# InterPlanetary Wayback

## The Next Step Towards Decentralized Web Archiving

**Sawood Alam**, Mat Kelly, Michele C. Weigle, Michael L. Nelson
Web Science and Digital Libraries Research Group
Old Dominion University
Norfolk, Virginia, USA
**@WebSciDL**

`http://github.com/oduwsdl/ipwb`

IPFS Lab Day, Decentralized Web Summit, 2018
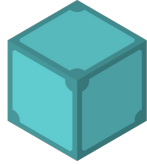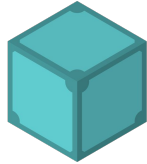San Francisco, CA (USA)
August 3, 2018

# Content Addressing



http://foo.com/spaceDog.jpg

IPFS

QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4

===

QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4



http://example.org/yuri.jpg

IPFS

$ ipfs cat QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4 > doge.jpg

# Rendered HTML vs. Source Code

# HTTP Response vs. WARC Record

```
                    less hello-dweb.warc 63x33
WARC/1.0
WARC-Type: response
WARC-Target-URI: https://www.cs.odu.edu/~salam/dweb/
WARC-Date: 2018-08-02T15:12:34Z
WARC-Record-ID: <urn:uuid:6e829124-b8a3-02ae-4912-f9be5978e09a>
Content-Type: application/http; msgtype=response
Content-Length: 654

HTTP/1.1 200 OK
Server: nginx
Date: Thu, 02 Aug 2018 15:06:34 GMT
Content-Type: text/html
Transfer-Encoding: chunked
Connection: keep-alive
Vary: Accept-Encoding
Front-End-Https: on

<!DOCTYPE html>
<html lang="en" dir="ltr">
  <head>
    <meta charset="utf-8">
    <title>Hello Decentralized Web!</title>
    <link rel="stylesheet" href="style.css">
  </head>
  <body>
    <h1>Hello Decentralized Web!</h1>
    <img src="wsdl-logo.png" alt="">
    <p>We, at <a href="https://ws-dl.cs.odu.edu/">Web Science a
nd Digital Libraries Research Group</a>, tweet from the <a href
="https://twitter.com/WebSciDL">@WebSciDL</a> handle.</p>
  </body>
</html>
hello-dweb.warc
```
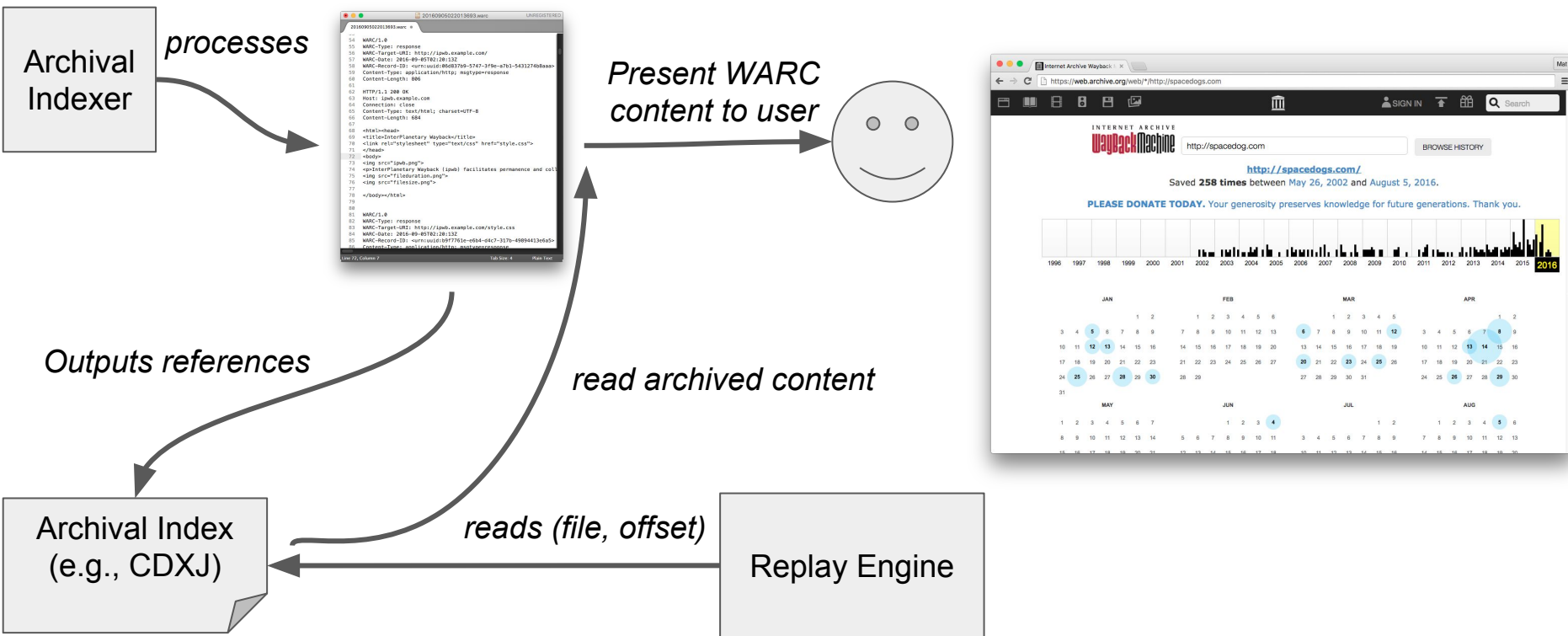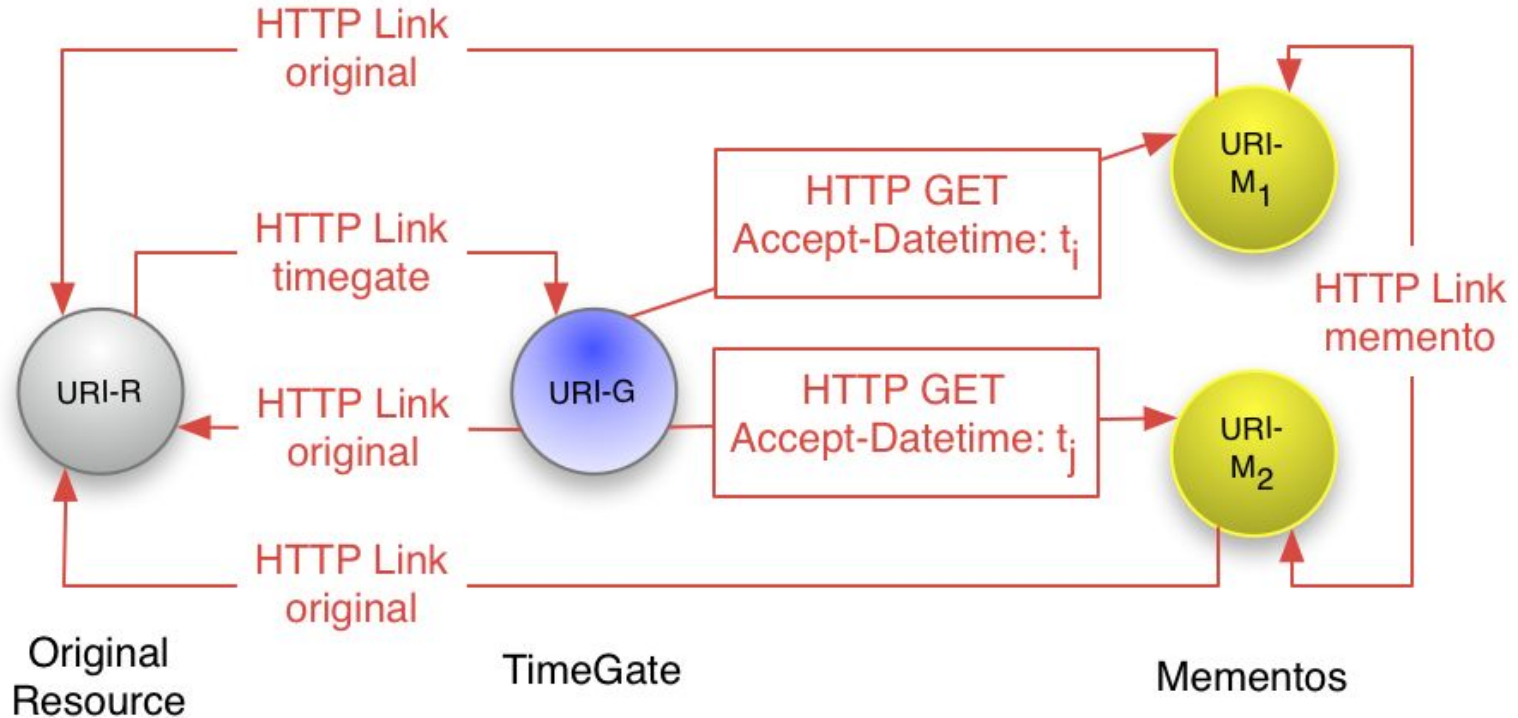
WARC headers

```
        curl -i https://www.cs.odu.edu/~salam/dweb/ 63x25
$ curl -i https://www.cs.odu.edu/~salam/dweb/
HTTP/1.1 200 OK
Server: nginx
Date: Fri, 03 Aug 2018 19:56:33 GMT
Content-Type: text/html
Transfer-Encoding: chunked
Connection: keep-alive
Vary: Accept-Encoding
Front-End-Https: on

<!DOCTYPE html>
<html lang="en" dir="ltr">
  <head>
    <meta charset="utf-8">
    <title>Hello Decentralized Web!</title>
    <link rel="stylesheet" href="style.css">
  </head>
  <body>
    <h1>Hello Decentralized Web!</h1>
    <img src="wsdl-logo.png" alt="">
    <p>We, at <a href="https://ws-dl.cs.odu.edu/">Web Science a
nd Digital Libraries Research Group</a>, tweet from the <a href
="https://twitter.com/WebSciDL">@WebSciDL</a> handle.</p>
  </body>
</html>
```

HTTP headers

Payload

4

# How Wayback Machine Works?



@ibnesayeed

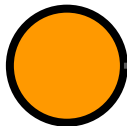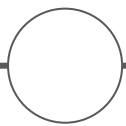# Memento: Time Dimension to the Web

# Why IPWB?



- Persistence of archived web dependent on resilience of organizations
- Availability of data is subject to censorship
- Redundancy in web archive files of exact duplicate content
- Lack of public participation in web archiving
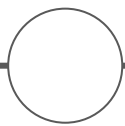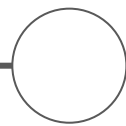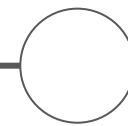- Discoverability issue of small web archives

# Indexing

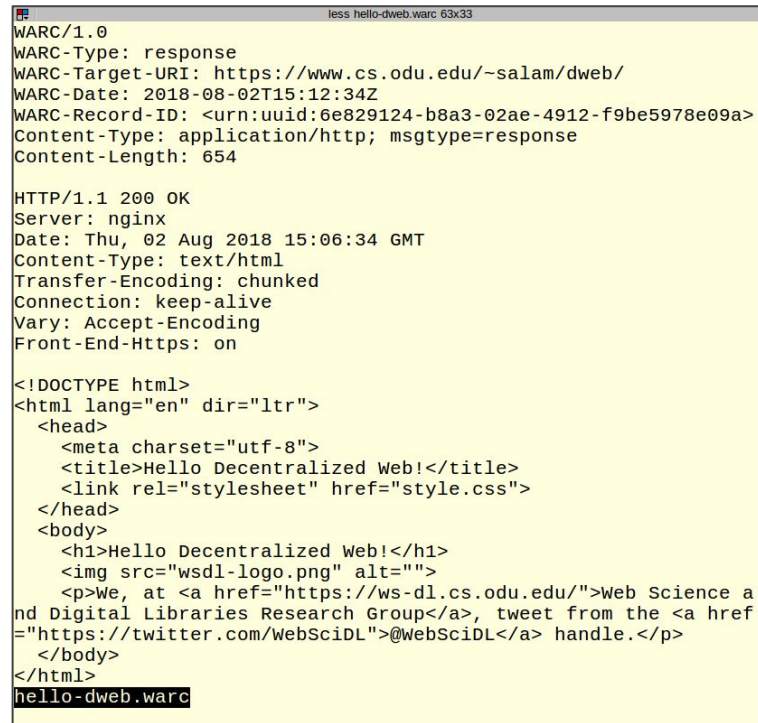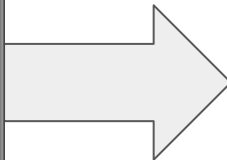WARC Creation — HTTP Header & Payload Extraction — Push to IPFS — Generate CDXJ — WARC-CDXJ Correspondence

Browser showing https://www.cs.odu.edu/~salam/dweb/ with "Hello Decentralized Web!" and WSDL Old Dominion University logo. Text: "We, at Web Science and Digital Libraries Research Group, tweet from the @WebSciDL handle."

```
                     less hello-dweb.warc 63x33
WARC/1.0
WARC-Type: response
WARC-Target-URI: https://www.cs.odu.edu/~salam/dweb/
WARC-Date: 2018-08-02T15:12:34Z
WARC-Record-ID: <urn:uuid:6e829124-b8a3-02ae-4912-f9be5978e09a>
Content-Type: application/http; msgtype=response
Content-Length: 654

HTTP/1.1 200 OK
Server: nginx
Date: Thu, 02 Aug 2018 15:06:34 GMT
Content-Type: text/html
Transfer-Encoding: chunked
Connection: keep-alive
Vary: Accept-Encoding
Front-End-Https: on

<!DOCTYPE html>
<html lang="en" dir="ltr">
  <head>
    <meta charset="utf-8">
    <title>Hello Decentralized Web!</title>
    <link rel="stylesheet" href="style.css">
  </head>
  <body>
    <h1>Hello Decentralized Web!</h1>
    <img src="wsdl-logo.png" alt="">
    <p>We, at <a href="https://ws-dl.cs.odu.edu/">Web Science a
nd Digital Libraries Research Group</a>, tweet from the <a href
="https://twitter.com/WebSciDL">@WebSciDL</a> handle.</p>
  </body>
</html>
hello-dweb.warc
```

9

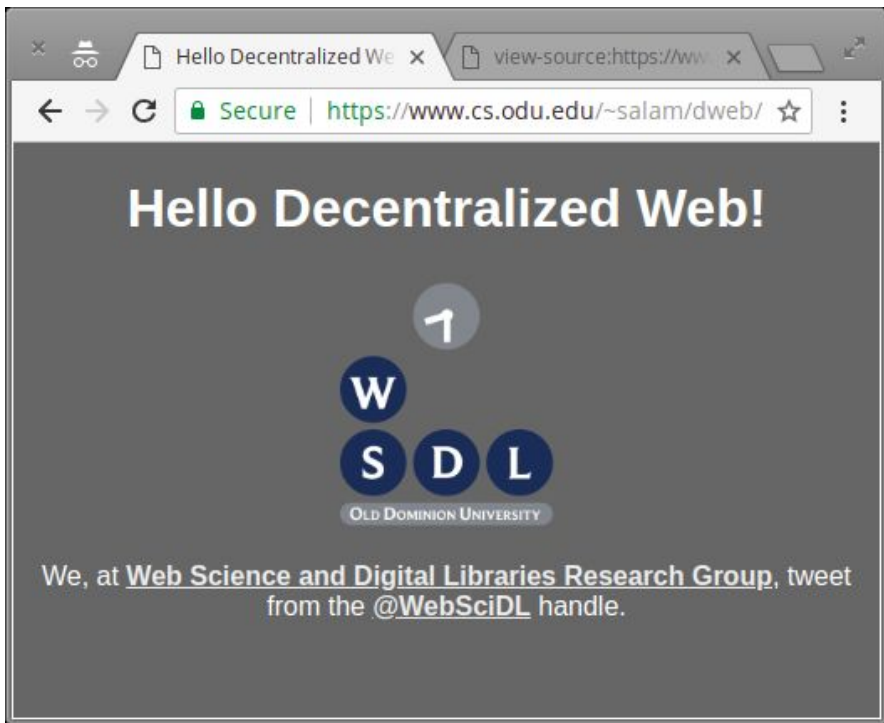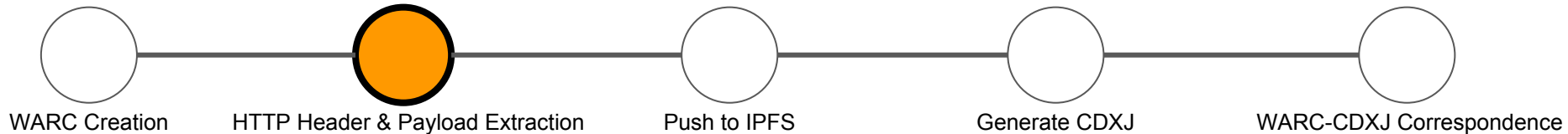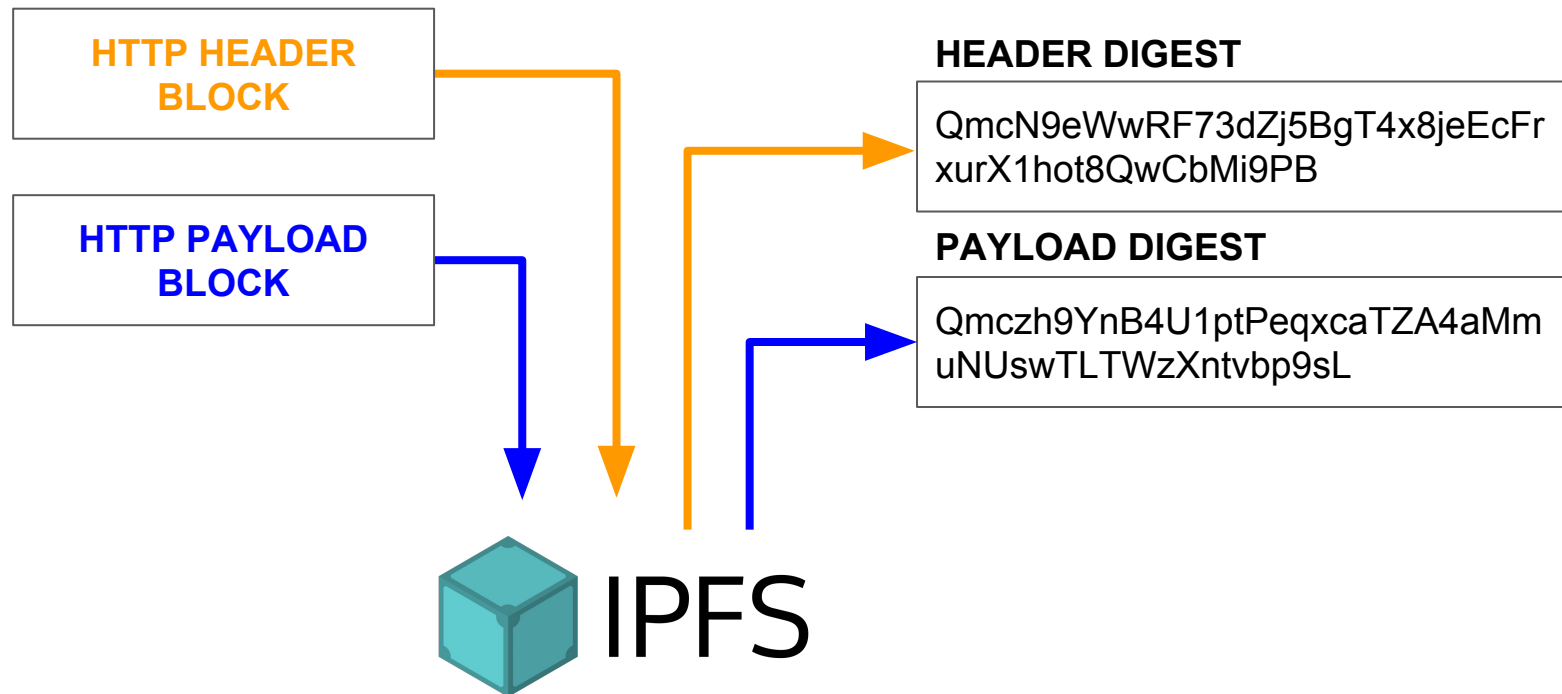WARC Creation · **HTTP Header & Payload Extraction** · Push to IPFS · Generate CDXJ · WARC-CDXJ Correspondence
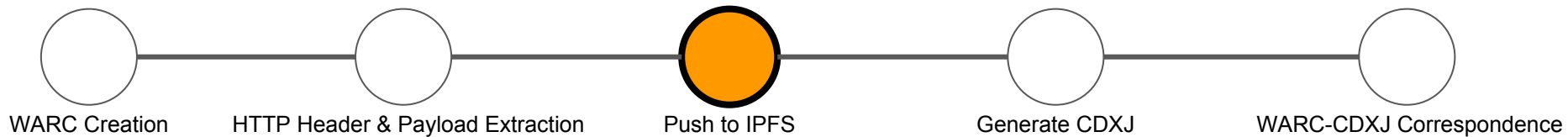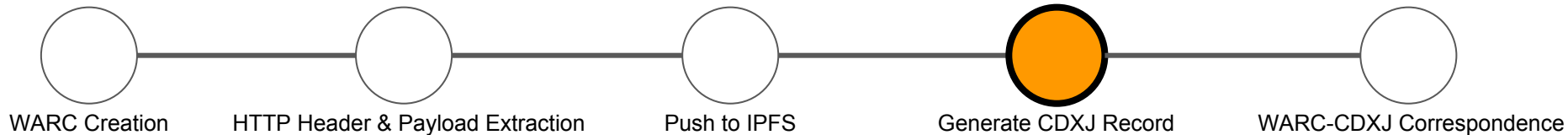
```
                    less hello-dweb.warc 63x33
WARC/1.0
WARC-Type: response
WARC-Target-URI: https://www.cs.odu.edu/~salam/dweb/
WARC-Date: 2018-08-02T15:12:34Z
WARC-Record-ID: <urn:uuid:6e829124-b8a3-02ae-4912-f9be5978e09a>
Content-Type: application/http; msgtype=response
Content-Length: 654

HTTP/1.1 200 OK
Server: nginx
Date: Thu, 02 Aug 2018 15:06:34 GMT
Content-Type: text/html
Transfer-Encoding: chunked
Connection: keep-alive
Vary: Accept-Encoding
Front-End-Https: on

<!DOCTYPE html>
<html lang="en" dir="ltr">
  <head>
    <meta charset="utf-8">
    <title>Hello Decentralized Web!</title>
    <link rel="stylesheet" href="style.css">
  </head>
  <body>
    <h1>Hello Decentralized Web!</h1>
    <img src="wsdl-logo.png" alt="">
    <p>We, at <a href="https://ws-dl.cs.odu.edu/">Web Science a
nd Digital Libraries Research Group</a>, tweet from the <a href
="https://twitter.com/WebSciDL">@WebSciDL</a> handle.</p>
  </body>
</html>
hello-dweb.warc
```

**HTTP HEADER BLOCK**

**HTTP PAYLOAD BLOCK**

*@ibnesayeed*

10

**HTTP HEADER BLOCK**

**HTTP PAYLOAD BLOCK**

**HEADER DIGEST**

QmcN9eWwRF73dZj5BgT4x8jeEcFr
xurX1hot8QwCbMi9PB

**PAYLOAD DIGEST**
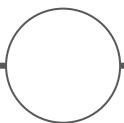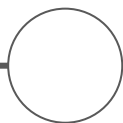
Qmczh9YnB4U1ptPeqxcaTZA4aMm
uNUswTLTWzXntvbp9sL

IPFS

*@ibnesayeed*

11

WARC Creation — HTTP Header & Payload Extraction — Push to IPFS — **Generate CDXJ Record** — WARC-CDXJ Correspondence

## HEADER DIGEST

QmcN9eWwRF73dZj5BgT4x8jeEcFrxurX1hot8QwCbMi9PB

## PAYLOAD DIGEST

Qmczh9YnB4U1ptPeqxcaTZA4aMmuNUswTLTWzXntvbp9sL

```
com,example,ipwb)/ 20180802012013 {
     "locator":"urn:ipfs/QmcN9eWwRF73dZj5BgT4x8jeEcFrxurX1hot8QwCbMi9PB/Qmczh9YnB4U1ptPeqxca
     TZA4aMmuNUswTLTWzXntvbp9sL",
     "mime_type": "text/html",
     "status_code": 200,
     "other_fields": "other values..."
}
// * This is a single-line record, line breaks and indentation are added for readability only
```
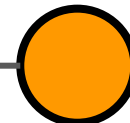
**CDXJ:** http://ws-dl.blogspot.com/2015/09/2015-09-10-cdxj-object-resource-stream.html

*@ibnesayeed*  12

WARC Creation · HTTP Header & Payload Extraction · Push to IPFS · Generate CDXJ · WARC-CDXJ Correspondence

hello-dweb.warc

```
edu,odu,cs)/~salam/dweb/ 20180802012013 {
    "locator": "urn:ipfs/QmcN9eWwRF73dZj5.../Qmczh9YnB4U1ptPe...",
    "mime_type": "text/html",
    "status_code": "200"
}

edu,odu,cs)/~salam/dweb/style.css 20180802012013 {
    "locator": "urn:ipfs/QmU1k71bT6ibZBSd.../QmbvUAo9U31wSdvA...",
    "mime_type": "text/css",
    "status_code": "200"
}

edu,odu,cs)/~salam/dweb/wsdl-logo.png 20180802012013 {
    "locator": "urn:ipfs/QmTjfMxFGvbP4nwF.../QmYMKZbnk53kuPJi...",
    "mime_type": "image/png",
    "status_code": "200"
}
```
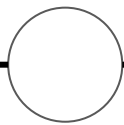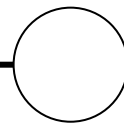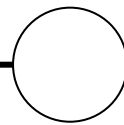
13

# Replay

Lookup in CDXJ  —  Fetch from IPFS  —  Reconstruct Response  —  Reroute Resources

https://www.cs.odu.edu/~salam/dweb/

SEP

|    | 1  | 2  | 3  | 4  | 5  | 6  |
|    | 7  | 8  | 9  | 10 | 11 | 12 | 13 |
|    | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|    | 21 | 22 | 23 | 24 | 25 | 26 | 27 |
|    | 28 | 29 | 30 | 31 |

```
edu,odu,cs)/~salam/dweb/ 20180802012013 {
    "locator": "urn:ipfs/QmcN9eWwRF.../Qmczh9YnB4...",
    "mime_type": "text/html",
    "status_code": "200"
}

edu,odu,cs)/~salam/dweb/style.css 20180802012013 {
    "locator": "urn:ipfs/QmU1k71bT6.../QmbvUAo9U3...",
    "mime_type": "text/css",
    "status_code": "200"
}

edu,odu,cs)/~salam/dweb/wsdl-logo.png 20180802012013 {
    "locator": "urn:ipfs/QmTjfMxFGv.../QmYMKZbnk5...",
    "mime_type": "image/png",
    "status_code": "200"
}
```

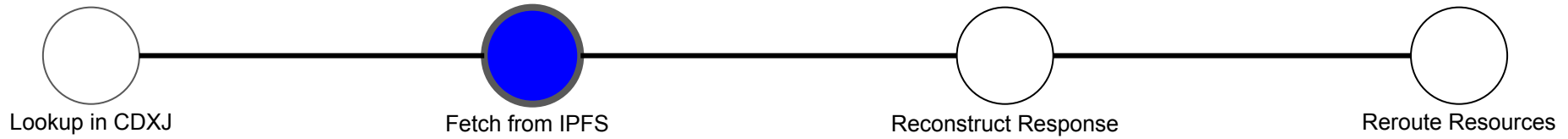Lookup in CDXJ     Fetch from IPFS     Reconstruct Response     Reroute Resources
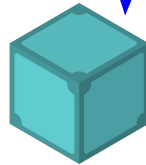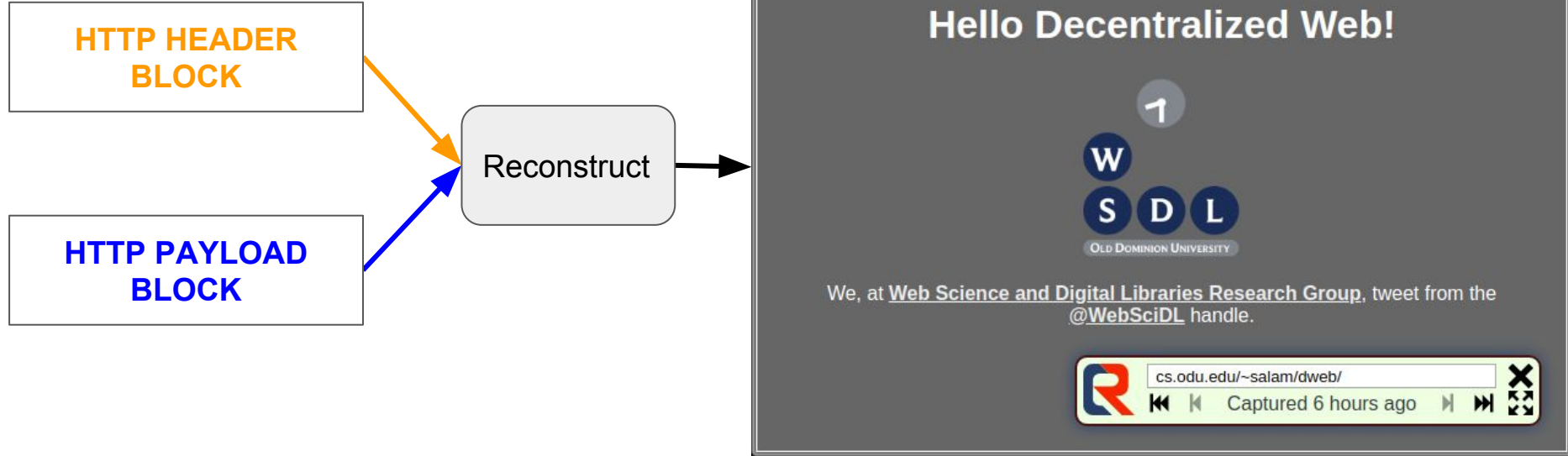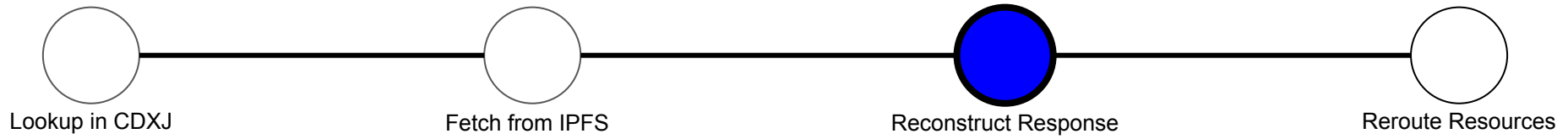
```
edu,odu,cs)/~salam/dweb/ 20180802012013 {
    "locator": "urn:ipfs/QmcN9eWwRF73dZj5.../Qmczh9YnB4U1ptPe...",
    "mime_type": "text/html",
    "status_code": "200"
}
```

IPFS

HTTP HEADER BLOCK

HTTP PAYLOAD BLOCK

Lookup in CDXJ

Fetch from IPFS

Reconstruct Response

Reroute Resources

**HTTP HEADER BLOCK**

**HTTP PAYLOAD BLOCK**

Reconstruct

Lookup in CDXJ     Fetch from IPFS     Reconstruct Response     Reroute Resources

- Avoids zombies (live-leakage)
- Adds an unobtrusive archival banner (Custom HTML Element)

- https://oduwsdl.github.io/Reconstructive/
- http://ws-dl.blogspot.com/2018/01/2018-01-08-introducing-reconstructive.html

*@ibnesayeed*     18

# IPWB Indexing and Replay

# Decentralization

# Current Issues

- IPFS is permanent, but not persistent
- DHT-based IPNS is history-unaware
- CDXJ index, a critical piece of replay, is centralized

# Persistence

- Data persistence is critical for web archiving
- A decentralized storage with sufficient replication is needed
- Memory organizations should contribute storage infrastructure
- Qri, Filecoin, IPFS-Cluster, IPFS-Sync etc. can be helpful

# IPNS: InterPlanetary Naming System

| URI | IPFS Hash |
|---|---|
| http://example.org/yuri.jpg | QmZAD4xeeNeYF3TmwWgBXypLKTiCGwGRMXHW7MtheWKtw4 |
| http://example.com/style.css | QmbvUAo9U31wSdvARjvbPeVBTAwCjN1kyPhQ4ho3n8TAZo |
| http://example.com/logo.png | QmYMKZbnk53kuPJirahJHGevCCy2afLyePRdX38TukFUwd |

How about changes and history?

| http://example.com/style.css | QmAAykNAehCjbn43wvUQUd1ovhTPjRo39VvoT8bSA1ZBwP |
|---|---|

# IPNS Blockchain

- URI $\rightarrow$ Latest hash
- URI + DateTime $\rightarrow$ A historical hash
- URI $\rightarrow$ List historical hashes with times

IPNS + Blockchain + Memento

https://github.com/oduwsdl/IPNS-Blockchain

| Owner | URI | Time | Hash | PrevBlock |
|-------|-----|------|------|-----------|
| Pub $K_1$ | $URI_1$ | $T_1$ | $H_1$ | 1234567... |
| Pub $K_2$ | $URI_2$ | $T_2$ | $H_2$ | 0000000... |
| Pub $K_3$ | $URI_3$ | $T_3$ | $H_3$ | 9876543... |

| Owner | URI | Time | Hash | PrevBlock |
|-------|-----|------|------|-----------|
| Pub $K_1$ | $URI_1$ | $T_4$ | $H_5$ | 5463728... |
| Pub $K_3$ | $URI_4$ | $T_5$ | $H_6$ | 0000000... |

# Lazy Relationship Evaluation



/<namespace>/about/<URI>

MementoOf
(Active Relation)

HasMemento
(Lazy Evaluation)

Memento

Memento

Memento

IP-LD to the rescue?

https://github.com/oduwsdl/ipwb/issues/61

# Evaluation

# Storage Space and Time Overhead



- Reported IPFS slowness https://github.com/ipfs/go-ipfs/issues/1216
  - Has since been fixed, but we did not evaluate again

# Replay Time

- 600 requests in 222 seconds
- Slower than PyWB (which took 5.26 seconds)
- File vs. rich object based retrieval
- Never expiring cache

https://github.com/ibnesayeed/ipfsapi-concurrency-test

# Future Works

- Evaluate the improved IPFS on large dataset
- Evaluate deduplication
- Implement an index-free collaborative archiving system
- Utilize IPNS to reference URI-Rs with datetime
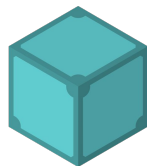
*@ibnesayeed*

# Conclusions

- A proof of concept system to leverage a novel approach to archiving and retrieval
- Storage and time costs evaluation and qualitative analysis
- It can only work for small archives in its current state
- A path to answer "who will archive the archives?"
- More work to be done to make it a truly decentralized archiving system

# InterPlanetary Wayback

## The Next Step Towards Decentralized Web Archiving

**Sawood Alam**, Mat Kelly, Michele C. Weigle, Michael L. Nelson

`http://github.com/oduwsdl/ipwb`

**@WebSciDL**