

DATA QUALITY, LANGUAGE COVERAGE

Common Crawl ~ 9PiB

↳ only TEXT (HTML or PDF)

Domain level sampling → "BUDGET" for DOMAIN

↳ GRAPH-BASED LINKING

→ Below rank ~ 80M → SPAR

IN TIME...

• List based → page rank → harmonic centrality (now)

REPRESENTATIVE? → New
→ Duplicate
→ Depth of date

+ new & REVISITS
of URL

Text quality > Regional coverage

DOWNLOADING Common Crawl

→ DOWNLOAD-PATH

CC-downloader download-paths CC-MAIN-2024-46 wet
path To Folder

Formats

- 1) WARC - web page captures
- 2) ARC - older, before 2012
- 3) WAT - metadata and links
- 4) WET - plain text extracted from HTML
- 5) Index → CDX format or Parquet
- 6) Webgraphs & RANKS
- 7) Crawl metrics

WARC → • Freezes internet traffic between a client

• Successful | Not successful collected separately

• robots.txt

• Shuffled Files

• 1MiB max content → otherwise TQNCABET

Metadata → replicability of data
- WARC & HTTP headers
- robots.txt, h04

Applications

• Before only Wik. & Books were used

• But GPT-BERT performed well on CC data !!

Not representative
↳ Bias
↳ Code
↳ Unbalanced domain & language

Text extraction

• justText

• TiFiLatura → Different or

• Resiliparse

• Readability

Heuristic Filtering

- Length
- symbol-to-word-ratio
- Remove lists
- Require alphabetic character
- Stop words lists
- Repetitions

→ Discourage blocklist of words

Deduplication

• xxhash64 (not GPU!! and FAST)

• MinHash (fuzzy and GPU)

Quality Filter (Acplexy Filtering)

- KenLM trained on Wikipedia
 - KenLM trained on adult content
 - Quality classifiers
 - Education classifiers
- } n-gram models

[LangID project → Language identification
Web Language Project