INTERNET ARCHIVE

# The Annual Wikimania Update

Mark Graham
September 18th, 2023

# Team "Turn All References Blue" (TARB)



Stephen Balbach



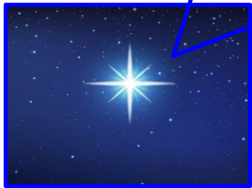Mark Graham



Max Doerr



Jake Orlowitz



James Hare



Dr. Sawood Alam

# Turn All [References](#) Blue

In the future nearly everything ever written or recorded, and all related resources, will be a click, gesture or a thought away.
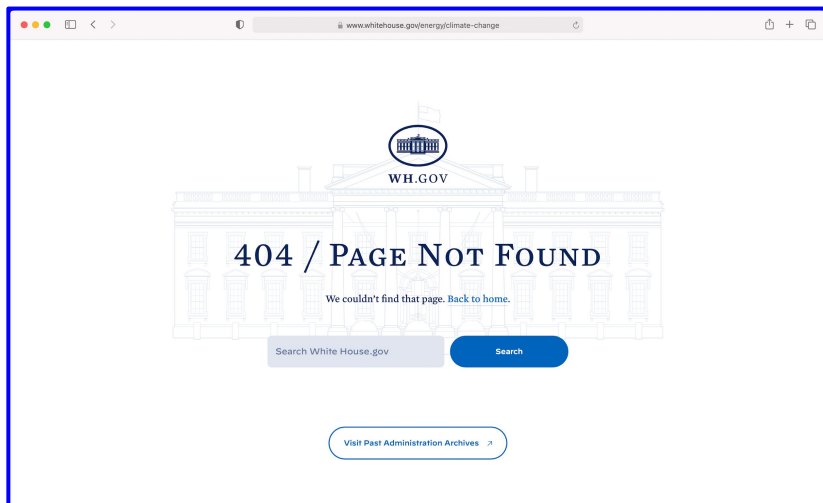
The goal of the "Turn All References Blue" project is to contribute to making that a reality, **starting with Wikipedia**.
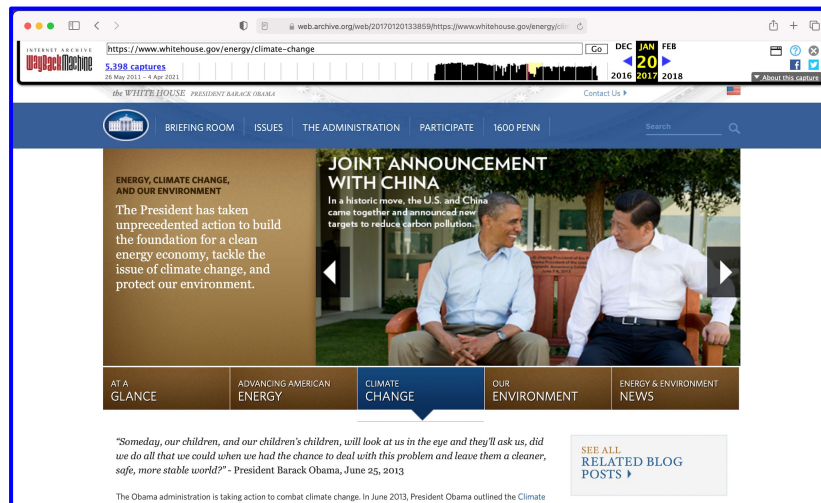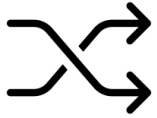
# Link Rot: When Good URLs go Bad

## Live Web (404)



## Archived Version

# Content Drift: Change Happens

## April 2, 2020

Strategic National Stockpile

Strategic National Stockpile is the nation's largest supply of life-saving pharmaceuticals and medical supplies for use in a public health emergency severe enough to cause local supplies to run out.

When state, local, tribal, and territorial responders request federal assistance to support their response efforts, the stockpile ensures that the right medicines and supplies get to those who need them most during an emergency. Organized for scalable response to a variety of public health threats, this repository contains enough supplies to respond to multiple large-scale emergencies simultaneously.

About the Stockpile

Sustaining the Stockpile

Stockpile Products

## April 6, 2020

Strategic National Stockpile

The Strategic National Stockpile's role is to supplement state and local supplies during public health emergencies. Many states have products stockpiled, as well. The supplies, medicines, and devices for life-saving care contained in the stockpile can be used as a short-term stopgap buffer when the immediate supply of adequate amounts of these materials may not be immediately available.

About the Stockpile

Sustaining the Stockpile

Stockpile Products

# We get URLs cited in Wikipedia via the EventStreamAPI

# We Archive Millions of Wikipedia Sourced URLs/day

cawg 🤖 | web | **top** **space** **dash** | **dump** scan **reduce** **viz**

From | 2020 ⇕ | 4 ⇕ | 1 ⇕ | To | 2023 ⇕ | 8 ⇕ | 7 ⇕ | **Set**

Days | 1224 | **collection** wikipedia-eventstream | ☆ Meta-counts | 1146/1220 (94%) | ⋯

| | wikipedia-eventstream | | | | | |
|---|---|---|---|---|---|---|
| | date | items | count | | size | | B/uri |
| 1 | 2023-08-07⭐ | 7 | 1,298,950 | 0% | 82,731,858,081 | 0% | 63.7KB |
| 2 | 2023-08-06⭐ | 8 | 1,410,651 | 0% | 81,116,676,520 | 0% | 57.5KB |
| 3 | 2023-08-05⭐ | 6 | 824,281 | 0% | 60,862,226,293 | 0% | 73.8KB |
| 4 | 2023-08-04⭐ | 6 | 960,376 | 0% | 62,092,605,426 | 0% | 64.7KB |
| 5 | 2023-08-03⭐ | 7 | 1,076,207 | 0% | 82,631,879,715 | 0% | 76.8KB |
| 6 | 2023-08-02⭐ | 6 | 890,475 | 0% | 81,443,293,754 | 0% | 91.5KB |
| 7 | 2023-08-01⭐ | 8 | 1,070,500 | 0% | 81,339,625,595 | 0% | 76.0KB |
| 8 | 2023-07-31⭐ | 8 | 1,405,654 | 0% | 83,193,245,880 | 0% | 59.2KB |
| 9 | 2023-07-30⭐ | 12 | 2,305,861 | 0% | 132,761,583,245 | 0% | 57.6KB |
| 10 | 2023-07-29⭐ | 9 | 1,390,846 | 0% | 93,290,117,260 | 0% | 67.1KB |

| | ALL collections | | | | |
|---|---|---|---|---|---|
| | ☆ | items | count | size | B/uri |
| 1 | ⭐ | 5,934 | 8.6B | 57.7TB | 6.7KB |
| 2 | ⭐ | 6,337 | 9.4B | 61.6TB | 6.5KB |
| 3 | ⭐ | 6,241 | 9.6B | 61.4TB | 6.4KB |
| 4 | ⭐ | 6,863 | 7.5B | 65.5TB | 8.8KB |
| 5 | ⭐ | 5,741 | 3.4B | 55.0TB | 16.2KB |
| 6 | ⭐ | 5,888 | 8.2B | 56.1TB | 6.8KB |
| 7 | ⭐ | 4,816 | 6.0B | 45.7TB | 7.6KB |
| 8 | ⭐ | 4,317 | 3.1B | 40.4TB | 13.0KB |
| 9 | ⭐ | 4,253 | 2.8B | 37.6TB | 13.2KB |
| 10 | ⭐ | 4,372 | 2.4B | 41.3TB | 17.1KB |

# How 1 seed URL can result in 31,259 Archived URLs



## 1 seed page



## 174 outlinks



## 338 embeds



## 30,746 outlink embeds



Results from cnn.com on June 17, 2021

ARCHIVE INTERNET

# We also Archive Millions of Outlinked URLs

**cawg** 🤖 | web | top | space | dash | dump | scan | reduce | viz

From | 2020 ⇕ | 4 ⇕ | 1 ⇕ | To | 2023 ⇕ | 8 ⇕ | 7 ⇕ | Set

Days | 1224 | collection | wikipedia-eventstream-outlinks | ☆ Meta-counts | 100% | ...

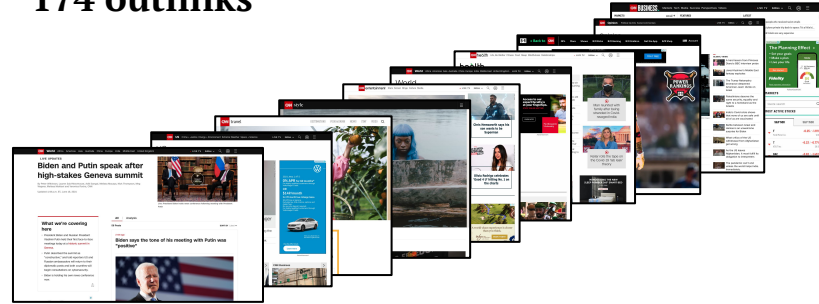| | wikipedia-eventstream-outlinks | | | | |
|---|---|---|---|---|---|
| | date | items | count | size | B/uri |
| 1 | 2023-08-07⭐ | 46 | 6,846,451 0% | 467,531,356,333 1% | 68.3KB |
| 2 | 2023-08-06⭐ | 62 | 9,797,050 0% | 639,846,297,269 1% | 65.3KB |
| 3 | 2023-08-05⭐ | 53 | 7,779,631 0% | 561,287,474,335 1% | 72.1KB |
| 4 | 2023-08-04⭐ | 41 | 5,246,811 0% | 416,772,334,461 1% | 79.4KB |
| 5 | 2023-08-03⭐ | 48 | 5,908,604 0% | 496,407,096,697 1% | 84.0KB |
| 6 | 2023-08-02⭐ | 52 | 6,007,106 0% | 544,033,016,699 1% | 90.6KB |
| 7 | 2023-08-01⭐ | 64 | 6,799,250 0% | 661,515,880,081 1% | 97.3KB |
| 8 | 2023-07-31⭐ | 64 | 8,472,577 0% | 672,326,227,733 2% | 79.4KB |
| 9 | 2023-07-30⭐ | 95 | 15,519,149 1% | 973,588,838,302 3% | 62.7KB |
| 10 | 2023-07-29⭐ | 64 | 9,759,397 0% | 673,523,724,314 2% | 69.0KB |

| | | ALL collections | | | |
|---|---|---|---|---|---|
| | ☆ | items | count | size | B/uri |
| 1 | ⭐ | 5,934 | 8.6B | 57.7TB | 6.7KB |
| 2 | ⭐ | 6,337 | 9.4B | 61.6TB | 6.5KB |
| 3 | ⭐ | 6,241 | 9.6B | 61.4TB | 6.4KB |
| 4 | ⭐ | 6,863 | 7.5B | 65.5TB | 8.8KB |
| 5 | ⭐ | 5,741 | 3.4B | 55.0TB | 16.2KB |
| 6 | ⭐ | 5,888 | 8.2B | 56.1TB | 6.8KB |
| 7 | ⭐ | 4,816 | 6.0B | 45.7TB | 7.6KB |
| 8 | ⭐ | 4,317 | 3.1B | 40.4TB | 13.0KB |
| 9 | ⭐ | 4,253 | 2.8B | 37.6TB | 13.2KB |
| 10 | ⭐ | 4,372 | 2.4B | 41.3TB | 17.1KB |

INTERNET ARCHIVE

# An Interactive Page from the 117th Congress

IABot has run on **215** out of **334**
Wikipedia language editions, and **58** additional wikis

# Progress Fixing Broken Links

- Fixed more than **19 million** formerly broken links across 305 Wikipedia languages

- More than **8 million** fixed in the last 2 years alone!

**404**
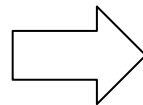
# Turning Book **Citations** Blue

# Better World Books

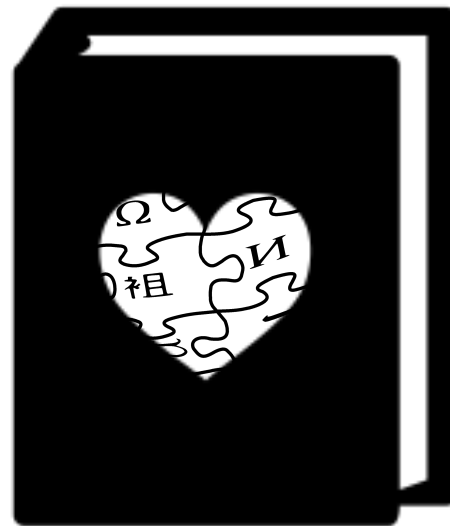The Internet Archive creates a wish list of books based on Wikipedia citations

# Progress Adding Links to Books

More than **1.8 million links** to books from [archive.org](archive.org) across have been added to ~**60 Wikipedia language** editions

(1 million added by our software and 800,000 added by Wikipedia Editors)
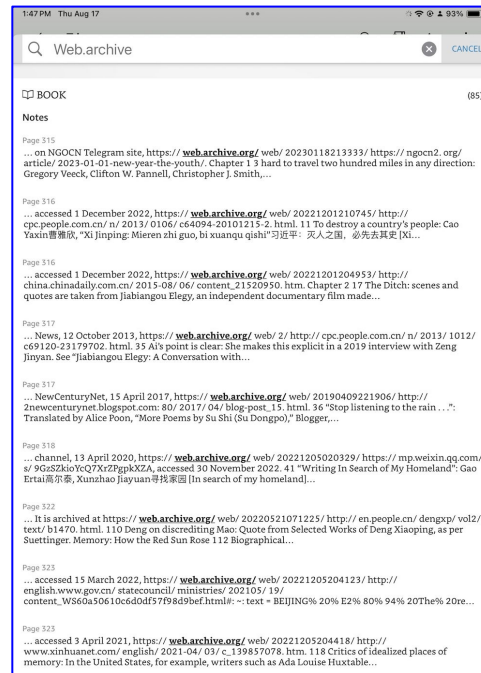
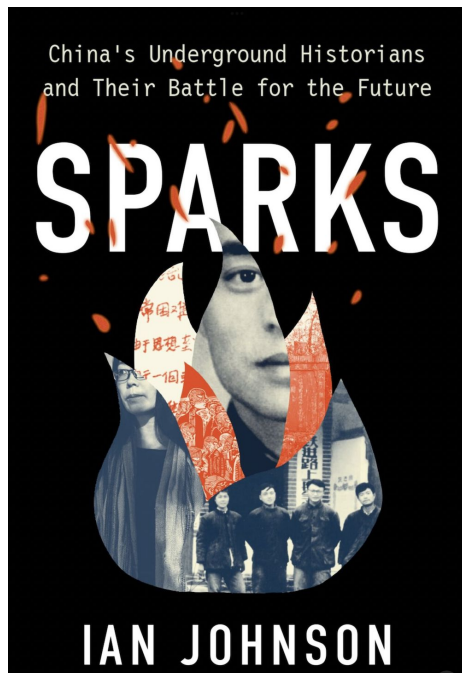# Progress Adding Links to Papers

Our software has added more than **141,000** links to academic papers and articles across ~**60 Wikipedia language** editions

# 85 Wayback Machine Links in this book

INTERNET ARCHIVE

## SPARKS

China's Underground Historians and Their Battle for the Future

**IAN JOHNSON**

---

1:47 PM  Thu Aug 17 · · · · · · 93%

Web.archive ⊗ CANCEL

📖 BOOK (85)

**Notes**

Page 315

… on NGOCN Telegram site, https:// **web.archive.org**/ web/ 20230118213333/ https:// ngocn2. org/ article/ 2023-01-01-new-year-the-youth/. Chapter 1 3 hard to travel two hundred miles in any direction: Gregory Veeck, Clifton W. Pannell, Christopher J. Smith,…

Page 316

… accessed 1 December 2022, https:// **web.archive.org**/ web/ 20221201210745/ http:// cpc.people.com.cn/ n/ 2013/ 0106/ c64094-20101215-2. html. 11 To destroy a country's people: Cao Yaxin曹雅欣, "Xi Jinping: Mieren zhi guo, bi xuanqu qishi"习近平: 灭人之国，必先去其史 [Xi…

Page 316

… accessed 1 December 2022, https:// **web.archive.org**/ web/ 20221201204953/ http:// china.chinadaily.com.cn/ 2015-08/ 06/ content_21520950. htm. Chapter 2 17 The Ditch: scenes and quotes are taken from Jiabiangou Elegy, an independent documentary film made…

Page 317

… News, 12 October 2013, https:// **web.archive.org**/ web/ 2/ http:// cpc.people.com.cn/ n/ 2013/ 1012/ c69120-23179702. html. 35 Ai's point is clear: She makes this explicit in a 2019 interview with Zeng Jinyan. See "Jiabiangou Elegy: A Conversation with…

Page 317

… NewCenturyNet, 15 April 2017, https:// **web.archive.org**/ web/ 20190409221906/ http:// 2newcenturynet.blogspot.com: 80/ 2017/ 04/ blog-post_15. html. 36 "Stop listening to the rain . . .": Translated by Alice Poon, "More Poems by Su Shi (Su Dongpo)," Blogger,…

Page 318

… channel, 13 April 2020, https:// **web.archive.org**/ web/ 20221205020329/ https:// mp.weixin.qq.com/ s/ 9GzSZkioYcQ7XrZPgpkXZA, accessed 30 November 2022. 41 "Writing In Search of My Homeland": Gao Ertai高尔泰, Xunzhao Jiayuan寻找家园 [In search of my homeland]…

Page 322

… It is archived at https:// **web.archive.org**/ web/ 20220521071225/ http:// en.people.com.cn/ dengxp/ vol2/ text/ b1470. html. 110 Deng on discrediting Mao: Quote from Selected Works of Deng Xiaoping, as per Suettinger. Memory: How the Red Sun Rose 112 Biographical…

Page 323

… accessed 15 March 2022, https:// **web.archive.org**/ web/ 20221205204123/ http:// english.www.gov.cn/ statecouncil/ ministries/ 202105/ 19/ content_WS60a50610c6d0df57f98d9bef.html#: ~: text = BEIJING% 20% E2% 80% 94% 20The% 20re…

Page 323

… accessed 3 April 2021, https:// **web.archive.org**/ web/ 20221205204418/ http:// www.xinhuanet.com/ english/ 2021-04/ 03/ c_139857078. htm. 118 Critics of idealized places of memory: In the United States, for example, writers such as Ada Louise Huxtable…

# In the Lab: Reference Explorer

# Next Up: Finding/Fixing "Soft 404s"

# Publishers Lawsuit (blog.archive.org)