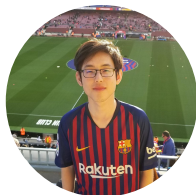# Jawa:
# Web Archival in the Era of JavaScript

**Ayush Goel**
**University of Michigan**

Jingyuan Zhu
*University of Michigan*

Ravi Netravali
*Princeton University*

Harsha V. Madhyastha
*University of Michigan*

How to **reduce the storage overhead** of web archives and **improve the quality** of archived pages?

# A Couple of Days After the Conference..

# A Couple of Days After the Conference..

# 10 Years From Now (2032)....



Ayush Goel

Publications     News     Misc

**About me**

Ayush Goel is a (full) Professor at [REDACTED]. His work on networking systems has won award papers at premier systems conferences like NSDI, SIGCOMM, OSDI and SOSP. ......

**Ayush Goel**

Professor in the Computer
Science Divison

📍 Ann Arbor, Michigan
✉ Email
🐦 Twitter

**Recent News (all)**

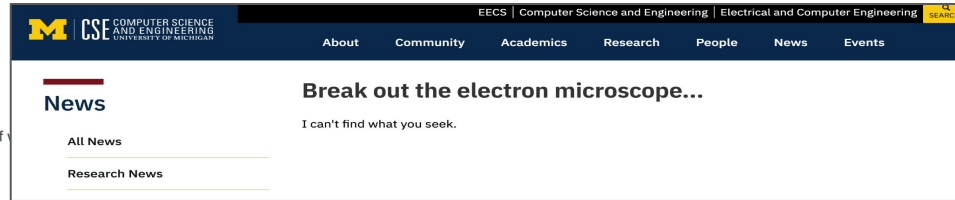**July 2032** I got awarded the 2032 ACM Turing Award for my pioneering work in the area of ... systems and PL.
........
........
........
**July 2022** Our work Jawa got featured in www.nytimes.com, hackernews, UMICH CSE
**March 2022** Our work Jawa got accepted at OSDI'22.

The New York Times

**Page Not Found**

We're sorry, we seem to have lost this page,
but we don't want to lose you.

Search NYTimes.com     GO

Report a broken link | Go to Home Page

M | CSE COMPUTER SCIENCE AND ENGINEERING UNIVERSITY OF MICHIGAN

EECS | Computer Science and Engineering | Electrical and Computer Engineering

About     Community     Academics     Research     People     News     Events

**News**

All News

Research News

**Break out the electron microscope...**

I can't find what you seek.

This site can't be reached

Check if there is a typo in news.ycombinator.com.

DNS_PROBE_FINISHED_NXDOMAIN

Reload

# High Prevalence of Link Rot on the Web

# Web Archives to the Rescue

# Web Archives to the Rescue



~700 Billion pages, ~100 Petabytes

# How Modern Web Archives Operate?

Crawler

Fetching HTML, JS, CSS, Images

User

Web pages on live web

Apply deduplication, compression

Fetch archived page HTML, CSS, JS, Images

Storage

# Problems with Web Archives: **No Snapshots for Many Pages**

# Problems with Web Archives: **Poor Page Fidelity**

# Outline

★ **Challenges faced by web archives**

★ Understanding the root cause

★ Our insights and design of a new crawler

★ Evaluation

# Outline

★ Challenges faced by web archives

★ **Understanding the root cause**

★ Our insights and design of a new crawler

★ Evaluation

# Root Cause: **Increasing JavaScript** on Web Pages



*Corpus*: Landing pages of 300 Alexa sites

# Root Cause: **Increasing JavaScript** on Web Pages

Per page snapshot **expensive** ➜ **Fewer snapshots**



*Corpus*: Landing pages of 300 Alexa sites

# Root Cause: **JavaScript Induced Non-Determinism**

**Resources fetched** different from **crawled → Poor page fidelity**

# Outline

★ Challenges faced by web archives

★ Understanding the root cause

★ **Our insights and design of a new crawler**

★ Evaluation

# Key Insight: Archived Page Differs from Live Page

A.  No back-end origin server

*Account login and subscription*

*Search query to www.nytimes.com*

The New York Times

PLAY THE CROSSWORD   Account ⌄

**JavaScript that interacts with server will not work ➔ Removing it will not impact fidelity**

storage overhead while ensuring high page fidelity

Give this article   💬 22

*Add comments to the article*

By **Kenneth Chang**

July 20, 2022

# Key Insight: Archived Page Differs from Live Page

B. Certain sources of non-determinism are absent

Reads client-side cookie:
**Always returns "no"**

Good evening.

Subscription overview  >

**YOUR CONTENT**

Saved articles  >

**GET SUPPORT**

Help Center  >

Log out  >

**JavaScript in certain control flows will never be executed ➜ Removing it will not impact fidelity**

✕

Print login button

SUBSCRIBE FOR $1/WEEK    LOG IN

19

# **Jawa**: A New Web Archive Crawler

**Improve page fidelity**



**Reduce storage cost**

➔ **Eliminate impact of non-determinism**

➔ Identify and remove non-functional code

➔ **Identify and remove unreachable code**

# **Jawa**: A New Web Archive Crawler

**Improve page fidelity**

**Reduce storage cost**

➔ **Eliminate impact of non-determinism**

→ Identify and remove non-functional code

→ **Identify and remove unreachable code**

# Make JavaScript Execution **Completely** Deterministic

Original page

Deterministic page (broken)

# Understand How Non-determinism Impacts Resources Fetched

Date, Math.random, Performance (DRP)

Client characteristics



Analyze impact on control flow on 3000 pages

**No impact on control flow**
➡ retain the non-determinism

**Influences control flow impacting resource fetches**
➡ Eliminate non-determinism

# **Jawa**: A New Web Archive Crawler

**Improve page fidelity**

**Reduce storage cost**

➡ **Eliminate impact of non-determinism**

➡ Identify and remove non-functional code

➡ **Identify and remove unreachable code**

# Insufficient to Save Only JS Executed While Crawling

# Challenge: Code Executed Depends on How User Interacts

# Challenge: Code Executed Depends on **Order**

# Challenge: Code Executed Depends on **Order**

➔ **Read-write dependencies** between different events on page

  ◆ **JavaScript/DOM state**

➔ **Can not predict** order of events

# Challenge: Code Executed Depends on **Input**



1000s of result

One

**Exploring all possible orderings/inputs impractical**

# Analysis Framework: Understand State Dependencies

Inject instrumentation code

```
function foo(a,b){
    var p = Proxy(win,handler)
    p.c = _(a)+_(b);
    return _ret(true);
}

foo(1,2);
```

Trigger each interaction

Monitor state accesses



*Corpus*: 3000 pages from 300 Alexa sites

# Findings From the Analysis

➔ Order **mostly** influences coverage for **analytics events**

**Both categories of events irrelevant for archived pages**

➔ Input **mostly** influences coverage for **text-based events**

Implication:

Estimate code coverage with **single execution of events with default inputs**

# Outline

★  Challenges faced by web archives

★  Understanding the root cause

★  Our insights and design of a new crawler

★  **Evaluation**

# Evaluation

## Storage

- Storage for web resources
- Storage for crawling/serving indices

## Fidelity

- Failed Resource fetches
- Visual comparison
- Functional interactions

## Throughput

- Crawling throughput
- Overhead of each technique

# Evaluation

## Storage

- **Storage for web resources**
- Storage for crawling/serving indices

## Fidelity

- **Failed Resource fetches**
- Visual comparison
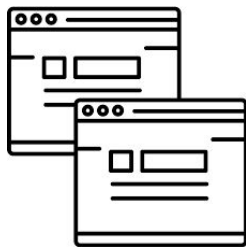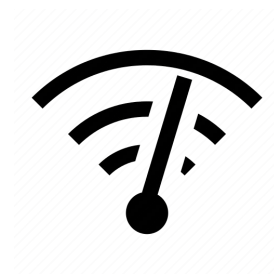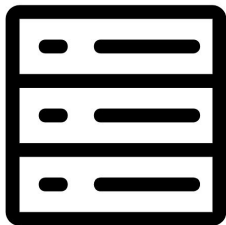- Functional interactions

## Throughput

- Crawling throughput
- Overhead of each technique

# Storage: Jawa Reduces Storage Overhead by **41%**

*Corpus*:

<u>1 million snapshots</u> of pages on 300 sites archived by Internet Archive

### JavaScript

**84%**

Total Storage (GB)

200

100

0

IA *            Jawa

### Total

**41%**

Total Storage (GB)

400

200

0

IA *            Jawa

# Fidelity: Jawa Eliminates Almost All Failed Resource Fetches



**On 10% of pages, >30% of resource fetches failed**

*Corpus*: 3000 web pages from 300 sites

# Fidelity: Jawa Eliminates Almost All Failed Resource Fetches



*Corpus*: 3000 web pages from 300 sites

# Conclusion

1.  JavaScript on web pages negatively impacts web archival

2.  Fundamental differences b/w live and archived pages can be exploited to overcome such negative impacts
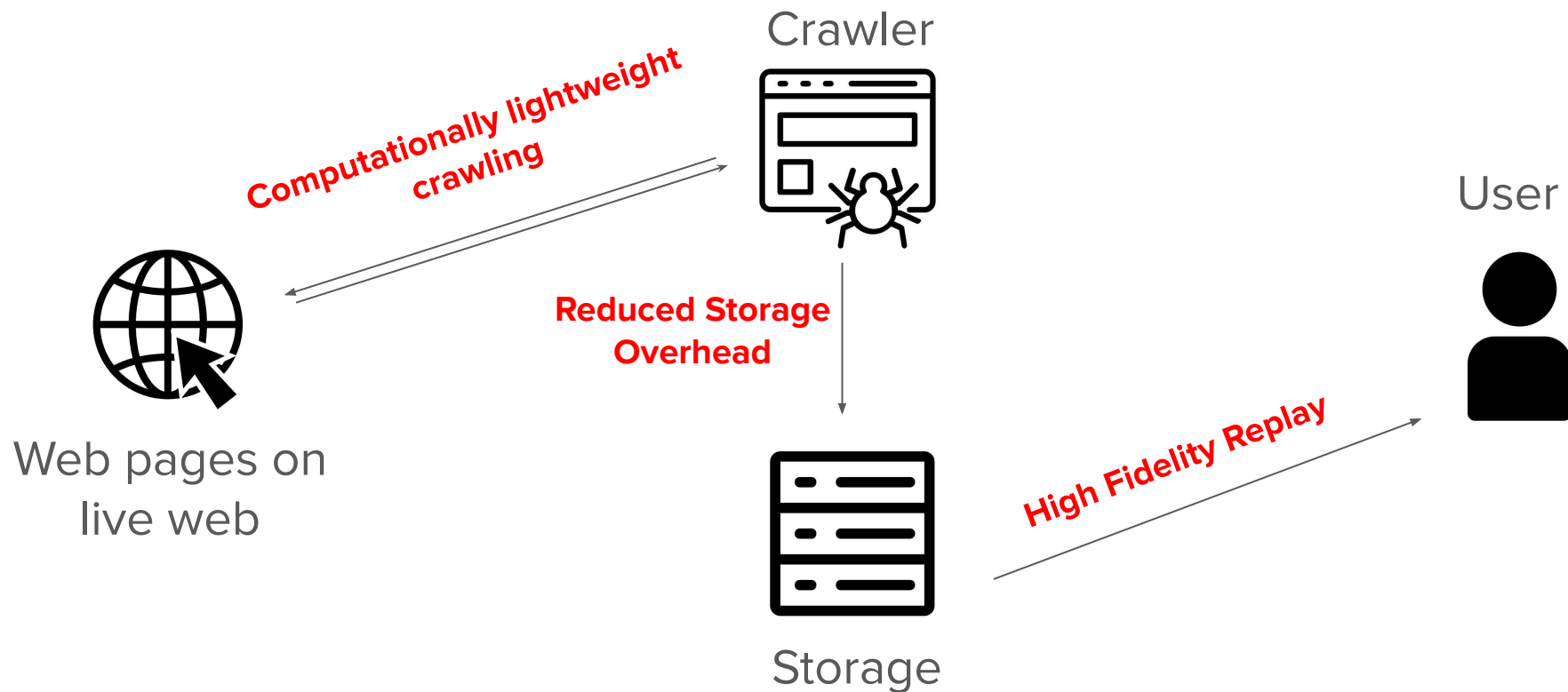
https://github.com/goelayu/Jawa

goelayu@umich.edu

ARTIFACT
EVALUATED
usenix
ASSOCIATION
AVAILABLE

ARTIFACT
EVALUATED
usenix
ASSOCIATION
FUNCTIONAL

ARTIFACT
EVALUATED
usenix
ASSOCIATION
REPRODUCED

# Evaluation.

# Backup

# Jawa: A New Web Archival Crawler



Crawler

Computationally lightweight crawling

Web pages on live web

Reduced Storage Overhead

User

Storage

High Fidelity Replay

# Improve Page Fidelity: Strawman #1

Screen capture



❌ No client interactions

# Jawa: A New Web Archival Crawler

**High Fidelity**
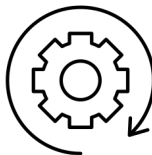
**Low Cost**



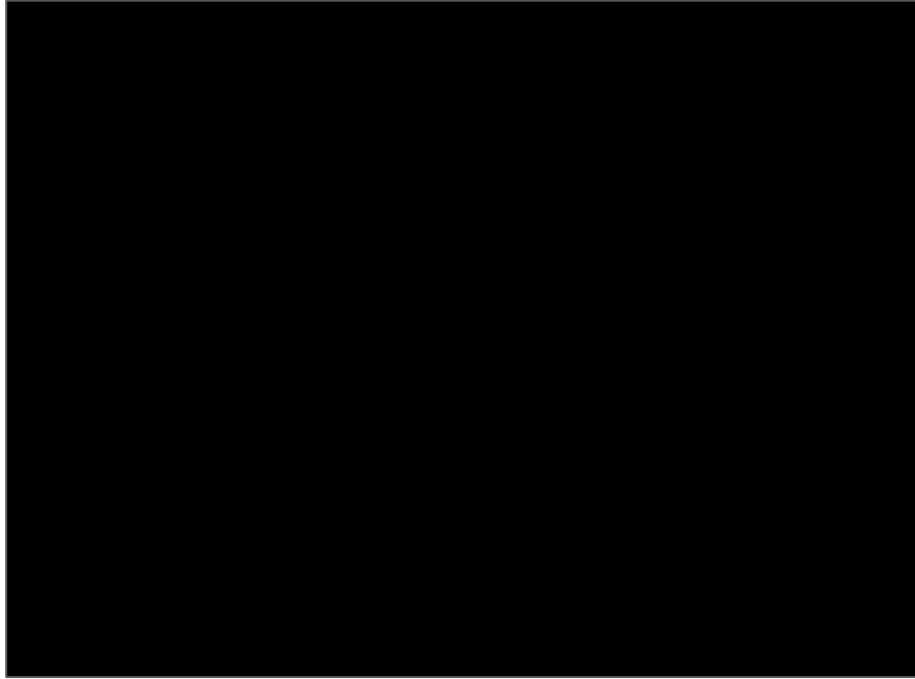Archived page

Visual ← Functionality (vertical double arrow) → Original page


Reduced storage


Computationally lightweight

# Web Archives to the Rescue
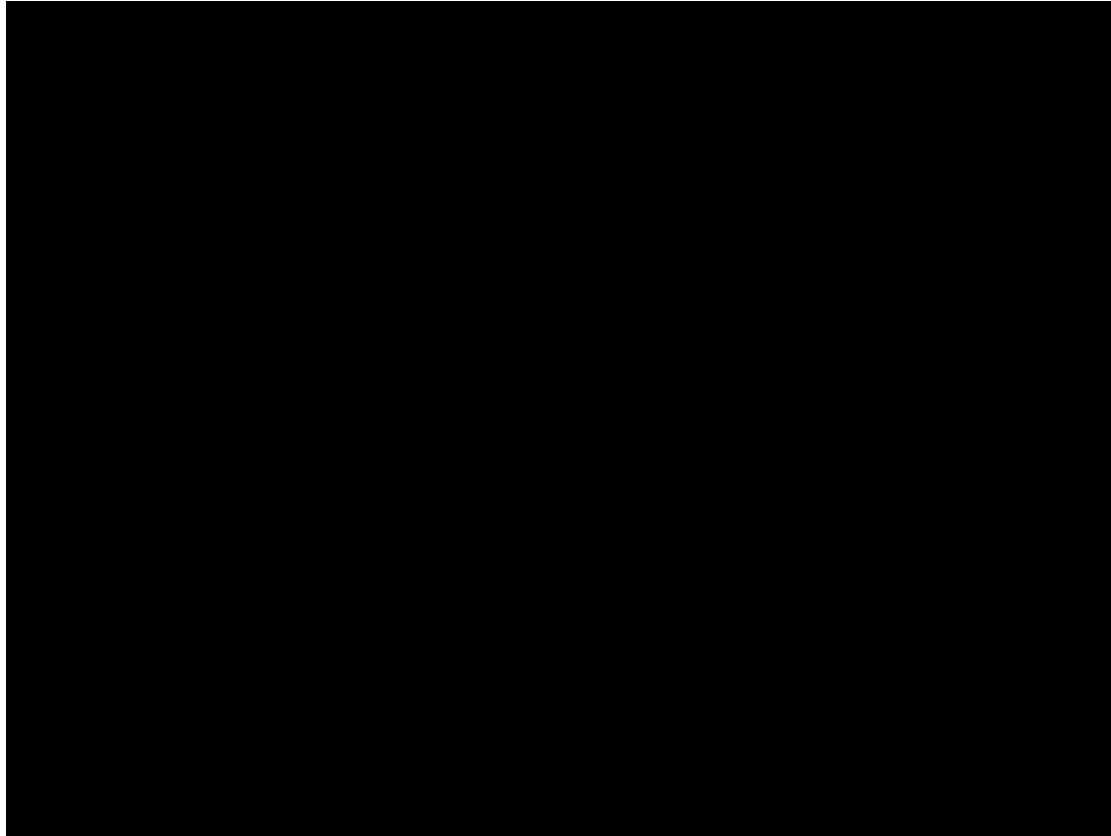


~700 Billion pages, ~100 Petabytes

# Web Archives to the Rescue

To tackle this problem of ephemeral web, various web archives have been established
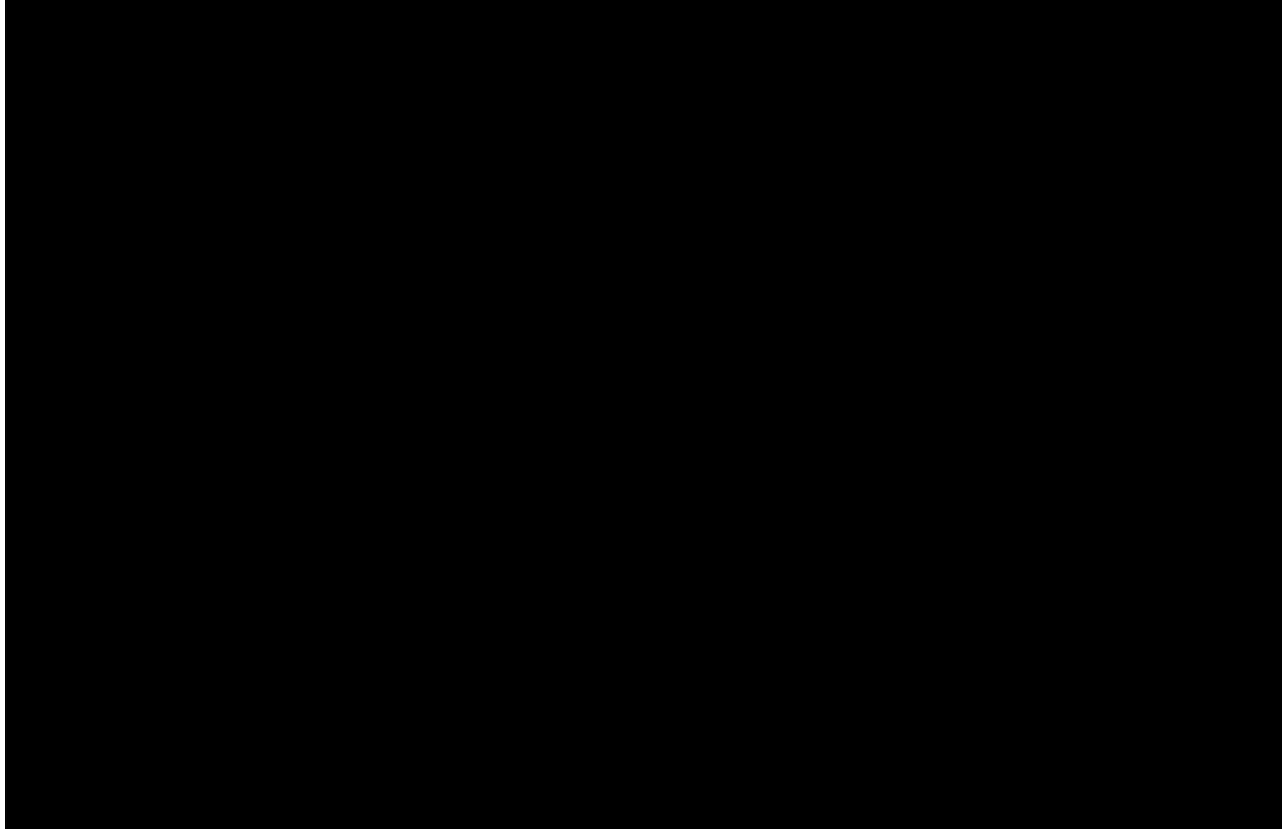
# Web Archives to the Rescue
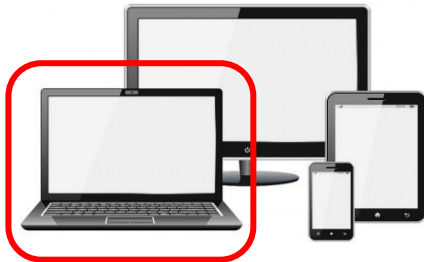
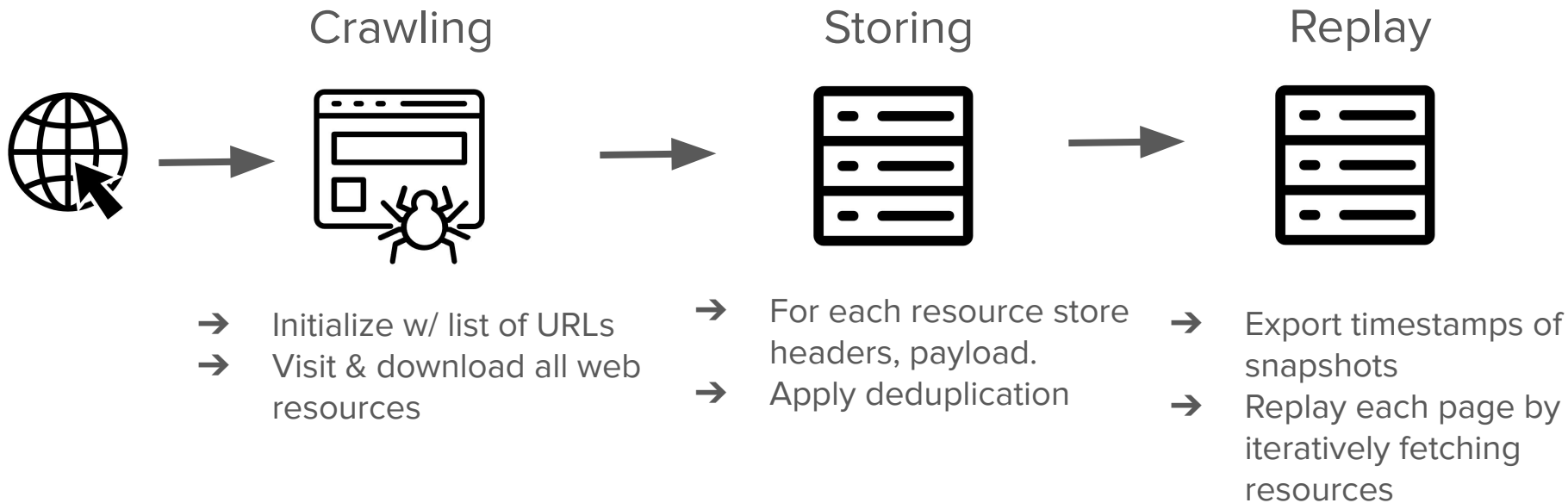# Web Archives to the Rescue

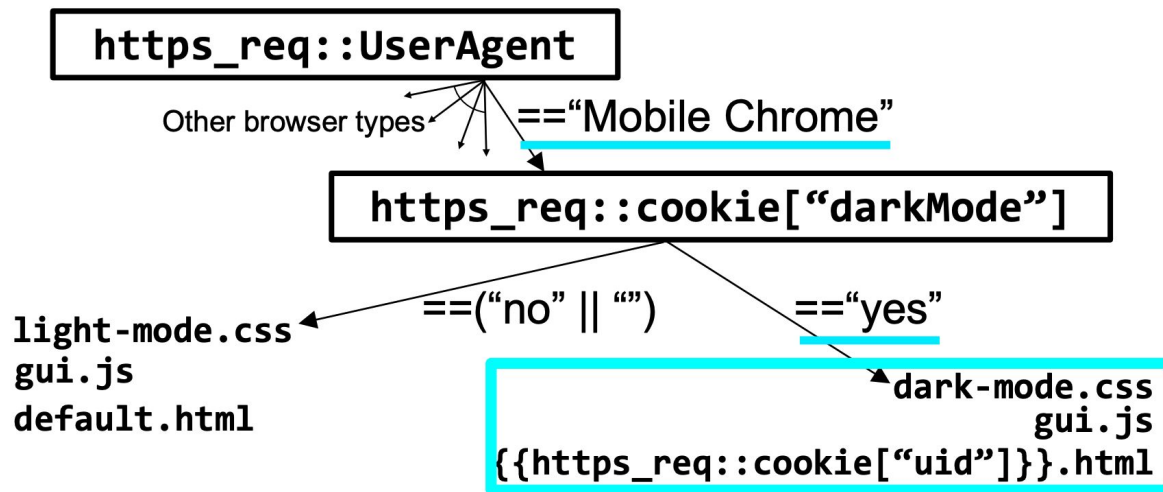# Web Archives to the Rescue

# Approach to Fix Fidelity

➔ Do not make DRP APIs deterministic

➔ **Leverage server-side matching** techniques to account for such URL differences

➔ Enforce same client characteristics for loads as the ones used for crawling

# How Modern Web Archives Operate

## Crawling

➔ Initialize w/ list of URLs
➔ Visit & download all web resources

## Storing

➔ For each resource store headers, payload.
➔ Apply deduplication

## Replay

➔ Export timestamps of snapshots
➔ Replay each page by iteratively fetching resources

# **Non-determinism** Resulting in Poor **Fidelity**



```
https_req::UserAgent
```

Other browser types        =="Mobile Chrome"

```
https_req::cookie["darkMode"]
```

```
light-mode.css
gui.js
default.html
```

==("no" || "")        =="yes"

```
                    dark-mode.css
                            gui.js
{{https_req::cookie["uid"]}}.html
```

# Strawman: Preserve Code Executed **During Page Load**
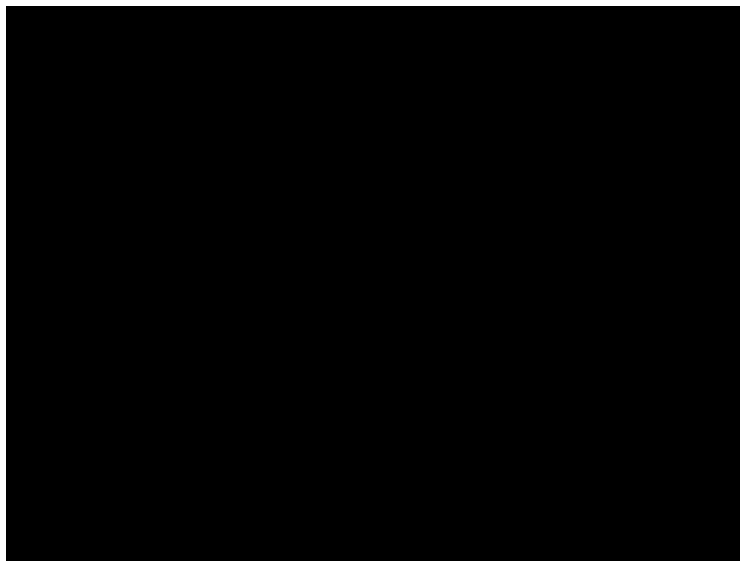
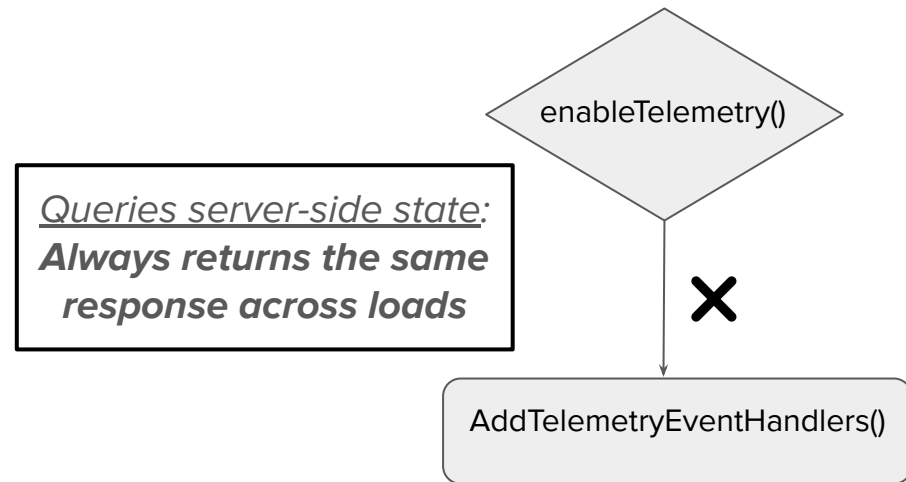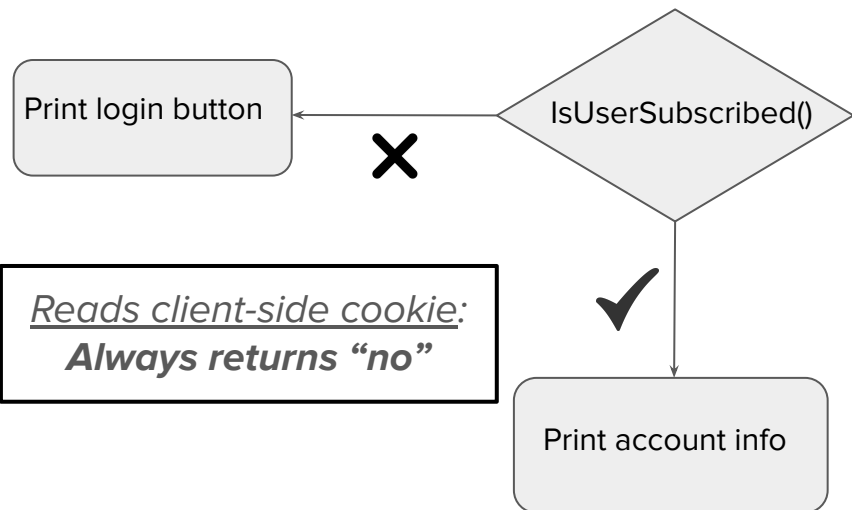# Make JavaScript Execution **Completely** Deterministic

Original page

Deterministic page (broken)

# Key Insight #1: Archived Page ≠ Live Page

B. Certain sources of non-determinism are absent



Print login button

IsUserSubscribed()

✗

*Reads client-side cookie:*
**Always returns "no"**

✓

Print account info

enableTelemetry()

*Queries server-side state:*
**Always returns the same response across loads**

✗

AddTelemetryEventHandlers()

≠