

SWEB: A LARGE WEB DATASET FOR THE SCANDINAVIAN LANGUAGES

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper presents the hitherto largest pretraining dataset for the Scandinavian languages: the Scandinavian WEb (SWEb), comprising over one trillion tokens. The paper details the collection and processing pipeline, and introduces a novel model-based text extractor that significantly reduces complexity in comparison with rule-based approaches. We also introduce a new cloze-style benchmark for evaluating language models in Swedish, and use this test to compare models trained on the SWEb data to models trained on FineWeb, with competitive results. All data, models and code are shared openly.

1 INTRODUCTION

Large language models have made significant strides in recent years due to their general capabilities in language-processing tasks. This progress has been largely driven by the development of extensive and high-quality pretraining datasets sourced from open web data (Wenzek et al., 2020; Brown et al., 2020; Abadji et al., 2022; Penedo et al., 2023; 2024). However, the majority of research aimed at improving pretraining data focuses on high-resource languages such as English. Our goal is to create a large-scale and high-performing open pretraining dataset specifically for the Scandinavian (*north-germanic*) languages: Swedish, Danish, Norwegian, and Icelandic.

Existing large-scale datasets for these languages primarily include mC4 (Xue et al., 2021), OSCAR (Abadji et al., 2022), and HPLT Datasets 1.2 (de Gibert et al., 2024). The Scandinavian portion of mC4 comprises approximately 100B tokens, 10B tokens for OSCAR 23.01, and 35B tokens for HPLT, which are all relatively small numbers considering that state-of-the-art large language models today are trained on trillions of high-quality tokens.

In this paper we make the following contributions:

- We release¹ the largest to date pretraining dataset for the Scandinavian languages: Scandinavian **WEb (SWEb)**. SWEb is the result of running our proposed pipeline on 98 Common Crawl snapshots. SWEb contains 1.01 trillion tokens in the Scandinavian languages, approximately an order of magnitude more than other available open alternatives.
- We introduce a new cloze-style benchmark for evaluating language models in Swedish, **HP-MEK**, a subset of the Swedish Scholastic Aptitude Test (Högskoleprovet) used for university admissions in Sweden. Using HP-MEK, we show our data performs on-par with data from the recently proposed FineWeb (Penedo et al., 2024) pipeline.
- We propose a new comprehensive pipeline for curating pretraining data for large language models, built around a model-based text extractor that significantly reduces complexity and is easily adaptable through rapid data annotation². Most notably, we demonstrate that our pipeline returns about +60% more high quality tokens than FineWeb on the same input data.

¹Data available here: <https://huggingface.co/datasets/...>

²Code and extractor model is available here: <https://github.com/...>

054 **2 BACKGROUND AND RELATED WORK**

055

056

057 Early efforts to extract massive amounts of text from the open internet for LLM training start from
 058 WebText (Radford et al., 2019), developed for training GPT-2. In this case, outbound links from
 059 Reddit with a certain number of upvotes were used as the content selection criterion. Text was
 060 extracted using Dragnet (Peters et al., 2018) and Newspaper³ and filtered with several heuristics,
 061 resulting in a dataset of 40GB after deduplication. Soon after, CCNet (Wenzek et al., 2020) and C4
 062 (Roberts et al., 2019) were proposed, both based on open web data from Common Crawl. C4 was
 063 initially developed exclusively for English but was later followed by a multilingual version, mC4
 064 (Xue et al., 2021). CCNet, on the other hand, was multilingual from the outset.

065 Both CCNet and C4 are based on the WET archives from Common Crawl, where all HTML format-
 066 ting has been stripped, leaving only the text. However, this text still contains a significant amount
 067 of noise in the form of menu and ad text, headers, footers, and sidebars, which are irrelevant to the
 068 page’s primary content. A successful method for extracting primary content from WET archives
 069 is to deduplicate the documents at the line level. C4 globally deduplicates all lines, while CCNet
 070 deduplicates over a subset of documents from the same Common Crawl dump. Line-by-line dedu-
 071 plication is the primary extraction method in CCNet, whereas C4 additionally employs a range of
 072 English-specific cleaning heuristics.

073 Following extraction comes a language detection and filtering step. Whilst more computationally
 074 expensive, performing language detection post extraction been shown to achieve better detection
 075 accuracy than filtering pre extraction (especially for low-resource languages) (Wenzek et al., 2020).
 076 Quality filtering differs slightly between the two, with C4 filtering using several heuristics, a bad
 077 words filter, and URL deduplication. In contrast, CCNet employs a model-based filter, using per-
 078 plexity as a quality measure with a KenLM model trained on Wikipedia.

079 CCNet has since been utilized in subsequent works such as RedPajama (v1 and v2) (Together Com-
 080 puter, 2023) and Dolma (Soldaini et al., 2024). RedPajama-Data v2 runs CCNet on an expanded
 081 number of Common Crawl snapshots and filters for five high-resource languages (none of which are
 082 Scandinavian, however). They also extend CCNet’s quality filtering by pre-computing a larger set
 083 of popular quality signals but leave the thresholding and filtering to the user.

084 Recently, several works have moved away from Common Crawl’s WET archives in favor of pro-
 085 cessing the raw HTML of webpages found in the WARC archives. Utilizing mor sophisticated text
 086 extraction turns out to be critical for the improving quality of the resulting data (Penedo et al., 2024).
 087 In MassiveWeb (Rae et al., 2021), the tree structure of HTML is utilized to more easily group and
 088 identify the primary content of pages. Some formatting is also retained, with the argument that
 089 this “diversity in formatting style translates effectively to the generative capabilities of the Gopher
 090 models.”

091 A similar approach is developed in NeuScraper (Xu et al., 2024), where a model is trained to – on
 092 an element level – decide whether it should be extracted or not. Both RefinedWeb and FineWeb use
 093 the open-source framework Trafilatura (Barbaresi, 2021) to extract text from HTML. Trafilatura is
 094 based on rules and heuristics on the DOM tree to identify primary content and has been shown to
 095 be the best non-commercial extractor for certain domains (Lopuhin, 2019). However, quality issues
 096 are still prevalent, and in RefinedWeb (Penedo et al., 2023) further (line-based) filters are added in
 097 an attempt to address these.

098 MassiveWeb introduce what they call “repetition filters” to remove documents with repetitive text,
 099 that is found beneficial with their extractor. These are also sucessfully reused in both RefinedWeb
 100 and later FineWeb. Through a systematic analysis, FineWeb further adds a small set of quality filters,
 101 that is shown through ablation experiments to yet increase quality. For a state of the art pipeline like
 102 FineWeb, the filtering can add up to 30 or more quantities and rules that might be difficult to oversee
 103 and adapt to new languages.

³<https://github.com/codelucas/newspaper>

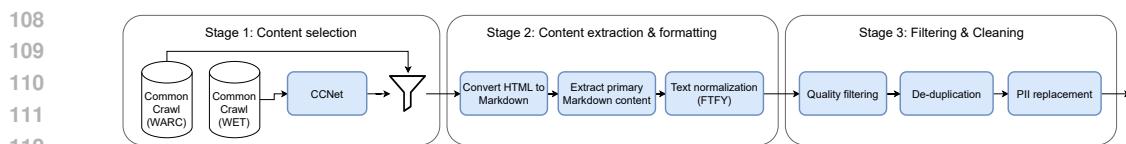


Figure 1: The SWEb pipeline. We use Common Crawl’s preprocessed WET archives for content selection, and WARC for extraction. At the center stage sits our model based Markdown extractor, that is the primary workhorse to produce our dataset.

3 THE SWEB PIPELINE

As evident by the previous section, much focus has been placed on the development of heuristics and filters to enhance the quality of the resulting data.

To move away from the extensive number of manual thresholds and complex extraction rules, we propose a more data-driven alternative. By learning a model for extraction, this complexity can be significantly reduced.

We begin by describing our pipeline that, like existing approaches, consists of the overarching steps of content selection, extraction, quality filtering, and deduplication (Figure 1).

3.1 STAGE 1: CONTENT SELECTION

Our pipeline begins with content selection, which aims to identify and select source documents from Common Crawl that are likely to be in one of the Scandinavian languages. Since the Scandinavian languages make up a very small portion of the entire Common Crawl, we want to implement this step early to filter out all non-relevant content.

We use CCNet to identify Scandinavian documents within the entire Common Crawl dataset. CCNet processes the WET archives, and after line-by-line deduplication, language detection is performed using fastText (Joulin et al., 2016b). Documents with a detected language score above 0.2 for any of the four languages are selected for the next stage.

3.2 STAGE 2: CONTENT EXTRACTION AND FORMATTING

In Stage 2, we start from the documents identified in Stage 1 but discard their content and instead use Common Crawl’s index to download their original HTML from the WARC archives. This means we use CCNet and the WET documents solely for content selection, but not for extraction. In the WET archives, all formatting and structure, such as header information, tables, text styles, bullet lists, and images, have been removed. We believe it is useful for language models to also model such structural information, in addition to plain text. Therefore, we aim to extract also this information from the webpages, and retain it in Markdown format.

We propose a new method for extracting primary content from the webpages, consisting of two steps: 1) Convert HTML to Markdown, 2) Extract primary content from the resulting Markdown through line-by-line filtering with a trained model.

3.2.1 CONVERT HTML TO MARKDOWN

Since we want to preserve basic textual formatting, we choose to convert from HTML to Markdown with its very lightweight markup, thus does not add many extra tokens. We convert all incoming HTML documents to Markdown using Pandoc, stripping links and images. See Listing 1 for an example.

No extraction has yet taken place, so these documents are still full of noise from menus, advertisements, and other extraneous content. We address this in the next step.

162

163 Listing 1: A webpage converted to markdown (translated, originally in Swedish), including title, top
 164 menu, headings and primary content. The document is truncated for brevity.

```

1 My Life, My Thoughts & My Training
2
3 ## The Blog
4 - The Blog
5 - Running Times Over the Years
6 - My Education
7 - Personal Training
8
9 ## Wednesday, December 14, 2011
10
11 ### The Tough Week Continues...
12
13 ...but tomorrow is a rest day.
14
15 I can feel in my body that I am right in the middle of a tough week *(I periodize my training, among other
16 things, by alternating between heavy, medium, and light weeks.)* and running was not exactly the first
17 thing I thought about when I woke up this morning. But after a nap together, sleep?\!
18
19 Posted by
20 Running & Life at
21 ...

```

178

179

180

181 3.2.2 MODEL-BASED CONTENT EXTRACTION

182 We observe that individual lines in the Markdown documents often correspond to specific elements
 183 such as headers, paragraphs, or navigation links. This makes lines an appropriate level for extraction.
 184 Therefore, we develop a custom annotation tool (details in Appendix B) to annotate which lines in
 185 these Markdown documents should be extracted and which should not. We ask annotators to mark
 186 what is considered the “main content” on the current webpage, and make some principled decisions
 187 for quality and consistency:

188

- 189 1. We do not extract navigation text such as menus or buttons.
- 190 2. A significant portion of the webpages are product pages. We decide to extract these only if
 191 there is a product description consisting of at least three complete sentences.
- 192 3. We extract tables if they are well-formatted and their content is tightly coupled to the main
 193 content.
- 194 4. On blogs or article pages that include user comments, we extract such comments in addition
 195 to the main content.
- 196 5. We do not extract information from sidebars unless it clearly constitutes main content.

197

198 While not explicitly excluded as per our guidelines, advertisement text isn’t considered to be main
 199 content and is thus implicitly excluded. The full annotation guidelines can be found in Appendix C.
 200 In total, we annotate 1,380 webpages, using 100 of these for validation and the remainder as training
 201 data for our extraction model.

202

203

Line Extraction Model Our dataset consists
 204 of Markdown documents with corresponding
 205 binary line annotations, see Figure 11. We
 206 aim to train a model to predict this label for
 207 each purpose, we choose to use
 208 a transformer encoder, where each newline is
 209 replaced with a special token [SEP]. We then
 210 feed the entire document through the encoder,
 211 with each [SEP] token representing the preceding
 212 line. This way, each line classification is
 213 contextualized by (theoretically) the full docu-
 214 ment context.

215

$$h_{0:n} = \text{Encoder}(x_{0:n}) \quad (1)$$

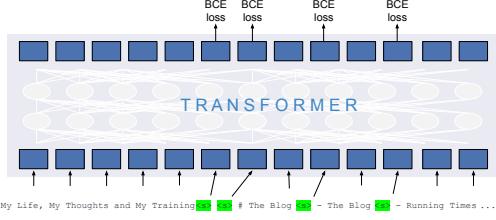


Figure 2: Illustration of our proposed line classification model. Each newline is replaced by a special $\langle s \rangle$ token, and the corresponding embeddings are used for classification

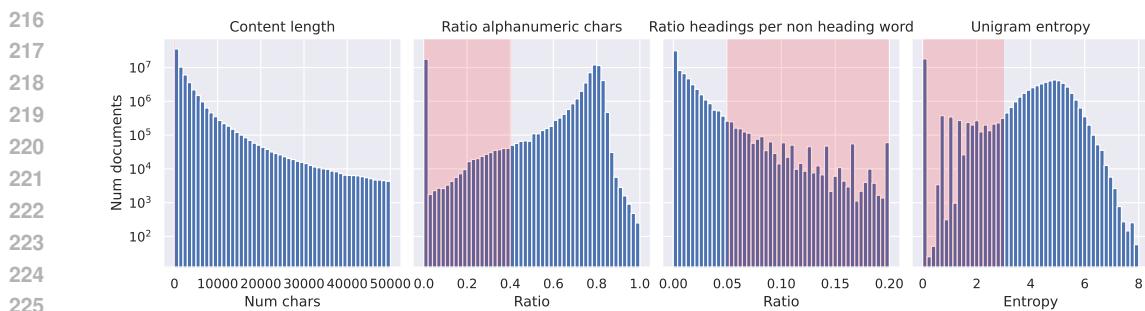


Figure 4: Filtering distributions on two Common Crawl dumps, and exclude regions marked in red. We exclude documents whose content length is shorter than 100 chars (invisible in the chart).

Through a linear projection of the output hidden state of each [SEP] token, we obtain logits for predicting the binary label of the current line. Let j denote the token index corresponding to each [SEP] token in the document. We then get the predicted probability for the line as:

$$p_j = \sigma(Wh_j + b) \quad (2)$$

where σ is the sigmoid function. The model is trained using binary cross-entropy loss between each p_j and an annotated line label. See Figure 2 for an illustration. We apply a fixed threshold to p_j to determine whether to include or exclude the line.

The Markdown documents can be very long, so we use the Longformer (Beltagy et al., 2020) architecture. Specifically, we use a pre-trained model that supports up to 16k tokens and has been trained for representation learning using masked language modeling⁴. The Longformer is a linear complexity transformer, thanks to its local self-attention, where each token only attends to a fixed number of nearby tokens. We use a local attention window size of 256 tokens and no global attention, as this turned out to only impair generalization.

We fine-tune the entire model on our training set of 1,280 documents, and the results on the validation set can be seen in Figure 3. We use the Adam optimizer with a constant learning rate of 1e-5. The results show that despite our small-scale training data, we achieve an F1 score of 87%.

Finally, we normalize the text using Fix Text For You (Speer, 2019).

3.3 STAGE 3: QUALITY FILTERING AND CLEANING

The third stage aims to filter for quality, reduce duplicate content and remove personally identifiable information (PII).

Quality Filtering A significant advantage of our model-based extraction is that it also implicitly performs much of the quality filtering. The extractor effectively learns to exclude content that is not of sufficient quality, such as spam and advertisements. This allows us to use only a minimal set

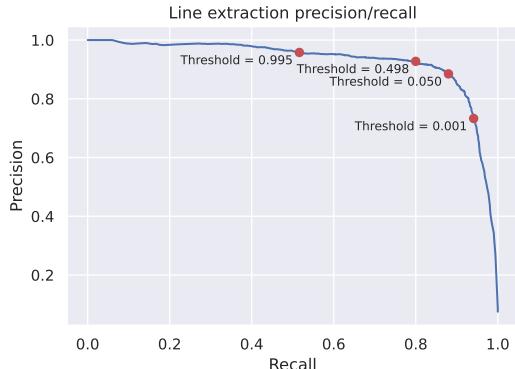


Figure 3: Precision/recall of our final line extraction model. We pick a threshold of 0.05 at inference, e.g. when applying the model for extraction.

⁴<https://huggingface.co/severinsimmler/xlm-roberta-longformer-base-16384>

270 of simple filters to remove edge cases where the extractor fails. Through qualitative analysis, we
 271 developed four filters to exclude such edge cases:
 272

- 273 1. **Content length:** We exclude cleaned documents that are shorter than 100 characters.
- 274 2. **Ratio of alphanumeric characters:** We exclude cleaned documents whose ratio of al-
 275 phanumeric characters is lower than 0.4. These documents primarily consist of data tables
 276 and are not relevant without additional context.
- 277 3. **Headings per non-heading word:** We note that in some documents, only headings are
 278 extracted with little or no accompanying text. We compute the ratio of the number of
 279 headings to the total number of words from non-heading lines. If the ratio is greater than
 280 0.05, we exclude the document.
- 281 4. **Unigram entropy:** Also used in Together Computer (2023), this measures the diversity
 282 of the content and is computed using $\sum -x/\text{total} * \log(x/\text{total})$ where the sum is taken
 283 over counts of unique words in the normalised content. By manual inspection, we found a
 284 threshold value of 3.0 to be reasonable, and exclude all documents below it.
 285

286 In Figure 4, we show the distributions of these four quantities, and in Appendix D, we provide
 287 examples of documents that are filtered out.
 288

289 **De-duplication** We used MinHashLSH (Leskovec et al., 2020) for document level near duplicate
 290 removal. The MinHash signatures were computed using unicode code point-level shingles of size
 291 16⁵, 14 bands, and 8 hashes per band (a total of 112 Hashes). Deduplication was done per band
 292 in an iterative fashion: For each band in order, we grouped documents by their hashes within that
 293 band, and kept only one document per group. Following FineWeb (Penedo et al., 2024), we only
 294 performed deduplication within snapshots, and not between them, as this was shown to increase
 295 downstream performance.
 296

297 **PII Replacement** As a final processing step, we make a best effort at removing personally identifi-
 298 able information from our data. To this end, we use regular expressions to replace email addresses
 299 and publicly facing IP-adresses with one of a few samples. This follows what has been done in
 300 previous works (Penedo et al., 2024; Soldaini et al., 2024).

301 4 EXPERIMENTS

304 How good is the data produced by our pipeline? To assess this question we conduct experiments
 305 against the recently proposed FineWeb pipeline (Penedo et al., 2024). We do this by performing a
 306 data ablation experiment. Here, we train two language models on data produced by 1) our pipeline
 307 and 2) the FineWeb pipeline respectively. We then evaluate the language models as a proxy for
 308 evaluating the datasets and, in turn, the pipelines.

309 FineWeb uses trafilatura (Barbaresi, 2021) as
 310 HTML extractor and relies on quality filter sets
 311 from both C4 and Gopher, as well as some
 312 novel additions. A notable difference is the fact
 313 that trafilatura (in the setting used by FineWeb)
 314 produces plain text content, while SWEb for-
 315 mat as Markdown. As mentioned in Section
 316 3.2, we primarily retain formatting via Mark-
 317 down *as a feature*, but note that this may also
 318 affect the learning behavior of the model. In
 319 this work however, we do not perform spe-
 320 cific ablations to single out this particular fac-
 321 tor. Please see Appendix E where we show
 322 side-by-side comparisons of trafilatura vs our
 323 extractor outputs.

How will Sweden be able to ____ itself in the international competition and strengthen its position as a leading knowledge nation? A first step is to look at the ____ that govern the allocation of state research funds.

- A activate – knowledge
- B mark – needs
- C assert – criteria**
- D entrust – institutions

Proper shoes are on the way out, while sneakers are spreading. The following ____ no longer causes any sensation: blazer, pleated trousers, and white sneakers.

- A manner
- B propriety
- C ensemble
- D attire**

Figure 5: Two examples from the HP-MEK task. Translated to English (originally in Swedish).

⁵We lowercased the text and removed non-alphabetic characters before creating shingles.

324
325
326
327
328

Exp. Dataset	#Docs	#Tokens	Tokens/doc
SWEb	32.3M	25.2B	779.7
FineWeb	19.2M	15.8B	820.3

329
330
331
332
333
334
335

Table 1: Stats of experimental datasets
SWEb and FineWeb

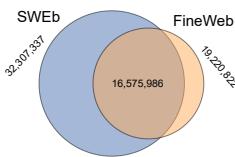


Figure 6: Venn diagram of documents in experimental SWEb and FineWeb datasets

4.1 BENCHMARK: HP-MEK

We investigated different benchmarks to evaluate the language models on. An appropriate benchmark should give good “early signals” of performance, in small model and data scales. For the Scandinavian benchmarks, the Scandeval suite (Nielsen, 2023) is commonly used. However, we found neither of its subtasks to be appropriate for this study, as the models didn’t reach good enough performance.

Instead, we chose to develop an alternative benchmark based on the Swedish Scholastic Aptitude Test (Högskoleprovet), that we denote **HP-MEK**⁶. We download and extract the MEK (sentence completion) section of all available historical tests, and end up with a total of 460 examples. HP-MEK is a cloze style test, with masked portions of a provided passage. For each passage, four alternatives of the masked portions are available, see Figure 5. We evaluate a model by inserting each of the four alternatives into the passage, and pick the alternative with the highest joint log likelihood.

In our experiments, we see early and consistently increased performance as we train on successively more data, which speaks for it being a suitable indicator for performance at larger scales.

4.2 EXPERIMENTAL SETUP

We extract, filter and deduplicate the 2024-10 and 2024-18 Common Crawl snapshots using our pipeline to form an experimental dataset (SWEb). We also run the FineWeb pipeline on the same input documents (selected from Stage 1) to form a competing dataset (FineWeb). Table 1 compares the two and Figure 6 shows a Venn diagram of their document (url) sets.

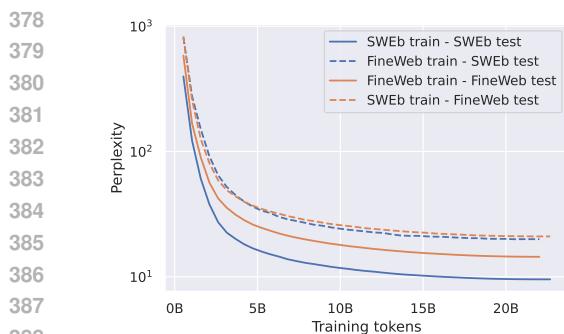
We note that the SWEb pipeline extracts significantly more documents (+62%) and tokens (+60%) than FineWeb’s pipeline. Most of FineWeb’s documents are contained in SWEb, while relatively few are uniquely selected by FineWeb. Interestingly, FineWeb extracts slightly more tokens per document on average, despite SWEb containing additional Markdown formating tokens.

We split the two datasets in 90/10 train/test splits and tokenize using the GPT-SW3 tokenizer (Ekgren et al., 2024). Then, we train small language models on each training set respectively (M_{SW} for SWEb and M_{FW} for FineWeb), and use the Llama architecture with 1.82B parameters (including embeddings) with a 2048 sequence length, a global batch size of 2 million tokens and a cosine decay learning rate schedule. Each model is trained for 10,811 steps, which corresponds to one full epoch for SWEb, and 1.6 epochs for FineWeb. We checkpoint every 250 steps to evaluate progression throughout training.

4.3 RESULTS

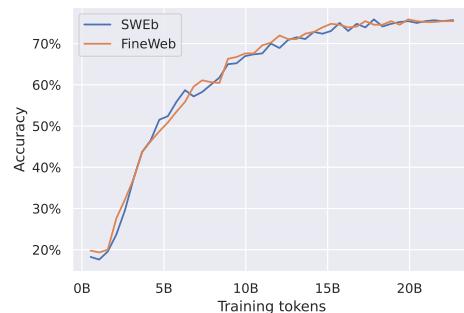
In Figure 7, we show perplexity plots where each model is evaluated on the each of the two test sets. We can first note that M_{SW} achieves lower perplexity on its own data than M_{FW} , i.e. SWEb seems “easier” to fit despite it being trained on more unique tokens. This could for example be due to the markdown formating, where markup tokens might be easier to predict. Secondly, M_{SW} performs relatively better on FineWeb than M_{FW} on SWEb. We speculate this could also be due to the markdown, where M_{FW} gets more confused not having seen Markdown during training.

⁶Available at <https://huggingface.co/datasets/>...



378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Figure 7: Perplexity cross-evaluation. The two models are evaluated on both SWEb and FineWeb test sets.



392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Figure 8: Learning curves. Performance of our two ablation models on HP-MEK throughout training.

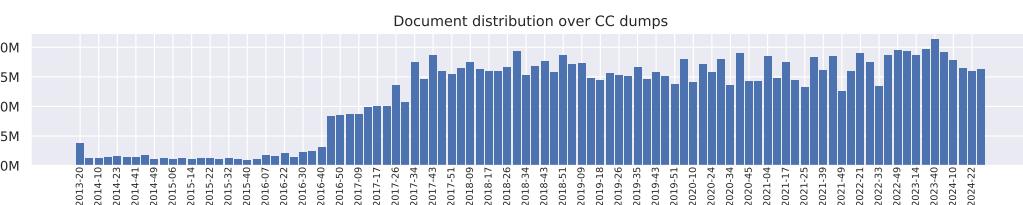


Figure 9: SWEb distribution over the Common Crawl snapshots.

Next, we evaluate M_{SW} and M_{FW} on HP-MEK, and plot learning curves in Figure 8. We can see that M_{SW} closely matches M_{FW} throughout the training, suggesting the two datasets are on-par with each other with regards to this task. This suggests that we are able to match the trafiletura extractor with just 1,380 annotated extraction samples, and at the same time reduce the complex filtering to only four simple quantities.

5 THE SWEB DATASET

We run our pipeline on 98 Common Crawl dumps, starting from 2013-20 until 2024-26, to produce the **Scandinavian Web (SWEb)** dataset. SWEb comprises a total of **1.01 trillion tokens⁷**, distributed over **1.2 billion documents**, resulting in **3.6TB** of raw (UTF-8) text.

This makes SWEb the largest open Scandinavian dataset to date, an order of magnitude larger than the (to our knowledge) previously largest mC4 dataset. In Figure 9, we show the document distribution across the Common Crawl dumps. As we can see, the amount of Scandinavian content has been steady since around 2017, averaging about 50M documents per dump.

To investigate the language distribution of SWEb, we use the fastText language identification classifier by Joulin et al. (2016a;b). Among the four Scandinavian languages, Swedish is the dominating one with 48% of documents classified as Swedish, 26% as Danish and 20% as Norwegian, see Figure 10. Only 2.3% are classified as Ice-

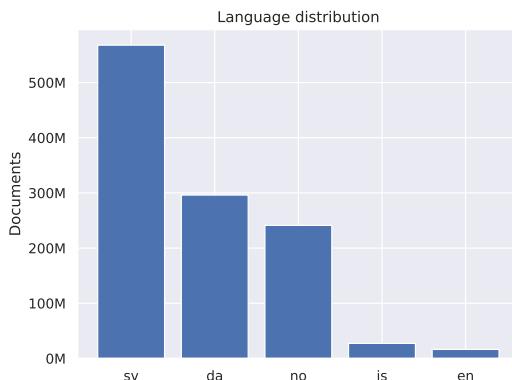


Figure 10: Language distribution over the SWEb dataset

⁷Using the GPT-SW3 (Ekgren et al., 2024) tokenizer

landic. A small portion of documents are classified as non-scandinavian after our content extraction, of which a majority is classified as English.

We release the SWeB dataset, the pipeline code, as well as our trained extractor model open source license, and hope this will further research and development of high performant Scandinavian LLMs. We also provide a datasheet detailing the dataset further in Appendix A.

6 DISCUSSION AND FUTURE WORK

Comparing to rule-based extractors such as *trafilatura*, our model based extractor offers greater flexibility as the desired extraction output is *demonstrated* instead of encoded as heuristics. Our work also highlights the data efficiency with which this can be done, i.e just 1,380 annotated examples in our case. However, this also comes with a cost. Running our model extractor for each document increases the compute required substantially over rule-based alternatives, which adds to these already compute-intensive workloads. In extracting SWeB, we consumed 20k AMD MI250X GPU-hours which is a significant amount, but comparing to the budgets required for training the downstream LLMs it is still negligible.

While training LLMs on larger datasets have shown to yield higher performance, a hypothesis is that there is only a subset of high quality documents that are behind the performance boosts. For example, in FineWeb-Edu, further filtering web data towards “educational content” is shown to significantly boosts performance in reasoning- and knowledge-intensive benchmarks. We see work on topic and content based filtering as a promising avenue for further refinement of SWeB towards particular LLM capabilities. This could potentially even be built into the extractor for more fine-grained control instead of as a binary post-hoc filter.

7 CONCLUSION

A major bottleneck for pre-training LLMs for smaller languages is the lack of large and high-quality open datasets. In this paper, we have presented the thus far largest open dataset for pre-training LLMs for the Scandinavian languages (Swedish, Danish, Norwegian and Icelandic). The dataset, which we call SWeB, comprises 1 trillion high-quality tokens in said four languages, and is openly shared in order to promote the development of LLMs for the Scandinavian languages. In creating SWeB, we have also developed a pipeline with a novel model-based text extractor that offers greater flexibility over the extraction process versus rule-based alternatives. We share both code and models for the novel text extractor openly. This paper has introduced a new benchmark for Swedish, which we use to compare models trained using our data with models trained using FineWeb, and we demonstrate that our data leads to models with performance on par with models trained from data using the state-of-the-art pipeline FineWeb.

ACKNOWLEDGMENTS

REFERENCES

- Julien Abadji, Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. Towards a Cleaner Document-Oriented Multilingual Crawled Corpus. *arXiv e-prints*, art. arXiv:2201.06642, January 2022.
- Adrien Barbaresi. Trafilatura: A web scraping library and command-line tool for text discovery and extraction. In Heng Ji, Jong C. Park, and Rui Xia (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pp. 122–131, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-demo.15. URL <https://aclanthology.org/2021.acl-demo.15>.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv:2004.05150*, 2020.

- 486 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
 487 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
 488 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
 489 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz
 490 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
 491 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In
 492 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-
 493 ral Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc.,
 494 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/
 495 file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf).
- 496 Ona de Gibert, Graeme Nail, Nikolay Arefyev, Marta Bañón, Jelmer van der Linde, Shaoxiong Ji,
 497 Jaume Zaragoza-Bernabeu, Mikko Aulamo, Gema Ramírez-Sánchez, Andrey Kutuzov, Sampo
 498 Pyysalo, Stephan Oepen, and Jörg Tiedemann. A new massive multilingual dataset for high-
 499 performance language technologies. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste,
 500 Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024 Joint
 501 International Conference on Computational Linguistics, Language Resources and Evaluation
 502 (LREC-COLING 2024)*, pp. 1116–1128, Torino, Italia, May 2024. ELRA and ICCL. URL
 503 <https://aclanthology.org/2024.lrec-main.100>.
- 504 Ariel Ekgren, Amaru Cuba Gyllenstein, Felix Stollenwerk, Joey Öhman, Tim Isbister, Evangelia
 505 Gogoulou, Fredrik Carlsson, Judit Casademet, and Magnus Sahlgren. GPT-SW3: An autore-
 506 gressive language model for the Scandinavian languages. In Nicoletta Calzolari, Min-Yen Kan,
 507 Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of
 508 the 2024 Joint International Conference on Computational Linguistics, Language Resources and
 509 Evaluation (LREC-COLING 2024)*, pp. 7886–7900, Torino, Italia, May 2024. ELRA and ICCL.
 510 URL <https://aclanthology.org/2024.lrec-main.695>.
- 511 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,
 512 Hal Daumé III au2, and Kate Crawford. Datasheets for datasets, 2021. URL <https://arxiv.org/abs/1803.09010>.
- 513 Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hervé Jégou, and Tomas
 514 Mikolov. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*,
 515 2016a.
- 516 Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient
 517 text classification. *arXiv preprint arXiv:1607.01759*, 2016b.
- 518 Jure Leskovec, Anand Rajaraman, and Jeffrey David Ullman. *Mining of massive data sets*. Cam-
 519 bridge university press, 2020.
- 520 Konstantin Lopuhin. Evaluating quality of article body extraction for commercial ser-
 521 vices and open-source libraries, 2019. [https://github.com/scrapinghub/
 522 article-extraction-benchmark](https://github.com/scrapinghub/article-extraction-benchmark).
- 523 Dan Nielsen. ScandEval: A benchmark for Scandinavian natural language processing. In Tanel
 524 Alumäe and Mark Fishel (eds.), *Proceedings of the 24th Nordic Conference on Computational
 525 Linguistics (NoDaLiDa)*, pp. 185–201, Tórshavn, Faroe Islands, May 2023. University of Tartu
 526 Library. URL <https://aclanthology.org/2023.nodalida-1.20>.
- 527 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,
 528 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The RefinedWeb
 529 dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv
 530 preprint arXiv:2306.01116*, 2023. URL <https://arxiv.org/abs/2306.01116>.
- 531 Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro
 532 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data
 533 at scale. *arXiv preprint arXiv:2406.17557*, 2024.

- 540 Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee,
 541 and Luke Zettlemoyer. Deep contextualized word representations. In Marilyn Walker, Heng Ji,
 542 and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of
 543 the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long
 544 Papers)*, pp. 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Lin-
 545 guistics. doi: 10.18653/v1/N18-1202. URL <https://aclanthology.org/N18-1202>.
- 546 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
 547 models are unsupervised multitask learners. 2019.
- 548 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, H. Francis
 549 Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford,
 550 Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche,
 551 Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth
 552 Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat
 553 McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden,
 554 Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lor-
 555 raine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Ange-
 556 liki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev,
 557 Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cy-
 558 prien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan
 559 Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew J. Johnson,
 560 Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Edward Lockhart, Simon
 561 Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne
 562 Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language mod-
 563 els: Methods, analysis & insights from training gopher. *CoRR*, abs/2112.11446, 2021. URL
 564 <https://arxiv.org/abs/2112.11446>.
- 565 Adam Roberts, Colin Raffel, Katherine Lee, Michael Matena, Noam Shazeer, Peter J Liu, Sharan
 566 Narang, Wei Li, and Yanqi Zhou. Exploring the limits of transfer learning with a unified text-to-
 567 text transformer. *Google, Tech. Rep.*, 2019.
- 568 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Author,
 569 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh
 570 Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas
 571 Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle
 572 Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke
 573 Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge,
 574 and Kyle Lo. Dolma: An Open Corpus of Three Trillion Tokens for Language Model Pretraining
 575 Research. *arXiv preprint*, 2024. URL <https://arxiv.org/abs/2402.00159>.
- 576 Robyn Speer. ftfy. Zenodo, 2019. URL <https://doi.org/10.5281/zenodo.2591652>.
 577 Version 5.5.
- 578 A Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, April
 579 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- 580 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán,
 581 Armand Joulin, and Édouard Grave. Ccnet: Extracting high quality monolingual datasets from
 582 web crawl data. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pp.
 583 4003–4012, 2020.
- 584 Zhipeng Xu, Zhenghao Liu, Yukun Yan, Zhiyuan Liu, Chenyan Xiong, and Ge Yu. Cleaner pre-
 585 training corpus curation with neural web scraping. In *Proceedings of the 62nd Annual Meeting of
 586 the Association for Computational Linguistics*, 2024.
- 587 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
 588 Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer.
 589 In *Proceedings of the 2021 Conference of the North American Chapter of the Association for*

594 *Computational Linguistics: Human Language Technologies*, pp. 483–498, Online, June 2021.
595 Association for Computational Linguistics. doi: 10.18653/v1/2021.nacl-main.41. URL <https://aclanthology.org/2021.nacl-main.41>.
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 **A SWEB DATASHEET**
 649

650 We provide a datasheet inspired by Gebru et al. (2021):
 651

Motivation	
Purpose of the dataset	We want to encourage open research and development of LLMs in the Swedish, Danish, Norwegian and Icelandic languages. We build and release SWEb to promote this objective and to address the linguistic challenges specific to underrepresented Scandinavian languages, improving access to language technology in these regions.
Curated by	XX
Funded by	XX
Composition	
Data Fields	<p>Each data instance contains:</p> <ol style="list-style-type: none"> 1. The source URL 2. The original Common Crawl WARC file path 3. The WARC date 4. The extracted text content, in markdown format 5. The detected language (using fastText classifier)
Data Splits	We split SWEb based on Common Crawl dump, to allow for download based on time of crawl. We also include a default split containing the entire dataset.
Errors and noise	As outlined in this paper, we propose a novel model based approach to extract text from websites. However, the model is not perfect and non-relevant content as well as noise are sometimes also erroneously extracted. We try to filter such examples in our third pipeline stage, but despite our best effort such examples may sometimes slip through.
Offensive and toxic content	As we don't attempt to filter based on content or topic in this work, SWEb might contain content that can be perceived as offensive, threatening or otherwise toxic. When considering using this dataset, it is important to be aware of this and that further processing might be necessary depending on use case.
Dataset Curation	
Curation rationale	We use Common Crawl as it is the (to our knowledge) largest and most diverse open corpus available in the Scandinavian languages.
Source data	The Common Crawl source data consist of large amounts of webpages crawled from the open web. Common Crawl's crawlers has always respected <code>nofollow</code> and <code>robots.txt</code> policies.
Time frames of collected data	We use all Common Crawl scraped dating back to week 50 of 2013 and up to week 26 of 2024.
Data processing steps	See Section 3.

700
 701

702		
703		
704		
705		
706		
707		
708	Considerations for using the data	
709	Social impact of dataset	With SWEb, our goal is to make LLM training more accessible to the machine learning community by: (a) making the dataset creation process more transparent, by sharing our entire processing setup including the codebase used (b) helping alleviate the costs of dataset curation, both in time and in compute, for model creators by publicly releasing our dataset with the community.
710		
711		
712		
713		
714		
715		
716	Bias and Representation	While the Common Crawl data gathers diverse text sources, biases present in the original content may still exist. Users should critically assess how these biases may affect model training and outcomes, especially in sensitive applications. It is recommended to implement bias-mitigation techniques during training and model development.
717		
718		
719		
720		
721		
722	Model Misuse	When training models with this dataset, it is crucial to prevent harmful uses of the resulting models, such as generating misleading or dangerous content (e.g., disinformation, hate speech). Always consider the societal impact of deploying models trained on this data, and take precautions to implement appropriate safeguards.
723		
724		
725		
726		
727		
728	Distribution	
729		
730	Distribution platform	The dataset will be distributed on the Huggingface Hub
731	License	<p>The data is released under the CC0 Creative Commons License. We make the following clarifications:</p> <ol style="list-style-type: none"> 1. We do not warrant or guarantee any rights to the underlying data contained within this dataset. Users are solely responsible for validating and securing the appropriate rights and licenses for their specific intended uses. 2. This license applies only to the structure and compilation of the dataset as provided by us. We do not claim any database rights or ownership over the underlying data itself. Users must ensure compliance with any legal obligations, including those related to third-party content, copyrighted material, or personal information (PII) that may be contained in the underlying data. 3. With the release of this dataset, our goal is to promote and advance open research and the development of Scandinavian language models, showcase research outcomes as well as enable research validation. Open datasets are essential to fostering innovation and expanding knowledge in AI. We disclaim any responsibility for other uses, including commercial applications. Users are responsible for ensuring the legality of their usage, especially in cases involving copyrighted material.
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

Notice and take-down policy	<p>Should you consider that our data contains material that is owned by you and should therefore not be reproduced here, please:</p> <ol style="list-style-type: none">1. Clearly identify yourself, with detailed contact data such as an address, telephone number or email address at which you can be contacted.2. Clearly identify the copyrighted work claimed to be infringed.3. Clearly identify the material that is claimed to be infringing and information reasonably sufficient to allow us to locate the material.4. You can reach us at XX <p>We will comply to legitimate requests by removing the affected sources from the next release of the corpus.</p>
-----------------------------	---

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810 B MARKDOWN ANNOTATION DETAILS

813 Text Annotation Tool

814 Previous 13 / 500 Next Ignored Annotations

815 http://psykopedia.org/wiki/Telefon

816

817 # Telefon

818 Från Psykopedia

819 Hoppa till navigering Hoppa till sök

820 ***Telefonen** uppfanns av Jan Telefon.

821 Telefon är en uppfinning som uppfanns i förrgår, tekniken går ut på att en sändare som sitter i marken läser av teckenspråk som mobilen på ett telepatiskt sätt framför via "ETERN" (Se Fear and Loathing in Las Vegas om du vill veta vad Eter är).

822 En gammal version av telefonen är telegrafen som man använde med hjälp av klossar som man ristade in bokstäver baklänges på. August Strindberg var den första som använde en telefon som toalett.

823 Hämtad från "<https://psykopedia.org/index.php?title=Telefon&oldid=48402>"

824 Kategorier:

- Teknik

825 ## Navigeringsmeny

826

827

828

829 Figure 11: Our web based annotation tool. On the right side the original web page is displayed

830 and on the left the corresponding markdown. Annotation is performed by selecting individual lines

831 (marked green) that constitute the main content of the page.

832

833 We develop a web based tool that we use to annotate markdown documents, see Figure 11. The tool

834 is used to annotate data for training and evaluating our text extractor model (Section 3.2.2).

835 The annotation was performed by the authors as well as additional lab colleagues, in total a group

836 of 12 people. We started by jointly annotating a gold standard test set of 100 examples (web pages).

837 This was useful to align and develop our annotation guidelines.

838 Next, we annotated a first set of 400 training examples and trained a first extractor model. This

839 model served as a first baseline. We then iteratively annotated additional training data in batches of

840 300-500 examples, re-trained and re-evaluated after each iteration.

841 Judging what is “main content” in web pages is not always obvious however. When the evaluation

842 didn’t improve after a new batch of annotations, we developed a method for discovering “confusing”

843 training examples in the new batch that we could jointly discuss and align on. For each example x

844 in the new training batch, we compute the loss $l^{M_n}(x, y) = \mathcal{L}(M_n(x), y)$, where \mathcal{L} is the average

845 over all BCE losses in the example and M_n is the model trained on all batches including iteration n .

846

847 By comparing this loss to the corresponding loss under the *previous* model M_{n-1} , we get a measure

848 of how “surprising” this example is:

849

$$\delta = l^{M_{n-1}}(x, y) - l^{M_n}(x, y) \quad (3)$$

850 Using this quantity, we could easily identify outliers and correct inconsistent annotations. By per-

851 forming this post-annotation fix-up, we were able to improve performance on our test set, for each

852 annotated batch of data.

853

864 C CONTENT EXTRACTION ANNOTATION GUIDELINES 865

866 *The following description was provided to our annotators*
867

868 In the provided annotation tool, please select individual lines by clicking and dragging across the
869 lines you want to select.
870

- 871 • Please look at the rendered web page on the right. We want to extract the “**main textual**
872 **content**” of the current page.
- 873 • **No navigation** (menus, links, button texts etc) should be selected, except well formatted
874 tables of content that link within the same page
- 875 • Include **headers** of main content
- 876 • If duplicate header, select the one closest to the main content
- 877 • Include **well formatted tables**
- 878 • Don’t include content of **sidebars** that is unrelated to the main content
- 879 • It is OK to leave the whole document unselected if there is no relevant content
- 880 • If there are many very **similar-looking pages**, they can be marked *Ignored* if they have al-
881 ready been annotated. Bad pages without any good content should **not** be ignored however.
- 882 • Include **comment sections** if there are any, but exclude navigation associated with those,
883 e.g. *Svara / Rapportera inlägg* or similar.
- 884 • Keep **comment headings**
- 885 • If text is broken with e.g. “...”, don’t include
- 886 • Select **top heading** if it exists
- 887 • Keep **at most 2 consecutive newlines**
- 888 • Remove empty formatting lines (e.g **), except for dividers (-----)
- 889 • Pages that are primarily “data” (e.g. tables, numbers) without much text should be unse-
890 lected. There should be at least **three consecutive sentences** of text. This puts a somewhat
891 high bar for product pages
- 892 • No HTML should be selected

893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918 **D FILTERED EXAMPLES**
919

920 We show examples of extracted documents that are *filtered out* by each of our four quality filters.
921

922 **D.1 CONTENT LENGTH < 100 CHARS**
923

924 <https://www.buskerudmynt.no/produkt/norske-mynter-etter-1874/norske-argangsmynter/50-ore/olav-v-1974-1991/50-ore-1977-kv.-0>

925
1 # 50 øre 1977 kv. 0
2
3 Tatt fra rull. Litt skjoldete mynt.
4
5 NOK5,00 inkl. mva.

930 <https://www.ovedanielsson.se/2021/08/30/ohrmans-fick-inte-bygga-nytt-mot-torget/embed/>

931
1 Öhrmans fick inte bygga nytt mot torget

934 <https://jesper.nu/spel/achtung-die-kurve>

935
1 # Achtung Die Kurve

937 **D.2 RATIO OF ALPHANUMERIC CHARACTERS < 0.4**

938 <https://www.innebandystats.se/statistik/219645/kevin-sandeback>

939
1 | | CL98IC | Juniorallsvenskan HJ18 | 14 | 16

940
941
942
943
944 https://nn.wikipedia.org/wiki/Kategori:Deltakarar_under_vinter-OL_1984_etter_Åvvинг

945
1 1896 ** Å... 1900 ** Å... 1904 ** Å... 1906 ** Å... 1908 ** Å... 1912 ** Å... ~ (1916) ~ ** Å... 1920 ** Å...
946 ** 1924 ** Å... 1928 ** Å... 1932 ** Å... 1936 ** Å... ~ (1940) ~ ** Å... ~ (1944) ~ ** Å... 1948 **
947 Å... 1952 ** Å... 1956 ** Å... 1960 ** Å... 1964 ** Å... 1968 ** Å... 1972 ** Å... 1976 ** Å... 1980
** Å... 1984 ** Å... 1988 ** Å... 1992 ** Å... 1996 ** Å... 2000 ** Å... 2004 ** Å... 2008 ** Å...
2012 ** Å... 2016** Å... 2020
2 **Vinter-OL**
3
4 Deltakarar etter **nasjon:***
5
6 1924 ** Å... 1928 ** Å... 1932 ** Å... 1936 ** Å... ~ (1940) ~ ** Å... ~ (1944) ~ ** Å... 1948 ** Å... 1952
** Å... 1956 ** Å... 1960 ** Å... 1964 ** Å... 1968 ** Å... 1972 ** Å... 1976 ** Å... 1980 ** Å...
1984 ** Å... 1988 ** Å... 1992 ** Å... 1994 ** Å... 1998 ** Å... 2002 ** Å... 2006 ** Å... 2010 ** Å...
2014 ** Å... 2018 ** Å... 2022
7
8 Deltakarar etter **øving:***

946 <https://historik.val.se/val/val2010/alkon/K/valdistrikt/12/80/0102/alderkon.html>

947
1 | -----: | -----: | -----: | -----: | -----: | -----: | -----: | -----: | --: |
948 -----: | -----: | -----: | ---: | -----: | -----: | -----: | -----
2 | Gamla staden, Stortorget | 1533 | 24,8% | 380 | 43,0% | 659 | 20,4% | 312 | 11,9% | 182 |
949 4,4% | 68 | 52,6% | 806 | 47,4% | 727 | 13,8%
212 | | | | | |
3 | Summa | 1533 | 24,8% | 380 | 43,0% | 659 | 20,4% | 312 | 11,9% | 182 |
950 4,4% | 68 | 52,6% | 806 | 47,4% | 727 | 13,8%
212 | | | | |
4
5 http://www.val.se

944 **D.3 HEADINGS PER NON-HEADING WORD > 0.05**

945 <https://www.sahlgrensk.se/for-dig-som-ar/vardgivare/laboratoriemedicin/analyslistan/specialprover-cytologisk-diagnostik/16648.html/>

946
1 # Glaskropp
2
3 # Glaskropp cytologisk diagnos
4
5 ### Synonymer
6

972
973 7 Specialprover, cytologisk diagnostik
974 8
974 9 ## Provtagningsanvisning
975 10
975 11 ### Provmaterial
976 12
976 13 ### Rör el. motsv
977 14
977 15 10 ml rör med gul kork, konisk botten, steril (för mindre mängder material) eller Burk ca 40 ml m tä
978 tslutande lock, sterilt
978 16
979 17 ### Provtagning
980 18
980 19 Enligt inremitterande kliniks regler Provet skall snarast efter sköljningen transporteras till Cytologen.
981 20 Ofixerade vätskor ska föranmälas och lämnas direkt till lab. personal före kl 14.00.
981 21
982 22 ### Transport
982 23
983 24 Transport ska ske omgående till Laboratoriet för klinisk patologi och där lämnas direkt till provinlä
984 mningen.

986 <https://folk.su.se/björnlund/publications/211851>

```
987 1 # Olive Buchvold Juvik
988 2
989 3 Gullet vårt Olive «snat 2 år» «Ja, vennen, på lørdag 18. nov, fyller du 2 år» Me gratulerer så masse!\n        Kjempe gla' i deg. Klem fra tanter, onkler, besteforeldre og oldeforeldre.
990 4
991 5
992 6
993 7 ## Go'ungen (0-12 år)
994 8
995 9
996 10
997 11 ### Nora Silden Fredheim
```

[https://start.arcada.fi/sv/kurser/303000/2021-2022/IA-2-004/0](https://start.arcada.fi/sv/kurser/303000/2021-2022/IA-2-004/)

```
997 1 ## Kursens undervisningsperiod
998 2
999 3 3 (2022-01-01 till 2022-03-13)
999 4
1000 5 ## Nivå/kategori
1000 6 ## Cykel/nivå
1000 7 Yrkeshögskoleexamen
1001 8
1002 9 ## Rekommenderat studieår
1002 10
1003 11 1
1003 12 ## Omfattning
1004 13
1004 14 5 sp
1005 15
1006 16 ## Kompetensmål
1006 17
1007 18 I denna studieenhet står följande kompetenser i
1007 19 fokus:
1008 20 \- Kompetens inom serverprogrammering
1009 21 \- Kompetens inom databashantering och lagring av
1009 22 data
1010 23 \- Kompetensen att skapa dynamiska applikationer
1010 24
1011 25 ## Läranderesultat
1011 26
1012 27 Efter avlagd studieenhet:
1012 28 \- Du behärskar programmering med
1013 29 PHP (Kunskap)
1014 30 \- Du ser skillnaden mellan statiska, interaktiva
1014 31 och dynamiska webbsidor (Kunskap)
1015 32 \- Du kan hantera filer från klienten och på
1015 33 servern (Kunskap)
1016 34 \- Du kan bygga dynamiska webbappar (Färdighet)
1016 35 \- Du kan lagra data säkert i en databas
1016 36 (Färdighet)
1017 37 \- Du insrer problematik och lösningar kring att
1017 38 lagra känslig information om en användare
1017 39 (Förhållningssätt)
1018 40 \- Du uppfattar olika sätt att överföra och lagra
1018 41 data och dess koppling till säkerhet och
1018 42 prestanda (Förhållningssätt)
1019 43 \- Du uppfattar din makt och ditt ansvar som
```

D.4 UNIGRAM ENTROPY < 3.0

<https://hastkatalogen.se/content/horse/info.php?id=31999>

1026
1027 1 # Catinkaox
1028 2
1029 3 ## Arabiskt Fullblod
1030 4
1031 5 Catinka är ett sto som föddes 1976 i Sverige.
1032 6
1033 7
1034 8
1035 9 - Ras: Arabiskt Fullblod
1036 10 - KÖN: Sto
1037 11 - Färg: brun
1038 12 - Stofamilj:
1039 13 |

1040
1041
1042
1043
1044
1045
1046
1047
1048
1049 <https://www.nilssonsilammhult.se/hallmobler/ida-skohorn-ek/>
1050 1 # Ida skohorn ek
1051 2
1052 3 430 kr
1053 4
1054 5 Ida skohorn i oljad ek från småländska Westroth. Tillverkad i formpressat trä. En fin detalj till hallen!\!
1055 6
1056 7
1057 8
1058 9 ## Veisludagur runninn upp
1059 10 -
1060 11 -
1061 12
1062 13 ## Dvalarflokkur
1063 14
1064 15 Höfundur: Heiðbjört Arney|2017-07-12T10:01:09+00:0012. júlí 2017|
1065 16
1066 17 -
1067 18
1068 19 ## Leikjanámskeið 2
1069 20
1070 21 Höfundur: Heiðbjört Arney|2017-06-28T10:15:51+00:0028. júní 2017|
1071
1072
1073
1074
1075
1076
1077
1078
1079

1080 E COMPARING OUR MODEL EXTRACTOR VS TRAFILATURA 1081

1082 We compare our model based extractor vs traflatura (in the settings used by FineWeb).
1083

1084 <https://www.ark.no/produkt/boker/dokumentar-og-faktabøker/eksil-9788202253912>
1085

1086 Trafilatura

1088 1 Innbundet
2 2005
1089 3 Norsk, Bokmål
4 «Denne boken dreier seg om eksil og dannelsen.
1090 5 Lesning av Dante ga meg en italiensk regel:
Dannelsen oppstår alltid og bare i eksil.
1091 Det vesle som fins av dannelsen i Norge,
dannes ut fra evnen til distanse i et
livsnødvendig indre eller ytre eksil.
1092 Dannelsen er det motsatte av turisme. Slik
førte min selvomsorg meg stadig mer inn
og ut av norsk kultur (og underholdning)
1093 til jeg ble uhelbredelig gudløs og partiløs
øs i en vag, men livslang interesse for
1094 eksemplariske flyktninger og forrædere
fra Klosterlassen til Asbjørn Sunde..»
1095 6 (fra Georg Johannsenes forord)
1096 7 Klikk&Hent
1097 8 På lager i 8 butikker
1098 9 Nettlager Sendes normalt innen 1-2 virkedager
1099 10 Bytt i alle våre butikker
1100 11 –
1101 12 Klikk og hent
13 –

1102 Trafilatura

1103 1 Hur dom än
1104 2 Färgerna är blekare än igår
1105 3 tiden är för mörk för att vi ska kunna le
1106 4 jag vill inte höra deras röst mer
1107 5 Illusioner av tröst som drar mig ner
1108 6 Hur de än sargar oss så ska vi hålla hand
1109 7 Halva jävla världen är i brand
1110 8 O hur dom än sänker oss så ska vi skrika högst
1111 9 ett nej är alltid ett nej
1112 10 Vart vi än går ser vi ner
1113 11 aldrig mer igen, aldrig mer
1114 12 hela tiden får vi säga till
1115 13 ljusen runtomkring står bara still
1116 14 Hur de än sargar oss så ska vi hålla hand
1117 15 Halva jävla världen är i brand
1118 16 O hur dom än sänker oss så ska vi skrika högst
1119 17 ett nej är alltid ett nej
1120 18 En vacker stråle som försann
1121 19 innan det blev mörkt
1122 20 innan det blev kallt
1123 21 Och om det var dina skrik som inte hördes
22 eller var din dotter som fördes iväg
1124 23 hur skulle det kännas, hur skulle dääll
1125 24 Hur de än sargar oss så ska vi hålla hand
1126 25 Halva jävla världen är i brand
1127 26 O hur dom än sänker oss så ska vi skrika högst
1128 27 ett nej är alltid ett nej

Model extractor (ours)

1 1 # Eksil - om klosterlassen og andre eksempler
2 2
3 3 Av Georg Johannesen
4 4
5 5 «Denne boken dreier seg om eksil og dannelsen.
Lesning av Dante ga meg en italiensk
regel: Dannelsen oppstår alltid og bare i
eksil. Det vesle som fins av dannelsen i
Norge, dannes ut fra evnen til distansen i
et livsnødvendig indre eller ytre eksil.
Dannelsen er det motsatte av turisme.
Slik førte min selvomsorg meg stadig mer
inn og ut av norsk kultur (og
underholdning) til jeg ble uhelbredelig
gudløs og partiløs i en vag, men livslang
interesse for eksemplariske flyktninger
og forrædere fra Klosterlassen til Asbjørn
Sunde..» (fra Georg Johannsenes forord)

Model extractor (ours)

1 1 # Vad dom än tror
2 2
3 3 Text: Clara Rudelius
4 4
5 5 https://gipomusic.se/wp-content/uploads/2014/10/04_Vad-dom-än-tror.mp3
6 6
7 7 **Hur dom än**
8 8
9 9 Färgerna är blekare än igår
10 10 tiden är för mörk för att vi ska kunna le
11 11 jag vill inte höra deras röst mer
12 12 Illusioner av tröst som drar mig ner
13 13
14 14 Hur de än sargar oss så ska vi hålla hand
15 15 Halva jävla världen är i brand
16 16 O hur dom än sänker oss så ska vi skrika högst
17 17 ett nej är alltid ett nej
18 18
19 19 Vart vi än går ser vi ner
20 20 aldrig mer igen, aldrig mer
21 21 hela tiden får vi säga till
22 22 ljusen runtomkring står bara still
23 23
24 24 Hur de än sargar oss så ska vi hålla hand
25 25 Halva jävla världen är i brand
26 26 O hur dom än sänker oss så ska vi skrika högst
27 27 ett nej är alltid ett nej
28 28
29 29 En vacker stråle som försann
30 30 innan det blev mörkt
31 31 innan det blev kallt
32 32
33 33 Och om det var dina skrik som inte hördes
34 34 eller var din dotter som fördes iväg
35 35 hur skulle det kännas, hur skulle dääll
36 36
37 37 Hur de än sargar oss så ska vi hålla hand
38 38 Halva jävla världen är i brand
39 39 O hur dom än sänker oss så ska vi skrika högst
40 40 ett nej är alltid ett nej
41 41
42 42 ## Albumspår

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

Trafilatura

1149

1 Turinformasjon

2 Tur fra Vågstrand til Norvika i Eidsbygda
Rauma i kjølvannet av bilfergen Vaage-

Norvik som gikk der fra 1930 til 1945.

3 Vei til Åndalsnes ble til stor del bygget

ferdig av okkupasjonsmakten under andre

verdenskrig og veien åpnet rundt

tidspunktet for freden i 1945.

4 Denne fergen ble bygget av samme båtbygger som

båten vi går turen med og det blir

fortalt historie rundt dette samt

hendelsene rundt den tragiske ulykken i

oktober 1942 hvor Kultur og

Propagandaminister i Quisling regjeringen

Gulbrand Lunde m/frae omkom ved

fergekaien på Vaage.

5 Turprisen er oppgitt pr passasjer basert på

max antall. Ta kontakt for alternativer

og evt allergier.

6 Eventuelt servering ombord!

7 1. Rik tomat/chili basert kremet fiskesuppe

servert m/nybakt brød, dessert (Tilslørte

bondepiker) og kokekaffe. Kr. 350.-

8 Lunsjpakke fra Braud Håndverksbakeri Vestnes:

9 2. Påsmurt bagett med ost & skinke +

kanelbolle alt. solskinnsbolle. Kr. 110.-

10 3. Påsmurt bagett med kylling & karri +

kanelbolle alt. solskinnsbolle. Kr. 120.-

11 4. Pastasalat med kylling og karri. Kr. 175.-

12 Mineralvann og annen drikke fås kjøpt separat

om bord.

13 5 Timer

14 -

15 Maks. Passasjerer: 12

16 -

17 Vestnes

18 -

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

Model extractor (ours)

1 # I kjølvannet av Bilfergen Vaage-Norvik

2

3 ### 1 100 NOK pr passasjer

4

5 ## Turinformasjon

6

7 Tur fra Vågstrand til Norvika i Eidsbygda

Rauma i kjølvannet av bilfergen Vaage-

Norvik som gikk der fra 1930 til 1945.

8

9 Vei til Åndalsnes ble til stor del bygget

ferdig av okkupasjonsmakten under andre

verdenskrig og veien åpnet rundt

tidspunktet for freden i 1945. Denne

fergen ble bygget av samme båtbygger som

båten vi går turen med og det blir

fortalt historie rundt dette samt

hendelsene rundt den tragiske ulykken i

oktober 1942 hvor Kultur og

Propagandaminister i Quisling regjeringen

Gulbrand Lunde m/frae omkom ved

fergekaien på Vaage.

10

11 Turprisen er oppgitt pr passasjer basert på

max antall. Ta kontakt for alternativer

og evt allergier.

12

13 **Eventuelt servering ombord\!**

14

15 1\. Rik tomat/chili basert kremet fiskesuppe

servert m/nybakt brød, dessert (Tilslørte

bondepicker) og kokekaffe. Kr. 350.-

16

17 Lunsjpakke fra Braud Håndverksbakeri Vestnes:

18 2\. Påsmurt bagett med ost & skinke +

kanelbolle alt. solskinnsbolle. Kr. 110.-

19 3\. Påsmurt bagett med kylling & karri +

kanelbolle alt. solskinnsbolle. Kr. 120.-

20 4\. Pastasalat med kylling og karri. Kr. 175.-

21

22 Mineralvann og annen drikke fås kjøpt separat

om bord.

23

24 - **5 Timer

25 - **Maks. Passasjerer: 12

26 - Avgang:Vestnes

27 - Turspråk:Engelsk, Norsk