | [2]               |                   |
| Aleph Alpha*       | core                | [71]              | ✓                 |
| Cerebras-GPT*      | core                | [48]              | ✓                 |
| FairSeq*           | core                | [70]              | ✓                 |
| Pythia(-Dedup)*    | core                | [48]              | ✓                 |
| OPT*               | core                | [48]              | ✓                 |
| GPT-J*             | core                | [70]              | ✓                 |
| GPT-NeoX 20B*      | core                | [70]              | ✓                 |

## 827 H Tasks, models, and datasets from the state-of-the-art

### 828 H.1 Task aggregates

829 We average zero-shot performance over diverse task aggregates, outlined in Table 3:

- 830 • **small**: small-scale ablation studies, tasks with non-zero performance for 1B parameters
- 831 models trained on 30GT;
- 832 • **core**: comparisons with a wide range of models, based on the tasks reported in [48];
- 833 • **main**: tasks available in the GPT-3 and PaLM papers [2, 22];
- 834 • **ext**: tasks available in the work of the BigScience Architecture and Scaling group [11].

835 Detailed evaluation results are available in this dedicated spreadsheet: [https://docs.google.com/spreadsheets/d/1u0HqZVtNxe2bYmF\\_1lQneR0FH-s6T0njiRE0bV1LtEA/](https://docs.google.com/spreadsheets/d/1u0HqZVtNxe2bYmF_1lQneR0FH-s6T0njiRE0bV1LtEA/). When comparing

836 with other models, we source results from papers detailed in Table 10.

837

### 838 H.2 Models

839 We compare against 10 series of models trained on a variety of curated corpora, presented in Table 11.

840

841 **Cerebras-GPT with  $\mu$ -parametrization.** The Cerebras-GPT series [48] also comes in a smaller

842 series, up to 2.7B parameters, using  $\mu$ -parametrization [88]. As we found the performance of this

843 smaller series to be close to the main series of models (see Figure 8), and as it does not include

844 models of a similar compute scale as the ones we compare to, we do not report it in our main figures.

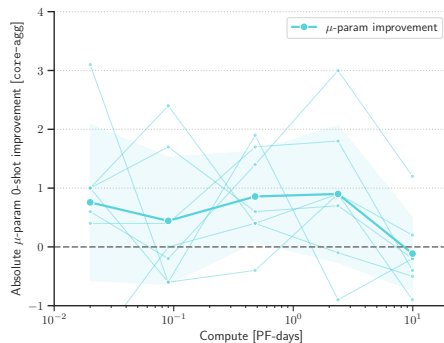


Figure 8:  $\mu$ -parametrization [88] slightly improves performance in the Cerebras-GPT series [48]. Zero-shot performance on our core aggregate, gap between Cerebras-GPT with  $\mu$ -param and without. Individual curves for per-task results and 1- $\sigma$  standard deviation across all tasks in shade.

845 **Pythia and deduplication.** The Pythia series of models is available in two flavours: one trained  
 846 on the vanilla version of The Pile, and another trained on a version deduplicated with MinHash.  
 847 Performance between these two flavours was noted to minimally differ [42]; in Figure 9, we find  
 848 the deduplicated version may be slightly ahead of the non-deduplicated one under our aggregate.  
 849 The higher end of this improvement is broadly in line with our findings in Table 5. Nevertheless, a  
 850 difference in our findings and theirs remain. We posit a few possible hypotheses:

- 851 • **Differences between curated and web data.** It is possible that web data is more sensitive to  
 852 duplicates. For instance, the most common duplicates in web data (e.g., spam) may be more  
 853 detrimental than the most common duplicates in curated data. This suggests a qualitative  
 854 component to deduplication that we have not studied in this work.
- 855 • **Differences in deduplication pipeline.** Because [42] uses the MinHash settings from [30],  
 856 they are mostly identical to ours. However, we also apply exact deduplication: while their  
 857 deduplication incurs a 30% reduction in size, our deduplication is more aggressive, resulting  
 858 in a 45% reduction in size. This may explain why our results in Table 5 show a stronger  
 859 gain from deduplication than theirs in Figure 9.
- 860 • **Differences in pretraining.** Finally, we note that [42] chooses to perform a partial extra  
 861 epoch on the deduplicated data to reach 300GT, while we always perform a single epoch.  
 862 Their setting corresponds to a data-constrained scenario, which is more realistic for the  
 863 curated data they study; for us, web data is plentiful, so deduplication never truly limits the  
 864 size of the datasets we can use.

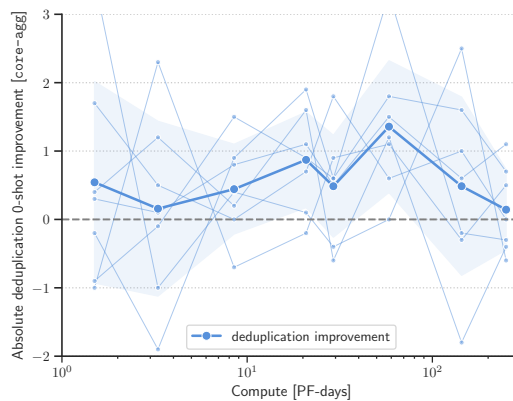


Figure 9: **In our core aggregate, deduplication brings a small improvement to the Pythia suite [42].** Zero-shot performance on our core aggregate, gap between Pythia trained on the deduplicated and vanilla Pile. Individual curves for per-task results and  $1-\sigma$  standard deviation across all tasks in the aggregate in transparent.

### 865 H.3 Datasets

866 We extend on Table 1 in Table 12, providing details on the filtering and deduplication strategies used  
 867 across the literature.

Table 11: **Full-scale models trained on RefinedWeb (Falcon-RW) and other models from the state-of-the-art.** Across models trained on The Pile, the Pythia models are the closest to our achitecture: they use FlashAttention with rotary embeddings—with for only notably exception the use of parallel attention and feedforward for their models. Training budget  $C$  in PF-days calculated using  $C = 6ND$ , with  $N$  the number of parameters, and  $D$  the pretraining dataset size [5].

| Series             | GPT-3 (paper) <sup>†</sup> |       | GPT-3 (API)* |       | BigScience* | PaLM <sup>†</sup> | Ours        |           |       |
|--------------------|----------------------------|-------|--------------|-------|-------------|-------------------|-------------|-----------|-------|
| <b>Model</b>       | XL                         | XXL   | babbage      | curie | BS-A&S      | PaLM-8B           | Ours (Pile) | Ours (RW) |       |
| <b>Dataset</b>     | GPT-3                      | GPT-3 | GPT-3        | GPT-3 | Pile        | PaLM              | Pile        | RW        | RW    |
| <b>Params.</b>     | 1.3B                       | 6.7B  | 1.3B         | 6.7B  | 1.3B        | 8.6B              | 1.3B        | 1.3B      | 7.5B  |
| <b>Pretraining</b> | 300GT                      | 300GT | 300GT        | 300GT | 300GT       | 780GT             | 350GT       | 350GT     | 350GT |
| <b>PF-days</b>     | 27                         | 140   | 27           | 140   | 27          | 466               | 32          | 32        | 182   |
| <b>Citation</b>    |                            |       | [2]          |       | [11]        | [22]              | This paper  |           |       |

| Series             | EleutherAI* |       |              | Pythia*        |
|--------------------|-------------|-------|--------------|----------------|
| <b>Model</b>       | GPT-Neo     | GPT-J | GPT-NeoX 20B | Pythia(-Dedup) |
| <b>Dataset</b>     | Pile        | Pile  | Pile         | Pile (dedup)   |
| <b>Params.</b>     | 1.3B        | 6.7B  | 20B          | 70M-12B        |
| <b>Pretraining</b> | 380GT       | 402GT | 472GT        | 300GT          |
| <b>PF-days</b>     | 34          | 187   | 656          | 1.5 - 250      |
| <b>Citation</b>    | [77]        | [74]  | [70]         | [42]           |

| Series             | Aleph Alpha*       | Cerebras-GPT* | OPT*                    | FairSeq*  |
|--------------------|--------------------|---------------|-------------------------|-----------|
| <b>Model</b>       | Luminous           | Cerebras-GPT  | OPT                     | FairSeq   |
| <b>Dataset</b>     | <i>undisclosed</i> | Pile          | Pile (subset) + curated | curated   |
| <b>Params.</b>     | 13B                | 111M-13B      | 125M - 175B             | 1.3 - 13B |
| <b>Pretraining</b> | 400GT              | 2 - 257GT     | 300GT                   | 300GT     |
| <b>PF-days</b>     | 361                | 0.02 - 232    | 3 - 3646                | 27 - 271  |
| <b>Citation</b>    | [71]               | [48]          | [43]                    | [76]      |

Table 12: **Common massive web-scrape and LLM English datasets.** Datasets such as OSCAR and C4 also have significant multilingual versions, which have enjoyed wide adoption [89]. For OSCAR, the size corresponds to the non-deduplicated version, and is estimated from the number of words x0,75 (average number of words per tokens).

| General information         |   | Web data |              |      |                  | Language ID                     | Heuristics                          | Content filtering           | Deduplication   |
|-----------------------------|---|----------|--------------|------|------------------|---------------------------------|-------------------------------------|-----------------------------|---|
| Dataset                     | Notable models  | Size     | Availability | Web  | HTML extraction  |                                 |                                     |                             |   |
| <b>MASSIVE WEB DATASETS</b> |   |          |              |      |                  |                                 |                                     |                             |   |
| <b>C4</b> [9]               | T5 [9]  | ~ 360GT  | Public       | 100% | .WET files       | Document-level w/ langdetect    | Document and line-level             | Rules-based: code, NSFW     | <b>Exact:</b> three sentences span  |
| <b>OSCAR 21.09</b> [8]      |   | ~ 370GT  | Public       | 100% | .WET files       | Line-level w/ fastText [28]     | Line < 100 characters               | None                        | (optional) <b>Exact:</b> per line (~ 55% removed)                         |
| <b>OSCAR 22.01</b> [90]     |   | ~ 283GT  | Public       | 100% | .WET files       | Document-level w/ fastText [28] | Line-level, optional document-level | Optional NSFW blacklist     | (optional) <b>Exact:</b> per line   |
| <b>CURATED DATASETS</b>     |   |          |              |      |                  |                                 |                                     |                             |   |
| <b>■ GPT-3</b> [2]          |   | 300GT    | Private      | 60%  | Unknown          | Unknown                         | Unknown                             | fastText trained on HQ-data | <b>Fuzzy:</b> min-hash with 10 hashes (~ 10% removed)                     |
| <b>▼ The Pile</b> [10]      | GPT-J [74], GPT-NeoX-20B [70], Pythia [42], Cerebras-GPT [48] | ~ 340GT  | Public       | 18%  | jusText [91]     | Document-level pycld2 [92]      | None                                | fastText on curated crawl   | <b>Fuzzy:</b> min-hash with 10 hashes, sim. threshold 0.5 (~ 26% removed) |
| <b>MassiveWeb</b> [12]      | Gopher [12], Chinchilla [4]                                   | 1, 400GT | Private      | 48%  | Custom           | Unknown                         | Document-level                      | SafeSearch                  | <b>Exact &amp; fuzzy:</b> exact documents, minhash w/ sim. threshold 0.8  |
| <b>★ PaLM</b> [22]          |   | 780GT    | Private      | 27%  | Unknown          | Unknown                         | Document-level                      | ML-based filter on HQ data  | Unknown   |
| <b>OURS</b>                 |   |          |              |      |                  |                                 |                                     |                             |   |
| <b>● REFINEDWEB</b>         | Falcon-RW   | 5,000GT  | 500GT Public | 100% | trafilatura [46] | From [27]                       | Document and line-level             | URL blacklist               | <b>Exact &amp; fuzzy</b>  |

## 868 I Details of the Macrodata Refinement pipeline

### 869 I.1 URL filtering

870 As discussed in Section 3.1, we base our filtering of adult documents only on the URL itself, and not  
871 on the content of the documents. This design choice was motivated by: (1) challenges in avoiding  
872 overfiltering content from minorities when using ML-based classifiers on the content of documents  
873 [44]; (2) NSFW words block-list applied on content (such as the one used in C4) also resulting in  
874 overfiltering of legal and medical content [31].

875 Our URL filtering focuses on finding domains that are related to adult content, that may be harmful  
876 to users, or that are very likely to contain mostly unstructured text/spam (e.g., file hosting websites).  
877 First, we aggregated a list of 4.6M domains, detailed in Appendix I.1.1, that we explicitly ban; then,  
878 we built a simple URL scoring system, based on matching subwords in the URL against a list of  
879 words we curated (see Appendix I.1.2). We curated this list of words based on manual inspection,  
880 cross-referencing results with pages surfaced by ToxicBERT as being outliers in toxicity [93].

#### 881 I.1.1 URL Blocklist

882 **Origin of the list.** We use an aggregated list of about 4.6M URLs that we explicitly ban. This  
883 list is broken in categories (e.g. pornography, gambling); we outline the categories we selected in  
884 Table 13. The list is regularly updated, with an original intended usage as a blocklist for universities.

885 **Curation.** We noticed the list blocked a number of domains inappropriately; while these domains  
886 were few (<100), they accounted for a significant portion of the data filtered by the list, as these were  
887 rather prolific domains, with thousands of pages of content. To identify these false positive domains,  
888 we applied the blocklist to a subset of 832M pages. 6.04M (0.73%) pages matched with the blocklist,  
889 and the number of occurrences per URL ranged from 1 to 79k. We manually inspected all URLs  
890 matched more than 4k times, which represented an appreciable portion of the dataset. We found a  
891 number of benign domains, such as pop culture news websites, or blogging platforms, which we  
removed from the list.

Table 13: We select categories likely to contain adult or malicious content, as well as spam or unstructured text.

| Category    | Description  | Number of links |
|-------------|--|-----------------|
| adult       | adult websites: from eroticism to hard pornography | 4516478         |
| phishing    | phishing websites, malwares, etc.                  | 42445           |
| dating      | dating websites                                    | 3829            |
| gambling    | online casino                                      | 1365            |
| filehosting | websites hosting files, videos, pictures, music    | 909             |
| ddos        | websites related to ddos attacks                   | 421             |
| agressif    | hate, racism, etc                                  | 390             |
| chat        | online chat websites                               | 244             |
| mixed adult | websites with some adult content                   | 153             |
| arjel       | French regulated gambling websites                 | 69              |

892

#### 893 I.1.2 URL Scoring with a Word-List

894 To score URLs, we used three matching patterns based on a soft, hard, and strict violation word-list:

- 895 • **Strict subword matching:** `http://foobann.edsub-wo.rdbar.com/any/bar`, matching words  
896 such as `xvideos`, `groupsex`;
- 897 • **Hard whole word matching:** `http://www.foo.bannedword-bar.com`, with words such as  
898 `porn`, `xxx`, `orgy`;

<https://dsi.ut-capitole.fr/blacklists/>



899 • **Soft words matching:** <http://www.foo.soft1-bar-soft2.com>, with "softer" words such as  
900 `sex`, `webcam`, `escort`.

901 Each list is associated with a different level of severity: for the strictest one (strict subword matching),  
902 we ban any URL matching a banned word in its substrings (as fraudulent websites may attempt  
903 to escape similar recognition schemes by breaking-up adult keywords); for the hard whole word  
904 matching, we ban URLs with a whole word matching in the list; finally, a minimum of two matches  
905 are required with the soft word matching.

906 We curated the lists based on manual inspection of the data, informed by top hits reported by  
907 ToxicBERT. For the strict subword matching, we included words that were unequivocally related to  
908 adult content (e.g., `groupsex`). We avoided partial unclear matches (e.g., `ass`), that may be part of  
909 neutral words (e.g., `massachusetts`). In the soft word list, we included words that do not constitute  
910 a sufficient reason to discard the document on their own, but which are suspicious when multiple  
911 words from the list result in a match. This helped with keeping medical or legal content unaffected  
912 (e.g., a single match of `dick`).

### 913 I.1.3 Excluded High Quality Sources

914 Since our paper focuses on the study of RefinedWeb alone, we chose to exclude common online  
915 sources of curated data from it. This serves two objectives: (1) it strengthens our results, by ensuring  
916 that RefinedWeb doesn't end-up actually being made mostly of known high-quality sources (e.g.,  
917 Wikipedia represents a significant portion of C4); (2) future works may be interested in combining  
918 RefinedWeb with existing curated corpora, which would require further deduplication if they are  
919 included in RefinedWeb. Accordingly, we remove common sources used in The Pile [10] from  
920 RefinedWeb. The full list of curated data sources domains that we blocked is in Table 14.

Table 14: **RefinedWeb is stripped from common so-called high-quality sources to simplify combining it with existing curated corpora.** This blocklist is applied at the URL filtering stage, along with the adult content blocklist.

| Curated data source   | Domain name blocked  |
|-----------------------|----------------------|
| arxiv                 | arxiv.org            |
| AskUbuntu             | askubuntu.com        |
| StackOverflow         | stackoverflow.com    |
|                       | stackapps.com        |
|                       | stackexchange.com    |
|                       | mathoverflow.net     |
| NIH Abstracts         | exporter.nih.gov     |
|                       | ncbi.nlm.nih.gov     |
| Github                | github.com           |
| Ubuntu IRC            | irclogs.ubuntu.com   |
| HackerNews            | news.ycombinator.com |
| FreeLaw               | courtlistener.com    |
| Reddit                | reddit.com           |
| Europarl              | statmt.org           |
| United States Patents | uspto.gov            |
| Wikipedia             | wikipedia.org        |

### 921 I.2 Line-wise filtering

922 Despite the improvements brought forth by running text extraction with Trafilatura, we found that a  
923 number of irrelevant lines still seeped through. These lines are usually related to navigation menus,  
924 call to actions, or social media counters. Following manual inspection of the data, we devised a  
925 line-wise filtering strategy. We analyse documents line-by-line, and discard or edit the lines based on  
926 the following rules:

- 927 • If it is mainly composed of uppercase characters (discard);
- 928 • If it is only composed of numerical characters (discard);
- 929 • If it is a counter (e.g. 3 likes) (discard);
- 930 • If it only contains one word (discard);
- 931 • If it is short ( $\leq 10$  words) and matches a pattern (edit):
  - 932 – At the beginning of the line (e.g. sign-in);
  - 933 – At the end of the line (e.g. Read more...);
  - 934 – Anywhere in the line (e.g. items in cart).

935 Finally, if the words in the flagged lines represent more than 5% of the total document words, the  
 936 document is discarded. We derived these filters through manual inspection of the data, and note that  
 937 they require adaptation across languages.

### 938 I.3 Deduplication

939 We make use of the two deduplication methods described in [30]: EXACTSUBSTR and NEARDEDUP  
 940 (detailed in Appendix I.3.1 and Appendix I.3.2; see Appendix I for samples of duplicates).

941 We start with the most scalable approach, NEARDEDUP. We remove similar documents by applying  
 942 MinHash [34], whereby a signature/sketch supporting efficient approximate similarity queries is  
 943 computed for each document in the dataset, and document pairs with a high  $n$ -gram overlap are  
 944 identified.

945 We then use EXACTSUBSTR, leveraging the implementation from [30] to identify ranges of  
 946 exact duplicate text of at least 50 tokens. We experiment with three different approaches for these  
 947 ranges: EXACTSUBSTR-CUT, where we remove them from the original text, as done in the original  
 948 implementation; EXACTSUBSTR-MASK, where the dataset is unchanged but we do not compute the  
 949 loss on the duplicated ranges; and EXACTSUBSTR-DROP, where we simply drop an entire document  
 950 if the duplicated ranges make up more than a certain percentage of its content.

951 We present small-scale ablations around these different approaches in Appendix G.1.

#### 952 I.3.1 MinHash Approximate Matching

953 We employ MinHash to find approximate duplicate documents in our web corpora at a very large  
 954 scale. This technique allows us to identify templated pages or otherwise very similar content where  
 955 most of the interspersed duplicated sections are small enough to not be identified by exact matching  
 956 methods (anything smaller than 50 tokens).

957 **Signing.** We start by normalizing the content to increase recall: punctuation is removed, text is  
 958 lowercased, NFD Unicode normalization is applied, accents are removed, and all whitespace is  
 959 normalized. We tokenize the resulting text using the GPT-2 tokenizer [17] and obtain the set of  
 960 unique  $n$ -grams for each document. Hash functions are used to obtain a signature for each document:  
 961 for each hash function, the smallest value is kept from hashing every unique  $n$ -gram in the document.  
 962 If two documents are similar, then there is a high probability that they will have the same minimum  
 963 hash (MinHash) for at least some of the hash functions used [34]. The ratio of matching hashes  
 964 between two documents approximates the Jaccard Similarity [94] of the sets of their unique  $n$ -grams  
 965 (the sets being  $d_i$  and  $d_j$ ):

$$J(d_i, d_j) = \frac{|d_i \cap d_j|}{|d_i \cup d_j|} \quad (1)$$

966 **Matching.** Since comparing MinHash signatures between every possible document pair is compu-  
 967 tationally expensive, we apply a locality sensitive hashing version of MinHash, MinHash LSH. A  
 968 document signature is split into  $r$  buckets, each with  $b$  minhashes. Documents are indexed by these  $b$

---

<sup>1</sup><https://github.com/google-research/deduplicate-text-datasets>

969 minhashes on each of the  $r$  buckets, and we mark two documents as duplicates if their  $b$  minhashes  
970 are exactly the same on at least one of the buckets. These two parameters,  $b$  and  $r$ , will determine  
971 the probability that similar documents will be detected. For two documents  $i$  and  $j$  whose ratio of  
972 matching hashes between their MinHash signatures is  $s_{i,j}$ , the probability that there is a match in a  
973 given bucket is  $s_{i,j}^b$ ; the probability that there isn't a match in any of the buckets is  $(1 - s_{i,j}^b)^r$ ; and  
974 finally that there is a match in at least one of the buckets:

$$P = 1 - (1 - s_{i,j}^b)^r \quad (2)$$

975 We use the same parameters as [30]:  $n = 5$  (5-grams);  $b = 20$  and  $r = 450$ . This means that for each  
976 document, we compute a total of 9000 minhashes, and that the probability that a document pair with  
977 similarity 0.75 or 0.8 will be marked as duplicates will be 76% and 99.4% (respectively), diminishing  
978 rapidly for smaller similarity values.

979 Finally, we cluster documents across all buckets — if documents A and B match in one bucket and B  
980 and C in another, A-B-C becomes a cluster. We randomly remove all but one of the documents in  
981 each cluster.

982 [30] also proposed filtering down on false positives by computing the real Jaccard similarity, or other  
983 metrics such as the edit similarity between identified document pairs. Given the large amount of data  
984 we have available across all of CommonCrawl, and that our main concern is improving recall, we  
985 decided to skip this additional step.

### 986 I.3.2 Exact substring deduplication

987 We make use of the EXACTSUBSTR implementation publicly released by [30] for exact text matching.  
988 We apply exact substring deduplication to data that has already been deduplicated by MinHash,  
989 reducing by nearly 40% size of the dataset on which we have to operate. EXACTSUBSTR will find  
990 long strings of text that are present, character for character, across multiple documents. Some of  
991 these may have escaped the earlier stage of approximate deduplication: they might not constitute  
992 a big enough portion of the document; one document might have repeated sections sourced across  
993 many different documents; or they may simply not have been found due to the approximate nature of  
994 MinHash.

995 **Finding duplicates.** EXACTSUBSTR concatenates all the documents in the dataset to create a single  
996 long text sequence; then, it builds a suffix array [33] in linear time—an array of the indexes to a  
997 lexicographical ordering of all the suffixes in the sequence. Finally, duplicate sequences can also  
998 be found in linear time using the suffix array, by simply traversing the ordered list of suffixes and  
999 comparing the beginning of each pair of two consecutive suffixes.

1000 We apply the same normalization and tokenization as for MinHash to the content of our documents  
1001 before concatenating them. One important difference is that reversibility is important: for MinHash,  
1002 we were discarding entire documents, and thus never relying on the normalized+tokenized repre-  
1003 sentation for downstream use. Here, once we have identified duplicate normalized+tokenized spans,  
1004 we need to revert to the original span to remove it. Accordingly, we include normalization in the  
1005 tokenization process, and validate that the process is reversible.

1006 If a match is longer than 50 tokens, there will be multiple overlapping duplicated ranges. These  
1007 overlapping duplicated ranges in the concatenated dataset sequence are merged before we save them  
1008 to a file. We then take these ranges and retrieve the original document that produced them, obtaining  
1009 the character substrings corresponding to the duplicated token ranges.

1010 **Removing duplicates.** We considered applying the following transformations to the duplicate  
1011 spans:

- 1012 • EXACTSUBSTR-CUT: we remove the duplicated spans, and discard documents where there  
1013 are fewer than 20 non-duplicated characters left—this is the vanilla setting used by [30];

- 1014 • EXACTSUBSTR-MASK: we loss-mask the duplicated spans, preventing a loss from being  
1015 computed on the duplicated text during pretraining, and discard documents where there are  
1016 fewer than 20 non-masked characters left.
- 1017 • EXACTSUBSTR-DROPPARTIAL: if more than 20% of the document is duplicated, we  
1018 remove the entire document;
- 1019 • EXACTSUBSTR-DROPANY: we drop any document with a duplicated span in it.

1020 Broadly speaking, EXACTSUBSTR-CUT might remove text mid-sentence resulting in disconnected  
1021 text; EXACTSUBSTR-MASK does not have this issue, but might be less efficient as a significant portion  
1022 of the training tokens will not directly contribute to updating the model’s weights; EXACTSUBSTR-  
1023 DROP might still keep considerable duplicated sections in its PARTIAL version, especially on larger  
1024 documents, while the ANY version might be overly aggressive. Following ablations in Appendix [G.1](#),  
1025 we choose to stick with the vanilla approach, EXACTSUBSTR-CUT.

1026 Note that in all cases, while MinHash keeps one copy of the duplicated documents, our exact  
1027 deduplication removes all copies of the duplicated span.

#### 1028 **I.4 Execution environment**

1029 Most data processing took place in large CPU clusters, with 100-250 AWS c5.18xlarge instances;  
1030 each instance has 72 vCPUs and 144 GiB of memory. We usually run with 10,000-20,000 vCPUs in  
1031 the cluster, enabling rapid parallel processing.

1032 For EXACTSUBSTR, the entire dataset being deduplicated needs to be loaded onto memory: we  
1033 leveraged the AWS x2iedn instances, which come with up to 2 TiB of memory in a single instance.

1034 **J Deduplication samples from RefinedWeb**

1035 **J.1 MinHash clusters**

1036 We report the 8 largest duplicate clusters found by MinHash in Table 15 – each spanning hundreds  
 1037 of thousands of documents. We also found a large number of duplicate document pairs to be due to  
 1038 different URL GET parameters not resulting in significantly different content. An example of this  
 1039 behaviour can be seen in the URLs presented in Table 16

Table 15: **Top-8 largest MinHash clusters found when building RefinedWeb.** We cut some of the longest samples in the interest of readability, only keeping a brief description.

| Description  | Example document   |
|--|--|
| Wordpress sitemap notice generated by the Google Sitemap Generator Plugin                  | This is a XML Sitemap which is supposed to be processed by search engines which follow the XML Sitemap standard like Ask.com, Bing, Google and Yahoo. It was generated using the WordPress content management system and the Google Sitemap Generator Plugin by Arne Brachhold. You can find more information about XML sitemaps on sitemaps.org and Google’s list of sitemap programs. This file contains links to sub-sitemaps, follow them to see the actual sitemap content.   |
| Cloudflare notice to enable Javascript   |  |
| Templated disability notice, with different phone numbers across pages                     | Welcome to our website! As we have the ability to list over one million items on our website (our selection changes all of the time), it is not feasible for a company our size to record and playback the descriptions on every item on our website. However, if you are an American with a disability we are here to help you. Please call our disability services phone line at [redacted] or [redacted] during regular business hours and one of our kind and friendly personal shoppers will help you navigate through our website, help conduct advanced searches, help you choose the item you are looking for with the specifications you are seeking, read you the specifications of any item and consult with you about the products themselves. There is no charge for the help of this personal shopper for any American with a disability. Finally, your personal shopper will explain our Privacy Policy and Terms of Service, and help you place an order if you so desire. |
| Templated cookies notice   |  |
| Templated domain name for sale page  |  |
| www.metoperashop.org and sub-URLs, with content changes but always the same (large) footer |  |
| Different pages across more than 80 different domain names but with a common section       | DC Customers also liked: Special event items are produced by manufacturers only after the outcome of a game or event. These are advanced sale items and will ship immediately after they are received in our warehouse. Manufacturer direct items are shipped directly from the manufacturer. These items are not available for international or expedited shipping. Customized items can be personalized with options such as your name, your favorite number, and/or designs. Some options may be limited by league rules.   |
| http://www.boxofficemojo.com/daily and sub-URLs  |  |

Table 16: URL with different GET parameters don't always result in significantly different page content.

|   |   |
|---|---|
| <pre>http://gamesandbiz.blogspot.com/2010/07/bad-reviews-can-hurt-game-sales.html?showComment=1278486430242</pre>   | <pre>http://gamesandbiz.blogspot.com/2010/07/bad-reviews-can-hurt-game-sales.html?showComment=1278499674195</pre>   |
| <pre>https://www.ocean-oxygen.org/home;jsessionid=1E3290E84F668552FAC643D0A8F81BEC?p_p_id=122_INSTANCE_Zy6zjkRLAg7v&amp;p_p_lifecycle=0&amp;p_p_state=normal&amp;p_p_mode=view&amp;p_p_col_id=column-2&amp;p_p_col_pos=1&amp;p_p_col_count=6&amp;p_r_p_564233524_resetCur=true&amp;p_r_p_564233524_categoryId=1346016</pre> | <pre>https://www.ocean-oxygen.org/home?p_p_id=122_INSTANCE_Zy6zjkRLAg7v&amp;p_p_lifecycle=0&amp;p_p_state=normal&amp;p_p_mode=view&amp;p_p_col_id=column-2&amp;p_p_col_pos=1&amp;p_p_col_count=6&amp;p_r_p_564233524_resetCur=true&amp;p_r_p_564233524_categoryId=1346016</pre> |

1040 **J.2 Exact substring matches**

1041 Examples of exact matches found by exact substring deduplication can be seen in Table [17](#).

Table 17: **Matches found by exact substring deduplication** (in *italics*).

|   |   |
|---|---|
| <p>it appears there is a transfer of ranking signals in this relationship. Supporting this finding is a quote from Google's guidelines: <i>Using JavaScript to redirect users can be a legitimate practice. For example, if you redirect users to an internal page once they're logged in, you can use JavaScript to do so. When examining JavaScript or other redirect methods to ensure your site adheres to our guidelines, consider the intent. Keep in mind that 301 redirects are best when moving your site, but you could use a JavaScript redirect for this purpose if you don't have access to your website's server.</i> NOTE: Their experiment is based on a live page with status code 200 and NOT an inactive page. So if you want to implement this for legacy</p> | <p>Some examples of sneaky redirects include: - Search engines shown one type of content while users are redirected to something significantly different. - Desktop users receive a normal page, while mobile users are redirected to a completely different spam domain. <i>Using JavaScript to redirect users can be a legitimate practice. For example, if you redirect users to an internal page once they're logged in, you can use JavaScript to do so. When examining JavaScript or other redirect methods to ensure your site adheres to our guidelines, consider the intent. Keep in mind that 301 redirects are best when moving your site, but you could use a JavaScript redirect for this purpose if you don't have access to your website's server.</i></p> |
| <p>Find Palm Beach FL homes for sale and other Palm Beach real estate on homesofthepalmbeaches.com. Browse and search Palm Beach houses, condos, townhomes and single-family homes by community , building, or location. <i>Our extensive database of real estate listings provide the most comprehensive property details including home values, features and local school and neighborhood info so you can be sure that you have nearly all the facts you need upfront. Search homesofthepalmbeaches.com today! Want a closer look at what other Palm Beach properties are available?</i></p>   | <p>Search Stuart houses, condos, townhomes and single-family homes by price and location. <i>Our extensive database of real estate listings provide the most comprehensive property details including home values, features and local school and neighborhood info so you can be sure that you have nearly all the facts you need upfront. Search Stuart Listings today! Want a closer look at what other Stuart properties are available? Also search our listings for the Newest Stuart Listings and Stuart Homes with Price Reductions now. Stuart FL Homes for Sale - Stuart Real Estate Listings FREE to search Stuart Property</i></p>  |
| <p><i>To find the correct size you should measure your foot from the heel to the toe point. Add approximately 1 - 1,5cm to get the actual inner sole length. Measure both feet and fit shoes to the larger foot. Measure feet at the end of the day, when your feet are at their largest.</i> Lente shoes are women's easy slip-on leisure shoes for everyday use. These lightweight shoes have a breathable textile mesh upper made of recycled PET bottles and cool Lycra lining.</p>   | <p><i>To find the correct size you should measure your foot from the heel to the toe point. Add approximately 1 - 1,5cm to get the actual inner sole length. Measure both feet and fit shoes to the larger foot. Measure feet at the end of the day, when your feet are at their largest.</i> Enjoy your summer days with Masera leisure sneakers. These low-cut women's sneakers are extremely lightweight thanks to phy-lon midsole and breathable textile mesh upper</p>   |
| <p>This bandana makes the perfect addition to every fur babies birthday collection! With its sparkly crown pattern, your pup will be ready for every birthday celebration! <i>With snaps for security, this bandana is made with love, down to the very last stitch ! Fabric: cotton Care Instructions: Hand wash only, iron as needed, on low heat Always supervise your pup while wearing Faithful Paws Co. accessories, as it could become a choking hazard if consumed.</i></p>   | <p>This bandana makes the perfect addition to every fur babies summer collection! With its vibrant watercolor popsicle pattern, your pup will be ready for every summer cookout! <i>With snaps for security, this bandana is made with love, down to the very last stitch ! Fabric: cotton Care Instructions: Hand wash only, iron as needed, on low heat Always supervise your pup while wearing Faithful Paws Co. accessories, as it could become a choking hazard if consumed.</i></p>   |