

# FastWARC: Optimizing Large-Scale Web Archive Analytics

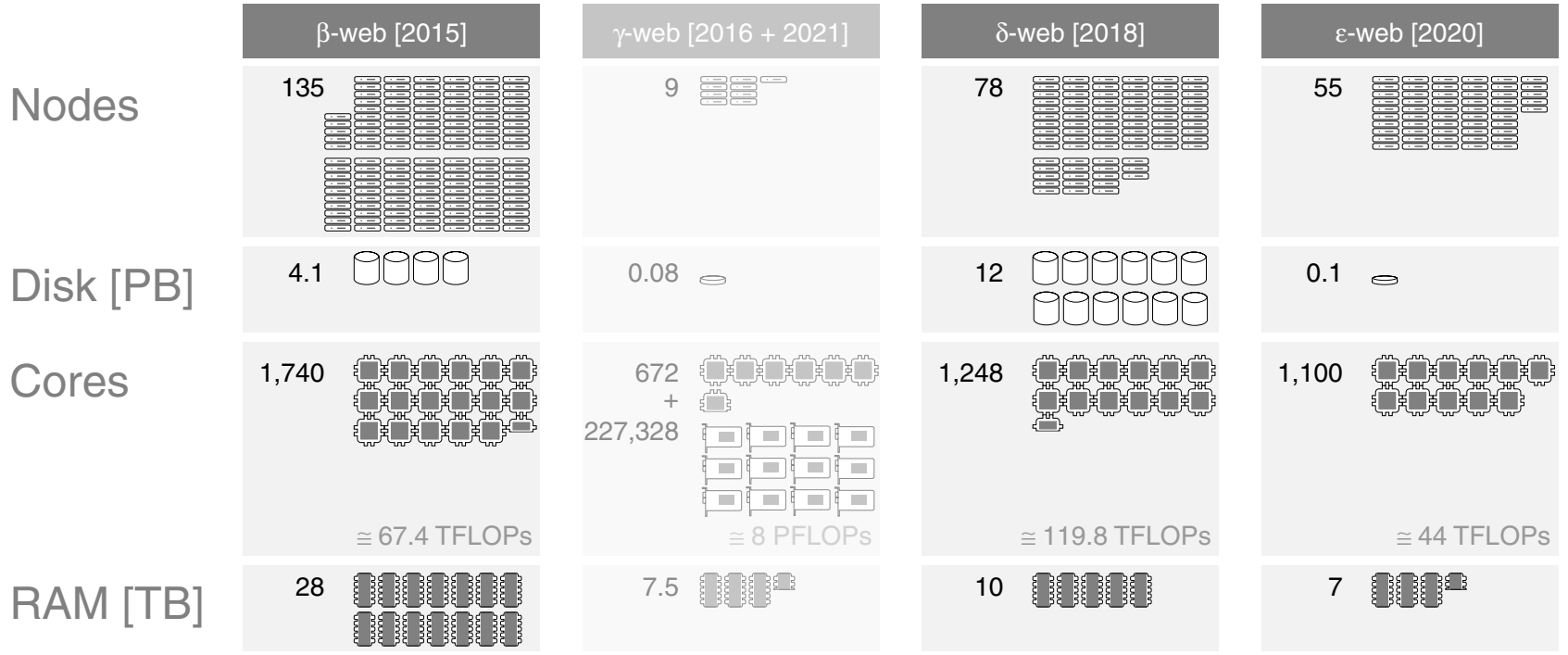
---

Janek Bevendorff, Martin Potthast, Benno Stein

Webis Group  
Bauhaus-Universität Weimar, Leipzig University  
[webis.de](http://webis.de) | [resiliparse.chatnoir.eu](http://resiliparse.chatnoir.eu)

OSSYM 2021 – 3rd International Open Search Symposium, October 12, 2021

# Webis BigData Cluster Infrastructure



# Webis Web Archive Analytics Stack

	Task stack	Technology stack	Vendor stack
<b>Data Consumption Layer</b>	<ul style="list-style-type: none"> <li>· Query and explore</li> <li>· Visualize and interact</li> <li>· Explain and justify</li> </ul>	<ul style="list-style-type: none"> <li>· Visual analytics</li> <li>· Immersive technologies</li> <li>· Intelligent agents</li> </ul>	
<b>Data Analytics Layer</b>	<ul style="list-style-type: none"> <li>· Diagnose and reason</li> <li>· Structure identification</li> <li>· Structure verification</li> </ul>	<ul style="list-style-type: none"> <li>· Distributed learning</li> <li>· State-space search</li> <li>· Symbolic inference</li> </ul>	
<b>Data Management + Hardware Layer</b>	<ul style="list-style-type: none"> <li>· Provenance tracking</li> <li>· Normalization</li> <li>· Cleansing</li> <li>· Monitoring</li> <li>· Replication</li> </ul>	<ul style="list-style-type: none"> <li>· Key-value store</li> <li>· RDF triple store</li> <li>· Graph store</li> <li>· Object store</li> <li>· Orchestration</li> <li>· Parallelization</li> <li>· Virtualization</li> </ul>	
<b>Data Acquisition Layer</b>	<ul style="list-style-type: none"> <li>· Replay</li> <li>· Collect</li> <li>· Log</li> </ul>	<ul style="list-style-type: none"> <li>· Distant supervision</li> <li>· Crowdsourcing</li> <li>· Crawling and archiving</li> </ul>	

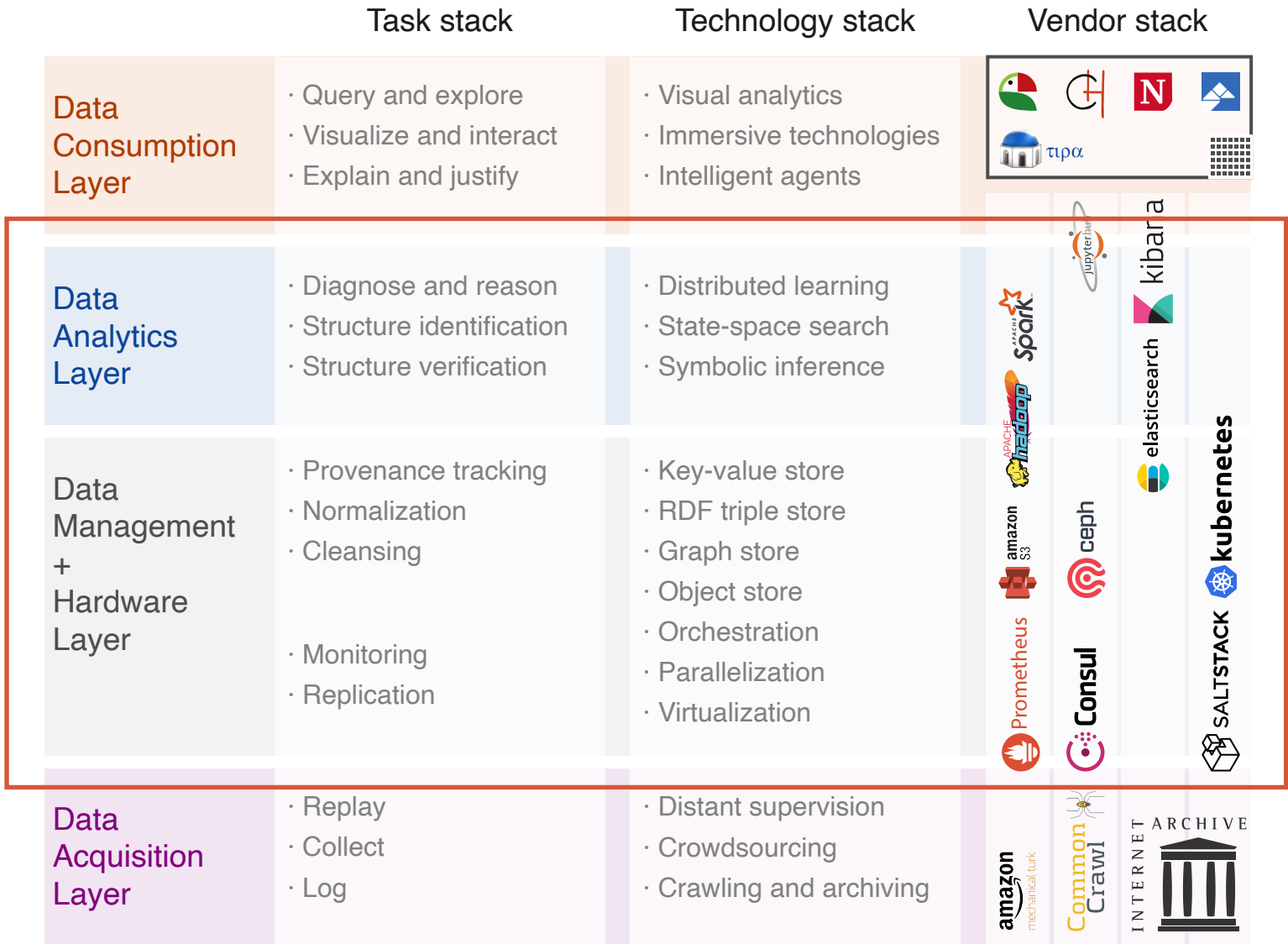
# Webis Web Archive Analytics Stack



## The ChatNoir Web Search Engine

- ❑ Indexes Clueweb09, Clueweb12, Common Crawl 11/2015, 04/2017
- ❑ Approx. 5 billion documents
- ❑ 87 TiB of index data
- ❑ 1 920 shards
- ❑ Service provided free of charge to the IR community.

# Webis Web Archive Analytics Stack



# Webis Web Archive Analytics Stack

	Task stack	Technology stack	Vendor stack
<b>Data Consumption Layer</b>	<ul style="list-style-type: none"> <li>· Query and explore</li> <li>· Visualize and interact</li> <li>· Explain and justify</li> </ul>	<ul style="list-style-type: none"> <li>· Visual analytics</li> <li>· Immersive technologies</li> <li>· Intelligent agents</li> </ul>	
<b>Data Analytics Layer</b>	<ul style="list-style-type: none"> <li>· Diagnose and reason</li> <li>· Structure identification</li> <li>· Structure verification</li> </ul>	<ul style="list-style-type: none"> <li>· Distributed learning</li> <li>· State-space search</li> <li>· Symbolic inference</li> </ul>	
<b>Data Management + Hardware Layer</b>	<ul style="list-style-type: none"> <li>· Provenance tracking</li> <li>· Normalization</li> <li>· Cleansing</li> <li>· Monitoring</li> <li>· Replication</li> </ul>	<ul style="list-style-type: none"> <li>· Key-value store</li> <li>· RDF triple store</li> <li>· Graph store</li> <li>· Object store</li> <li>· Orchestration</li> <li>· Parallelization</li> <li>· Virtualization</li> </ul>	
<b>Data Acquisition Layer</b>	<ul style="list-style-type: none"> <li>· Replay</li> <li>· Collect</li> <li>· Log</li> </ul>	<ul style="list-style-type: none"> <li>· Distant supervision</li> <li>· Crowdsourcing</li> <li>· Crawling and archiving</li> </ul>	

# WARC Basics

WARC/1.0

WARC-Type: response

WARC-Date: 2021-10-08T17:27:41Z

WARC-Record-ID: <urn:uuid:25a91250-dffe-4cab-a56b-2ffa07c21ec7>

Content-Length: 36578

Content-Type: application/http; msgtype=response

WARC-Warcinfo-ID: <urn:uuid:dba95eaf-b2ef-4453-82b8-21d106ec3463>

WARC-Concurrent-To: <urn:uuid:3ea5bff2-8f7f-4d80-9c11-320393c638ad>

WARC-Target-URI: https://example.com

WARC-Payload-Digest: sha1:CUJAG44UJDP64LTXHBIFKFXD54UHBP4S

WARC-Block-Digest: sha1:QXLZPJ4UUKOBZFWZWSAGHPHBT7X7UD2Q

WARC-Identified-Payload-Type: text/html

HTTP/1.1 200 OK

Content-Type: text/html

X-Crawler-Content-Encoding: gzip

Vary: Accept-Encoding

Server: Microsoft-IIS/8.5

X-Powered-By: PHP/5.2.17

...

# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.

<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>



# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.
- ❑ That's **345 hours** of compute! (5.5 hours / TiB)

<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>

# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.
- ❑ That's **345 hours** of compute! (5.5 hours / TiB)
- ❑ Time per PiB: **5 670 hours**

<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>

# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.
- ❑ That's **345 hours** of compute! (5.5 hours / TiB)
- ❑ Time per PiB: **5 670 hours**

Can we speed this up?



<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>

# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.
- ❑ That's **345 hours** of compute! (5.5 hours / TiB)
- ❑ Time per PiB: **5 670 hours**

Can we speed this up?

A reduction to only **12 s / file** (speedup: 1.6x) would save. . .



<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>

# Sample Calculation

- ❑ A somewhat recent Common Crawl<sup>1</sup>:
  - 64 000** WARC files, totalling **62.5 TiB** GZip-compressed
- ❑ WARCIO<sup>2</sup> needs approx. **19.5 s / file**.
- ❑ That's **345 hours** of compute! (5.5 hours / TiB)
- ❑ Time per PiB: **5 670 hours**

## Can we speed this up?

A reduction to only **12 s / file** (speedup: 1.6x) would save...

- ❑ **132 hours** on a Common Crawl,
- ❑ **2 157 hours per PiB!**



<sup>1</sup><https://commoncrawl.org/>

<sup>2</sup><https://github.com/webrecorder/warcio>

# Well... We can!

We introduce: **FastWARC**, a new WARC parsing library:

- ❑ Written from the ground up in optimized C/C++ using Cython
- ❑ API familiar to WARCIO users
- ❑ Compatible with Python 3.7+
- ❑ Pre-built binaries for Linux, macOS, Windows on PyPi
- ❑ Open Source under the Apache 2.0 license

```
from fastwarc.warc import ArchiveIterator
```

```
for record in ArchiveIterator(open('warcfile.warc.gz', 'rb')):  
    pass
```

# FastWARC Benchmarks (Then)

<b>Comp.</b>	<b>Parser</b>	<b>Records/s</b>	<b>Speedup</b>
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	13 971.5	–
None	FastWARC	64 698.0	4.6
None	WARCIO+HTTP	13 570.2	–
None	FastWARC+HTTP	58 354.0	4.3
None	WARCIO+HTTP+Checksum	7 890.9	–
None	FastWARC+HTTP+Checksum	11 528.6	1.5
GZip	WARCIO	5 898.8	–
GZip	FastWARC	8 899.1	1.5
GZip	WARCIO+HTTP	5 986.1	–
GZip	FastWARC+HTTP	8 659.0	1.4
GZip	WARCIO+HTTP+Checksum	4 544.7	–
GZip	FastWARC+HTTP+Checksum	5 022.6	1.1
LZ4	FastWARC	36 862.8	6.2*
LZ4	FastWARC+HTTP	36 327.9	6.2*
LZ4	FastWARC+HTTP+Checksum	10 110.0	2.2*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 865.7	–
None	FastWARC	29 307.7	3.7
GZip	WARCIO	3 438.4	–
GZip	FastWARC	4 583.3	1.3
LZ4	FastWARC	18 337.0	5.3*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.

# FastWARC Benchmarks (Now)

<b>Comp.</b>	<b>Parser</b>	<b>Records/s</b>	<b>Speedup</b>
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	16 945.5	—
None	FastWARC	108 488.0	6.4 (4.6)
None	WARCIO+HTTP	11 661.6	—
None	FastWARC+HTTP	79 297.0	6.8 (4.3)
None	WARCIO+HTTP+Checksum	6 986.7	—
None	FastWARC+HTTP+Checksum	21 320.9	3.1 (1.5)
GZip	WARCIO	6 460.1	—
GZip	FastWARC	10 413.4	1.6 (1.5)
GZip	WARCIO+HTTP	5 435.6	—
GZip	FastWARC+HTTP	10 101.5	1.9 (1.4)
GZip	WARCIO+HTTP+Checksum	4 121.6	—
GZip	FastWARC+HTTP+Checksum	7 433.0	1.8 (1.1)
LZ4	FastWARC	49 825.4	7.7 (6.2)*
LZ4	FastWARC+HTTP	42 394.5	7.8 (6.2)*
LZ4	FastWARC+HTTP+Checksum	16 992.2	4.1 (2.2)*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 969.1	—
None	FastWARC	49 396.5	6.2 (3.7)
GZip	WARCIO	3 555.7	—
GZip	FastWARC	6 335.1	1.8 (1.3)
LZ4	FastWARC	28 313.8	8.0 (5.5)*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.



# FastWARC Benchmarks (Now)

Comp.	Parser	Records/s	Speedup
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	16 945.5	—
None	FastWARC	108 488.0	6.4 (4.6)
None	WARCIO+HTTP	11 661.6	—
None	FastWARC+HTTP	79 297.0	6.8 (4.3)
None	WARCIO+HTTP+Checksum	6 986.7	—
None	FastWARC+HTTP+Checksum	21 320.9	3.1 (1.5)
GZip	WARCIO	6 460.1	—
GZip	FastWARC	10 413.4	1.6 (1.5)
GZip	WARCIO+HTTP	5 435.6	—
GZip	FastWARC+HTTP	10 101.5	1.9 (1.4)
GZip	WARCIO+HTTP+Checksum	4 121.6	—
GZip	FastWARC+HTTP+Checksum	7 433.0	1.8 (1.1)
LZ4	FastWARC	49 825.4	7.7 (6.2)*
LZ4	FastWARC+HTTP	42 394.5	7.8 (6.2)*
LZ4	FastWARC+HTTP+Checksum	16 992.2	4.1 (2.2)*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 969.1	—
None	FastWARC	49 396.5	6.2 (3.7)
GZip	WARCIO	3 555.7	—
GZip	FastWARC	6 335.1	1.8 (1.3)
LZ4	FastWARC	28 313.8	8.0 (5.5)*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.

# FastWARC Benchmarks (Now)

Comp.	Parser	Records/s	Speedup
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	16 945.5	—
None	FastWARC	108 488.0	6.4 (4.6)
None	WARCIO+HTTP	11 661.6	—
None	FastWARC+HTTP	79 297.0	6.8 (4.3)
None	WARCIO+HTTP+Checksum	6 986.7	—
None	FastWARC+HTTP+Checksum	21 320.9	3.1 (1.5)
GZip	WARCIO	6 460.1	—
GZip	FastWARC	10 413.4	1.6 (1.5)
GZip	WARCIO+HTTP	5 435.6	—
GZip	FastWARC+HTTP	10 101.5	1.9 (1.4)
GZip	WARCIO+HTTP+Checksum	4 121.6	—
GZip	FastWARC+HTTP+Checksum	7 433.0	1.8 (1.1)
LZ4	FastWARC	49 825.4	7.7 (6.2)*
LZ4	FastWARC+HTTP	42 394.5	7.8 (6.2)*
LZ4	FastWARC+HTTP+Checksum	16 992.2	4.1 (2.2)*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 969.1	—
None	FastWARC	49 396.5	6.2 (3.7)
GZip	WARCIO	3 555.7	—
GZip	FastWARC	6 335.1	1.8 (1.3)
LZ4	FastWARC	28 313.8	8.0 (5.5)*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.

# FastWARC Benchmarks (Now)

Comp.	Parser	Records/s	Speedup
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	16 945.5	—
None	FastWARC	108 488.0	6.4 (4.6)
None	WARCIO+HTTP	11 661.6	—
None	FastWARC+HTTP	79 297.0	6.8 (4.3)
None	WARCIO+HTTP+Checksum	6 986.7	—
None	FastWARC+HTTP+Checksum	21 320.9	3.1 (1.5)
GZip	WARCIO	6 460.1	—
GZip	FastWARC	10 413.4	1.6 (1.5)
GZip	WARCIO+HTTP	5 435.6	—
GZip	FastWARC+HTTP	10 101.5	1.9 (1.4)
GZip	WARCIO+HTTP+Checksum	4 121.6	—
GZip	FastWARC+HTTP+Checksum	7 433.0	1.8 (1.1)
LZ4	FastWARC	49 825.4	7.7 (6.2)*
LZ4	FastWARC+HTTP	42 394.5	7.8 (6.2)*
LZ4	FastWARC+HTTP+Checksum	16 992.2	4.1 (2.2)*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 969.1	—
None	FastWARC	49 396.5	6.2 (3.7)
GZip	WARCIO	3 555.7	—
GZip	FastWARC	6 335.1	1.8 (1.3)
LZ4	FastWARC	28 313.8	8.0 (5.5)*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.





# FastWARC Benchmarks (Now)

Comp.	Parser	Records/s	Speedup
<i>AMD Ryzen Threadripper 2920X (NVMe SSD)</i>			
None	WARCIO	16 945.5	—
None	FastWARC	108 488.0	6.4 (4.6)
None	WARCIO+HTTP	11 661.6	—
None	FastWARC+HTTP	79 297.0	6.8 (4.3)
None	WARCIO+HTTP+Checksum	6 986.7	—
None	FastWARC+HTTP+Checksum	21 320.9	3.1 (1.5)
GZip	WARCIO	6 460.1	—
GZip	FastWARC	10 413.4	1.6 (1.5)
GZip	WARCIO+HTTP	5 435.6	—
GZip	FastWARC+HTTP	10 101.5	1.9 (1.4)
GZip	WARCIO+HTTP+Checksum	4 121.6	—
GZip	FastWARC+HTTP+Checksum	7 433.0	1.8 (1.1)
LZ4	FastWARC	49 825.4	7.7 (6.2)*
LZ4	FastWARC+HTTP	42 394.5	7.8 (6.2)*
LZ4	FastWARC+HTTP+Checksum	16 992.2	4.1 (2.2)*
<i>Intel(R) Xeon(R) CPU E5-2620 v2 (remote Ceph storage)</i>			
None	WARCIO	7 969.1	—
None	FastWARC	49 396.5	6.2 (3.7)
GZip	WARCIO	3 555.7	—
GZip	FastWARC	6 335.1	1.8 (1.3)
LZ4	FastWARC	28 313.8	8.0 (5.5)*

\*LZ4 speedup is over WARCIO with GZip, since WARCIO does not support LZ4.





# ChatNoir Resiliparse

**FastWARC** is part of the **ChatNoir Resiliparse** web archive analytics library:

-  ChatNoir: **<https://chatnoir.eu>**
-  GitHub: **<https://git.io/resiliparse>**
-  (Extensive) documentation: **<https://resiliparse.chatnoir.eu>**
-  PyPi package: `pip install fastwarc`

# ChatNoir Resiliparse

**FastWARC** is part of the **ChatNoir Resiliparse** web archive analytics library:

-  ChatNoir: <https://chatnoir.eu>
-  GitHub: <https://git.io/resiliparse>
-  (Extensive) documentation: <https://resiliparse.chatnoir.eu>
-  PyPi package: `pip install fastwarc`

Thank you!