# Detecting and Diagnosing Errors in Replaying Archived Web Pages

Jingyuan Zhu          Huanchen Sun          Harsha V. Madhyastha
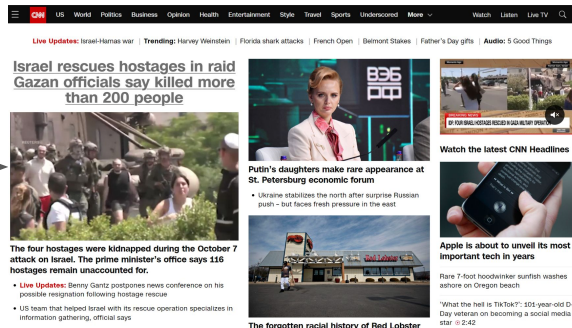
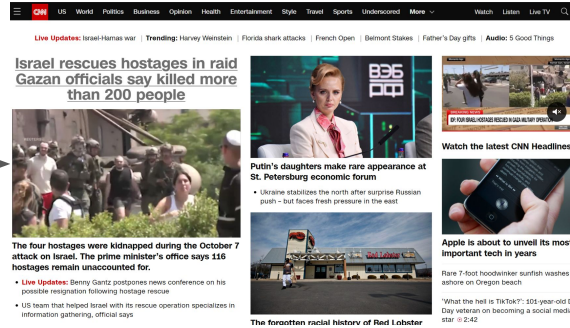UNIVERSITY OF MICHIGAN

USC University of Southern California

# Web archives capture and serve snapshots of webpages



Record: Crawler — Crawl → [webpage snapshot] — Generate → .warc

# Web archives capture and serve snapshots of webpages



Record:

Crawler → Crawl → [webpage snapshot] → Generate → .warc

Replay:

.warc → Input → INTERNET ARCHIVE WayBackMachine / pywb → [webpage snapshot]
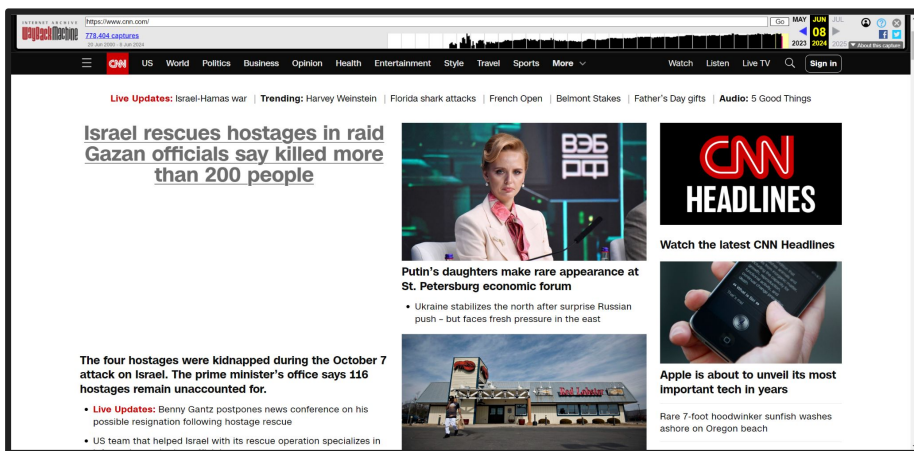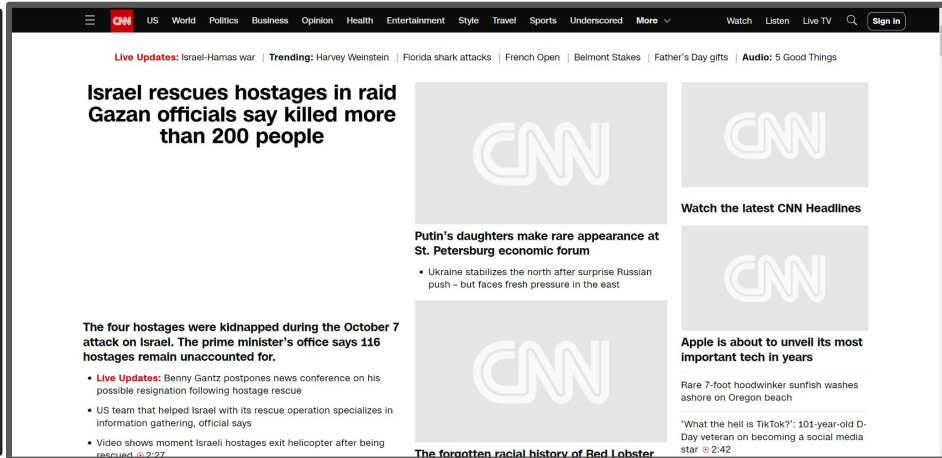
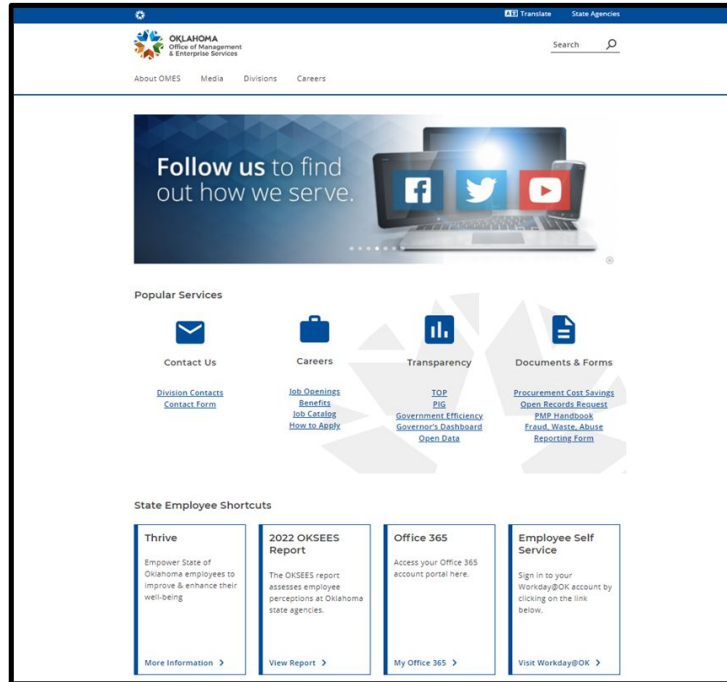# Need to rewrite URLs for crawled resources

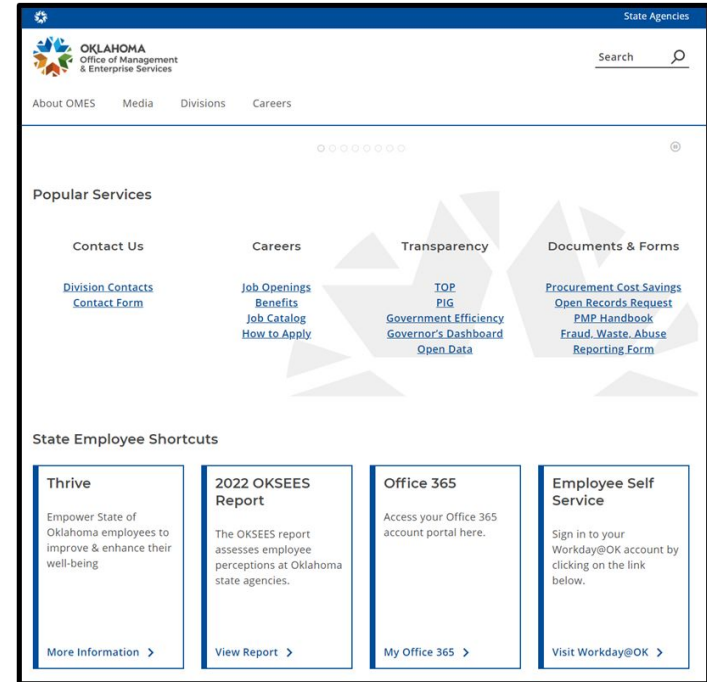## With rewriting



## Without rewriting

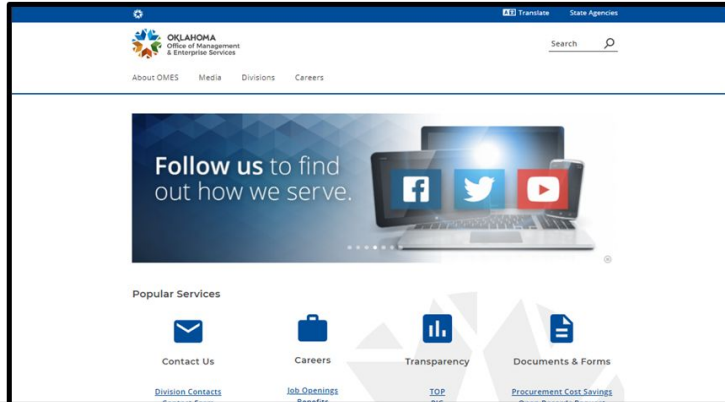# However, JS rewrites can lead to fidelity issues

## Live web page

## Archived copy

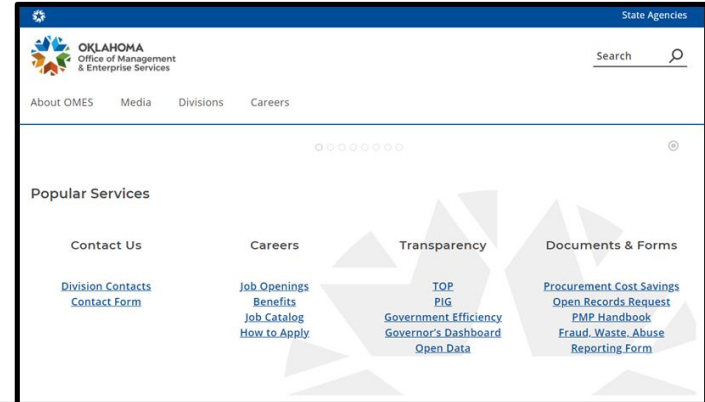# However, JS rewrites can lead to fidelity issues

Live web page

Archived copy



```
❌ Uncaught TypeError: string.indexOf is not a function          wombat.js:21
      at Wombat.startsWithOneOf (wombat.js:21:21211)
      at Object.get (wombat.js:21:61588)
      at readData (clientlib-base.min.js:2450:22)
      at HTMLDocument.onDocumentReady (clientlib-base.min.js:2695:56)
```
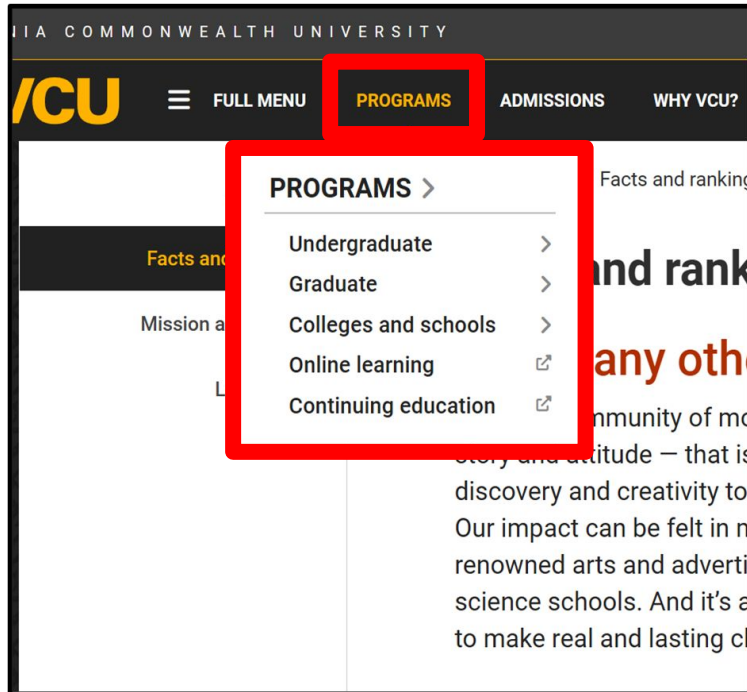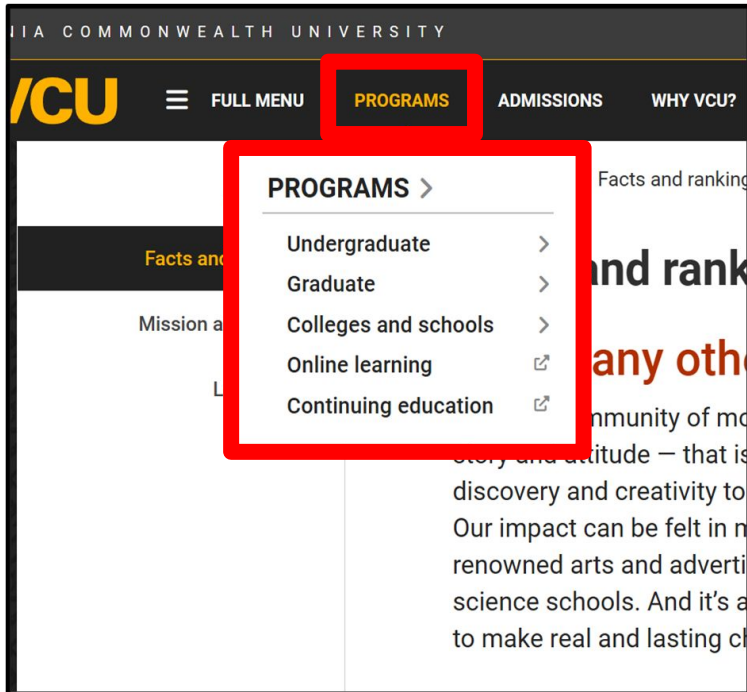
# However, JS rewrites can lead to fidelity issues

## Live web page

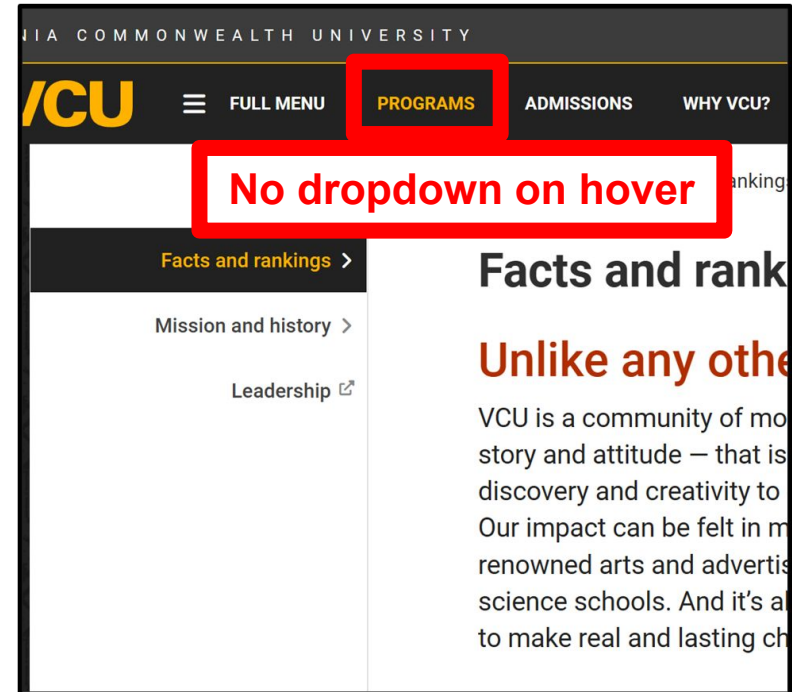# However, JS rewrites can lead to fidelity issues

## Live web page

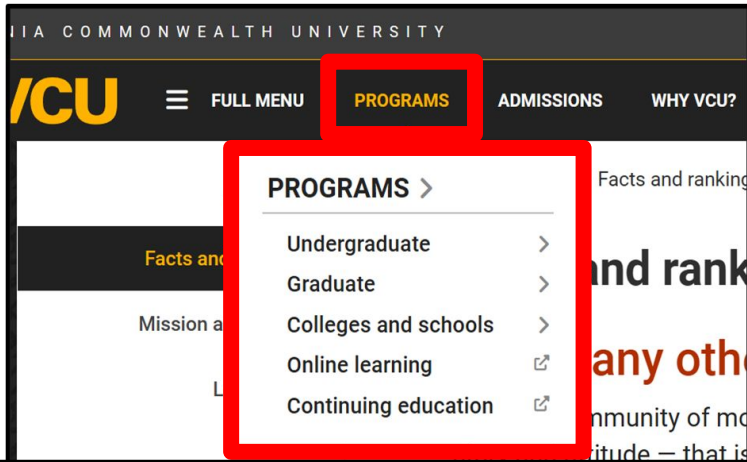# However, JS rewrites can lead to fidelity issues
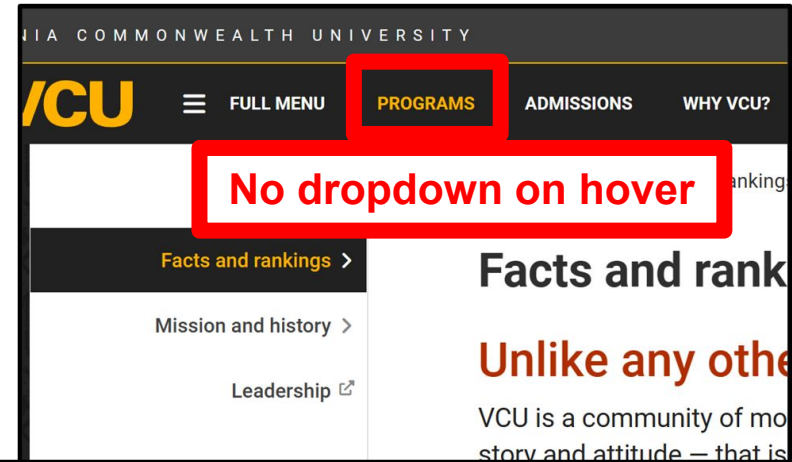
## Live web page



## Archived copy



**No dropdown on hover**

# However, JS rewrites can lead to fidelity issues

Live web page

Archived copy



No dropdown on hover

```
Uncaught ReferenceError: Logger is not defined                facts-and-rankings/:4486
    at new navFullMenu (facts-and-rankings/:4486:23)
    at facts-and-rankings/:4890:13
```

# Goals of our project

- Detect fidelity violations when replaying page snapshots

    - Focus only on fidelity violations caused due to rewriting

    - We do not address incomplete/inaccurate recording

# Goals of our project

- Detect fidelity violations when replaying page snapshots

    - Focus only on fidelity violations caused due to rewriting

    - We do not address incomplete/inaccurate recording

- Diagnose erroneous rewrites which lead to fidelity issues

    - Not trying to *fix* fidelity issues

    - Once diagnosed, developers will need to identify necessary fixes

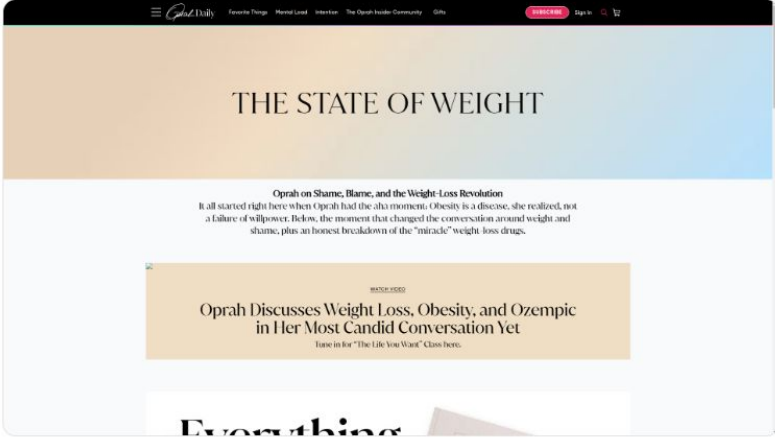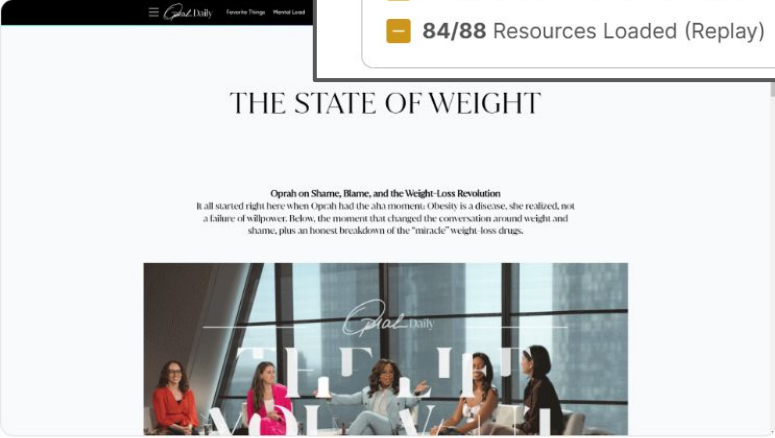# Existing Quality Assurance in Web Archiving

# Screenshots Differ across Multiple Loads of Same Page

# Runtime Error ⇏ Fidelity Violation

Hard to analyze if an error will cause violations
by manually inspecting its stack trace



15

# Fidelity Violation ⇏ Runtime Error

Errors can be <span style="color:red">hard to detect</span>

- Semantic errors, no error messages

# Fidelity Violation ⇏ Runtime Error

Errors can be hard to detect

- Semantic errors, no error messages

# Fidelity Violation ⇏ Runtime Error

Errors can be hard to detect

● Semantic errors, no error messages



Fail to load

FidEx (**Fid**elity violation **Ex**poser)



Live web
server





Replay server

FidEx (**Fid**elity violation **Ex**poser)



1. Crawl

Live web server

WARC files

Replay server

# FidEx (**Fid**elity violation **Ex**poser)



1. Crawl

Live web server

2. Load Archived Page

WARC files

Replay server

# FidEx (**Fid**elity violation **Ex**poser)

# FidEx (**Fid**elity violation **Ex**poser)



Instrumentation

Live web server

1. Crawl

WARC files

2. Load Archived Page

Replay server

Instrumented data

Detection

3. Pages with poor fidelity

# FidEx (**Fid**elity violation **Ex**poser)

**Instrumentation**

Live web server

1. Crawl

WARC files

2. Load Archived Page

Replay server

**Instrumented data**

**Detection**

**3. Pages with poor fidelity**

**Pinpointing**

**4. Problematic rewrites**

# Fidelity violation detection



**Instrumentation**

1. Crawl

Live web

**Instrumented data**

**Detection**

**3. Pages with poor fidelity**

**4. Problematic rewrites**

Replay server

- Detect fidelity violations accurately
  - **Collect better representations for matching pages**
  - Collect representations at the right timing

# How a browser renders pages

# How a browser renders pages

## 1. Fetch

```
          ┌──────────┐
     ┌───▶│   HTML   │
     │    └──────────┘
┌─────────┐    ┌────────────┐
│ Network │───▶│ JavaScript │
└─────────┘    └────────────┘
     │    ┌──────────┐
     └───▶│   CSS    │
          └──────────┘
```

# How a browser renders pages

**1. Fetch**     **2. Build**

Network → HTML → DOM

Network → JavaScript

Network → CSS → CSSOM

# How a browser renders pages

**1. Fetch**          **2. Build**

Network → HTML → DOM

Network → JavaScript

Write (JavaScript → DOM)

Write (JavaScript → CSSOM)

Network → CSS → CSSOM

# How a browser renders pages

**1. Fetch** | **2. Build** | **3. Layout**

Network → HTML → DOM

Network → JavaScript

Network → CSS → CSSOM

Write (DOM ↔ JavaScript)

Write (JavaScript ↔ CSSOM)

DOM → Layout

CSSOM → Layout

JavaScript → Layout

# How a browser renders pages

**1. Fetch** | **2. Build** | **3. Layout** | **4. Paint**

Network → HTML → DOM

Network → JavaScript

Write (DOM ↑)

Write (↓ CSSOM)

Network → CSS → CSSOM

JavaScript → Layout → Paint

DOM, CSSOM → Layout

31

# How a browser renders pages

**1. Fetch**  **2. Build**  **3. Layout**  **4. Paint**

Network → HTML → DOM

JavaScript

Write

Write

CSS → CSSOM

DOM, CSSOM → Layout → Paint

Screenshot

Capture

Paint

32

# How a browser renders pages

**1. Fetch**

**2. Build**

**3. Layout**

**4. Paint**

Network

HTML

JavaScript

CSS

DOM

Write

Write

CSSOM

Layout

Screenshot

Capture
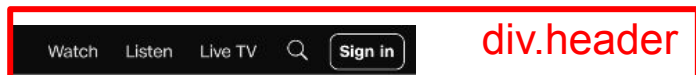
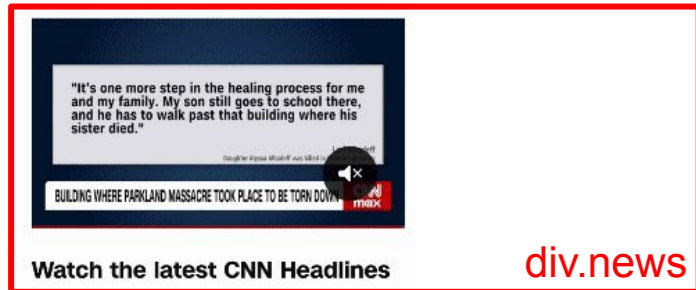Paint

# Fidelity violation detection: Match layout tree

# Fidelity violation detection: Match layout tree



div.header

div.links

div.news

div.news

# Fidelity violation detection: Match layout tree



div.header

div.links

video

div.news

img

div.news

# Fidelity violation detection: Match layout tree



div.header

div.links

video

div.news

img

div.news

**Layout tree unchanged**

# Fidelity violation detection: Match JS writes

# Fidelity violation detection: Match JS writes



```
<span>
{text}
</span>
```

# Fidelity violation detection: Match JS writes



```
<span>
{text}
</span>
```

```
for (var n = function(e) {
    setTimeout(function() {
        t.inputPlaceholder.textContent = t.searchTerm.substr(0, e)
    }, t.attractLoopTypeRate * e)
}, r = 1; r < this.searchTerm.length + 1; r += 1)
    n(r)
```

40

# Fidelity violation detection: Match JS writes



```
<span>
{text}
</span>
```

```
for (var n = function(e) {
    setTimeout(function() {
        t.inputPlaceholder.textContent
    }, t.attractLoopTypeRate * e)
}, r = 1; r < this.searchTerm.length +
    n(r)
```

| | |
|---|---|
| ➤ (anonymous) | main.js?siolrj:1 |
| **setTimeout** | |
| n | main.js?siolrj:1 |
| value | main.js?siolrj:1 |
| value | main.js?siolrj:1 |
| value | main.js?siolrj:1 |
| (anonymous) | main.js?siolrj:1 |
| **setInterval** | |
| (anonymous) | main.js?siolrj:1 |
| **setTimeout** | |
| (anonymous) | main.js?siolrj:1 |
| **setTimeout** | |
| value | main.js?siolrj:1 |
| t | main.js?siolrj:1 |
| (anonymous) | main.js?siolrj:1 |
| u | main.js?siolrj:1 |
| fireWith | main.js?siolrj:1 |
| r | main.js?siolrj:1 |
| (anonymous) | main.js?siolrj:1 |
| **load** | |
| send | main.js?siolrj:1 |
| ajax | main.js?siolrj:1 |

41

# Fidelity violation detection: Match JS writes

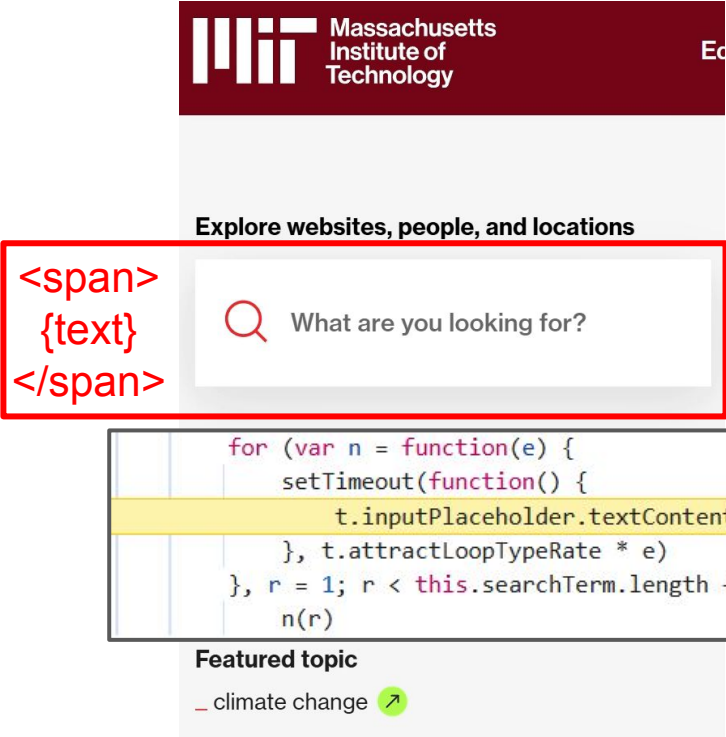<span>
{text}
</span>

```
for (var n = function(e) {
    setTimeout(function() {
        t.inputPlaceholder.textContent
    }, t.attractLoopTypeRate * e)
}, r = 1; r < this.searchTerm.length +
    n(r)
```

**Associated {JS writes} unchanged**

# Fidelity violation diagnosis

- Pinpoint causes of violations
  - **Discover** errors not reported by browser
  - Identify if an error truly matters to fidelity

3. Pages with poor fidelity

ection

WARC files

Archived Page

Replay server

Pinpointing

4. Problematic rewrites

# Evaluation

- **Datasets**
  - 80K pages sampled from top 1 million sites, 5 pages per site
  - Pages sampled from End-of-Term and CARTA collections

# Evaluation

- **Datasets**

  - 80K pages sampled from top 1 million sites, 5 pages per site

  - Pages sampled from End-of-Term and CARTA collections


- **Baselines for detecting fidelity violations**

  - Match screenshots, HTML extracted texts (Browsertrix's QA)

  - Check if archive has extra errors
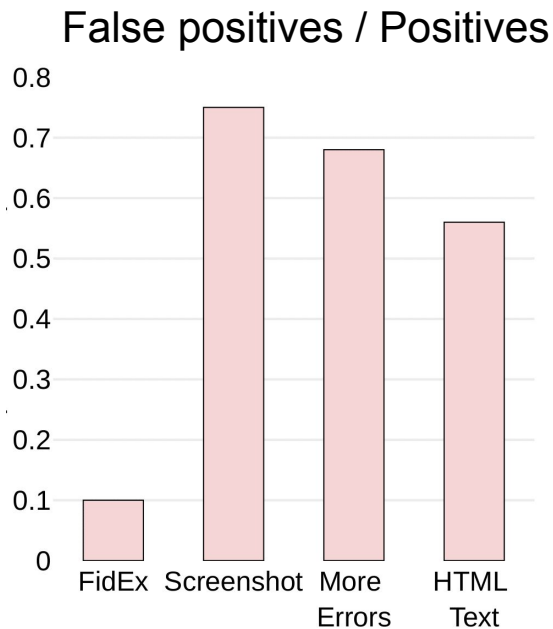
# FidEx improves accuracy of detecting fidelity violations

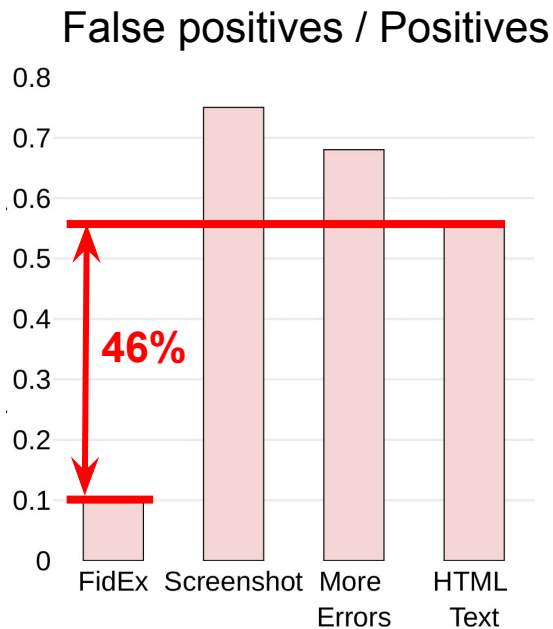| | FidEx | Screenshots | More Errors | HTML Text |
|---|---|---|---|---|
| Positive rate | 15.5% | 72.6% | 43.7% | 24.9% |

*Positive: detected fidelity violation(s)*

# FidEx improves accuracy of detecting fidelity violations

|  | FidEx | Screenshots | More Errors | HTML Text |
|---|---|---|---|---|
| Positive rate | 15.5% | 72.6% | 43.7% | 24.9% |

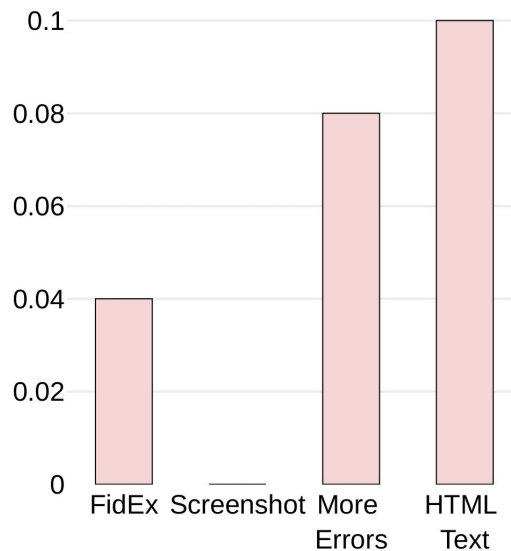*Positive: detected fidelity violation(s)*

### False positives / Positives

# FidEx improves accuracy of detecting fidelity violations

|  | FidEx | Screenshots | More Errors | HTML Text |
|---|---|---|---|---|
| Positive rate | 15.5% | 72.6% | 43.7% | 24.9% |

*Positive: detected fidelity violation(s)*

False positives / Positives

# FidEx improves accuracy of detecting fidelity violations

| | FidEx | Screenshots | More Errors | HTML Text |
|---|---|---|---|---|
| Positive rate | 15.5% | 72.6% | 43.7% | 24.9% |

*Positive: detected fidelity violation(s)*

### False positives / Positives



**46%**

### False negatives / Negatives



49

# FidEx improves accuracy of detecting fidelity violations

| | FidEx | Screenshots | More Errors | HTML Text |
|---|---|---|---|---|
| Positive rate | 15.5% | 72.6% | 43.7% | 24.9% |

*Positive: detected fidelity violation(s)*

## False positives / Positives



**46%**

## False negatives / Negatives



**If classify every page as negative, false negative rate will be ~15%**

# Fidelity Violations Detected by FidEx in EOT

## Live web page

## Archived copy

# Fidelity Violations Detected by FidEx in EOT

## Live web page



## Archived copy



**Unclickable**

# False Positives with Screenshot Comparison in EOT

# Fidelity Violations Detected by FidEx in CARTA

## Live web page



## Archived copy
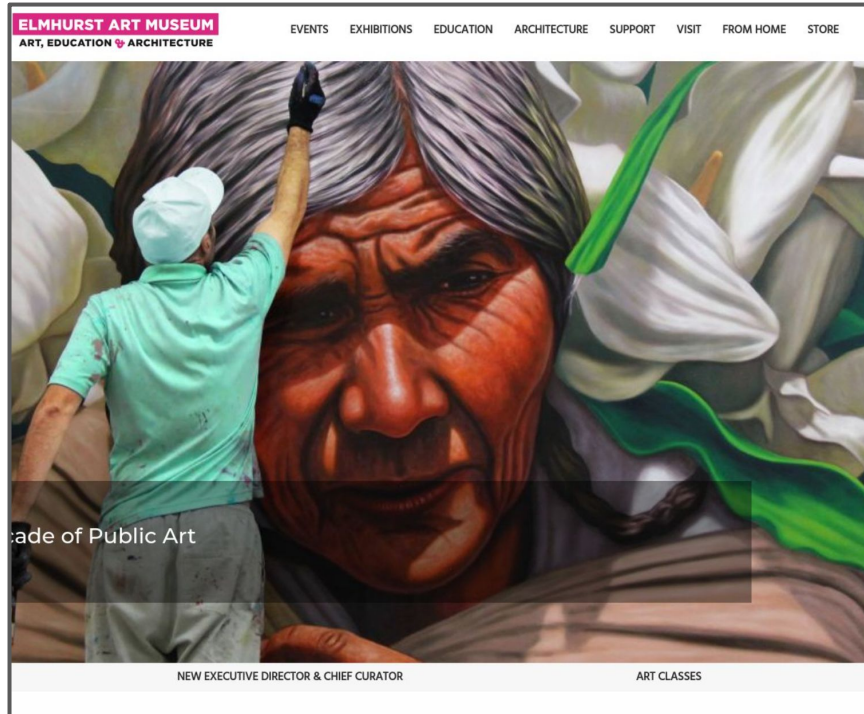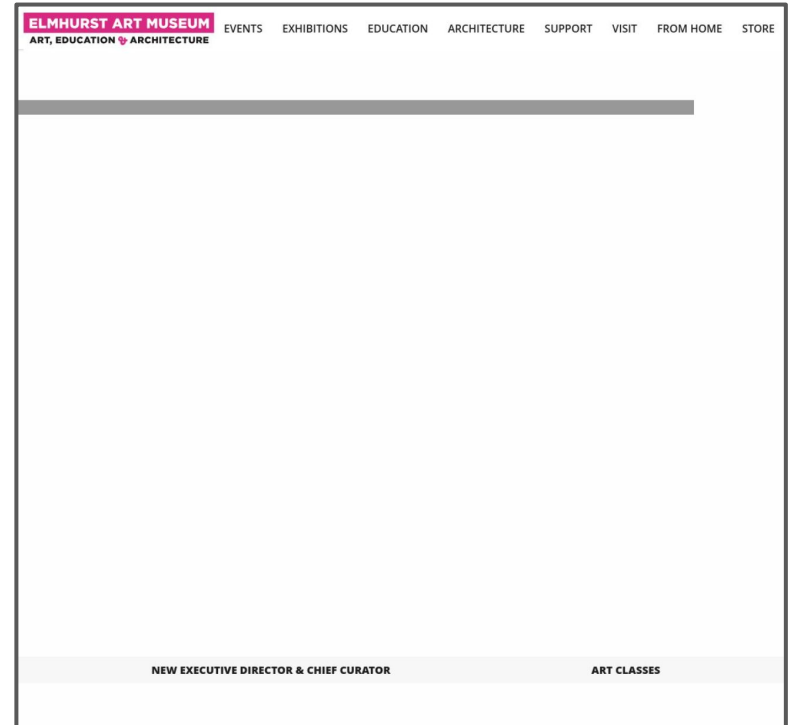
# Fidelity Violations Detected by FidEx in CARTA
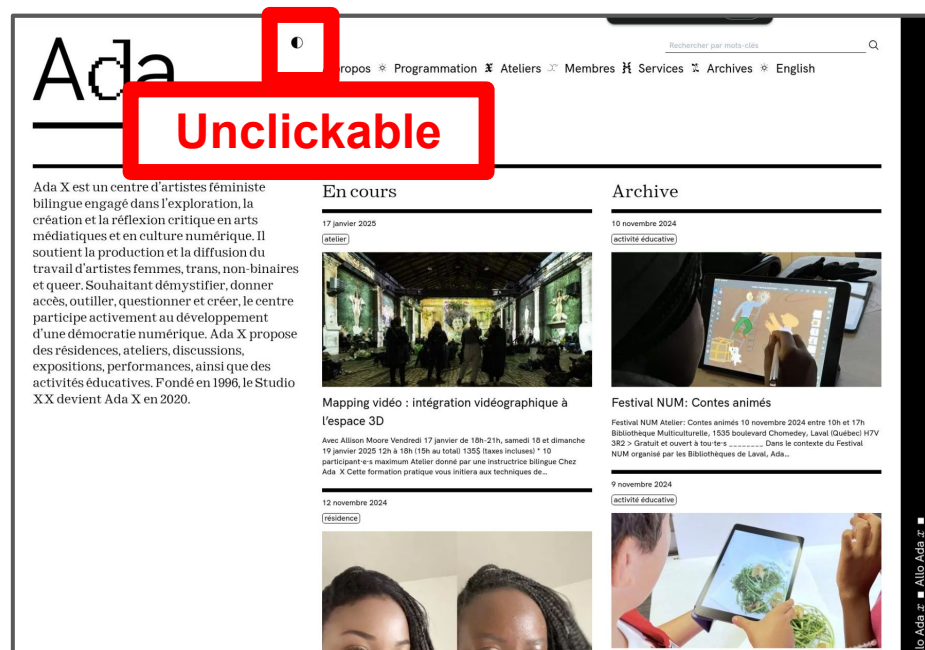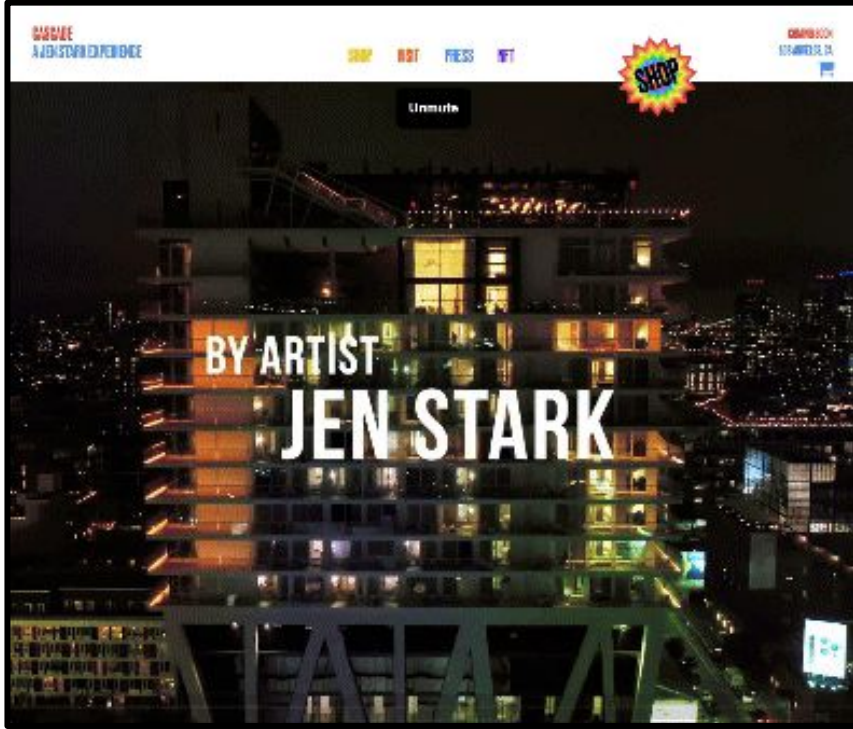
## Live web page

## Archived copy



Unclickable

# False Positives with Screenshot Comparison in CARTA

# Summary

- Bugs in JS rewrites → <span style="color:red">Poor fidelity replay</span> of archived pages

# Summary

- Bugs in JS rewrites → Poor fidelity replay of archived pages

**Detection**

- Screenshot comparison leads to false positives

- Instead, compare layout tree and JS writes

58

# Summary

- Bugs in JS rewrites → Poor fidelity replay of archived pages

| **Detection** | **Diagnosis** |
|---|---|
| ● Screenshot comparison leads to false positives<br><br>● Instead, compare layout tree and JS writes | ● Discover errors not reported by browser<br><br>● Confirm which errors are worth investigating |

Thank you!
Q & A