

Jawa: Web Archival in the Era of JavaScript

Ayush Goel, Jingyuan Zhu, Ravi Netravali, and Harsha V. Madhyastha



Objective: High fidelity and low cost web archive

<p>Ephemeral web and link rot → Web Archives</p> <p>RESEARCH ARTICLE Scholarly Context Adrift: Three out of Four URI References Lead to Changed Content</p> <p>Shawn M. Jones^{1*}, Herbert Van de Sompel^{1*}, Harihar Shankar^{1*}, Martin Klein^{1*}, Richard Tobin^{2*}, Claire Grover^{2*}</p> <p>DASH DIGITAL ACCESS to SCHOLARSHIP at HARVARD DASH.HARVARD.EDU HARVARD LIBRARY Office for Scholarly Communication</p> <p>The Paper of Record Meets an Ephemeral Web: An Examination of Linkrot and Content Drift within The New York Times</p>		<h3>Problems with Web Archives</h3> <table border="1"> <tr> <td data-bbox="968 674 1473 1004"> <h4>Limited page snapshots</h4> <p>2010-2020: 4x increase in JavaScript</p> <p>↓</p> <p>Per page snapshot expensive</p> </td> <td data-bbox="1473 674 2044 1004"> <h4>Poor Page Fidelity</h4> <p>JavaScript non-determinism: Date, math.random, user-agent</p> </td> </tr> </table>		<h4>Limited page snapshots</h4> <p>2010-2020: 4x increase in JavaScript</p> <p>↓</p> <p>Per page snapshot expensive</p>	<h4>Poor Page Fidelity</h4> <p>JavaScript non-determinism: Date, math.random, user-agent</p>
<h4>Limited page snapshots</h4> <p>2010-2020: 4x increase in JavaScript</p> <p>↓</p> <p>Per page snapshot expensive</p>	<h4>Poor Page Fidelity</h4> <p>JavaScript non-determinism: Date, math.random, user-agent</p>				

Improve fidelity: Fix sources of JavaScript variation

<h3>Randomness (date, math.random, performance)</h3>		<h3>Client characteristics</h3>	
<h4>Insight</h4>	<h4>Solution</h4>	<h4>Insight</h4>	<h4>Solution</h4>
<p>1-1 mapping between diverged URLs</p>	<p>Use server-side matching techniques</p> <p>a.com/b.js?ts=5467</p> <p>↕</p> <p>a.com/b.js?ts=8967</p>	<p>Key contributor to URL divergence.</p>	<p>Enforce same values across loads</p> <p>Crawling device</p> <p>↕</p> <p>Replay device</p> <p>APIs:</p> <ul style="list-style-type: none"> - User-agent - Screen size - Operating system - Geolocation -

Reduce storage overhead: Remove non-functional code

<h3>No back-end origin server</h3>	<h3>Insights</h3>		<h3>Solution</h3>
<p>Remove code that interacts with back-end server</p>	<p>1 Compartmentalized into files</p> <p>2 Hosted on third-party domains</p>		<p>Filter list to identify URLs of non-functional scripts</p> <p>Rules:</p> <ul style="list-style-type: none"> - "cdn.onesignal.com" - "cdn.pushly.com" - "cdn.parsley.com" - ...

Reduce storage overhead: Remove unused code

<h3>Remove (unused) alternate code flows</h3>	<h3>Preserve code for event handlers</h3>	
<p>Insight: Absent sources of non-determinism:</p> <ul style="list-style-type: none"> - Server-side state - Client APIs return values - Client-side state <p>Solution: Dynamically identify used code (during page load)</p>	<h4>Insight</h4> <p>Order of execution, and inputs to events do not impact code coverage</p>	<h4>Solution</h4> <p>Trigger each event once with default inputs</p> <p>+ Track executed lines of code</p>

Results

<h3>Storage</h3> <p>41% reduction in total storage overhead</p>	<h3>Fidelity</h3> <p>Eliminate almost all failed network fetches</p>	<h3>Crawling throughput</h3> <p>Improved throughput by 31%</p>
---	--	--