# Common Crawl:
# open web data for everybody

Laurie Burchell

laurie@commoncrawl.org

# Overview

- What is the Common Crawl Foundation?

- What is the data like?

- A representative sample?

- Ongoing research: language diversity

- Next steps

slides

**COMMON**
CRAWL

# What is the Common Crawl Foundation?
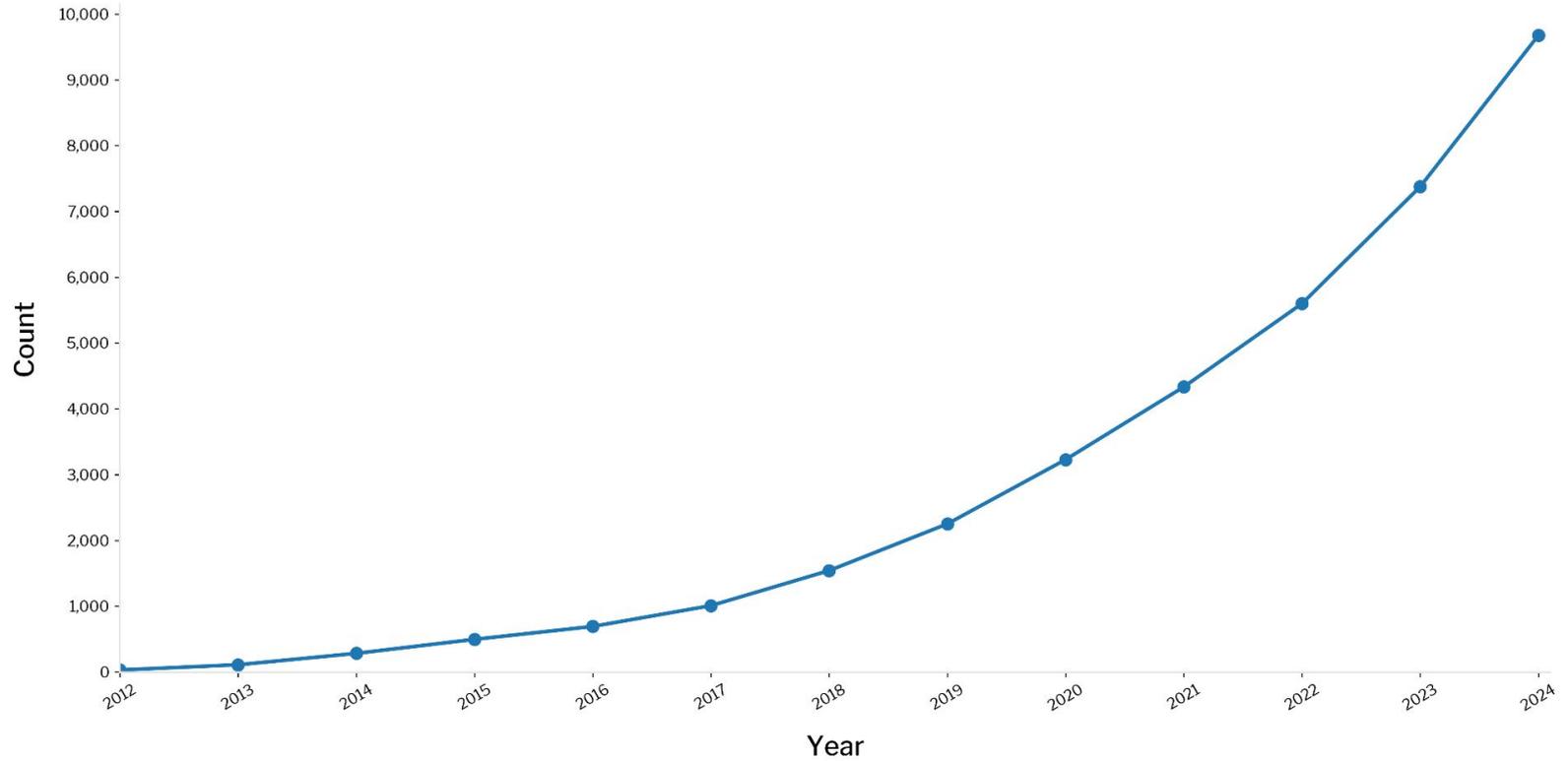
COMMON
CRAWL

# What is the Common Crawl Foundation?

- Non-profit founded in 2007 by Gil Elbaz

- Mission: **make web data accessible** to researchers, developers

- Small (but growing!) team of staff and volunteers

- Part of AWS Open Datasets Program: free to access!

**COMMON**
CRAWL

COMMON CRAWL

# 11.3 PB

…and growing by >4 TB each month

Dataset size as of September 2025.

Plot of Common Crawl citations (cumulative) in Google Scholar until January 2025

https://commoncrawl.org/research-papers
https://huggingface.co/datasets/commoncrawl/citations

# Projects with Common Crawl

- [Creating a large-scale, multilingual corpus](#)

- [Analysing disappearing links over time](#)

- [Detecting misinformation sources](#)

- [Censorship of Amazon products](#)

- [COVID-19 news mood map](#)

**COMMON**
CRAWL

# The data

# Overview

The Common Crawl corpus contains petabytes of data, regularly collected since 2008.

| Choose a crawl... ⌃ |
| --- |
| CC-MAIN-2025-33 |
| CC-MAIN-2025-30 |
| CC-MAIN-2025-26 |
| CC-MAIN-2025-21 |
| CC-MAIN-2025-18 |
| CC-MAIN-2025-13 |
| CC-MAIN-2025-08 |
| CC-MAIN-2025-05 |
| CC-MAIN-2024-51 |

Access our data via HTTP(S) or AWS:

see commoncrawl.org/get-started

(the most important) **Data products**

- **WARCs**: web page captures

- CDXJ and columnar **indices**

- **Web Graph**: structure and connectivity

**COMMON**
CRAWL

# WARCs: web page captures

- Web ARChive format

- Raw crawl data

  - Content payload

  - HTTP headers

  - Connection metadata (datetime, IP address)

- Derived data: WETs (text), WATs (metadata)

- Stats at [cc-crawl-statistics](cc-crawl-statistics)

**COMMON CRAWL**

# CDXJ and columnar indices

- Different indices for different kinds of access

- [CDXJ index](): single URLs or domains

- [Columnar index ]()(Parquet)

  - SQL queries

  - Big data toolkits

**COMMON**
CRAWL

# Web Graph

- Structure and connectivity of the web

- Two levels: host and domain

- Tools at cc-webgraph

- Stats at cc-webgraph-statistics

**COMMON CRAWL**

# There's more!

See our website:
[commoncrawl.org](commoncrawl.org)

Tools and examples:
[github.com/common crawl](github.com/commoncrawl)



## COMMON CRAWL

### Common Crawl August 2025 Crawl Archive (CC-MAIN-2025-33)

The August 2025 crawl archive contains 2.44 billion pages, see the [announcement](announcement) for details.

### Data Size and File Listings

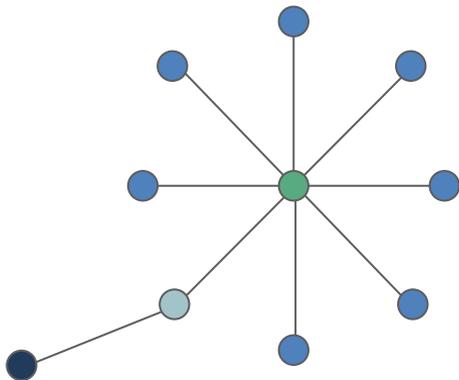| Data Type | File List | #Files | Total Size Compressed (TiB) |
|---|---|---|---|
| Segments | segment.paths.gz | 100 | |
| WARC | warc.paths.gz | 100000 | 88.24 |
| WAT | wat.paths.gz | 100000 | 16.71 |
| WET | wet.paths.gz | 100000 | 6.63 |
| Robots.txt files | robotstxt.paths.gz | 100000 | 0.15 |
| Non-200 responses | non200responses.paths.gz | 100000 | 2.97 |
| URL index files | cc-index.paths.gz | 302 | 0.19 |
| Columnar URL index files | cc-index-table.paths.gz | 900 | 0.21 |

A representative sample
of the web?

# How does crawling work?

- Start from seeds and spider out

- We crawl politely!

  - Slowly

  - Respecting robots.txt and opt out

  - No log ins

- NB: our crawl is **text-only**

COMMON CRAWL

# Why sample?

- TL;DR: **the web is big, our resources are finite**

- Per monthly crawl:

    - 2.5 billion page captures

    - 500+ billion links

    - 20+ billion unique URLs linked (excluding media links)

- Also: avoiding overload, spam traps!

$$H(v) = \sum_{u \neq v} \frac{1}{d(v, u)}$$

Where $H(v)$ is the **Harmonic Centrality** of vertex $v$,
and $d(v,u)$ is the shortest path distance between vertices $v$ and $u$.

Sample by budgeting domains,

guided by **harmonic centrality** rank
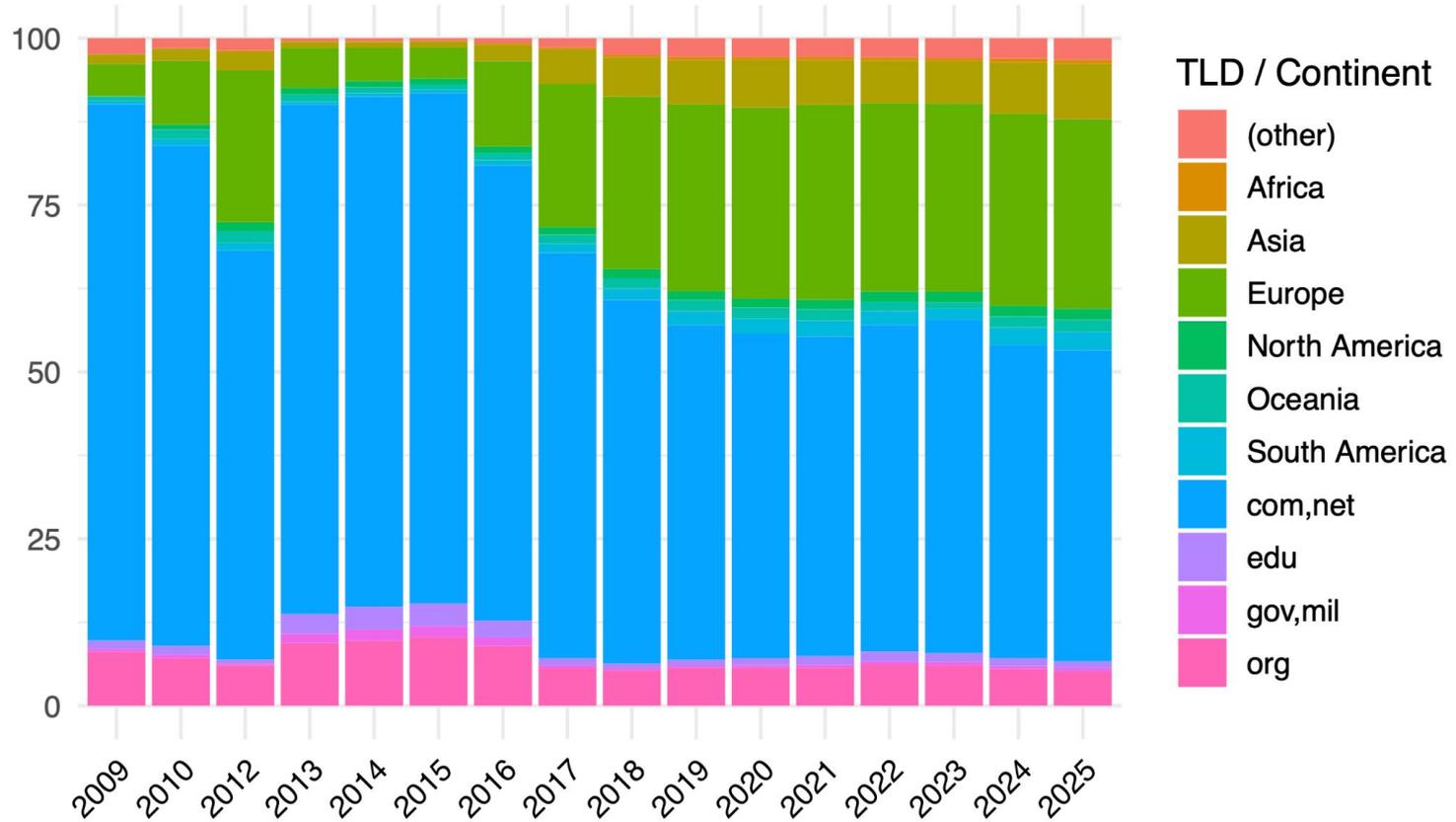
COMMON
CRAWL

# What is 'representative'?

- Breadth versus depth

- Freshness (new content)

- Amount of (near-)duplicates (per crawl and over multiple crawls)

- Regional coverage (top-level domains, content languages)

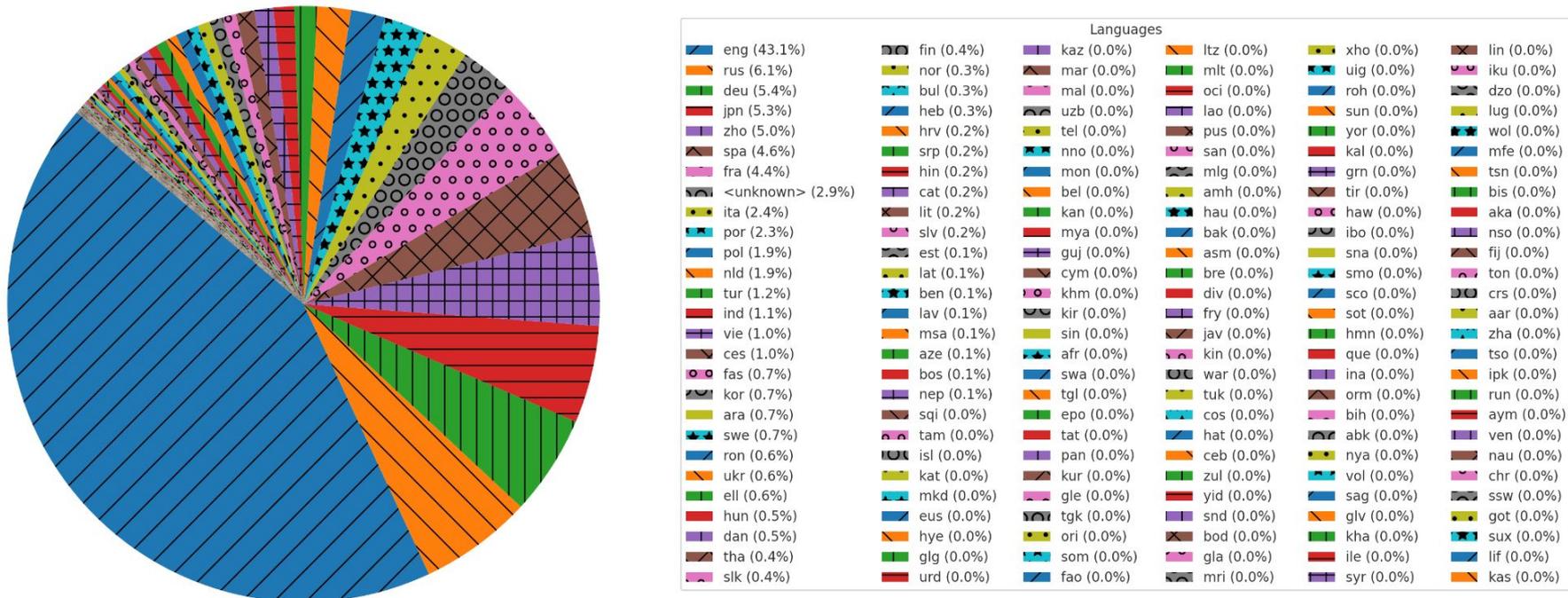- Content 'quality'

**COMMON CRAWL**

# What's the impact?

- 'Representative' depends on use case!

- We aim for compromise:
    - Breadth versus depth
    - Site categories, topics, languages, geographic regions...

- Focus on USA and Western world (for now...)

COMMON
CRAWL

# Top-Level Domains and Geographical Coverage

# Ongoing research: language diversity

COMMON
CRAWL

**Languages**

| | | | | | |
|---|---|---|---|---|---|
| eng (43.1%) | fin (0.4%) | kaz (0.0%) | ltz (0.0%) | xho (0.0%) | lin (0.0%) |
| rus (6.1%) | nor (0.3%) | mar (0.0%) | mlt (0.0%) | uig (0.0%) | iku (0.0%) |
| deu (5.4%) | bul (0.3%) | mal (0.0%) | oci (0.0%) | roh (0.0%) | dzo (0.0%) |
| jpn (5.3%) | heb (0.3%) | uzb (0.0%) | lao (0.0%) | sun (0.0%) | lug (0.0%) |
| zho (5.0%) | hrv (0.2%) | tel (0.0%) | pus (0.0%) | yor (0.0%) | wol (0.0%) |
| spa (4.6%) | srp (0.2%) | nno (0.0%) | san (0.0%) | kal (0.0%) | mfe (0.0%) |
| fra (4.4%) | hin (0.2%) | mon (0.0%) | mlg (0.0%) | grn (0.0%) | tsn (0.0%) |
| <unknown> (2.9%) | cat (0.2%) | bel (0.0%) | amh (0.0%) | tir (0.0%) | bis (0.0%) |
| ita (2.4%) | lit (0.2%) | kan (0.0%) | hau (0.0%) | haw (0.0%) | aka (0.0%) |
| por (2.3%) | slv (0.2%) | mya (0.0%) | bak (0.0%) | ibo (0.0%) | nso (0.0%) |
| pol (1.9%) | est (0.1%) | guj (0.0%) | asm (0.0%) | sna (0.0%) | fij (0.0%) |
| nld (1.9%) | lat (0.1%) | cym (0.0%) | bre (0.0%) | smo (0.0%) | ton (0.0%) |
| tur (1.2%) | ben (0.1%) | khm (0.0%) | div (0.0%) | sco (0.0%) | crs (0.0%) |
| ind (1.1%) | lav (0.1%) | kir (0.0%) | fry (0.0%) | sot (0.0%) | aar (0.0%) |
| vie (1.0%) | msa (0.1%) | sin (0.0%) | jav (0.0%) | hmn (0.0%) | zha (0.0%) |
| ces (1.0%) | aze (0.1%) | afr (0.0%) | kin (0.0%) | que (0.0%) | tso (0.0%) |
| fas (0.7%) | bos (0.1%) | swa (0.0%) | war (0.0%) | ina (0.0%) | ipk (0.0%) |
| kor (0.7%) | nep (0.1%) | tgl (0.0%) | tuk (0.0%) | orm (0.0%) | run (0.0%) |
| ara (0.7%) | sqi (0.0%) | epo (0.0%) | cos (0.0%) | bih (0.0%) | aym (0.0%) |
| swe (0.7%) | tam (0.0%) | tat (0.0%) | hat (0.0%) | abk (0.0%) | ven (0.0%) |
| ron (0.6%) | isl (0.0%) | pan (0.0%) | ceb (0.0%) | nya (0.0%) | nau (0.0%) |
| ukr (0.6%) | kat (0.0%) | kur (0.0%) | zul (0.0%) | vol (0.0%) | chr (0.0%) |
| ell (0.6%) | mkd (0.0%) | gle (0.0%) | yid (0.0%) | sag (0.0%) | ssw (0.0%) |
| hun (0.5%) | eus (0.0%) | tgk (0.0%) | snd (0.0%) | glv (0.0%) | got (0.0%) |
| dan (0.5%) | hye (0.0%) | ori (0.0%) | bod (0.0%) | kha (0.0%) | sux (0.0%) |
| tha (0.4%) | glg (0.0%) | som (0.0%) | gla (0.0%) | ile (0.0%) | lif (0.0%) |
| slk (0.4%) | urd (0.0%) | fao (0.0%) | mri (0.0%) | syr (0.0%) | kas (0.0%) |

Detected language distribution (averaged) in the last three crawls using CLD2 as the language identifier (**CC-MAIN-2024-46**, **CC-MAIN-2024-51**, and **CC-MAIN-2025-05**)

https://commoncrawl.github.io/cc-crawl-statistics/plots/languages
https://github.com/CLD2Owners/cld2

# Increasing language coverage

Two strategies:

1. More **diverse seeds**

2. Better **language detection**

COMMON
CRAWL

# Web Languages

Research workshop:
**WMDQS**

# Language Identification Shared Task

# Next steps with Common Crawl

# Getting started

- Visit our website: [Get Started guide](#)

- [Whirlwind Tour in Python](#)

- Join our community: [Discord](#), [contact us!](#)

**COMMON CRAWL**

# Thank you!

## Questions?

Laurie Burchell
laurie@commoncrawl.org