# WASP: Web Archiving and Search Personalized
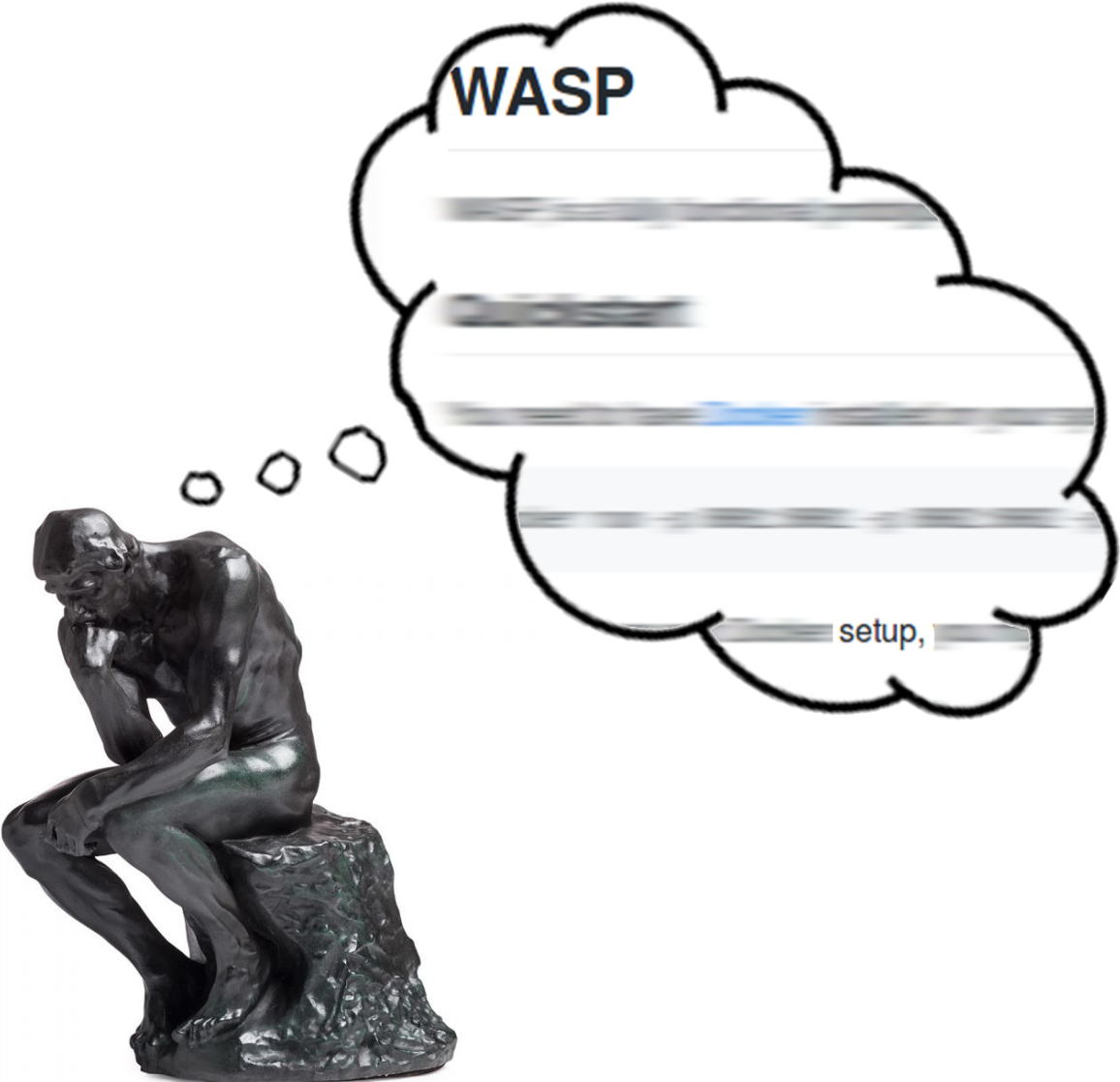
**Johannes Kiesel**, Arjen P. de Vries, Matthias Hagen, Benno Stein and Martin Potthast
@KieselJohannes, @arjenpdevries, @matthias_hagen, @bennostein, @martinpotthast

DESIRES, August 29[th] 2018

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation

🕐 History ⌄                    ✕
─────────────────────────────────

🔍 Search history          View⌄

▸ 🕐 Today
▸ 🕐 Yesterday
▸ 🕐 Last 7 days
▸ 🕐 This month
▸ 🕐 July
▸ 🕐 June
▸ 🕐 May
▸ 🕐 April
▸ 🕐 March
▸ 🕐 Older than 6 months

# The Personal Search Engine: Motivation



🕐 History ⌄                              ✕

| 🔍 Search history          | View ⌄ |

▸ 🕐 Today
▸ 🕐 Yesterday
▸ 🕐 Last 7 days
▸ 🕐 This month
▸ 🕐 July
▸ 🕐 June
▸ 🕐 May
▸ 🕐 April
▸ 🕐 March
▸ 🕐 Older than 6 months

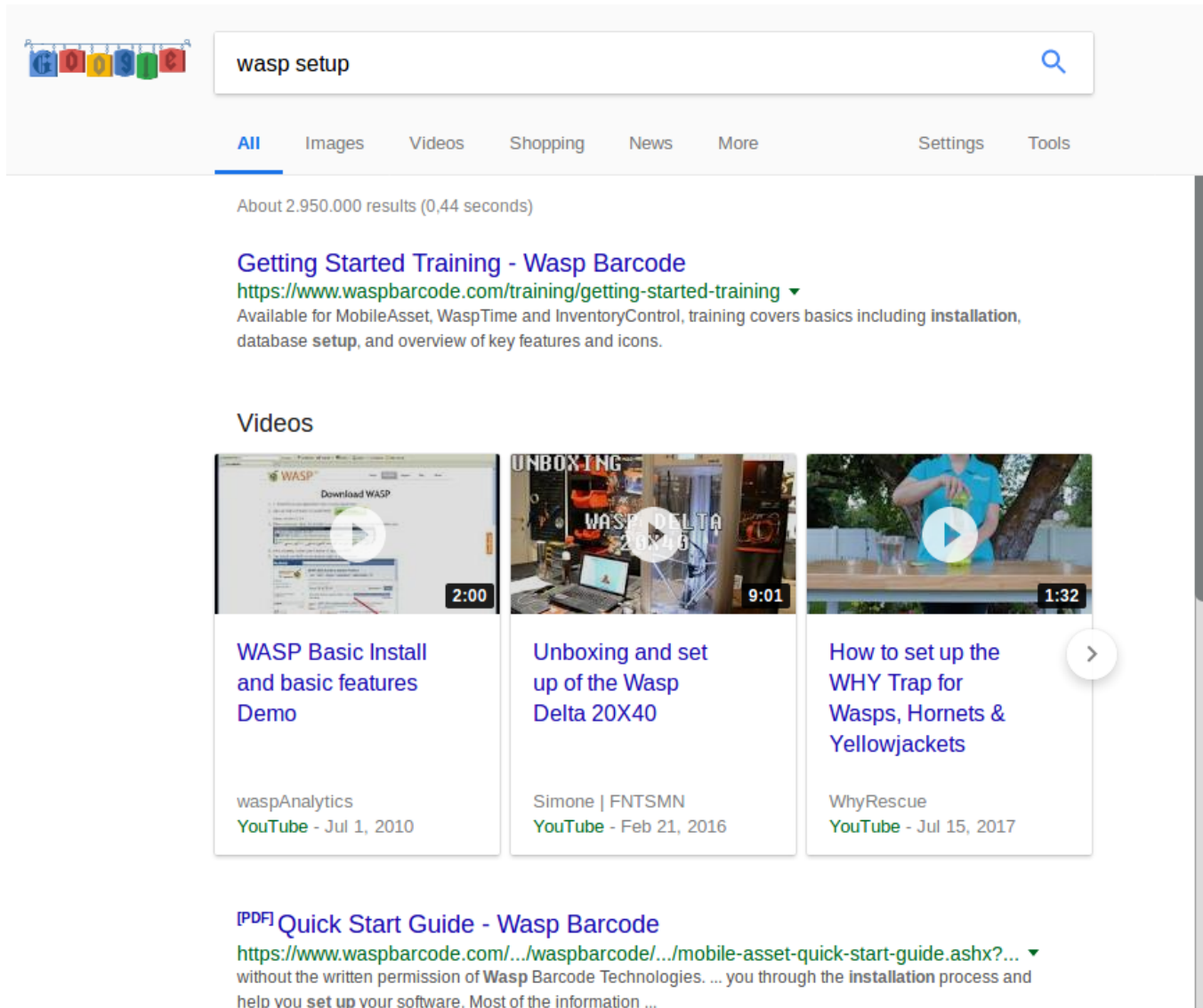🕐 History ⌄                              ✕

| 🔍 wasp              ⌫ | View ⌄ |

G webis wasp – Google Search
G wasp – Google Search
G wasp webis – Google-Suche
🐦 Johannes Kiesel auf Twitter:…
○ webis-de/wasp
G Google – My Activity
▦ WASP: Web Archiving and Sear…
G Google Advanced Search
G Search Settings
○ wasp/src/de/webis/warc at ma…
○ wasp/src at master · webis-d…
○ wasp/pom.xml at master · web…
○ wasp/Index.java at master · …
○ wasp/WarcIndexer.java at mas…
○ wasp/src/de/webis/warc/index…
🐳 webis/wasp – Docker Hub
   WASP: Web Archiving and Sear…
○ wasp/Warcs.java at master · …
G amber wasp – Google Search
▦ Wasp Insect Sting · Free ima…

# The Personal Search Engine: Motivation



🕐 History ⌄                    ✕

🔍 Search history       View⌄

▸ 🕐 Today
▸ 🕐 Yesterday
▸ 🕐 Last 7 days
▸ 🕐 This month
▸ 🕐 July
▸ 🕐 June
▸ 🕐 May
▸ 🕐 April
▸ 🕐 March
▸ 🕐 Older than 6 months

🕐 History ⌄                    ✕

🔍 wasp         ⌫    View⌄

G webis wasp – Google Search
G wasp – Google Search
G wasp webis – Google-Suche
🐦 Johannes Kiesel auf Twitter:…
⬤ webis-de/wasp
G Google – My Activity
▦ WASP: Web Archiving and Sear…
G Google Advanced Search
G Search Settings
⬤ wasp/src/de/webis/warc at ma…
⬤ wasp/src at master · webis-d…
⬤ wasp/pom.xml at master · web…
⬤ wasp/Index.java at master · …
⬤ wasp/WarcIndexer.java at mas…
⬤ wasp/src/de/webis/warc/index…
🐳 webis/wasp – Docker Hub
  WASP: Web Archiving and Sear…
⬤ wasp/Warcs.java at master · …
G amber wasp – Google Search
▦ Wasp Insect Sting - Free ima…

🕐 History ⌄                    ✕

🔍 wasp setup      ⌫    View⌄

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation



GOOGLE

**wasp setup visited:last-week**

All   Images   Videos   Shopping   News   More        Settings   Tools

About 1 results (0,44 seconds)

GitHub - webis-de/wasp
https://github.com/webis-de/wasp ▾
**WASP. WASP**, is a fully functional prototype of a personal web archive and search system. Quickstart. You need to have Docker installed on your system.

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation

# The Personal Search Engine: Motivation

# The Personal Search Engine: Inspiration

# The Personal Search Engine: Inspiration



Personal
search engine!

# WASP



## WASP

WASP, is a fully functional prototype of a personal web archive and search system.

## Quickstart
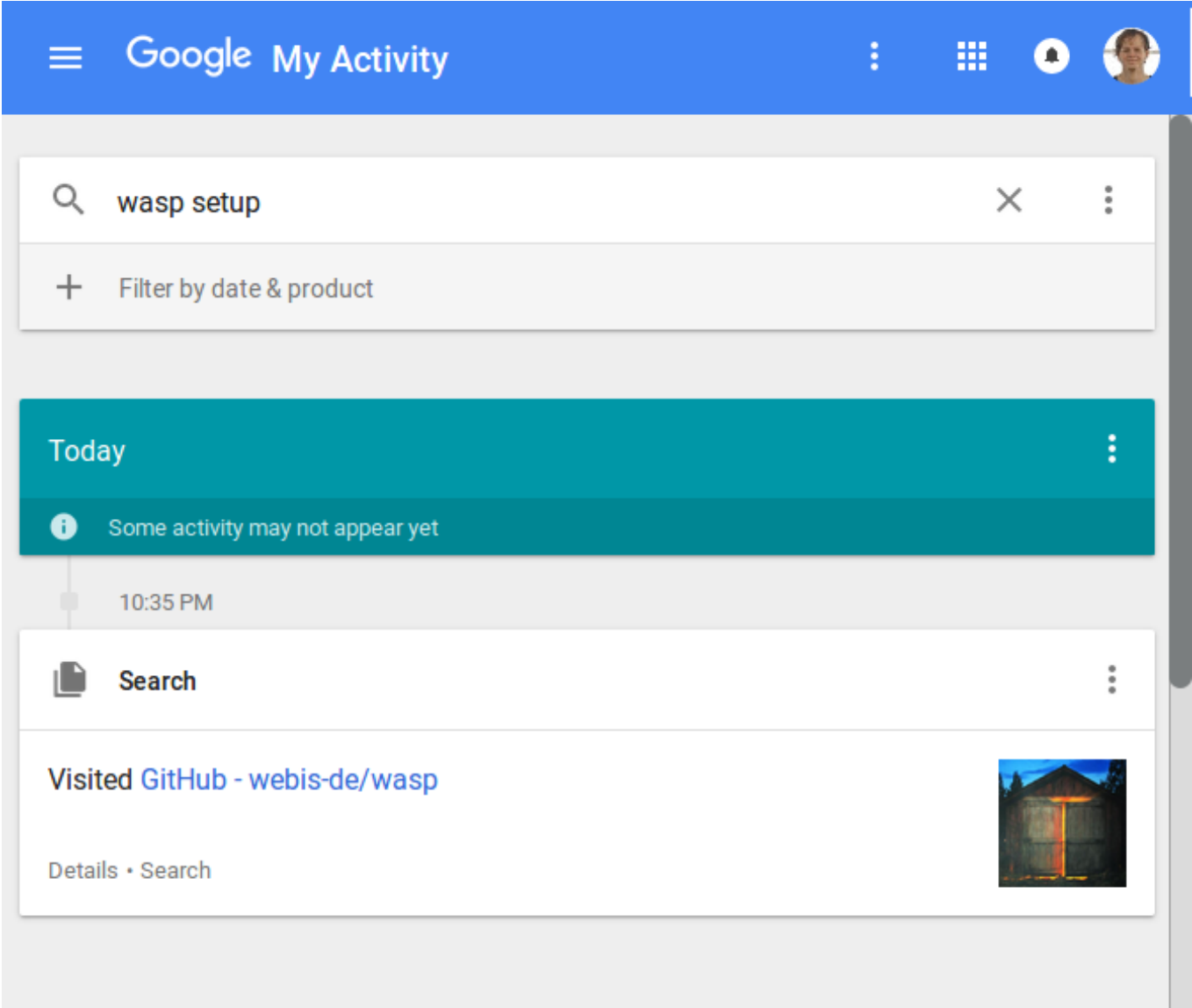
You need to have Docker installed on your system. You can then immediately start WASP like this:

```
docker run -p 8001:8001 -p 8002:8002 -p 8003:8003 --name wasp -d webis/wasp:0.2.0
```

Depending on your Docker setup, you might need to run this command as an administrator or add `sudo` in front of it.

After a few seconds, you should already see the search interface in your browser at (http://localhost:8003/search)

Next, you will have to change the proxy settings of your browser to route the requests and responses through it. How this can be done depends on the browser (Chrome, Edge, Firefox, Opera, Safari). However, also consider using a tailored proxy switching plugin for your browser, which allows you to disable/enable WASP easily. In all cases, you need to specify (localhost:8001) as proxy for HTTP and HTTPS and have "localhost" as an exception (often so by default).

Now it is time to clear your browser's cache (or use incognito mode). Otherwise you will be disappointed when some elements are missing in the archive.

# WASP

# WASP

# WASP

# WASP

---

localhost:8002/archive/20180826205834/https://githu

📖 **README.md**

# WASP

WASP, is a fully functional prototype of a personal web archive and search system.

## Quickstart

You need to have Docker installed on your system. You can then immediately start WASP like this:

```
docker run -p 8001:8001 -p 8002:8002 -p 8003:8003 --name wasp -d webis/wasp:0.2.0
```

Depending on your Docker setup, you might need to run this command as an administrator or add `sudo` in front of it.

After a few seconds, you should already see the search interface in your browser at (http://localhost:8003/search)

Next, you will have to change the proxy settings of your browser to route the requests and responses through it. How this can be done depends on the browser (Chrome, Edge, Firefox, Opera, Safari). However, also consider using a tailored proxy switching plugin for your browser, which allows you to disable/enable WASP easily. In all cases, you need to specify (localhost:8001) as proxy for HTTP and HTTPS and have "localhost" as an exception (often so by default).

Now it is time to clear your browser's cache (or use incognito mode). Otherwise you will be disappointed when some elements are missing in the archive.

Archived Version

# WASP



Browser

Search Interface

Index

pywb

WARCs

warcprox

World Wide Web

# WASP



❏ All requests (→) and responses (←) while browsing are stored and indexed

# WASP



- ❑ All requests (➤) and responses (◄) while browsing are stored and indexed

- ❑ Page on localhost allows to search. Result page links to archive where...

# WASP



- ❑ All requests (→) and responses (←) while browsing are stored and indexed

- ❑ Page on localhost allows to search. Result page links to archive where...

- ❑ Visited pages are reproduced for the corresponding time

# WASP



Personal
search engine!

# WASP

# Insight 1: Not Indexing Near-duplicates

Page 1 for `wasp setup` from 2018-08-19 23:43 until now

GitHub - webis-de/wasp [archive] [live]
https://github.com/webis-de/wasp                                              @2018-08-26 22:56
WASP, is a fully functional prototype of a personal web archive and search system. ... -name **wasp** -d webis/**wasp** 0.2.0 Depending on your Docker **setup**, you might need to run this command ... You are already set up to use **WASP** for HTTP connections. ... For HTTPS connections, you have to trust **WASP** to identify foreign web pages for you. ... **wasp**.

wasp/README.md at master · webis-de/wasp · GitHub [archive] [live]
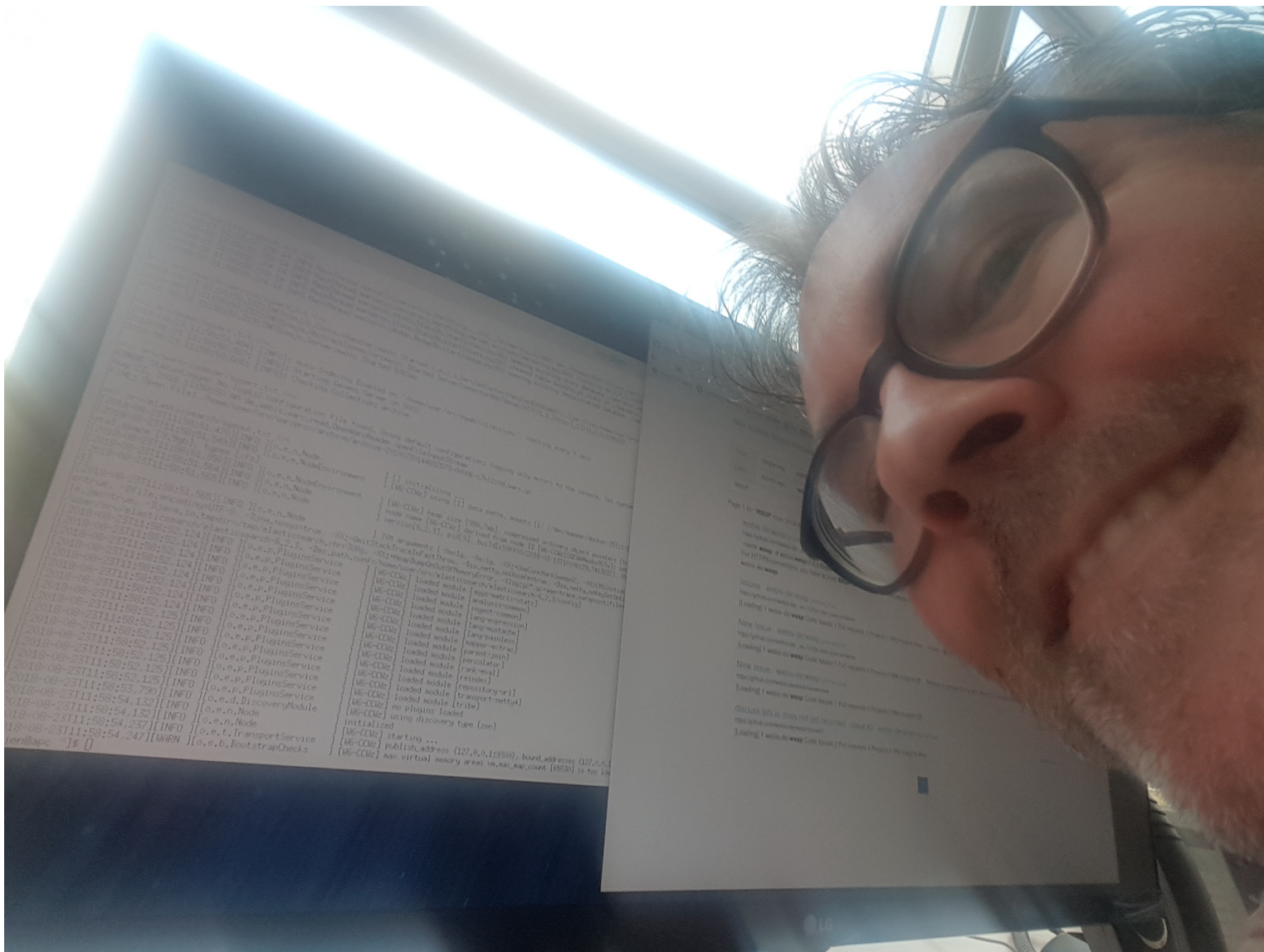https://github.com/webis-de/...ax=%23js-repo-pjax-container                   @2018-08-26 22:48
-name **wasp** -d webis/**wasp** 0.2.0 Depending on your Docker **setup**, you might need to run this command ... You are already set up to use **WASP** for HTTP connections. ... **wasp**. ... Docker A pre-build Docker image of **WASP** is available on dockerhub. ... **wasp**/README.md at master · webis-de/**wasp** · GitHub

wasp/README.md at master · webis-de/wasp · GitHub [archive] [live]
https://github.com/webis-de/...ax=%23js-repo-pjax-container                   @2018-08-26 23:00
**WASP**, is a fully functional prototype of a personal web archive and search system. ... -name **wasp** -d webis/**wasp** 0.2.0 Depending on your Docker **setup**, you might need to run this command ... You are already set up to use **WASP** for HTTP connections. ... For HTTPS connections, you have to trust **WASP** to identify foreign web pages for you. ... **wasp**.

❏ What changes warrant a re-archiving?

# Insight 2: Browsable and Deletable History

# Insight 3: Easy (De-)activation of Archiving

Easy activation, deactivation, and status-check



## Patterns

Add patterns to prevent this proxy being used for localhost and intranet/private IP addresses **Help**      Add

### White Patterns

| Name | Pattern | Type | http(s) | On/Off | |
|------|---------|------|---------|--------|--|
| all URLs | * | wildcard | both | on | 🗑 |

### Black Patterns

| Name | Pattern | Type | http(s) | On/Off | |
|------|---------|------|---------|--------|--|
| local hostnames (usually ... | ^(?:[^:@/]+(?::[^@/]+)?@)... | reg exp | both | on | 🗑 |
| local subnets (IANA reser... | ^(?:[^:@/]+(?::[^@/]+)?@)... | reg exp | both | on | 🗑 |
| localhost - matches the l... | ^(?:[^:@/]+(?::[^@/]+)?@)... | reg exp | both | on | 🗑 |

# Insight 4: Combined Indexing of Sub-pages

zuwandte, als Mitglied des Streitkräfteausschusses um die Welt reiste und auf
Veranstaltungen wie der Münchner Sicherheitskonferenz quasi zum Inventar
gehörte, zeigt die Blickrichtung: weg von den Mühen der
innenpolitischen Ebene.
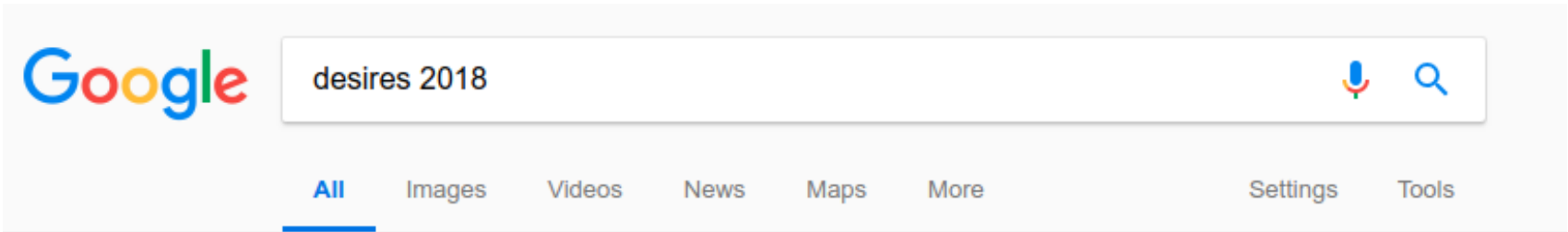
**Seite 1   Er war der Letzte seiner Art**
Seite 2   Warum es nicht zum US-Präsidenten reichte

nächste Seite ›

❏  Should visited sub-pages of a single article be indexed as one?

❏  If so, to which sub-page should be linked in the result list?

# Insight 5: Personalized Search



Google

desires 2018

All    Images    Videos    News    Maps    More         Settings    Tools

About 145.000.000 results (0,41 seconds)

## DESIRES 2018
desires.dei.unipd.it/ ▾
Vision. A systems-oriented biennial conference, complementary in its mission to the mainstream
Information Access and Retrieval conferences, emphasizing the ...
You've visited this page many times. Last visit: 8/23/18

# Insight 5: Personalized Search



**Google**

desires 2018

All   Images   Videos   News   Maps   More     Settings   Tools

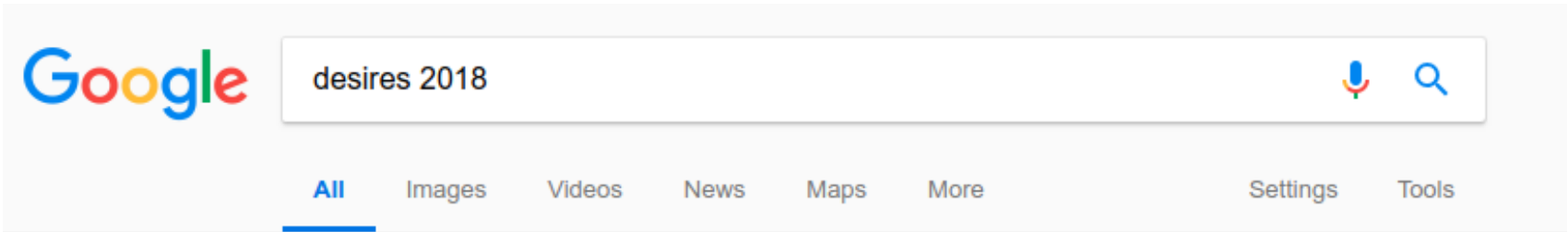About 145.000.000 results (0,41 seconds)

**DESIRES 2018**
desires.dei.unipd.it/ ▼
Vision. A systems-oriented biennial conference, complementary in its mission to the mainstream
Information Access and Retrieval conferences, emphasizing the ...
You've visited this page many times. Last visit: 8/23/18

## You've visited this page many times. Last visit: 8/23/18

# WASP: Web Archiving and Search Personalized

Insights overview

- ❑ Not Indexing Near-duplicates

- ❑ Browsable and Deletable History

- ❑ Easy (De-)activation of Archiving

- ❑ Combined Indexing of Sub-pages

- ❑ Personalized Search

Code and Instructions on Github



`github.com/webis-de/wasp`

# Thank you for your attention!