# PARENTAL
# ADVISORY
## EXPLICIT CONTENT

@badams

# What is Technical SEO?

# Google Processes



Crawler

Indexer

Ranker

@badams

POLEMIC
DIGITAL

# Google Processes



Crawler

Indexer

Ranker

Technical SEO

@badams

POLEMIC
DIGITAL

**Kevin_Indig** @Kevin_Indig · Mar 24

What are all the jobs of a technical SEO?

I got optimize
1. Crawling & rendering
2. Page experience (CWV, etc.)
3. Internal linking
4. SEO hygiene (solving problematic status codes, etc.)
5. Indexing
6. Mobile optimization
7. Structured Data/Rich Snippets

What else?

59     52     289

**Barry Adams** 🗂️
@badams

Replying to @Kevin_Indig

Anything that's not content or links.

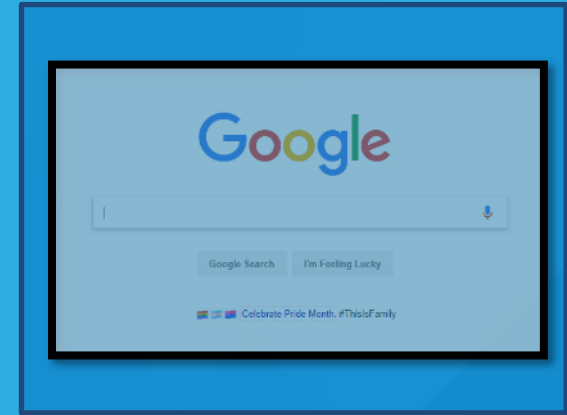10:28 PM · Mar 24, 2021 · Twitter for iPhone

**POLEMIC**
D I G I T A L

# 1. Crawler (Googlebot)



Crawler          Indexer          Ranker

**POLEMIC**
D I G I T A L

# Crawling: Googlebot

- URL discovery
  - ➢ <a href> tags in HTML
  - ➢ XML sitemaps
  - ➢ Other sources?
- Crawl queue management
  - ➢ De-duplication based on URL patterns
  - ➢ Crawl prioritisation & scheduling
- Crawling
  - ➢ Fetching raw HTML
  - ➢ Crawl 'politeness'

# Crawl Management

- Robots.txt Disallow
  - ➤ Strongest crawl management signal
  - ➤ Evaporates crawl budget



```
User-agent: *
Disallow: /*s=
Disallow: /search$

# Separate rules for Google Adsbot because it ignores blanket disallow rules.
# Fuck you, Adsbot.
User-agent: AdsBot-Google
Disallow: /

Sitemap: https://www.polemicdigital.com/sitemap_index.xml
```

Was this helpful? 👍 👎

# Large site owner's guide to managing your crawl budget 🔖

**Send feedback**

This guide describes how to optimize Google's crawling of very large and frequently updated sites.

If your site does not have a large number of pages that change rapidly, or if your pages seem to be crawled the same day that they are published, you don't need to read this guide; merely keeping your sitemap up to date and checking your index coverage 🔗 regularly is adequate.

If you have content that's been available for a while but has never been indexed, this is a different problem; use the URL Inspection tool 🔗 instead to find out why your page isn't being indexed.

@badams

**POLEMIC**
D I G I T A L

Was this helpful? 👍 👎

# Large site owner's guide to managing your crawl budget 🔖

Send feedback

This guide describes how to optimize Google's crawling of very large and frequently updated sites.

If your site does not have a large number of pages that change rapidly, or if your pages seem to be crawled the same day that they are published, you don't need to read this guide; merely keeping your sitemap up to date and checking your index coverage 🔗 regularly is adequate.
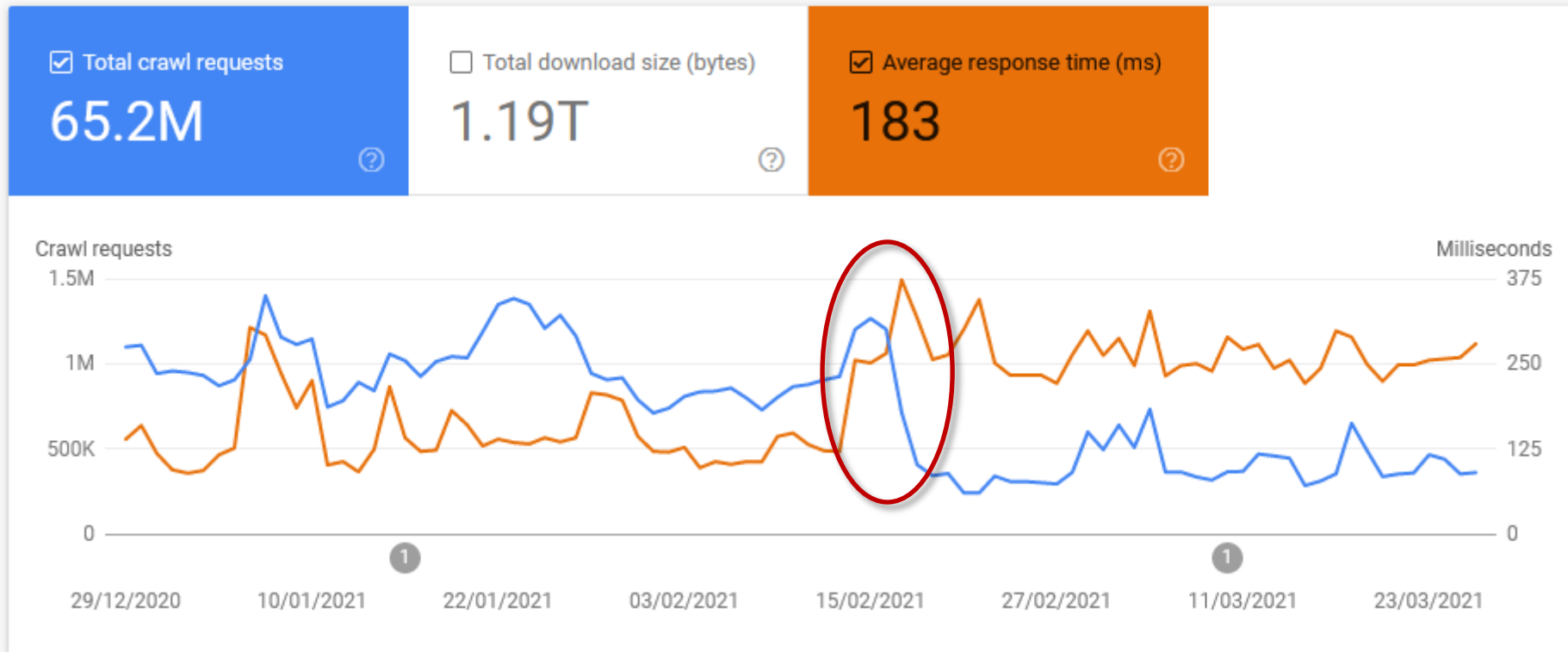
Don't use robots.txt to temporarily reallocate crawl budget for other pages; use robots.txt to block pages or resources that you don't want Google to crawl at all. **Google won't shift this newly available crawl budget to other pages** unless Google is already hitting your site's serving limit.

@badams

POLEMIC
D I G I T A L

# Crawl Management

- Robots.txt Disallow
  - ➢ Strongest crawl management signal
  - ➢ Evaporates crawl budget

- Canonicals & noindex are NOT crawl management
  - ➢ Google needs to **see** meta tags before it can act on them
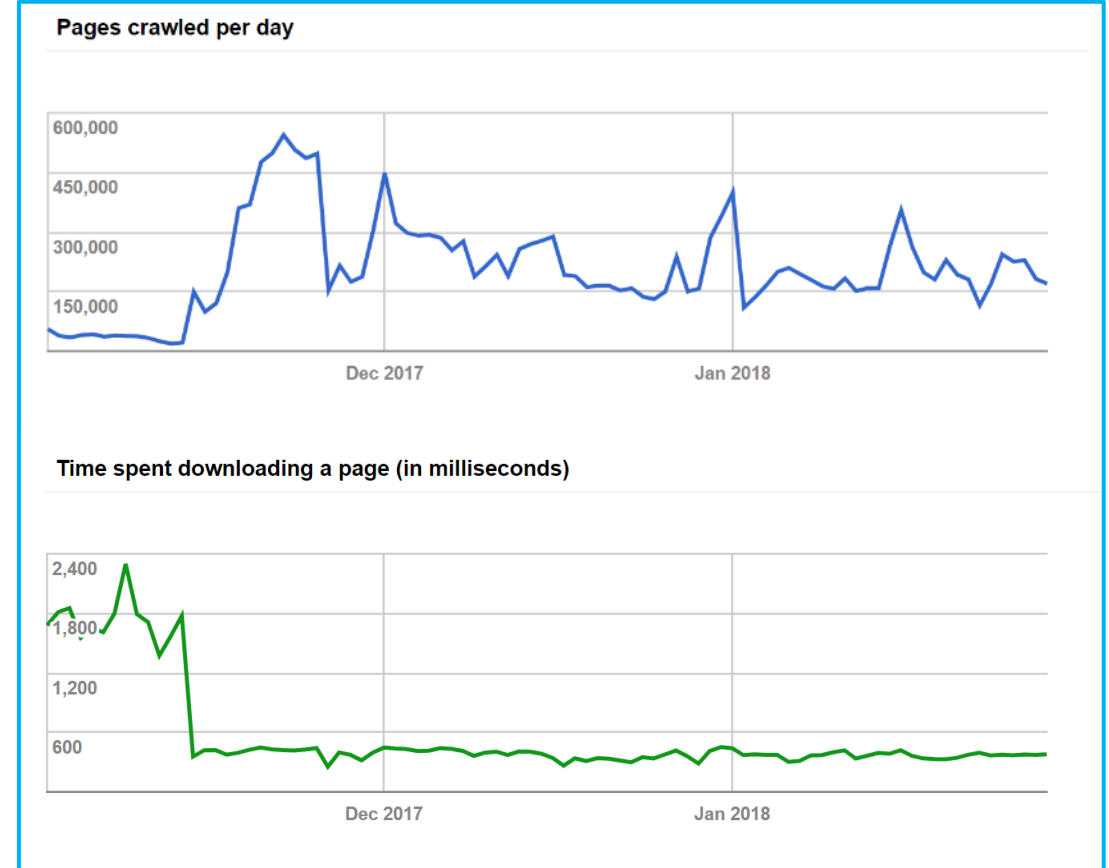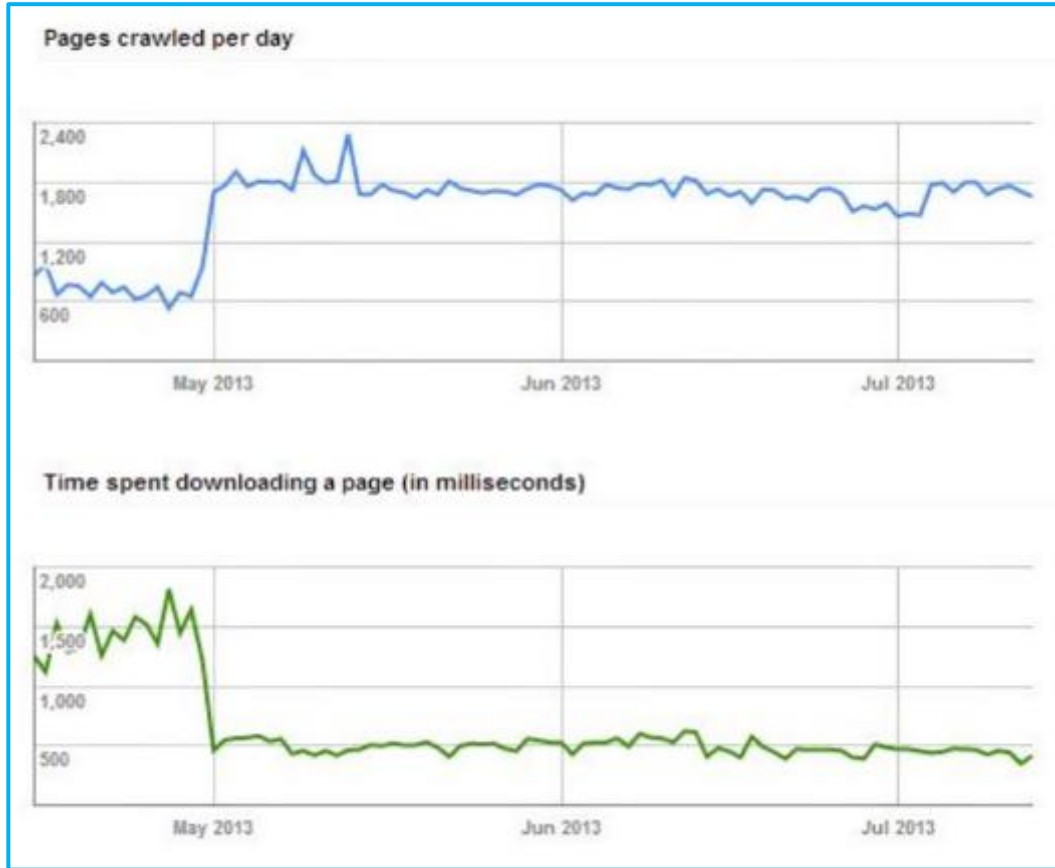  - ➢ That means Googlebot still crawls those URLs

# Optimise Crawling

- Server Response Time

POLEMIC
DIGITAL

# Load Speed

## Fast response time = optimal use of Googlebot

POLEMIC DIGITAL

# GSC Crawl Stats

## By response

| | | |
|---|---|---|
| OK (200) | 77% | |
| Not modified (304) | 21% | |
| Moved permanently (301) | 1% | |
| Not found (404) | < 1% | |
| Moved (other) | < 1% | |

Rows per page: 5 ▾    1-5 of 11    ‹ ›

## By file type

| | | |
|---|---|---|
| JavaScript | 60% | |
| HTML | 16% | |
| JSON | 1% | |
| CSS | < 1% | |
| Image | < 1% | |

Rows per page: 5 ▾    1-5 of 10    ‹ ›

## By purpose

| | | |
|---|---|---|
| Refresh | 99% | |
| Discovery | < 1% | |

Rows per page: 5 ▾    1-2 of 2    ‹ ›

## By Googlebot type

| | | |
|---|---|---|
| Page resource load | 63% | |
| Smartphone | 35% | |
| Desktop | 3% | |
| Image | < 1% | |
| AdsBot | < 1% | |

Rows per page: 5 ▾    1-5 of 7    ‹ ›

# Optimise Crawling

- Serve correct HTTP status codes

  - ➤ 200 OK
  - ➤ 301 / 302 Redirects
  - ➤ 304 Not Modified
  - ➤ 401 / 403 Permission Issues
  - ➤ 404 / 410 Not Found/Gone
  - ➤ 5xx Error

# HTTP Status Codes



- All redirects should be **301**

  ➢ Except geo-targeting redirects, those should be **302**

  ➢ **304** means the URL hasn't changed since the last crawl

  ➢ **307** relates to HSTS preload list

POLEMIC
D I G I T A L

# HTTP Status Codes

- Accidental not found: **404**

- Deliberate deletion:    **410**

POLEMIC
D I G I T A L

# HTTP Status Codes

- Accidental server error:    **500**

- Deliberate downtime:       **503**

**POLEMIC**
D I G I T A L

# Optimise Crawling

- ALL resources consume crawl budget;
  - ➢ Not just HTML pages
  - ➢ Reduce HTTP requests per page

**POLEMIC**
D I G I T A L

# Content breakdown by MIME type (First View)

## Requests



Legend:
- css
- font
- html
- image
- js
- other
- Other

Pie chart values: 30.6% other, 25.1% html, 30% image, 11.7% js

| MIME Type | Requests |
|---|---|
| other | 365 |
| image | 358 |
| html | 300 |
| js | 140 |
| css | 19 |
| font | 11 |
| video | 1 |
| flash | 0 |

## Bytes



Legend:
- css
- font
- html
- image
- js
- other
- video

Pie chart values: 11.6%, 9.6%, 44.5%, 21.6%

| MIME Type | Bytes | Uncompressed |
|---|---|---|
| js | 2,920,760 | 8,627,503 |
| image | 1,418,862 | 1,418,828 |
| video | 762,448 | 762,448 |
| html | 628,587 | 2,055,520 |
| font | 425,712 | 425,712 |
| other | 336,069 | 1,265,643 |
| css | 65,135 | 390,706 |
| flash | 0 | 0 |

## By response

| | | |
|---|---|---|
| OK (200) | 86% | |
| Moved temporarily (302) | 10% | |
| Not modified (304) | 2% | |
| Moved permanently (301) | 1% | |
| Not found (404) | < 1% | |
| Other client error (4XX) | < 1% | |
| Server error (5XX) | < 1% | |
| Page could not be reached | < 1% | |
| robots.txt not available | < 1% | |
| Unauthorised (401/407) | < 1% | |

Rows per page: 10    1-10 of 10    〈  〉

## By file type

| | | |
|---|---|---|
| HTML | 52% | |
| JSON | 1% | |
| CSS | < 1% | |
| JavaScript | < 1% | |
| Image | < 1% | |
| Syndication | < 1% | |
| PDF | < 1% | |
| Video | < 1% | |
| Other XML | < 1% | |
| Other file type | 45% | |
| Unknown (failed requests) | < 1% | |

Rows per page: 25    1-11 of 11    〈  〉

## By purpose

By Googlebot type

| | | |
|---|---|---|
| Smartphone | 42% | |
| Page resource load | 36% | |
| Desktop | 20% | |
| Image | 1% | |
| AdsBot | < 1% | |
| Video | < 1% | |
| Other agent type | < 1% | |

Rows per page: 10 ▾     1-7 of 7     < >

@badams

POLEMIC
DIGITAL

# Optimise Crawling

- ALL resources consume crawl budget;
  - ➤ Not just HTML pages
  - ➤ Reduce HTTP requests per page

- AdsBot can consume crawl budget;
  - ➤ Double-check your Google Ads campaigns

**POLEMIC**
D I G I T A L

By Googlebot type

| | | |
|---|---|---|
| AdsBot | 69% | |
| Smartphone | 21% | |
| Desktop | 10% | |
| Page resource load | < 1% | |
| Image | < 1% | |
| Other agent type | < 1% | |

Rows per page: 10 ▼    1-6 of 6    ‹ ›

POLEMIC
DIGITAL

# Optimise Crawling

- ALL resources consume crawl budget;
  - ➢ Not just HTML pages
  - ➢ Reduce HTTP requests per page

- AdsBot can consume crawl budget;
  - ➢ Double-check your Google Ads campaigns

- Link equity (PageRank) impacts crawl budget;
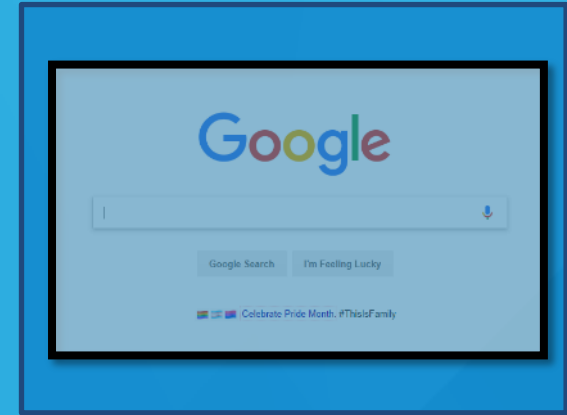  - ➢ More link equity = more crawl budget
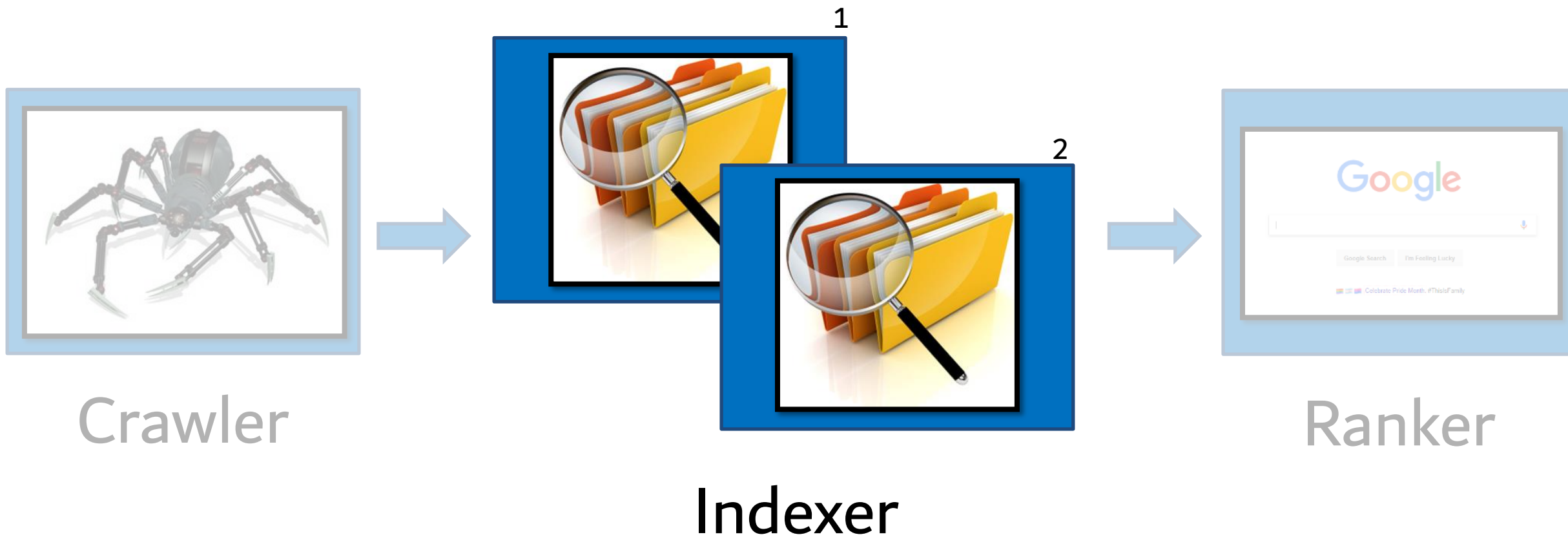
# 2. Indexer



Crawler

Indexer

Ranker

@badams

**POLEMIC**
D I G I T A L

# Two Stages* of Indexing

Crawler → Indexer (1, 2) → Ranker

*At least – indexing is a collection of interconnected processes

@badams

POLEMIC
DIGITAL

# Indexing

- HTML lexer
  - Cleaning & tokenising the HTML
- Index selection
  - De-duping prior to indexing
- Indexing
  - First-pass based on HTML
  - Potential rendering (not guaranteed)
- Index integrity
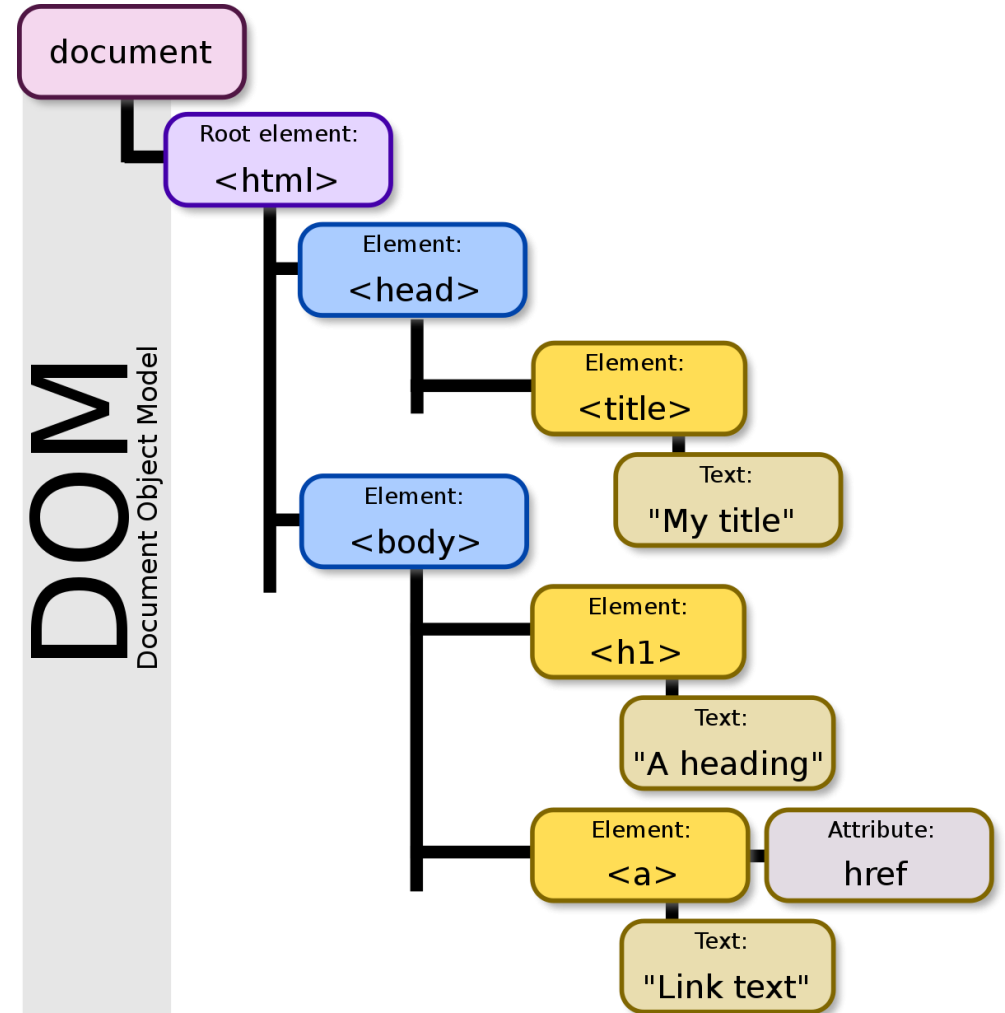  - Canonicalisation & de-duplication

# Extraction

Can Google easily extract a page's content from its DOM?

POLEMIC
D I G I T A L

# Optimise Extraction (1)

- Clean HTML;
  - ➤ Yes, really!

  - ➤ There is a max HTML size Google will parse
    - – Speculation: ~1 MB

  - ➤ Less clutter = easier parsing

# Optimise Extraction (2)
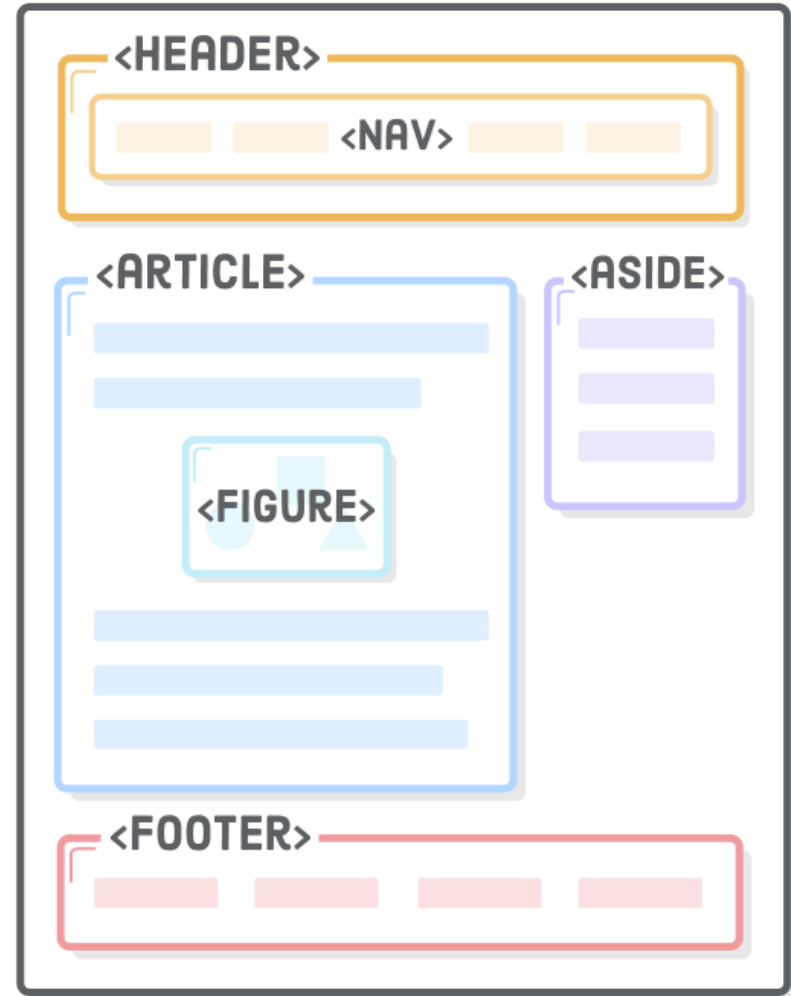
- Clean <head>;
  - ➢ Critical meta tags high in the <head>
    - Title & description
    - Open Graph
    - Canonical, hreflang & mobile alternate
    - Structured Data

  - ➢ Internal CSS & JS lower in the <head>

# Optimise Extraction (3)

- Minimise DOM-manipulation;

  ➢ Client-side JavaScript that manipulates the DOM can impact extraction

  ➢ Can also impact Core Web Vitals

POLEMIC

DIGITAL

# Semantics

Can Google understand
what the page is about?

POLEMIC
D I G I T A L

# Optimise Semantics

- Good content;
  - ➢ Easily identifiable entities and relationships


- Semantic HTML;
  - ➢ Enables Google to separate style & boilerplate from content


- Structured Data;
  - ➢ Makes page contents explicitly clear

# Test Entities in Content

Google NLP API: https://cloud.google.com/natural-language

POLEMIC
DIGITAL

# GSC: Mix of Crawling & Indexing Issues



@badams

POLEMIC
DIGITAL

# IndexNow



**Microsoft Bing** Blogs

This is a place devoted to giving you deeper insight into the news, trends, people and technology behind Bing.

⊕ Blogs   ⊕ Regions   ⊕ Skip to content

Follow us:   f 🐦 in 📌 📷   Subscribe RSS

**OCTOBER 18 2021**

## IndexNow - Instantly Index your web content in Search Engines

Ensuring timely information is available for searchers is critical. Yet historically one of the biggest pain points for website owners has been to have search engines quickly discover and consider their latest website changes. It can take days or even weeks for new URLs to be discovered and indexed in search engines, resulting in loss of potential traffic, customers, and even sales.

IndexNow is a new protocol created by Microsoft Bing and Yandex, allowing websites to easily notify search engines whenever their website content is created, updated, or deleted. Using an API, once search engines are notified of updates they quickly crawl and reflect website changes in their index and search results.

# Live Indexing API (?)



Google Search Central > Indexing API

🔍 Search    English ▼   ⚙️   POLEMIC

**Guides**    Reference

**Quickstart**

Prerequisites

Using the API

Errors

Client libraries

Authorize requests

Quota and pricing

Home > Search Central > Indexing API

Rate and review 👍 👎

## Indexing API Quickstart 🔖

**Send feedback**

The Indexing API allows any site owner to directly notify Google when pages are added or removed. This allows Google to schedule pages for a fresh crawl, which can lead to higher quality user traffic. Currently, the Indexing API can only be used to crawl pages with either `JobPosting` or `BroadcastEvent` embedded in a `VideoObject`. For websites with many short-lived pages like job postings or livestream videos, the Indexing API keeps content fresh in search results because it allows updates to be pushed individually.

Here are some of the things you can do with the Indexing API:

- **Update a URL**: Notify Google of a new URL to crawl or that content at a previously-submitted URL has been updated.

- **Remove a URL**: After you delete a page from your servers, notify Google so that we can remove the page from our index and so that we don't attempt to crawl the URL again.

- **Get the status of a request**: Check the last time Google received each kind of notification for a given URL.

- **Send batch indexing requests**: Reduce the number of HTTP connections your client has to make by combining up to 100 calls into a single HTTP request.

**Table of contents**

Sitemaps and the Indexing API

Get started

# Structured Data

Constantly evolving schemas

New rich snippets in SERPs

https://sitebulb.com/structured-data-history/

POLEMIC
DIGITAL

# Structured Data

- 'author.url' now recommended in Article SD

**POLEMIC**
D I G I T A L

# Edge SEO

https://dantaylor.online/edge-seo/

POLEMIC
D I G I T A L

https://www.searchpilot.com/resources/blog/edge-seo/

**POLEMIC**
D I G I T A L

# A/B Testing



**SearchPilot**

How it works    Features    Pricing    Resources    Newsletter    About    Jobs    Contact          Log in

- SEO A/B Testing
- Meta-CMS
- Full Funnel Testing
- Server-Side Testing
- Professional Services
- Enterprise Features

## Making SEO A/B Testing Easy

At SearchPilot, we are on a mission to prove the value of SEO for the world's biggest websites. It can be incredibly hard to connect specific changes to their associated SEO benefit without controlled testing. Our platform makes SEO A/B testing easy in three main ways:

1. By automatically splitting site sections into statistically-similar groups of pages, including or excluding any groups of pages we want
2. By making it easy to make the changes you want to test
3. By running advanced statistical models on your analytics data to understand the impact

When you have a winning test, you can deploy it to 100% of pages using the SearchPilot platform, or build it out yourself.

This page will focus on how SearchPilot makes testing easy. If you want to learn more about what SEO testing is and how it works, you can dive deeper into that in this blog post.

## Automatically splitting site sections

@badams          https://www.searchpilot.com/          **POLEMIC** DIGITAL

# Less hassle with JavaScript



Home > Search Central > Documentation > Advanced SEO
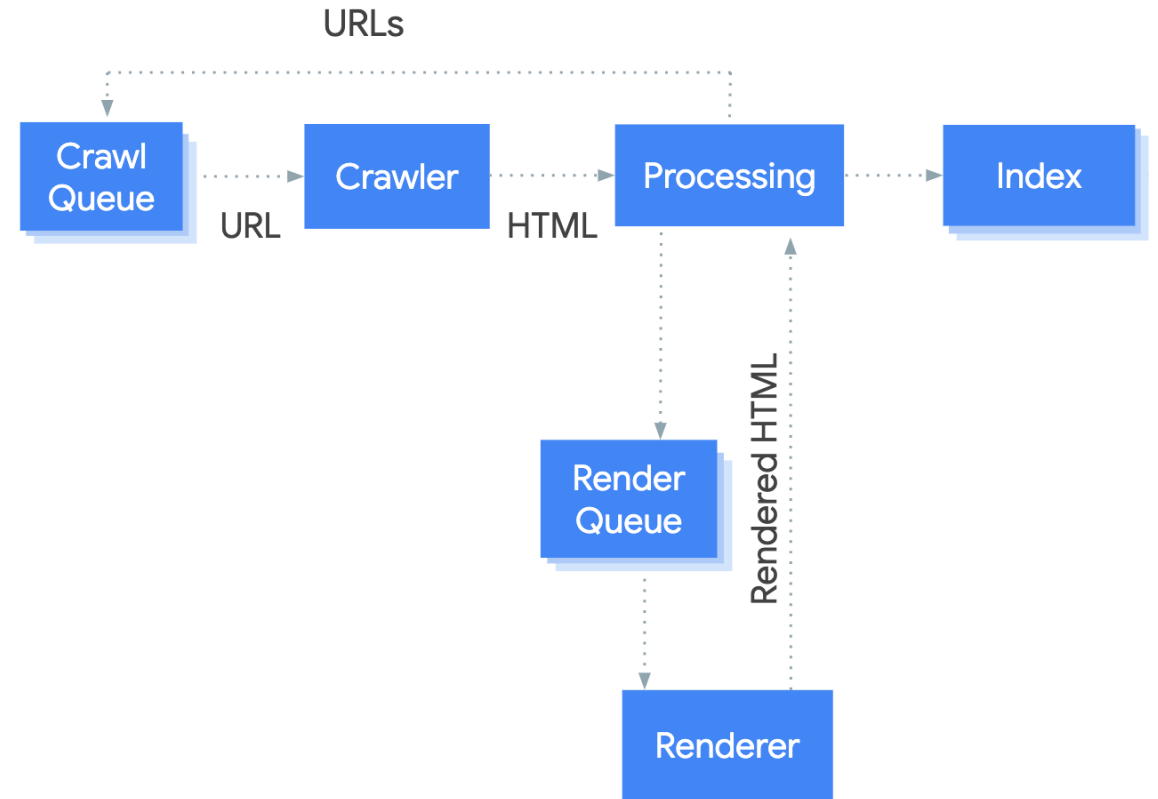
Rate and review 👍 👎

## Fix Search-related JavaScript problems  🔖

**Send feedback**

This guide helps you identify and fix JavaScript issues that may be blocking your page, or specific content o
powered pages, from showing up in Google Search. While Googlebot does run JavaScript, there are some d
and limitations that you need to account for when designing your pages and applications to accommodate
access and render your content.

💡 Our guide on JavaScript SEO basics has more information on how you can optimize your JavaScript site for Google S

Googlebot is designed to be a good citizen of the web. Crawling is its main priority, while making sure it doe
the experience of users visiting the site. Googlebot and its Web Rendering Service (WRS) component contin
analyze and identify resources that don't contribute to essential page content and may not fetch such resou
example, reporting and error requests that don't contribute to essential page content, and other similar type
are unused or unnecessary to extract essential page content.

URLs

Crawl Queue → URL → Crawler → HTML → Processing → Index

Render Queue

Rendered HTML

Renderer

POLEMIC
DIGITAL

# Better GSC Reports

More useful info to empower SEOs & Devs



Home > Search Central > Google Search Central Blog                    Rate and review  👍 👎

## Index Coverage Data Improvements  🔖                      [ Send feedback ]

*Monday, January 11, 2021*

Helping people understand how Google crawls and indexes their sites has been one of the main objectives of Search Console since its early days ⤢. When we launched the new Search Console, we also introduced the Index Coverage report ⤢, which shows the indexing state of URLs that Google has visited, or tried to visit, in your property.

Based on the feedback we got from the community, today we are rolling out significant improvements to this report so you're better informed on issues that might prevent Google from crawling and indexing your pages. The change is focused on providing a more accurate state to existing issues, which should help you solve them more easily. The list of changes include:

- Removal of the generic "crawl anomaly" issue type - all crawls errors should now be mapped to an issue with a finer resolution.

- Pages that were submitted but blocked by robots.txt and got indexed are now reported as "indexed but blocked" (warning) instead of "submitted but blocked" (error)

- Addition of a new issue: "indexed without content ⤢" (warning)

- Soft 404 reporting is now more accurate

POLEMIC
D I G I T A L

# Better Google Documentation



Home > Search Central > Documentation > Advanced SEO

Rate and review 👍 👎

## Large site owner's guide to managing your crawl budget 🔖

**Send feedback**

### Overview

This guide describes how to optimize Google's crawling of very large and

If your site does not have a large number of pages that change rapidly, or that they are published, you do not need to read this guide; merely keeping index coverage regularly should be adequate.

If you have content that's been available for a while but has never been in Inspection tool instead to find out why your page isn't being indexed.

Home > Search Central > Google Search Central Blog

Rate and review 👍 👎

## New resources for video SEO 🔖

**Send feedback**

*Wednesday, March 17, 2021*

As global online video consumption continues to grow, Google aims to surface video content from diverse sources across the web. We want to make it easy for site owners to get their videos indexed and surfaced on Google.

Today, we're excited to share two new resources to help you optimize your videos for Google Search and Discover.

### Search Central Lightning Talk

In this new lightning talk 🔗, we discuss how Google indexes videos, highlight features where videos appear on Google, and share five key tips to optimize your videos for Search and Discover:

@badams

POLEMIC
DIGITAL

# SEO Crawlers

- DeepCrawl
  https://www.deepcrawl.com/

- Sitebulb
  https://sitebulb.com/

- Screaming Frog
  https://www.screamingfrog.co.uk/seo-spider/

# SEO Review & Monitoring

- Little Warden
  https://littlewarden.com/

- ContentKing
  https://www.contentkingapp.com/

- SEO Info
  https://weeblr.com/doc/products.seoinfo/current/overview/

- SEOBrowse
  https://seobrowse.com/

# Performance Analysis

- PageSpeed Insights
  https://pagespeed.web.dev/

- WebPagetest.org
  https://www.webpagetest.org/

- GTmetrix
  https://gtmetrix.com/

# Barry Adams

➢ Doing SEO since 1998

➢ Specialist in Technical SEO & News SEO

➢ Newsletter: SEOforGoogleNews.com

POLEMIC
D I G I T A L