# To push, or not to push?!

A journey of resource loading in the browser

Fluent Conference, June 2018

Patrick Hamann

@patrickhamann

The New York Times

reddit

airbnb

deliveroo

ATRESMEDIA

stripe

lonely planet

Pinterest

La Redoute

wallapop

wayfair
a zillion things home

ticketmaster

the guardian

imgur

Rakuten.com

Boots

A&E

GitHub

SHAZAM

FAST COMPANY

BBC

CONDÉ NAST

BuzzFeed

RollingStone

BUSINESS INSIDER

KICKSTARTER

Hotel Tonight

FOURSQUARE

W
WENNER MEDIA

New Relic

imgIX

GOV.UK

vimeo

shopify

marfeel

KAYAK

Etsy

# Why?

# "HTTP/2 will solve this"
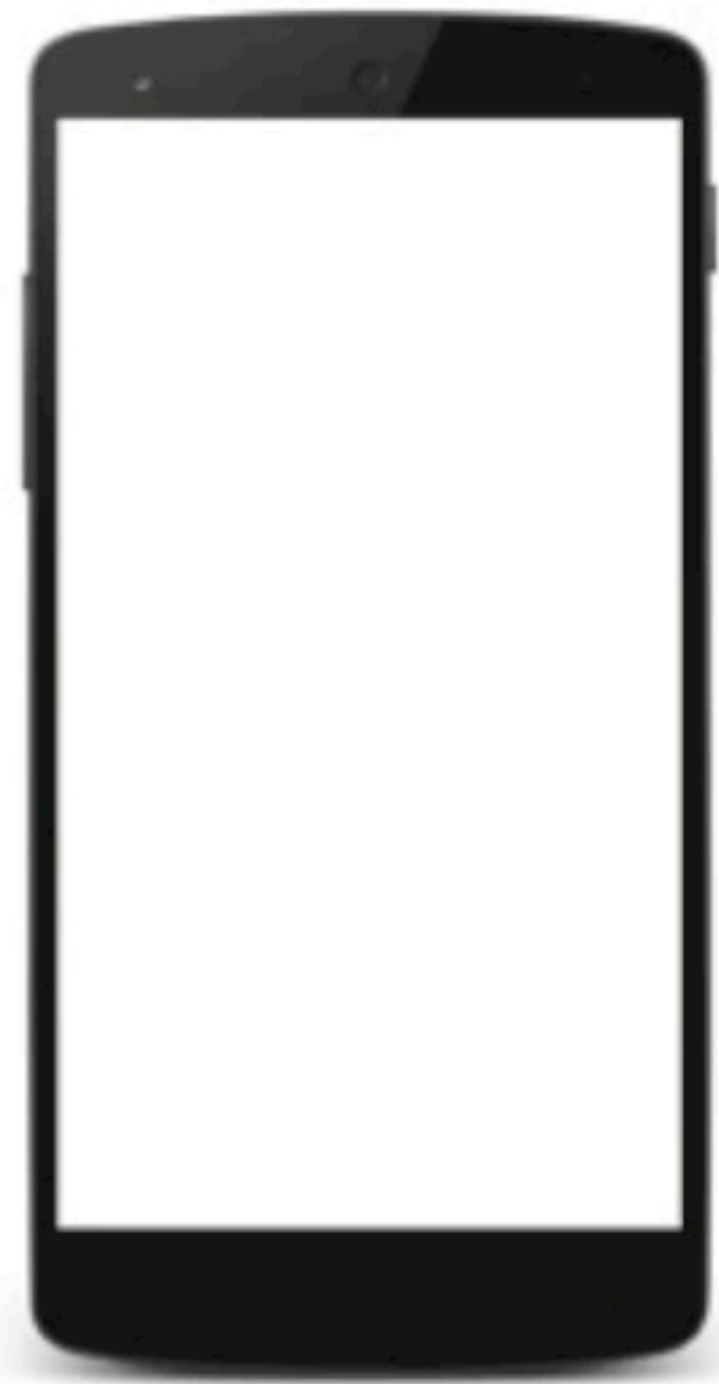
– Everybody

@patrickhamann

# Resource loading in the browser is hard.

# Resource loading is hard:

⏱️ Performance is tightly coupled to latency

🤝 Connection cost is high

📈 Congestion control is unavoidable

🙈 Critical resources can be hidden

💤 Bandwidth is often under-utilised

⚠️ Script execution is expensive

# How can we load our resources most efficiently?

@patrickhamann

A **critical request** is one that contains an asset that is **essential** to the **content** within the users viewport.

– Ben Schwarz, Calibre

# What are my critical resources?

✅ Critical CSS for current route

✅ Fonts

✅ Hero images

✅ Initial application route

✅ Application bootstrap data

**First Contentful Paint**

**Time to Interactive**

**User navigates**

**First Meaningful Paint**

**Fully loaded**

**First Contentful Paint**

**Time to Interactive**

**User navigates**

**First Meaningful Paint**

**Fully loaded**

# A good loading strategy:

✅ Prioritises above-the-fold rendering

✅ Prioritises interactivity

✅ Is easy to use

✅ Is measurable

# Preload

**Renderer**

| Request page | idle | Build DOM | idle | Build CSSOM | Render tree | First paint | Text paint |

**Network**

| GET html | response |

Render blocking

| GET css | response |

Text blocking

| GET font | response |

@patrickhamann

# What are my hidden sub-resources?

✅ Fonts

✅ Application data

✅ Application routes

✅ Async third parties

Provides a declarative fetch primitive that initiates an early fetch and separates fetching from resource execution.

Preload with HTTP header:

```
1  Link: <my-awesome-font.woff>; rel=preload; as=font; crossorigin
2  Link: <application-data.json>; rel=preload; as=fetch;
3  Link: <sub-module.mjs>; rel=modulepreload;
```

Preload with markup:

```
1  <!-- preload stylesheet resource via declarative markup -->
2  <link rel="preload" href="/styles.css" as="style">
3
4  <!-- or, preload stylesheet resource via JavaScript -->
5  <script>
6      const res = document.createElement("link");
7      res.rel = "preload";
8      res.as = "style";
9      res.href = "lazy-loaded-styles.css";
10     document.head.appendChild(res);
11 </script>
```

@patrickhamann

**Legend (top):**
dns · connect · ssl · html · js · css · image · flash · font · other · JS Execution

**Step_1**

| # | Resource | Time |
|---|----------|------|
| 1. | ttfmp.fastlylabs.com – index.html | 984 ms |
| 2. | ttfmp.fastlylabs.com – main.css | 411 ms |
| 3. | www.ft.com – h...-markets-data.png | 928 ms |
| 4. | www.ft.com – ftlogo:brand-myft | 210 ms |
| 5. | www.ft.com – fticon-v1:hamburger | 196 ms |
| 6. | www.ft.com – ftlogo:brand-ft | 215 ms |
| 7. | www.ft.com – f...brand-ft-masthead | 200 ms |
| 8. | www.ft.com – fticon-v1:search | 274 ms |
| 9. | www.ft.com – fticon-v1:arrow-up | 252 ms |
| 10. | www.ft.com – f...nd-nikkei-tagline | 259 ms |
| 11. | www.ft.com – fticon-v1:arrow-left | 246 ms |
| 12. | www.ft.com – fticon-v1:cross | 268 ms |
| 13. | www.ft.com – fticon-v1:arrow-right | 215 ms |
| 14. | www.ft.com – fticon-v1:play | 240 ms |
| 15. | www.ft.com – fticon-v1:play | 234 ms |
| 16. | www.ft.com – fticon-v1:arrow-right | 228 ms |
| 17. | www.ft.com – M...cWeb-Regular.woff | 752 ms |
| 18. | www.ft.com – F...yWeb-Regular.woff | 1522 ms |
| 19. | www.ft.com – M...Web-Semibold.woff | 931 ms |
| 20. | next-geebee.ft... – polyfill.min.js | 876 ms |
| ... | | |

**Before**

**After**

"

Shopify's switch to preloading fonts saw a 50% (1.2 second) improvement in time-to-text-paint. This removed their flash-of-invisible text completely.

– Shopify

@patrickhamann

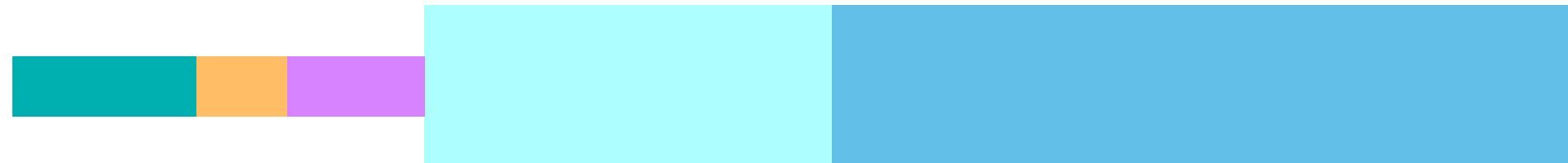# Preconnect

# No preconnect

Time

**index.html**

**main.css**
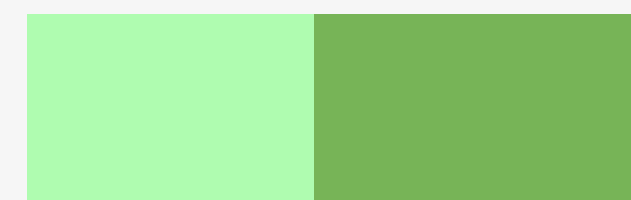
**app.js**

**font.woff**

Time

**index.html**

**main.css**

**app.js**

**font.woff**

# Preconnect

# Are indicating resource hints via the HTML response too late?

# Server push

Client      CDN/Surrogate/Server      Origin

Time

**GET** /index.html

**GET** /index.html

Server think time

200 /index.html

200 /index.html

**GET** main.css

**GET** /main.css

**Page**

Hey example.com, can I have your homepage please? 10:24

**Server**

Sure thing! Oh, but while I'm sending you that, here's a stylesheet, some images, some JavaScript, and some JSON. 10:24

**Page**

Uh, sure. 10:24

**Page**

I'm just reading the HTML here, and it looks like I'm going to need a stylesh… oh it's the one you're already sending me, cool! 10:25

**Stream**
A virtual channel within an established connection which carries bidirectional messages.

**Message**
A complete sequence of frames that map to a logical HTTP message, such as a request.

Connection

Stream

Message

Frame

```
:method: GET
:path: /image-2.jpg
```

Frame

```
:status 200
:version: HTTP/2.0
:vary: Accept-Encoding
```

Frame

```
… response payload
```

Client

Server

**Frame**
The smallest unit of communication, which carries a specific type of data—e.g., HTTP headers, payload, commands e.t.c.

# Multiplexing

@patrickhamann

**Client**

**CDN/Surrogate/Server**

**Origin**

Time

**GET** /index.html

**GET** /index.html

Server think time

200 /index.html
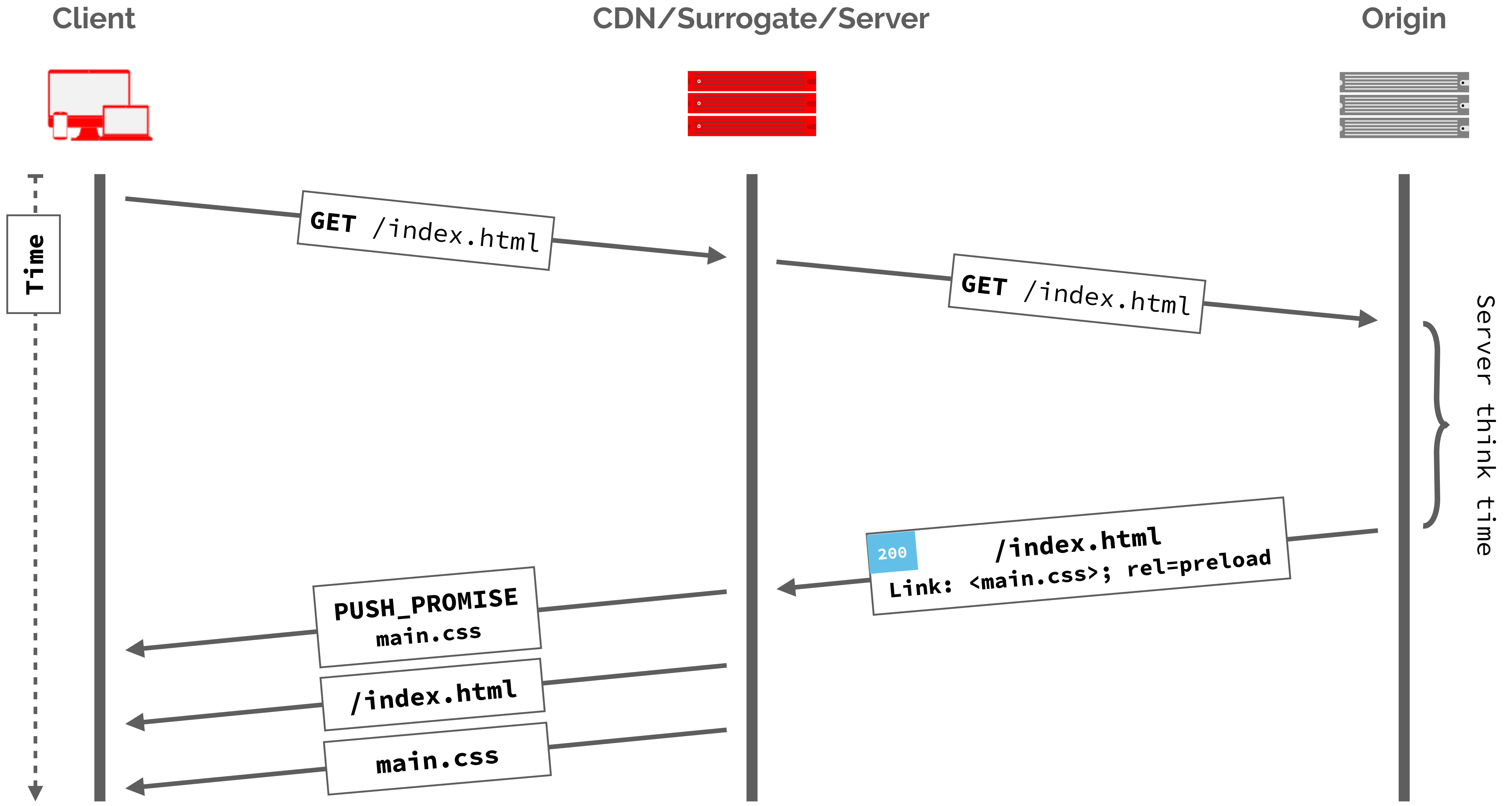Link: <main.css>; rel=preload

**PUSH_PROMISE**
main.css

/index.html

main.css

# So how can I push?

```
1   Link: <font.woff2>; rel=preload; as=font crossorigin
```

*Indicate push via preload Link header.*

```
1   Link: <main.css>; rel=preload; as=style; nopush
```

*Use no push attribute to disable push semantics and only use preload.*

```
1   Link: <application.js>; rel=preload; as=style; x-http2-push-only
```

*Fastly uses x-http2-push-only attribute to disable preload semantics*

@patrickhamann

Time

index.html

main.css

app.js

font.woff

**Before**

No Push

Time

index.html
main.css
app.js
font.woff

1 RTT saving!

Time

index.html
main.css
app.js
font.woff

Push

Time

index.html — Idle

main.css

app.js

font.woff

After

**Client**  **CDN/Surrogate/Server**  **Origin**

Time

GET /index.html

GET /index.html

Server think time

200 /index.html
Link: <main.css>; rel=preload

PUSH_PROMISE
main.css

/index.html

main.css

# Server push benefits:

✅ 1 RTT saving

✅ Useful for long server think time

✅ Useful for long RTT times

⚠️ Link header indication is too late

# Is indicating push via the HTML response too late?

# Async push

Client      CDN/Surrogate/Server      Origin

Time

**GET** /index.html

**GET** /index.html

**PUSH_PROMISE**
main.css

main.css

😎

Server think time

200 /index.html

/index.html

```javascript
const http2 = require('http2');

function handler(request, response) {
  if (request.url === "/index.html") {
    const push = response.push('/critical.css');
    push.writeHead(200);
    fs.createReadStream('/critical.css').pipe(push);
  }

  // Generate index response:
  // - Fetch data from DB
  // - Render template
  // etc ...

  response.end(data);
}

const server = http2.createServer(opts, handler);
server.listen(80);
```

@patrickhamann

```
1  sub vcl_recv {
2    if (fastly_info.is_h2 && req.url ~ "^/index.html") {
3      h2.push('/critical.css');
4    }
5
6    // etc ...
7
8  }
```

@patrickhamann

# Push

Time

index.html

main.css

app.js

font.woff

# Async push

Time

index.html

main.css

app.js

font.woff

Time

index.html

main.css

app.js

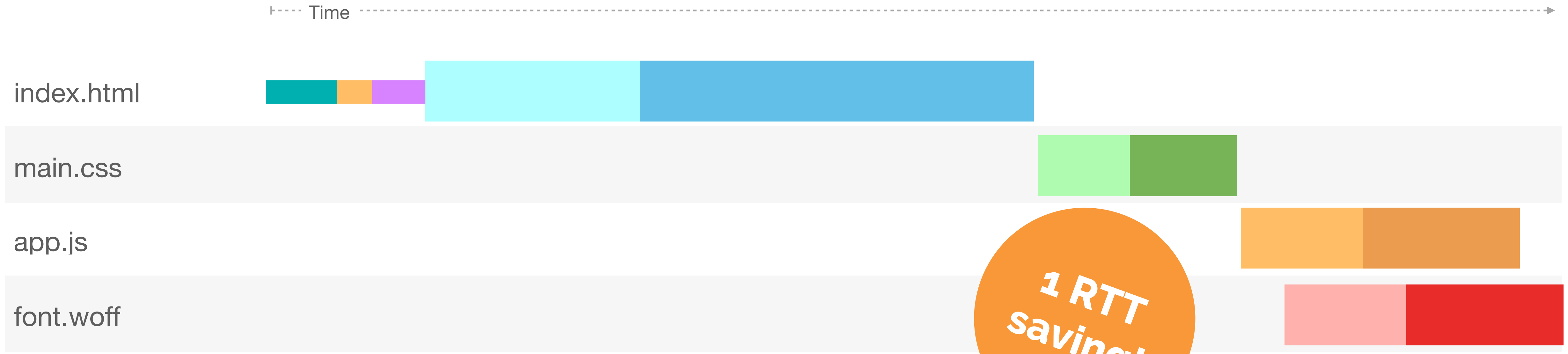font.woff

**Utilising idle network server think time == win!**

# What about the repeat view?

First view

Time

index.html

main.css

app.js

font.woff

Time

index.html

main.css

app.js          (from disk cache)

font.woff       (from disk cache)

Repeat view

# The server has no knowledge of client cache state.

# In the wild

"

Faster image loads times, **15%** reduction in time to first byte

– Facebook

HTTP2 server push - Facebook https://atscaleconference.com/videos/http2-server-push-lower-latencies-around-the-world/

# PRPL Pattern

# Poll?

# So what's the problem?

Client                          CDN/Surrogate/Server

Time

GET /index.html

PUSH_PROMISE
main.css

main.css

RST_STREAM
main.css

@patrickhamann

Client                          CDN/Surrogate/Server

Time

GET /index.html

PUSH_PROMISE
main.css

main.css

RST_STREAM
main.css

Separate push cache per
HTTP/2 connection.

Page

Memory cache

Service Worker

HTTP cache

Push cache

Push cache

Server

Page

Memory cache

@patrickhamann

# Push cache semantics

⚠️ Connection must be authoritative

⚠️ Cache per HTTP/2 connection

⚠️ Items can only be claimed once

⚠️ It's the last cache

⚠️ It's not spec'd

# HTTP/2 push is tougher than I thought

Posted 30 May 2017

"HTTP/2 push will solve that" is something I've heard a lot when it comes to page load performance problems, but I didn't know much about it, so I decided to dig in.

HTTP/2 push is more complicated and low-level than I initially thought, but what really caught me off-guard is how inconsistent it is between browsers – I'd assumed it was a done deal & totally ready for production.

This isn't an "HTTP/2 push is a douchebag" hatchet job – I think HTTP/2 push is really powerful and will improve over time, but I no longer think it's a silver bullet from a golden gun.

## Map of fetching

Between your page and the destination server there's a series of caches & things that can intercept the request:

Hello, I'm Jake and that is my face. I'm a developer advocate for Google Chrome.

**Elsewhere**

Twitter          Lanyrd

Github          Google+

Flickr

**Contact**

Feel free to throw me an email, unless you're a recruiter, in which case destroy every email-capable device you own to prevent this possibility.



HTTP/2 push is tougher than I thought – Jake Archibald https://jakearchibald.com/2017/h2-push-tougher-than-i-thought/

HTTP/2 Server Push - Browser inconsistencies

**0.008%** of requests on the Fastly network are push initiated.

# When should I push?

✅ You have long RTTs or server processing

✅ You can use async push

✅ You have a client-rendered app shell (PRPL)

✅ You control the client cache (SW, native, Electron etc)

# Is the 1 RTT saving worth the complexity?

# Are there other solutions?

@patrickhamann

# The future

# Can we fix the problems with push?

# Cache digests

Client

CDN/Surrogate/Server

Time

GET /index.html

CACHE_DIGEST

PUSH_PROMISE
main.css

main.css

/index.html

@patrickhamann

Time

index.html

main.css

app.js

font.woff

**First view**

Time

index.html

main.css     (from disk cache)

app.js       (from disk cache)

font.woff    (from disk cache)

**Repeat view**

Time

index.html

main.css

app.js

font.woff

Time

index.html

main.css        (from disk cache)

app.js          (from disk cache)

font.woff       (from disk cache)

Repeat view

### Cache Digests for HTTP/2
### draft-ietf-httpbis-cache-digest-04

Abstract

   This specification defines a HTTP/2 frame type to allow clients to
   inform the server of their cache's contents.  Servers can then use
   this to inform their choices of what to push to clients.

Note to Readers

   Discussion of this draft takes place on the HTTP working group
   mailing list (ietf-http-wg@w3.org), which is archived at
   https://lists.w3.org/Archives/Public/ietf-http-wg/ .

   Working Group information can be found at http://httpwg.github.io/ .

IETF Draft Cache Digests for HTTP/2 - K. Oku, Y. Weiss https://tools.ietf.org/html/draft-ietf-httpbis-cache-digest-04

# This still seems too complicated...

# 103 Early hints

The 103 (Early Hints) informational status code indicates to the client that the server is likely to send a final response with the header fields included in the informational response.

Client      CDN/Surrogate/Server      Origin

Time

**GET** /index.html

**GET** /index.html

**103** **Link:** <main.css>; rel=preload

**GET** /main.css

Server think time

main.css

**200** /index.html

/index.html

@patrickhamann

```
HTTP/1.1 103 Early Hints
Link: </style.css>; rel=preload; as=style
Link: </main.js>; rel=preload; as=script
Link: </application-data.json>; rel=preload; as=fetch

HTTP/1.1 200 OK
Date: Thu, 14 June 2018 16:10:00 PDT
Content-length: 1234
Content-type: text/html; charset=utf-8
Link: </main.css>; rel=preload; as=style
Link: </newstyle.css>; rel=preload; as=style
Link: </main.js>; rel=preload; as=script
```

@patrickhamann

## An HTTP Status Code for Indicating Hints

Abstract

   This memo introduces an informational HTTP status code that can be
   used to convey hints that help a client make preparations for
   processing the final response.

Status of This Memo

   This document is not an Internet Standards Track specification; it is
   published for examination, experimental implementation, and
   evaluation.

   This document defines an Experimental Protocol for the Internet
   community.  This document is a product of the Internet Engineering

IETF RFC8297 Early Hints - K. Oku https://tools.ietf.org/html/rfc8297

## 103 Early hints:

✅ Same benefits as push

✅ Much simpler

✅ Leverages the browser caches

✅ Allows client to initiate fetches

⚠️ No 1 RTT saving

# Priority hints

```html
<link rel="preload" as="script" href="critical-script.js">
<link rel="preload" as="style" href="theme.css" importance="low" onload="this.rel=stylesheet">

<style>/* critical-path styles */</style>
<img src="hero.jpg" importance="high">
<img src="meme.gif" importance="low">

<!-- superfluous fetch requests -->
<script>
  fetch('/api/related.json', { importance: 'low' });
</script>

<!-- scripts at the end of the document -->
<script src="critical-script.js"></script>
```

github.com/WICG/priority-hints

Features    Business    Explore    Marketplace    Pricing        This repository    Search        **Sign in** or **Sign up**

📖 **WICG** / **priority-hints**

👁 Watch    24        ⭐ Star    76        🍴 Fork    8

<> Code        ⓘ Issues **9**        Pull requests **0**        Projects **0**        Insights

Branch: **master** ▾    **priority-hints** / **EXPLAINER.md**        Find file    Copy path

**addyosmani** Explainer: drop requirement for custom importance groups        5d7720e on Nov 6, 2017

1 contributor

54 lines (42 sloc)    2.53 KB        Raw    Blame    History    ✏️    🗑

# Priority Hints

The browser's resource loading process is a complex one. Browsers discover needed resources and download them according to their heuristic priority. Browsers may also use this heuristic resource priority to delay sending certain requests in order to avoid bandwidth contention of these resources with more critical ones.

Currently web developers have very little control over the heuristic importance of loaded resources, other than speeding up their discovery using `<link rel=preload>`. Browsers make many assumptions on the importance of resources based on the

# Closing

# HTTP/2 doesn't solve everything.

# Resource loading is hard.

# Performance is for humans. Optimise for user experiences.

@patrickhamann

# The future is bright!

# Resource loading checklist:

✅ Identify your critical resources

✅ Preload hidden sub-resources

✅ Preconnect critical third-parties

❌ Avoid pushing with preload

⚠️ Use async push with care

🚀 Decorate HTML with priority hints

🚀 Use Early Hints when available

# Thanks!

Patrick Hamann
speakerdeck.com/patrickhamann

patrick@fastly.com
@patrickhamann