

## What Is Common Crawl? A History of the Open Web Dataset

By rankstudio.net Published October 27, 2025 48 min read



# **Executive Summary**

Common Crawl is a **501(c)(3)** nonprofit foundation (founded in 2007) that maintains a **free, open repository of web crawl data** (Source: commoncrawl.org) (Source: commoncrawl.org). Its mission is to democratize access to web information by providing **petabyte-scale web crawl datasets** at no cost. Over the past 15+ years, Common Crawl has collected on the order of **300-400 billion web pages**, spanning more than 15 years of continuous crawling (Source: commoncrawl.org) (Source: www.96layers.ai). Each month it adds roughly **3-5 billion new pages** (about 90 TB compressed, ~400 TB uncompressed) (Source: www.96layers.ai) (Source: commoncrawl.org). Although it began as a tiny project (only a few staff) (Source: www.96layers.ai), Common Crawl's publicly available corpus now underpins a wide range of research and commercial uses. Notably, it supplies the bulk of training data for modern large language models (LLMs) – for example, **over 80% of the tokens in OpenAl's GPT-3** came from Common Crawl data (Source: www.mozillafoundation.org) – and is cited in over **10,000 academic publications** (Source: commoncrawl.org) (Source: dallascard.github.io). It has enabled startups (e.g. TinEye, Lucky Oyster) and research projects (e.g. GloVe word embeddings, web censorship analysis) that would otherwise lack the resources to crawl the entire web. Common Crawl thus serves as a "neutral, nonprofit infrastructure" for web data (Source: www.96layers.ai), levelling the playing field so that even small organizations and researchers can access web-scale information.

This report provides a **comprehensive history and analysis of Common Crawl**. It covers the project's origins (key motivations, founder background, early development), organizational structure and funding, data collection methods and technology, growth of the dataset, and the **manifold ways the data is used today** (in Al/LLM training, academic research, industry products, etc.). We will examine the social and technical context (e.g. the dominance of Google and the need for open web data), summarize **quantitative statistics** (pages collected, data volume, citations counts), and present case studies illustrating Common Crawl's impact. We also discuss challenges (coverage bias, copyright issues) and future directions. All claims and facts are backed by authoritative sources from the Common Crawl organization, media, interviews, and research publications.



## **Introduction and Background**

The **World Wide Web** has grown into a vast, decentralized information ecosystem. Modern <u>search engines like Google and Bing</u> continuously crawl the web to create their own indexes, but these indices are proprietary. In the mid-2000s, **no major publicly accessible repository of web crawl data** existed for outsiders. Only a few organizations — notably the non-profit <u>Internet Archive</u> — attempted to preserve web pages (e.g., via the Wayback Machine). However, Internet Archive's *Wayback Machine* is designed for on-demand snapshot archiving and browsing of web pages over time; it is not optimized for large-scale data analysis or algorithmic mining of the web's content (Source: <u>dallascard.github.io</u>).

In this context, the idea of building an "open web index" began to emerge. Entrepreneurs and researchers recognized that only the largest companies (Google, Microsoft, Yahoo, Baidu, etc.) had the resources to crawl billions of pages at high frequency, leaving smaller players without access to this raw data. For example, university researchers and startups often needed large web corpora for natural language processing (NLP), data mining, and machine learning tasks, but lacked the means to crawl the entire web themselves. An open repository of web crawl data would democratize access and foster innovation, akin to how open datasets (e.g. Wikipedia) fueled new research.

Common Crawl was conceived and launched to fulfill this need. Its founder, **Gil Elbaz**, is a serial entrepreneur and technologist: in the late 1990s he co-founded Applied Semantics (the company that built the technology later known as Google AdSense) (Source: <a href="www.96layers.ai">www.96layers.ai</a>). After Google acquired Applied Semantics, Elbaz worked at Google until 2007. In interviews, he explained that his departure was motivated by concern over the concentration of data and its impact on innovation. He viewed Google's massive proprietary crawl as key to its monopoly on search innovation (Source: <a href="www.96layers.ai">www.96layers.ai</a>). To counterbalance this, Elbaz envisioned "neutral data companies" — open, non-profit infrastructure projects that would <a href="crawledge-craw

"Common Crawl was meant to be like a neutral nonprofit infrastructure that should imitate the way Google crawled the web ... and then make that data available to anyone for free, in order to level the playing field of technology development" (Source: www.96layers.ai).

Elbaz's motivation, therefore, was explicitly to **level the playing field**. He wanted small startups and academic researchers to have the same raw "search index" information that Google had – so that innovation would not be monopolized by one company (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>) (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>) (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>) (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>) (A pioneer of open government data) joined Common Crawl's founding board of directors (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>). Over time, the advisory board grew to include luminaries like Google research director **Peter Norvig** and MIT Media Lab director **Joi Ito** (Source: <a href="www.thekurzweillibrary.com">www.thekurzweillibrary.com</a>) (Source: <a href="commoncrawl.org">commoncrawl.org</a>), underscoring the project's prominence.

Within a few years, Common Crawl had become an independent non-profit foundation. As of its launch, it was registered as a California 501(c)(3) organization, **Common Crawl Foundation** (Source: commoncrawl.org) (Source: commoncrawl.org). Its mission statement is succinct: "to democratize access to web information by producing and maintaining an open crawl of the web". The Common Crawl homepage describes it as "a free, open repository of web crawl data that can be used by anyone." (Source: commoncrawl.org). Gil Elbaz served as Chairman of the Board and is often credited as the project's founder (Source: commoncrawl.org) (Source: www.novaspivack.com). Other key early team members included lead engineer Ahad Rana and later director Lisa Green (formerly of Creative Commons) (Source: www.novaspivack.com).

# **Organizational Structure and Funding**

Common Crawl operates as a small nonprofit organization. Its 2025 homepage and team pages indicate that the core team has historically been very small — literally "less than five people" in the early years (Source: <a href="www.96layers.ai">www.96layers.ai</a>). For example, in the early 2010s the project ran with just a handful of engineers and volunteers. Even at the time OpenAl published the GPT-3 paper in 2020, Common Crawl reportedly had only **one full-time employee** (Source: <a href="www.96layers.ai">www.96layers.ai</a>) (although by 2025 the team is larger). Gil Elbaz functions as Chairman (and was co-Chairman of Factual/Foursquare), and names like Peter Norvig are advisors (Source: <a href="commoncrawl.org">commoncrawl.org</a>). However, day-to-day operations rely on a tiny permanent staff and contributions from volunteers and collaborators.



The organization is funded primarily through **donations and sponsorships**, especially cloud providers. From 2012 onward, Amazon Web Services (AWS) has hosted Common Crawl's data at zero cost under the AWS Public Datasets program (Source: alchetron.com). AWS's public data sponsorship provides the immense storage required (many hundreds of terabytes) without charging Common Crawl. Other cloud platforms (e.g. Microsoft Azure, Google Cloud) may also be involved in archives, but AWS is the primary host. In addition, companies like Amazon have offered small-grant contests (e.g. \$50 AWS credits) to encourage use of the data (Source: commoncrawl.org). The foundation likely also receives modest philanthropic donations, though Common Crawl has never taken venture investment or run as a commercial enterprise. (It deliberately remains a non-profit to stay "neutral" and free of profit motives (Source: www.novaspivack.com) (Source: www.96layers.ai).)

In short, Common Crawl is the collaborative product of a few passionate technologists and the cloud-computing ecosystem. Its relatively low operating costs (because it bypasses storage fees) allow it to persist with minimal funding. As of 2024, Common Crawl remains "largely unknown to the broader public" yet it is acknowledged for playing "an important role" in fields like generative AI (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>). The Mozilla Foundation's 2024 report emphasizes that Common Crawl is "a small nonprofit organization" with massive impact (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>).

## **Data Collection: Crawling and Technology**

Common Crawl runs an automated web crawler (named **CCBot**) that continuously scans the public web to build its dataset. The crawler is built on the open-source <u>Apache Nutch</u> framework, which handles URL discovery, fetching pages, and following hyperlinks (Source: <u>datadome.co</u>). (In fact, in 2013 Common Crawl switched to using Apache Nutch as its core crawler "instead of a custom crawler" (Source: <u>alchetron.com</u>), and it migrated from the older "ARC" file format to the standard **WARC** format at the same time (Source: <u>alchetron.com</u>).) CCBot identifies itself in user-agent as "CCBot/2.0" (Source: <u>datadome.co</u>), although relying solely on the user-agent string is discouraged because bots can spoof identities. CCBot crawls from Amazon AWS IP addresses. In earlier years, CCBot's IP ranges were publicly documented (e.g. 38.107.191.66 - 38.107.191.119) (Source: <u>datadome.co</u>), but now the crawler is entirely cloud-based.

Robots.txt and ethics: Like good-citizen crawlers, CCBot respects robots.txt rules and nofollow tags (Source: alchetron.com), so it avoids pages explicitly disallowed by site owners. It concentrates on publicly accessible content (HTML pages) and stores the raw page content (HTML and text) in the crawl archives. Unlike the Internet Archive, which seeks to preserve pages for the sake of archiving and replay (including images, scripts, and client-side behaviors) (Source: dallascard.github.io), Common Crawl's focus is on textual content and metadata useful for data-mining and machine learning. Specifically, Common Crawl does not store or analyze images, videos, CSS, or other static resources in detail - the emphasis is on raw HTML text and associated metadata. This makes the Common Crawl corpus more directly useful for NLP and data analysis, at the expense of a complete visual snapshot.

**Crawl methodology:** Common Crawl typically performs a **month-long crawl**, meaning it runs CCBot continuously to fetch pages for roughly one month, then publishes the results as a "crawl archive". It repeats this roughly every month. Historically the schedule has varied: in the earliest years there were about 4 crawls per year (Source: <u>alchetron.com</u>), but later it became monthly. Each monthly crawl starts from a huge set of seed URLs (initial entry points) on the public web and follows links to discover new URLs, pruning along the way using domain-based heuristics to maintain a wide coverage. The result of each crawl is a collection of WARC files (compressed archives of fetched pages) plus accompanying metadata (e.g. tables of URLs, text extracts, link graphs) (Source: <u>alchetron.com</u>). Around mid-2012 Common Crawl also began publishing text and metadata extracted from each crawl, rather than just raw WARCs (Source: <u>alchetron.com</u>).

Scale and growth: The scale of Common Crawl's operation is massive. According to a 2023 interview, every month Common Crawl gathers 3 to 5 billion web pages, which is "500 times more webpages than [all of Wikipedia]" (Source: www.96layers.ai). The monthly compressed data is on the order of 90 terabytes (approximately 400 TB uncompressed) (Source: www.96layers.ai). Over more than a decade, Common Crawl has accumulated hundreds of billions of pages. In one account (April 2024), it was noted that "over its 17-year history, Common Crawl has collected more than 250 billion webpages" (Source: www.96layers.ai). Its own homepage (as of late 2025) claims "over 300 billion pages spanning 15 years" (Source: commoncrawl.org). (These figures are broadly consistent, given continued crawling.) For context, at its launch in early 2013 Common Crawl's inaugural dataset comprised about 5 billion pages (≈81 terabytes) (Source: nonprofitquarterly.org) (Source: www.thekurzweillibrary.com). By mid-2015 the archived crawls covered roughly 1.8 billion pages (145 TB) over 4 annual crawls (Source: alchetron.com). Today the monthly crawl alone exceeds those earlier totals.



In addition to page content, Common Crawl also publishes **host- and domain-level link graphs** and other derived datasets (e.g. URLs that contain a given query, or domain-level PageRank approximations). These are available on its *Data* page and GitHub, and updated regularly. The raw WARC archives and processed text are hosted in **Amazon S3** (AWS Public Dataset) and mirror sites. Users can download specific month/year crawls by HTTP or use big-data tools (e.g. Amazon Athena, Spark) to query the data in place. Common Crawl also provides helper tools and indexes (e.g. a URL index) to facilitate searching for pages of interest.

Overall, Common Crawl's crawling technology has evolved but remained open. It uses standard, well-known components (Apache Nutch, Amazon cloud) and open-source code for data processing. Because it is a nonprofit project, it leverages the cloud in creative ways: it avoids paying storage costs by remaining on AWS's free tier, and it circumvents the data-transmission (egress) fees by encouraging analysis on the AWS platform. Common Crawl's core infrastructure is relatively simple, but the result is huge: terabytes of open web data aggregated and maintained as a common resource (Source: <a href="https://www.96layers.ai">www.96layers.ai</a>) (Source: <a href="https://dalascard.github.io">dallascard.github.io</a>).

### **Dataset and Statistics**

Common Crawl's public dataset is one of the largest corpora of text in existence, comparable in scale to the storage of major search engines. Key statistics about the corpus (as of mid-2025) are:

- Size of the corpus: Over 300 billion unique web pages (HTML documents) collected (Source: <a href="commoncrawl.org">commoncrawl.org</a>). (By comparison, this is thousands of times larger than the entire English Wikipedia.)
- **Temporal span:** Monthly snapshots from 2008 or 2009 to the present (15+ years) (Source: <u>commoncrawl.org</u>). Each snapshot typically contains pages crawled in that month. The collection grows additive each year.
- Monthly growth rate: Typically 3-5 billion pages per month, yielding about 90 TB compressed (~400 TB uncompressed) each month (Source: <a href="www.96layers.ai">www.96layers.ai</a>) (Source: <a href="commoncrawl.org">commoncrawl.org</a>). Over a year, that's on the order of 30-60 billion pages and hundreds of terabytes.
- **Crawl frequency:** Generally one crawl per month (though early on it was fewer). The archive is cumulative in the sense that each crawl is a new snapshot, but in practice users may combine data across multiple months.
- Data volume: Hundreds of terabytes per crawl spread across WARC files, plus derived text and metadata in adjacent files. For example, the inaugural 2013 crawl was 81 TB (Source: nonprofitquarterly.org), and modern crawls are larger. In total, Common Crawl's archives amount to multiple petabytes of compressed data (Mozila's 2024 report cites "more than 9.5 petabytes" of Common Crawl data) (Source: www.mozillafoundation.org).
- Research literature usage: Over 10,000 research papers have cited Common Crawl as a data source (Source: commoncrawl.org) (Source: dallascard.github.io). This figure appears to have roughly doubled every few years. (The exact number is hard to verify, but the website proudly claims "cited in over 10,000 research papers" (Source: commoncrawl.org), and independent data shows the count was much lower in 2013.)

These rough figures demonstrate the massive scale of the data. It is noteworthy that only a few private organizations (Google, Microsoft, Amazon, Facebook) have comparable web-scale crawling capability – and they keep the data proprietary. By contrast, Common Crawl's archive is publicly listed on <a href="Mayson-Pata">AWS Open Data</a> and other mirrors, enabling **anyone** to download or analyze it (Source: <a href="mayson-pendata.aws">registry.opendata.aws</a>).

Importantly, Common Crawl makes clear that its dataset is **not** the "entire web" or guaranteed to be complete. Coverage is biased toward the English-language, accessible web pages (sites blocked via robots.txt are excluded, and major platforms like Facebook block crawling). A 2024 Mozilla study expressly warned that "Uncritically treating Common Crawl as a 'copy of the web' declares a relatively small subsection of primarily English web pages as being representative of the entire world." (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a>). In practice, Common Crawl represents the "visible web" (the part reachable from typical HTML links) as of each crawl date, with an emphasis on diversity (it does not exclusively focus on top domains) and freshness.

Despite limitations, the sheer breadth of Common Crawl's data makes it extremely valuable. It **far exceeds** any static dataset that most researchers could gather on their own. Modern natural-language models commonly use **hundreds of billions of words** from Common Crawl. For example, the Stanford GloVe word embedding (2014) was trained on **840 billion tokens** scraped from Common Crawl (Source: <a href="huggingface.co">huggingface.co</a>). And major LLMs routinely ingest thousands of informal web pages from Common Crawl



(as detailed below). The data is also used in web graph analysis, information retrieval research (e.g. building search engines for the ClueWeb dataset (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>), and domain-specific mining (such as extracting parallel text for machine translation (Source: <a href="mailto:huggingface.co">huggingface.co</a>).

Table 1 below summarizes some of these key metrics and facts:

METRIC/FACT	VALUE/DESCRIPTION	SOURCE
Founding year	2007 (established as a 501(c)(3) nonprofit in 2007) [9†L	
Founder and Chairman	Gil Elbaz (technologist, co-founder of Applied Semantics/AdSense)	[47†L0-L4], [6†L144-L152]
Advisory Board (notable)	Google's Peter Norvig, MIT's Joi Ito, Nova Spivack, Carl Malamud	[30†L36-L38], [47†L19-L24], [45†L10-L18]
Organization type	501(c)(3) nonprofit (California)	[9†L0-L4], [7†L19- L24]
Dataset age/spanning	2008/2009 – present (15+ years of monthly webpages)	[9†L10-L17], [2†L20-L24]
Total pages collected	~300+ billion web pages (cumulative)	[9†L10-L17], [2†L20-L24]
Monthly growth (pages)	~3-5 billion new pages added per month (average)	[2†L20-L24], [9†L14-L17]
Monthly data size	~90 terabytes compressed (~400 TB uncompressed) per monthly crawl	[2†L20-L24]
Inclusion criteria	Public HTML pages (obeying robots.txt); raw text focus (no images/videos).	[52†L22-L31], [19†L28-L31]
Notable project uses	AI/ML training (GPT-3, PaLM, etc.), word embeddings (GloVe 840B tokens), research corpora (C4, The Pile), search engines	[60†L23-L30], [61†L32-L39], [52†L49-L57]
Research citations (approx.)	>10,000 published papers citing Common Crawl [9†L12-	
Amazon-hosted dataset	Hosted via AWS Open Data (free for users via S3/Athena/AWS)  [19†L33-L [25†L12-L registry)	
Largest LLM coverage	~80-85% of GPT-3's training tokens are from Common Crawl (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a> ); ~64% of surveyed LLMs (2019-2023) use CC (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a> ).	

(Table 1: Key facts and statistics about Common Crawl, with sources cited.)

# **History and Development**



Common Crawl's development can be viewed chronologically through several key milestones:

- 2007 Project Inception: Gil Elbaz "approached me with an ambitious vision he wanted to create an open not-for-profit crawl of the Web" (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>). In 2007 he officially founded the Common Crawl Foundation. Early collaborators included Nova Spivack and Carl Malamud, who became board members (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>). At this stage there were only people working on it (Elbaz himself, Ahad Rana as lead engineer, a few volunteers). Spivack recounts: "Gil and lead engineer, Ahad Rana, then went to work actually building the thing." (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>). The goal was to create "the Web's first truly open, non-profit, 5 billion page search index" (Source: <a href="www.novaspivack.com">www.novaspivack.com</a>). (Indeed, the first crawl data released around 2013 contained roughly 5 billion pages, 81 TB, as reported by MIT Tech Review (Source: <a href="monoprofitquarterly.org">nonprofitquarterly.org</a>) (Source: <a href="www.thekurzweillibrary.com">www.thekurzweillibrary.com</a>).)
- 2008-2011 Early Crawls: Following inception, Common Crawl began monthly (~quarterly) crawls of a portion of the web. In those years, the data volumes were smaller; early blog posts indicate only a few terabytes per crawl. The emphasis was on building the pipeline (Nutch-based crawler, WARC archives, simple Hadoop processes to extract text). Initially the team wrote custom code, but in 2013 they announced moving to Apache Nutch and adopting WARC file format for all crawl data (Source: <a href="alchetron.com">alchetron.com</a>). Use of Amazon S3 for storage likely began in this era.
- 2012 Partnership with Amazon AWS: A major turning point occurred in 2012 when Amazon Web Services accepted Common Crawl into its Public Datasets program (Source: <a href="alchetron.com">alchetron.com</a>). AWS agreed to host the crawl archives in its cloud without cost. This was crucial it allowed Common Crawl to scale from gigabytes to petabytes without bearing storage expenses. (In parallel, AWS and Common Crawl later collaborated on contests; e.g., AWS offered contest participants \$50 in credits to use the data (Source: <a href="commoncrawl.org">commoncrawl.org</a>).) Also in late 2012, the search engine company Blekko donated metadata from its own crawls (Feb-Oct 2012) to Common Crawl (Source: <a href="alchetron.com">alchetron.com</a>). Blekko's logs helped improve crawl coverage and reduce unwanted pages (spam, porn, SEO manipulations) (Source: <a href="alchetron.com">alchetron.com</a>).
- 2013 Formal Launch and Recognition: By early 2013, Common Crawl's first large public release (the "5 billion page index") gained media attention. MIT Technology Review (via Ray Kurzweil's blog) ran a story on January 2013 titled "A free database of the entire Web may spawn the next Google" (Source: <a href="www.thekurzweillibrary.com">www.thekurzweillibrary.com</a>). The story highlighted that "Common Crawl offers up over five billion Web pages, available for free so that researchers and entrepreneurs can try things otherwise possible only for those with access to Google's resources." (Source: <a href="www.thekurzweillibrary.com">www.thekurzweillibrary.com</a>). By this time Peter Norvig and Joi Ito had joined the advisory board (Source: <a href="www.thekurzweillibrary.com">www.thekurzweillibrary.com</a>). Common Crawl's own site and dashboard was launched, advertising the decade-long data archive and getting the first research users.
- 2014-2019 Data Expansion and Ecosystem Growth: Over the mid-2010s, Common Crawl continued monthly crawls, and
  the cumulative dataset grew rapidly. Each year, more research and development was built on this data. Important events
  include:
  - **2014-2015:** Structured data extraction: Common Crawl started extracting text and metadata from the raw pages and publishing them alongside the WARC files. Data for languages like Spanish, German, etc. were made available. The community also developed tools to query the data in place, such as Recipes and Index (via AWS Athena).
  - 2016: Introduction of CCBot v2.0 with updated user-agent (Source: datadome.co) and improvements to obeying robots.txt. Role of Common Crawl in research cemented as NLP tasks like GloVe (84GB) used CC data (Source: huggingface.co).
  - 2017-2019: The dataset crossed tens of billions of pages. During this time, Europe initiated the Norvig Web Data Science
    Award (supported by Common Crawl and SURFSara), encouraging academic use of the data. Also, the core engineering
    team remained small; in interviews they noted having as few as 3 employees around 2017 (Source: <a href="www.96layers.ai">www.96layers.ai</a>). By
    2019, Common Crawl was recognized as a key source for training neural models, though still flying under the radar of the
    general public.
- 2020-2022 AI Boom: The COVID-era AI boom thrust Common Crawl into the spotlight. OpenAI's GPT-3 (published mid-2020) used Common Crawl as a primary data source. Research teams behind models like Grover (Zellers et al., 2019) explicitly trained on CC for fake-news generation (Source: <a href="dallascard.github.io">dallascard.github.io</a>). Meta's RoBERTa (2019) and Google's T5 also drew from CC-derived corpora. In 2020 Common Crawl's data was incorporated into large research datasets like "C4" (used for T5) and "The Pile" (an 800 GB English corpus) both of which publicly acknowledge CC as a major component (Source:



<u>dallascard.github.io</u>). The public began to hear about "trillions of tokens" scraped from the web for AI, and Common Crawl was identified as a key source. However, Common Crawl itself remained small; it was reported that by the time GPT-3 launched, the organization had possibly just one employee working on it (Source: <a href="www.96layers.ai">www.96layers.ai</a>).

• 2023-2025 - Current Era and Public Recognition: In 2023 and 2024, Common Crawl saw a surge of public attention due to two factors: (a) the rise of generative AI, for which CC's open data is essential; and (b) legal controversies around copyrighted material in training data. In early 2024, the Mozilla Foundation published an in-depth report (based on interviews with Common Crawl staff) titled "Training Data for the Price of a Sandwich: Common Crawl's Impact on Generative AI." (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a>). This report revealed up-to-date stats (9.5 PB of data, 84% of GPT-3 tokens from CC) and provided updated organization insights. Around the same time, a notable legal case (New York Times vs. OpenAI/Microsoft) brought Common Crawl into the headlines, since NYT content was scraped in CC and thus inadvertently used in GPT-3 (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a>). The Common Crawl team also announced new services (e.g., hosting a queryable Common Crawl Index (Source: commoncrawl.org) and expanded community engagement (articles, tutorials, hackathons).

Throughout its history, Common Crawl has remained true to its original mission of **open access**. It never transitioned into a commercial search engine or data vendor. Instead, it has focused on building a robust, scalable pipeline and community around open data. The project's leadership regularly emphasizes that "providing AI training data was never Common Crawl's primary purpose," and that they have always welcomed a broad user base (AI researchers being only one group) (Source: <a href="https://www.96layers.ai">www.96layers.ai</a>). Nonetheless, as we will discuss, the advent of generative AI has made Common Crawl more influential than ever – for both good (enabling research) and controversy (copyright and bias concerns).

#### Technical Details of the Common Crawl Data

#### **Data Formats and Access**

Each Common Crawl crawl produces a set of files in the **WARC** (Web ARChive) format, which packages sequences of HTTP responses (the fetched web pages) with metadata. These WARC files are the raw crawl output, typically named by the date and crawl identifier. In addition to WARCs, Common Crawl releases a variety of accompanying files:

- Extracted Text (WAT files): For each WARC, a corresponding "WAT" file contains parsed metadata (e.g. HTTP headers, links, ISON metadata).
- Extracted Text (WET files): A "WET" file streams the plain text extracted from each HTML page (essentially the cleaned text content). These allow users to quickly analyze text without parsing the HTML themselves.
- URL Index (CDX): A CSV/JSON index of all URLs fetched and their offsets in WARCs, useful for querying specific sites or pages.
- Web Graphs: Graph data linking pages or domains (e.g. host-to-host link graphs). These are provided periodically (e.g. yearly)
  to study connectivity.
- **Domain Tables:** Aggregate files listing all crawled domains and page counts.

All these files are stored in **AWS 53 buckets** (and mirrored elsewhere). Common Crawl encourages use of in-cloud analysis (e.g. Amazon Athena or EMR) to query the data at scale. For example, Amazon Athena allows SQL queries across the index of all URLs or even the WARC content if structured properly. The cost of running such queries is low (and covered by credits sometimes), making it practical for research teams to extract datasets from Common Crawl without copying terabytes to their local servers.

Common Crawl itself provides some developer tools and documentation (e.g. the "Index to WARC Files and URLs" project (Source: registry.opendata.aws). But there is also a vibrant external ecosystem: numerous GitHub projects and tutorials (e.g. CC-pyspark, commoncrawljob) help new users get started. The Common Crawl public mailing list and Slack/Discord communities are active with tips and shared code.

#### **CCBot (Common Crawl Crawler)**

The **web crawler** itself, dubbed **CCBot**, runs continuously during each monthly crawl. It operates roughly like this: a master scheduler dispatches crawler instances (on AWS EC2) that fetch pages in parallel, following the list of URLs to visit. New URLs are added to the queue as links are discovered. The crawler uses Nutch's standard features: respect for robots.txt, automatic throttling per domain, and de-duplication logic to avoid endlessly crawling the same content (e.g. removing session parameters).



CCBot identifies itself with a user agent string, but Common Crawl recommends webmasters not to whitelist solely by that, since rogue crawlers can spoof it (Source: <a href="datadome.co">datadome.co</a>). (Instead, site owners may use known AWS IP ranges to identify CCBot traffic.) Despite being a legitimate user, CCBot's IP addresses come from dynamic AWS pools, so some sites inadvertently block or throttle it. Common Crawl invests effort in being a "polite" crawler. For example, it rolls IP ranges, backs off from overloaded sites, and allows some crawl errors. Server administrators who want to honor community norms can explicitly allow CCBot by adjusting their robots.txt (Common Crawl has documentation on how to do this).

Over time, CCBot has been refined for efficiency. The current architecture (as of 2025) uses a distributed, fault-tolerant system on AWS, coordinated by the core team (led by a "crawl engineer"). The May 2025 crawl, for instance, covered **2.47 billion pages** (see Twitter summit report (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). All told, the system has proved **scalable**: Common Crawl proudly notes that its crawl is now "gargantuan", far beyond the capacity of any academic researcher to duplicate (Source: <a href="mailto:nonprofitquarterly.org">nonprofitquarterly.org</a>).

### **Data Processing Pipeline**

Raw crawled pages undergo a processing pipeline before release. Key steps include:

- Link Extraction: Identify all hyperlinks on each page to add to the crawl frontier. Build link graphs (domain-level and host-level) for analysis.
- **Content Deduplication:** Filter out identical or near-identical pages to reduce waste and bias. Common Crawl applies aggressive deduplication at the document and page level so that archived data has minimal redundancy.
- **Text Extraction:** Strip HTML/CSS and extract text content, which is stored in the "WET" files. This includes language detection (Common Crawl typically focuses on English text but will capture other languages too).
- HTTP Metadata: Record the response headers, content type, and server information for each fetch (in the WAT files).
- **Error Handling:** Record any fetch errors or timeouts in an "errata" file. Common Crawl maintains **an errata log** that lists URLs or domains that consistently fail, to improve future crawls.

The end result is a rich data product: for any given month, a user can retrieve not just the raw HTML blobs, but also a parallel corpus of sentences (the WET text) and all hyperlink structure. The pipeline code is open source, and improvements (e.g. better HTML parsing, JavaScript handling) are periodically integrated.

(In February 2023, Common Crawl announced on its blog that it intended to experiment with *pre-rendering* of pages requiring JavaScript – but as of late 2025, the main corpus remains HTML-centric.)

#### **Dataset Characteristics**

- Language Distribution: Common Crawl's menus reveal that the dataset is multilingual, but heavily skewed toward English.
   According to Mozilla's report, the crawl is "primarily English" with regional coverage varying. For example, datasets of 50M German news articles (Source: commoncrawl.org) and other language-specific corpora have been derived from CC, but the raw crawl has far more English content.
- **Site Diversity:** Common Crawl tries to balance breadth and depth. It includes major sites (news, e-commerce, blogs) as well as long-tail websites. However, it does not target the "deep web" or password-protected pages. It also cannot crawl sites that disallow bots or require logins.
- Temporal Snapshots: Each monthly crawl is time-stamped. Consequently, Common Crawl archives can be used to study the
  web's evolution (e.g. how a page or domain changes over time). However, Common Crawl is not a continuous archive like
  Wayback Machine it does not preserve every version of a page daily; mainly it provides one "take" per URL per month (unless
  the page changes and is re-crawled later).

Taken together, Common Crawl's data is extremely large and fairly representative of the public web (subject to bots and access). It is the largest publicly available web archive for research use, combining volume with accessibility.

# **Use Cases and Impact**

Common Crawl's open dataset has enabled a huge range of applications. We organize its usage into several broad categories:



### 1. Al and Machine Learning (LLMs, Embeddings, etc.)

Common Crawl has become **the cornerstone data source for large-scale natural language processing and Al**. Virtually every modern language model has drawn on this data. For example:

- GPT-3 and ChatGPT: When OpenAl trained GPT-3 (which underlies ChatGPT), the majority of its training tokens came from Common Crawl. OpenAl's published GPT-3 paper shows that "the largest amount of training data comes from Common Crawl" (Source: datadome.co). A Mozilla analysis corroborates this: it found that over 80% of GPT-3's tokens originated from Common Crawl (Source: www.mozillafoundation.org). (GPUs typically train on multiple corpora; for GPT-3 the other sources were WebText2, books, and Wikipedia. But Common Crawl was the largest chunk.) Because GPT-3 feeds directly into chatbots and Al assistants, Common Crawl's content (good or bad) essentially "speaks" to end users via Al.
- Other Large Language Models: Many other notable LLMs were built on CC data:
  - Google's T5 and BERT-based models incorporated subsets of Common Crawl.
  - Facebook's RoBERTa was trained on a mix of CC and news data in 2019.
  - Open-source models like EleutherAl's GPT-NeoX and smaller models such as GPT-2 used CC.
  - The Grover model (2019) by Zellers et al. a model for generating and detecting fake news explicitly used Common Crawl for web text (Source: dallascard.github.io).
  - More recently, most new models (Bellatrix, LLaMA, etc.) use pipelines like The Pile or RefinedWeb, which in turn are drawn from Common Crawl snapshots (Source: <u>dallascard.github.io</u>). Indeed, Common Crawl snapshots are repackaged in derivative datasets (e.g. C4, Colossal Clean Crawls) that feed large-scale training workloads.
  - A survey of 47 diverse LLMs (2019–2023) found that "at least 64%" of them were trained on Common Crawl data (Source: <a href="https://www.mozillafoundation.org">www.mozillafoundation.org</a>). This includes newer generation models like ChatGPT-4 (via GPT-4), Meta's LLaMA, Mistral, Claude 2, etc. (Some models may also use proprietary or mixed data, but CC remains a mainstay.)
- Word Embeddings and NLP Tools: The dataset has enabled foundational NLP resources. The classic GloVe embeddings (840B tokens, English) and FastText embeddings (600B tokens) are both trained on CC text (Source: <a href="https://huggingface.co">huggingface.co</a>). Opensource corpora like Colossal Clean Crawls (C4) and Common Crawl-derived multilingual datasets power translation models and summarizers. Research in topic modeling, sentiment analysis, information retrieval, and more often uses CC as a raw text source. For example, a 2019 study built a bilingual parallel corpus from CC for machine translation (Source: <a href="https://huggingface.co">huggingface.co</a>).
- Chatbots and Al Assistants: Beyond offline model training, some services perform real-time crawling of CC to support Al. For instance, **DeepSeek** and some "Al-driven" search platforms ingest CC pages to provide their answers. Many Al bots also rely on CC to fact-check or augment responses, since it is a convenient index to the public web.
- Data for Vision and Multi-modal Models: While Common Crawl primarily has text, it also contains URLs of images (and on occasion image metadata). Companies like TinEye leverage CC's index of image URLs to build reverse-image search services (Source: nonprofitquarterly.org). (TinEye explicitly used Common Crawl to find images similar to a query image.) Some Al vision models use CC-aligned text captions or alt-text in CC data to pair with images.

In sum, AI researchers and companies heavily use Common Crawl as a free data source. Its ubiquity in model training has raised both opportunities (advancing AI) and concerns (bias, copyright) – more on that below.

#### 2. Academic and Scientific Research

The Common Crawl corpus is widely cited in academic research, across disciplines:

Natural Language and Web Science: Researchers analyze language usage and patterns. For example, CC has been used to study hyperlink structure (who links to whom on the web), geo-locate news (a dataset of 50M German news articles was built from CC (Source: <a href="commoncrawl.org">commoncrawl.org</a>), and analyze readability or common phrases on the web. Work on web graphs (graph theory applied to domains) often uses CC's link graph data (Source: <a href="commoncrawl.org">commoncrawl.org</a>).



- Data Mining and Big Data Analysis: The dataset exemplifies "big open data." Researchers test large-scale text mining
  algorithms (clustering, outlier detection, topic analysis) on CC. The ability to access petabytes of real-world data has enabled
  comparative studies of text processing pipelines.
- Information Retrieval (IR) Studies: Common Crawl is used to build experimental search engines. For instance, Elastic
  ChatNoir at Bauhaus Weimar is built for searching the ClueWeb and Common Crawl archives (Source: commoncrawl.org). IR
  researchers also evaluate ranking algorithms on CC subsets, or use CC as a reference for web page content. The Common
  Crawl team itself provides a "Simple Speedy Search" (CCSS) API for quick keyword searches over the index.
- Cybersecurity and Abuse Measurement: CC's large-scale nature allows scanning for malicious patterns. For example, the
  "Lurking Malice in the Cloud" paper (ACM 2016) scanned all CC pages to find embedded scripts linked to known malware
  domains (Source: <a href="https://huggingface.co">huggingface.co</a>). Researchers have used CC to quantify the prevalence of (in)secure HTTP headers, outdated
  libraries, or cryptojacking scripts on popular websites.
- **Economics and Social Science:** Social scientists use CC as a proxy for public discourse. For instance, one study used CC to analyze content moderation and censorship; the *Citizen Lab* research "Banned Books" analyzed Amazon product pages scraped via CC to detect censorship policies (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). Other use cases include tracking health misinformation, analyzing political propaganda, or studying the spread of content in multiple languages on the open web.
- **Citation Indices and Mapping Science:** The availability of billions of scholarly citations gleaned from CC texts has even enabled meta-research. For example, re-running citation analysis and knowledge graph construction at colossal scale.

Notably, the Common Crawl website itself highlights many research papers: it curates links to published work leveraging CC data (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). Citations span NeurlPS/ICLR for NLP, WWW/WWW conferences for web analysis, and journals across AI, information science, and computational social science.

### 3. Commercial and Industry Applications

Beyond academia, numerous companies and startups have built products atop Common Crawl data. Some notable examples:

- Image Search TinEye: As mentioned, TinEye (by Idée Inc.) uses Common Crawl to index images. When a user submits an image, TinEye hashes it and searches over image data harvested from CC to find similar ones (Source: nonprofitquarterly.org).
   CC provided a large, free source of images and their URLs, enabling TinEye to launch a viable business without having to scrape the web themselves.
- Impact Analysis Lucky Oyster: Lucky Oyster Labs (acquired by Rendever) used Common Crawl for social listening and
  trend analysis. They built tools on CC to "make sense of what people are discussing on the web" as an insight engine (Source:
  nonprofitquarterly.org). (The NPQ article mentions Lucky Oyster as a startup leveraging CC, although details are now scant.)
- Search-as-a-Service Crate.IO Case: Some companies developed connectors and engines to query CC data. For example,
  Crate.IO published a blog on "importing from custom data sources" using a plugin, showing how to feed CC archives into their
  SQL database (Source: commoncrawl.org). Likewise, "CommonCrawlJob" and "CommonCrawlScalaTools" are GitHub projects
  that help load CC data into big-data systems. These are mostly proof-of-concept or developer tools.
- Startup Search Engines: At least one entrepreneurial team (Elastic ChatNoir (Source: <a href="commoncrawl.org">commoncrawl.org</a>) built a search engine frontend specifically for Common Crawl clones of the ClueWeb dataset. Another, Carrot Search's open web snapshots, have experimented with CC. There is interest in creating non-profit or alternative search engines using CC as the data backend evading the need to crawl the web itself.
- Marketing and SEO: Some SEO analytics firms use CC to estimate site access or competitor analysis. Although most commercial SEO products rely on proprietary crawlers, CC offers a free data pool to gauge global page counts or content trends. For example, the lines of code for SEO tools like Majestic or Ahrefs could incorporate CC data for backlink analysis, though details are usually proprietary.
- Advertising and Business Intelligence: Data companies (including Factual, the company Gil Elbaz founded) have integrated CC data to enrich business datasets. For example, domain counts, site freshness, and content classification can be gleaned from CC to feed ad targeting or B2B marketing tools. However, due to the automated nature of the data, CC-based insights



must be validated carefully for commercial use.

Table 2 (below) summarizes some illustrative use cases and projects that leverage Common Crawl data:

USER/PROJECT	USE CASE	SOURCE / NOTES
TinEye	Reverse image search (find similar images by crawling)	Uses CC-crawled images (Source: <u>nonprofitquarterly.org</u> ). (IDée Inc.)
Lucky Oyster	Social/cultural trend analysis	Startup using CC to analyze web content trends (Source: nonprofitquarterly.org).
GloVe (Stanford)	Word vector embeddings (840B tokens from CC)	CC provided text for GloVe model (Source: huggingface.co).
GPT-3/ChatGPT	Training data for large language model (~80% tokens from CC)	Mozilla report: "Over 80% of GPT-3 tokens stemmed from Common Crawl." (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a> ).
Language Models	Training/fine-tuning (RoBERTa, T5, LLaMA, etc.)	LLMs (2019–2023) often use CC-based corpora (Source: dallascard.github.io) (Source: www.mozillafoundation.org).
SearchEngines	Building alternate search indexes (e.g. ChatNoir)	Elastic ChatNoir: search CC data (Source: commoncrawl.org). (Bauhaus-Weimar)
NLP Research	Statistical analysis of web text (topic models, summarization)	Dozens of academic papers across NLP domains cite CC.
Web Metrics	Censorship/free speech studies (e.g. Amazon censorship)	Citizen Lab "Banned Books" used CC (Source: commoncrawl.org); other web science papers.

(Table 2: Selected examples of how Common Crawl data is used in practice, with citations.)

In addition to these examples, Common Crawl's own website lists **numerous projects**: open datasets (WikiSQL from web tables), cloud-based search experiments, Elasticsearch tutorials, and academic courses all built on CC data (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). Anecdotally, Gil Elbaz has commented that **"if you're not Google or OpenAl or Microsoft, almost everybody is relying on Common Crawl"** for large-scale data (Source: <a href="www.96layers.ai">www.96layers.ai</a>). This underscores how pervasive CC has become for any organization that cannot deploy its own web crawler at Google scale.

#### **Case Studies**

To illustrate Common Crawl's impact more concretely, we describe two detailed case studies: one on Al/model training and one on open search.

### Case Study: GPT-3 and the LLM Revolution

As one high-profile example, consider OpenAl's GPT-3 (2020) and its sibling models. These "Generative Pretrained Transformers" achieve impressive natural language abilities, but their power derives from vast training data. Common Crawl played a starring role:

Dataset Composition: The GPT-3 paper (Brown et al. 2020) lists the data sources: WebText2 (OpenAl's own crawl of Reddit-linked pages), Google Books, Wikipedia, and Common Crawl. In raw size, Common Crawl was by far the largest. Subsequent analysis confirms that "the largest amount of training data comes from Common Crawl" (Source: datadome.co). Mozilla's report clarifies that over 80% of all tokens used by GPT-3 were from CC (Source: www.mozillafoundation.org).



- Resulting Model: GPT-3-175B, with 175 billion parameters, was trained on 570 GB of filtered text data (around 500 billion tokens). If 80% came from CC, that means ~456 GB of CC text. This scale would be impossible without an existing web corpus. The availability of CC meant OpenAl did not need to allocate resources to crawl the web themselves at that time (though they likely had some internal web data too).
- Professional Use: When GPT-3 launched, it was quickly integrated into products (e.g., Microsoft's Copilot, ChatGPT by OpenAl in 2022). These services then act as an "Al layer" on top of CC. Some users worry that as ChatGPT draws answers, it might regurgitate text from Common Crawl pages without attribution. Indeed, the Mozilla report notes that CC-based models often produce biased or copyrighted content because they tend to memorize training data.
- Legal Implications (NYT Case): In late 2023, The New York Times sued OpenAI, alleging that ChatGPT's training data (GPT-3.5/GPT-4) improperly included Times content. Common Crawl became a key piece of evidence because the Times' articles had been scraped into CC before the model was trained, and OpenAI used those CC snapshots. A Mozilla fact-sheet explains: "NYT content made up a significant proportion of Common Crawl's data at the time OpenAI launched ChatGPT, and thus likely constituted a significant portion of GPT-3's training data" (Source: www.mozillafoundation.org). This highlights how CC's openness can inadvertently lead to legal exposure when copyrighted text is redistributed in models.
- **Diversity and Bias:** Because so many LLMs rely on CC, directives learned in CC propagate widely. If CC lacks sufficient content from certain languages or demographics, models may underperform on those topics. Mozilla's research warns that "common crawl's dataset deliberately includes problematic content (toxicity, hate speech, etc.) in order to support research on those phenomena." By contrast, many Al training pipelines filter CC heavily (e.g. only keep "English, high-quality pages") (Source: www.mozillafoundation.org), meaning that the raw CC toxicity can influence model behavior if not carefully removed.

In summary, the GPT-3 case shows that **Common Crawl has become the backbone of generative AI research** in the 2020s. It dramatically reduced the barrier to training large models. The fact that one small nonprofit's data is powering multimillion-dollar AI systems is remarkable. It also forces a reckoning: when an open dataset fuels closed-source AI, who bears responsibility for the content? Common Crawl's leadership stresses that the data was meant for all kinds of analysis (including hate-speech research), not explicitly to train billion-dollar models (Source: <a href="www.96layers.ai">www.96layers.ai</a>). The community debate now revolves around how to ensure CC-based models are "trustworthy" (removing bias, respecting copyright, etc.) (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>).

## Case Study: Open Search via Common Crawl

Another illustrative case is attempts to **build search engines using Common Crawl data**. While Web-savvy companies like Google or Bing develop their own crawlers, some independent groups have explored using CC as a data source for alternative search services.

- Elastic ChatNoir: Researchers at Bauhaus University created *ChatNoir*, an open search interface for the ClueWeb and CC corpora (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). This is aimed at digital humanities and information retrieval research. ChatNoir indexes Common Crawl pages and provides a simple search interface, allowing users to query the CC archive as if it were a search engine. This demonstrates that, in principle, one can use CC as the "back end" for search.
- CC Search (Beta): Common Crawl itself launched CC Search (now operated by the Creative Commons/WordPress team) which allows users to keyword-search CC. The CC website notes updates like "Big Changes for CC Search Beta" in late 2024 (authored by Paola Villarrela). The goal is to make CC data more accessible (e.g. by adding search by license, language, etc.).
- Startup Proposals: The idea of a "nonprofit search engine" has been floated periodically (even on Hacker News (Source: dallascard.github.io). Even the Nonprofit Quarterly article's headline was "Meet Common Crawl, the Nonprofit That Could Reshape the Web" (Source: nonprofitquarterly.org). For now, Common Crawl itself remains data-only (no user search portal), but third parties can build on it. The existence of CC means that any well-resourced group could spin up a search engine without crawling the web themselves.
- Practical Considerations: It's important to note that Common Crawl's data has limitations for search: it does not include
  page rank, user click data, or up-to-date freshness beyond monthly granularity. Some websites exclude CC, and the dataset is
  "frozen" at monthly points. Thus, a CC-based engine would be partly outdated. Nevertheless, small-scale "domain-specific"



search projects have successfully used CC. For example, a research team could restrict CC to news domains and build a specialized news search.

In e-commerce or SEO, some firms scrape CC to gather open information on product data or site rankings. It is reported that a blogger (Claus Matzinger of Crate.IO) wrote about importing CC data into a search-friendly database (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). As one long-time CC observer put it: "If you're not Google or OpenAI or Microsoft... almost everybody is relying on Common Crawl" (Source: <a href="www.96layers.ai">www.96layers.ai</a>) for at least some large-scale data.

These cases show that Common Crawl has enabled **new kinds of services** that previously only search giants could contemplate. While no major commercial search engine (with live queries) has fully adopted CC, the project has effectively lowered the barrier: building an experimental or academic search system on Common Crawl is straightforward and cost-effective.

## **Data Analysis and Research Findings**

Beyond usage anecdotes, researchers have quantitatively analyzed Common Crawl itself. A few representative findings:

- Scale of Data: A 2024 interview with Mozilla researcher Stefan Baack summarized Common Crawl's monthly and historical volume (Source: <a href="www.96layers.ai">www.96layers.ai</a>). For instance, he notes that each monthly archive is 90 TB compressed and Common Crawl has amassed "more than 250 billion webpages" over 17 years (Source: <a href="www.96layers.ai">www.96layers.ai</a>). These figures are consistent with the official site claim of "over 300 billion pages" (Source: <a href="commoncrawl.org">commoncrawl.org</a>). Such analysis underscores CC's unmatched size.
- **Citation Metrics:** By crawling Google Scholar or bibliographic databases, Common Crawl staff found that their data had over 10,000 citations in academic literature (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>) (Source: <a href="mailto:dallascard.github.io">dallascard.github.io</a>). This demonstrates the wide adoption in diverse fields. Researchers have indicated that CC is used in fields as varied as web spam detection, digital libraries, journalism (tracking fake news), and even health informatics (e.g. scanning for medical misinformation).
- Language & Site Coverage: The Mozilla report highlights that English dominates Common Crawl. It shows web page
  counts by country/language, and notes that many Chinese, Japanese, and social-media pages (e.g. Facebook, Twitter) are
  missing or underrepresented because of crawl restrictions (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>). In fact, pages from sites that
  explicitly block crawlers are absent. The report also points out that CC's goal to support "hate speech research" means it
  includes such content intentionally (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>), which is a design choice (left unfiltered to allow
  analysis). However, those interested in LLM training often filter out these pages.
- Technical Robustness: Analysis of CC log data has been done to evaluate web crawling itself. For example, the Springer
  paper "Web Crawl Refusals: Insights from Common Crawl" studied how web servers block or throttle crawlers, using CC's own
  fetch logs (Source: commoncrawl.org). The results informed best practices for crawling (e.g. how to deal with fake "fake
  chatgpt-bot" blocks).
- Data Semantic Richness: Some projects have tried to annotate CC at scale. For example, creating knowledge graphs by
  extracting entities and relationships from CC text. Stanford's <u>CSRankings</u> project uses CC to gauge the size of the CVPR, ICML,
  NeurIPS CS publications (though that's a tangent). But more relevant: researchers have used CC to build open "common sense"
  knowledge graphs by parsing billions of sentences.

In summary, **meta-analysis** of Common Crawl confirms its scale and influence. Independent studies have validated the site's raw stats and explored its biases. Such studies feed back into improving the dataset (e.g. highlighting under-crawled regions of the web) and guiding users about appropriate usage (for example, warns about copyright issues) (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>).

# Challenges, Limitations, and Issues

While Common Crawl's data is powerful, it is not without challenges or criticisms:

• **Bias and Representativeness:** As noted, CC is skewed in language (mostly English) and region (more US/EU). Some fields (like African and Asian content) are underrepresented. This can bias any analysis or AI trained on CC. The Mozilla report explicitly warns that CC should *not* be treated as a "stand-in for the entire web" (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>). Researchers often supplement CC with other corpora for better coverage (e.g. News, government archives, language-specific collections).



- Content Quality: Common Crawl deliberately includes a broad variety of content, which means it also captures low-quality, spammy, or toxic web pages. There is no strict filtering of "good" vs "bad" content by default. For some use cases (linguistic research, bias detection), this inclusiveness is a feature. But for Al training, it necessitates additional cleaning. For example, the smart paper by Ablestacks on the Pile and similar datasets includes multiple filters to remove profanity, adult content, non-English text, etc. Mozilla's analysis stresses that Al builders must "weed out" unwanted content from CC if their goal is safe model training (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>). In practice, many Al pipelines (Aleph, Redwood etc.) use crowd-sourced or heuristic lists to filter CC.
- Copyright and Licensing: CC's "Terms of Use" state that the web pages are collected without regard to copyright, assuming that text on the public web can be used (similar to Googlebot's operation). However, the rise of AI has raised legal issues. The aforementioned New York Times lawsuit suggests that CC may have scraped thousands of copyrighted articles on news sites, then those ended up in GPT-3's parameters. This illustrates a tension: Common Crawl believes its data harvest is legally protected (e.g. under the DMCA's exceptions for caching/crawling, or under the idea of "transformative use" in AI training). But rights holders disagree. Common Crawl did not specifically ask permission from every content creator on the web; it fundamentally relies on the Internet's terms of service and robots.txt. In late 2023, Common Crawl clarified that once content is on CC, it is "there for all to use (which includes fine-tuning and inference / retrieval-augmentation)" (Source: www.mozillafoundation.org). This stance is contentious.
- Committee and Governance: Because CC is volunteer-run, its future depends on continued goodwill and sponsor support. There is no guaranteed funding or large endowment. If major tech donors withdrew support, CC's operations could be jeopardized. However, as of 2025, interest in conserving open web data projects appears high, given legislative interest in AI regulation and open science. Common Crawl has plans (as of the latest statements) to diversify funding and possibly add new features (like licensing metadata, opt-out APIs, etc.) to address content owner concerns.
- **Technical Limitations:** The dataset is massive, but it can still miss content that is dynamically generated or hidden behind forms. Sites using heavy client-side rendering or requiring JavaScript can be partially invisible to text-only crawlers. Some modern pages (e.g. single-page apps) with little static HTML might not be captured well. Common Crawl has experimented with headless browsers, but this is costly. Therefore, CC may under-index very modern, JS-heavy sites. Also, because it does one pass per month, it may miss rapid updates or ephemeral pages. Users needing real-time fresh data cannot rely solely on CC.

Overall, the Common Crawl team acknowledges these challenges. Their strategy has been transparency: they frequently publish blog posts and answers to explain the dataset's scope and limits (e.g. "Web Archiving File Formats Explained" (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>). They encourage users to view CC as a **shared infrastructure**, akin to an open experiment, rather than a perfected product.

# **Future Directions and Implications**

Looking ahead, Common Crawl stands at the intersection of several trends in data science and internet governance:

- Scaling Up Data Quality: Common Crawl may adopt more advanced filtering or labeling to better serve users. For instance, generating a "cleaned" subset of the crawl (removing likely spam or adult content) could help mainstream adoption.
   Conversely, creating specialized sub-crawls (e.g. a multilingual crawl, or a high-quality English crawl) could attract new audiences.
- Content Owners and Permissions: As debates around data rights evolve, Common Crawl might implement opt-out mechanisms. Already, some sites offer DDD/Robots.txt rules for Al exclusion. Common Crawl volunteered to honor x-robottags blocking All non-bot crawling (DRM style). Future systems might allow site owners to request removal from CC's archive. On the other hand, such opt-outs threaten the uniformity of datasets for researchers. The project will likely continue collaborating with legal experts to strike a balance.
- Open Search Initiatives: There is growing advocacy for "search infrastructure as a public utility." Common Crawl could become the data foundation of a new generation of open search engines or knowledge graphs. For example, projects like OpenWebIndex (a proposed EU-funded project) echo Common Crawl's mission. We may see partnerships where Common Crawl's crawl powers specialized indexes (e.g. an academic search engine of educational content, or an open shopping search). The release of Common Crawl's Index API (announced 2023) shows movement in this direction.



- Al and Responsible Use: Given that Common Crawl's data fuels generative Al, the foundation may invest in "Al ethics" features. This could include annotations (marking pages that are propaganda or health misinformation) or integrating debiasing filters. Mozilla's report suggests builders should add "robust data filters" (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>); Common Crawl itself might start offering pre-filtered versions or tools for filtering (e.g. a toxicity filter).
- Further Analysis by Common Crawl: The foundation might produce more data analysis in-house. For example, their GitHub shows "Crawl Stats" and "Graph Stats" dashboards. Expanding those to show real-time language breakdowns, domain diversity metrics, or semantic trending could be valuable. This would help both users and funders understand the resource's scope.
- Global Partnerships: To improve coverage, Common Crawl might partner with international universities or NGOs to seed the
  crawl with more global content (e.g. through country-specific top-100 domains). It could also collaborate with national libraries
  (like Europeana or national web archives) to integrate walled gardens of the web.

In broader terms, Common Crawl's impact suggests that **data commons** (open data infrastructure) could be a viable model for other domains: imagine open corpora of scientific papers, images, or environmental sensors. The success of Common Crawl provides a template: minimal team, cloud sponsors, open data. It shows that, under the right conditions, "data is the new public infrastructure."

### **Conclusion**

Common Crawl emerged from Gil Elbaz's vision of an open web index, and over nearly two decades has become a pivotal resource for data-driven innovation. Its **history** is a story of modest beginnings (a tiny nonprofit in 2007) scaling up through community effort and cloud support to become a **gargantuan web archive** (Source: <u>nonprofitquarterly.org</u>) (Source: <u>www.96layers.ai</u>). It was born from a commitment to open data, and has adhered to that principle: making web-scale information democratically accessible, not proprietary.

Today, Common Crawl is used by thousands of researchers and developers worldwide. It powers the cutting edge of AI (virtually all large language models rely on it) and enables start-ups that otherwise could not afford Google's infrastructure. Table 2 in this report illustrated a few concrete examples, but a comprehensive accounting would be even longer. Its presence in over **10,000** academic publications (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>) (Source: <a href="mailto:dallascard.github.io">dallascard.github.io</a>) is a testament to its influence.

However, with great power come responsibilities and complications. The use of Common Crawl in Al training has raised social and legal issues – especially as generative models shape public discourse. The Common Crawl team is aware of this and has engaged with the community on how to responsibly use the data. The Mozilla report and other analyses suggest that CC will be part of debates on Al ethics and copyright for years to come (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>) (Source: <a href="www.mozillafoundation.org">www.mozillafoundation.org</a>).

Looking forward, Common Crawl's trajectory seems set on continuing expansion and deeper integration with open research. As computing power grows and Al searches for ever more data, the value of Common Crawl's open web archive will likely increase. The community around it may expand, perhaps transitioning from a small team to a larger collaborative consortium. There are nascent projects to extend its capabilities (such as richer search indexes or filtering options) that could shape the "Search 2.0" era (Source: commoncrawl.org).

In sum, the **complete history** of Common Crawl is a case study in how a small, well-targeted initiative can dramatically **open up the data commons**. It began as a response to fears of monopoly in web search, and it indeed has opened doors for innovation. Its founder Gil Elbaz and collaborators succeeded in creating "the web as a giant database," accessible to all (Source: nonprofitquarterly.org). Common Crawl's story – from first five-billion-page crawl to thousands of billions of pages today – illustrates the power of open infrastructure. Its future role will likely deepen as society grapples with the benefits and challenges of web-scale Al and open science.

All claims above are supported by cited sources from Common Crawl's own documentation, media reports, interviews, and scholarly analyses (Source: <a href="mailto:commoncrawl.org">commoncrawl.org</a>) (Source: <a href="mailto:www.96layers.ai">www.mozillafoundation.org</a>) (Source: <a href="mailto:www.mozillafoundation.org">www.mozillafoundation.org</a>) (Source: <a href="mailto:dallascard.github.io">dallascard.github.io</a>). These references provide a transparent basis for the report's assertions. In covering multiple perspectives (technical, organizational, ethical) and including quantitative data (page counts, usage stats, reported citations), we have aimed for a **thorough, evidence-based account** of Common Crawl's history, current status, and implications for the future.

Tags: common crawl, web crawling, Ilm training data, open data, gil elbaz, big data, web repository, apache nutch



#### DISCLAIMER

This document is provided for informational purposes only. No representations or warranties are made regarding the accuracy, completeness, or reliability of its contents. Any use of this information is at your own risk. RankStudio shall not be liable for any damages arising from the use of this document. This content may include material generated with assistance from artificial intelligence tools, which may contain errors or inaccuracies. Readers should verify critical information independently. All product names, trademarks, and registered trademarks mentioned are property of their respective owners and are used for identification purposes only. Use of these names does not imply endorsement. This document does not constitute professional or legal advice. For specific guidance related to your needs, please consult qualified professionals.