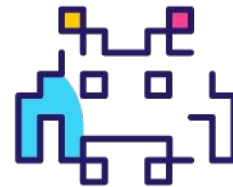THE SUPER
SEO GAME

**Webinar**

# Thinking **BIG**:
# Big Solutions for
# Big Sites

oncrawl

# GAMEPLAY & RULES

• Earn points by signing up, attending, and participating

• Unlock new levels, earn badges and check our leaderboard

• Use #SuperSEOGame to continue the conversation

• Have fun!

oncrawl

# SINGLE PLAYER of the day

**Rebecca Berbel**

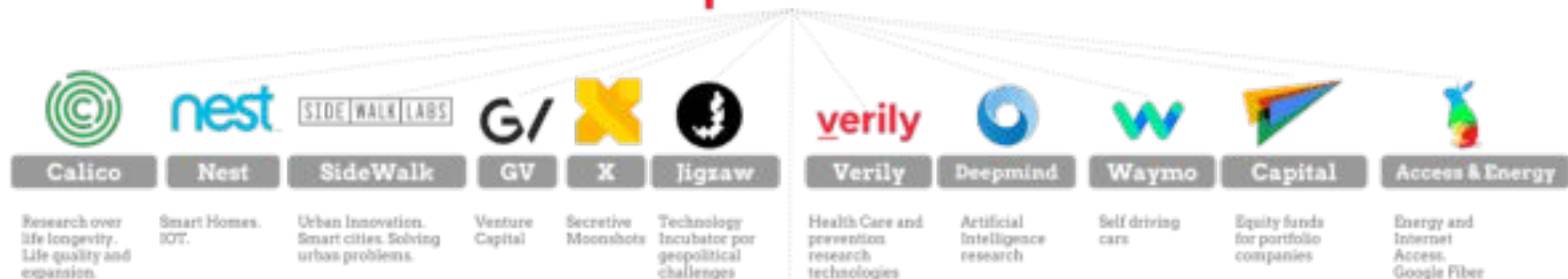Product Marketing Manager @Oncrawl

@rebberbel

# The BIG picture

1. Definitions & context

2. Technical constraints & solutions

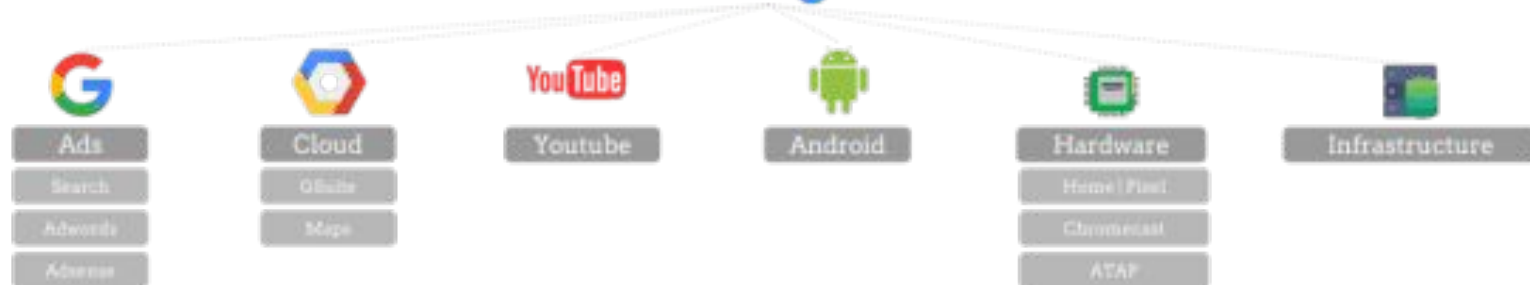3. Challenges in SEO strategy

4. Focus areas for successful SEO

# What is a **big** site?

# What do big sites have in **common**?

# Alphabet

| Calico | Nest | SideWalk | GV | X | Jigzaw | Verily | Deepmind | Waymo | Capital | Access & Energy |
|---|---|---|---|---|---|---|---|---|---|---|
| Research over life longevity. Life quality and expansion. | Smart Homes. IOT. | Urban Innovation. Smart cities. Solving urban problems. | Venture Capital | Secretive Moonshots | Technology Incubator por geopolitical challenges | Health Care and prevention research technologies | Artificial Intelligence research | Self driving cars | Equity funds for portfolio companies | Energy and Internet Access. Google Fiber |

# Google

| Ads | Cloud | Youtube | Android | Hardware | Infrastructure |
|---|---|---|---|---|---|
| Search | GSuite | | | Home | Pixel | |
| Adwords | Maps | | | Chromecast | |
| Adsense | | | | ATAP | |

# Links insight

| | | |
|---|---|---|
| Number of follow inlinks | **4,179,842,565** | N/A |
| Pages with 1 follow inlink | **22,848,953** | N/A |
| Pages with less than 10 follow inlinks | **49,858,530** | N/A |

[00:13:56] **Gary Illyes:** [00:13:56] Sounds awful. Let's do it. Okay, crawl budget. We published quite a bit about the crawl budget, I think, lately or the past couple of years.

[00:14:07] We've been pushing back on the crawl budget, historically, typically telling people that you don't have to care about it. And I stand my ground and I still say that most people don't have to care about it. We do think that there is a substantial segment of the ecosystem that has to care about it. That's why we publish more on that topic and talk a little bit more about it, but I still believe that - I'm trying to reinforce this here - that the vast majority of the people don't have to care about it. Everyone wants a number, basically how big your site has to be when you have to care about crawl budget. But I don't think it works like that. And I remember when we were writing one of the help center documentations, Josh, our help center tech writer, was also asking this question, okay, let's define this.

[00:14:59] Like how many pages do you think, or how many URLs on the site you have to have to start caring about crawl budget. And we were working with the Googlebot team on the documentation and both the Googlebot team and us - Search Relations or whatever we are called nowadays - were saying that well, it's not quite like that. It's like you can do stupid stuff on your site, and then Googlebot with start crawling like crazy. Or you can do other kinds of stupid stuff, and then Googlebot will just stop crawling altogether. And, eventually, he did convince us to give a number, which I don't remember, but I think it's around a million, I would say, URLs on the site and that's our baseline. So basically if you have fewer than a million URLs on the page, on your site, then you don't really have to care about crawl budget.
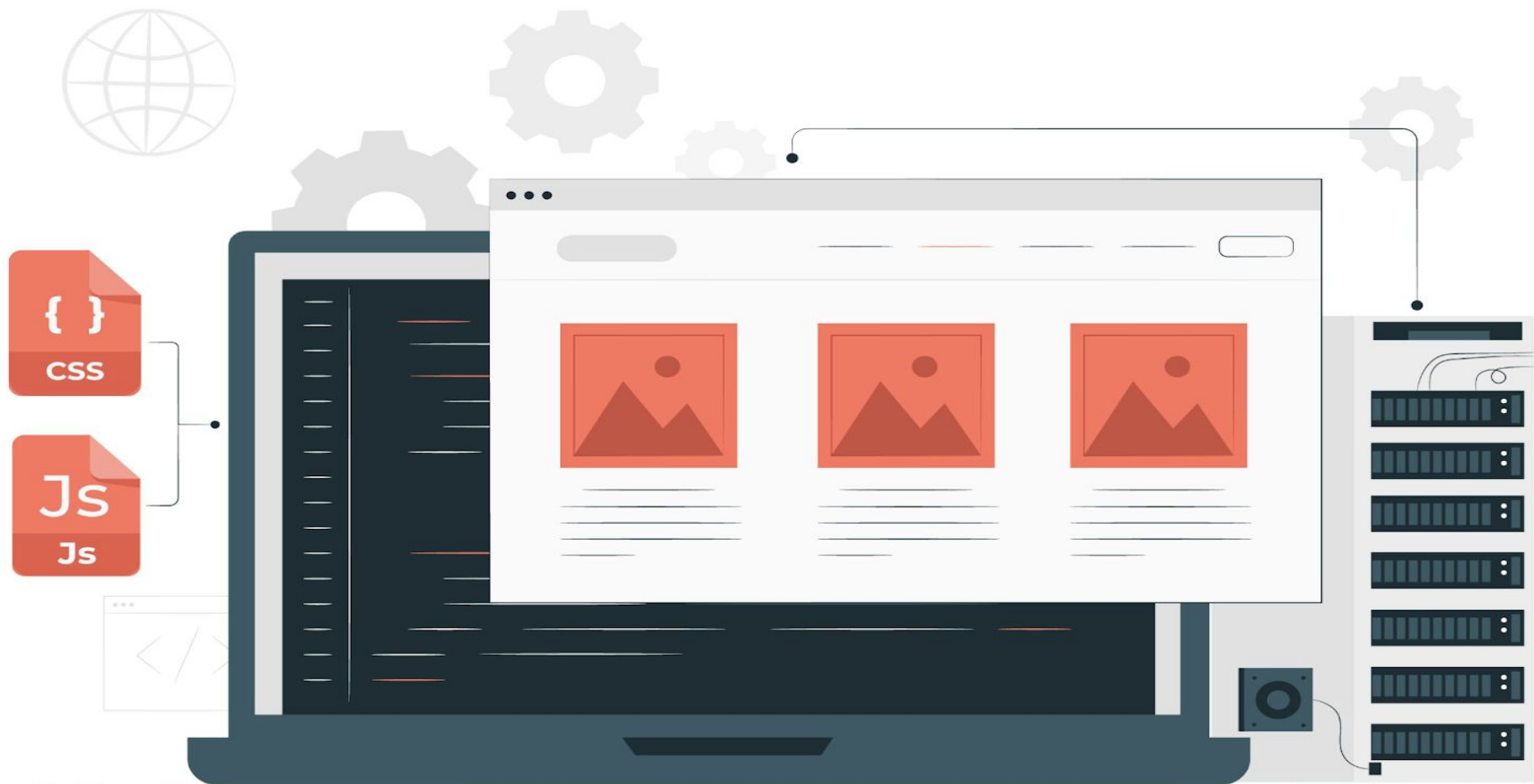
# Large site owner's guide to managing your crawl budget

## Who this guide is for

This is an advanced guide and is intended for:

- **Large sites** (1 million+ unique pages) with content that changes moderately often (once a week), or
- Medium or larger sites (10,000+ unique pages) with **very rapidly changing content** (daily).

Please note that the numbers given here are a **rough estimate** to help you classify your site. These are not exact thresholds.

oncrawl

Nov 16, 2020

2 weeks 5 days

default ⓘ

114,609,936*

Oct 28, 2020 17:07 / Nov 16, 2020 21:34

Sep 09,

1 week

default ⓘ

4,000,000+

□ △ ○ ✕

oncrawl

# How can see the **big** picture?

1 Collecting data

2 Preparing data

3 Analyzing data

4 Storing data

5 Retrieving data

6 Interpreting data

oncrawl

# Why is it hard to get a view of the whole site?

# How do you reduce the time it takes to crawl?

⚠️ **Site infrastructure**
⚠️ **Processing power**

oncrawl

# How do you reduce the time it takes to crawl?

*processing power*

limited by maximum available resources = requires more time

*time*

# How do you reduce the time it takes to crawl?

*processing power*

limited by maximum available
resources = requires more time

*time*

*processing power*

autoscaling crawl clusters as
required = requires more power

*time*

**oncrawl**

"Google Kubernetes Engine enables us to tackle gigantic projects immediately, allowing our customers to crawl even **hundreds of millions of JavaScript pages** without prior notice. We've seen our JavaScript crawler scale up from 10 to 750 pods in just a few minutes."

—*Philippe David, Chief Technology Officer, Oncrawl*

oncrawl

# Why is it hard to get a view of the whole site?

# Why is it hard to get a view of the whole site?

**LOGS**

Number of log files ingested per month

# Managing data for the full site

- **Extracting**

- **Structuring**

- **Sorting**

- **Filtering**

- **Grouping**

**100 000 000** (pages)

= (**100 000 000** × **99 999 999**) ÷ **2**

≅ **100 000 000**$^2$

# 10 000 000 000 000 000

(duplication analyses)

1 Collecting data

2 Preparing data

× 1    × 2    × 3

3 300 000 000 000 000

oncrawl
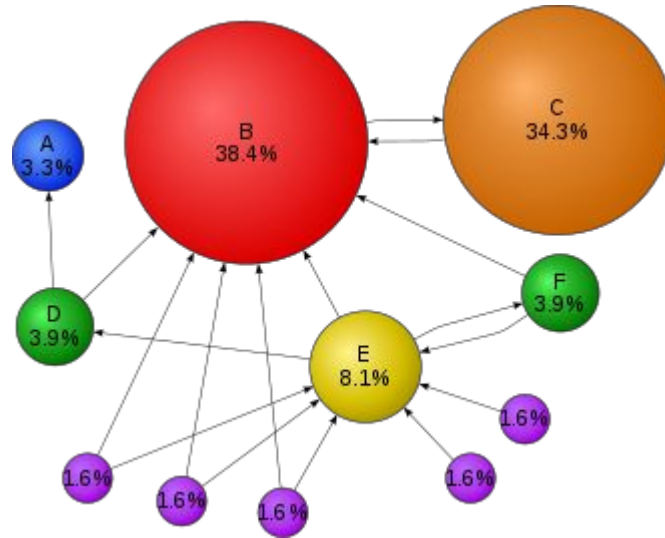
1 Collecting data

2 Preparing data

3 Analyzing data

$$PR(A) = \frac{1-d}{N} + d\left(\frac{PR(B)}{L(B)} + \frac{PR(C)}{L(C)} + \frac{PR(D)}{L(D)} + \cdots\right).$$

A
3.3%

B
38.4%

C
34.3%

D
3.9%

F
3.9%

E
8.1%

1.6%

1.6%

1.6%

1.6%

1.6%

oncrawl

## Collecting data

## Preparing data

## Analyzing data

A locality sensitive hashing scheme is a distribution on a family $\mathcal{F}$ of hash functions operating on a collection of objects, such that for two objects $x, y$,

$$\mathbf{Pr}_{h \in \mathcal{F}}[h(x) = h(y)] = sim(x, y),$$

where $sim(x, y) \in [0, 1]$ is some similarity function defined on the collection of objects. Such a scheme leads to a compact representation of objects so that similarity of objects can be estimated from their compact sketches, and also leads to efficient algorithms for approximate nearest neighbor search and clustering. Min-wise independent permutations provide an elegant construction of such a locality sensitive hashing scheme for a collection of subsets with the set similarity measure $sim(A, B) = \frac{|A \cap B|}{|A \cup B|}$.

We show that rounding algorithms for LPs and SDPs used in the context of approximation algorithms can be viewed as locality sensitive hashing schemes for several interesting collections of objects. Based on this insight, we construct new locality sensitive hashing schemes for:

1. A collection of vectors with the distance between $\vec{u}$ and $\vec{v}$ measured by $\theta(\vec{u}, \vec{v})/\pi$, where $\theta(\vec{u}, \vec{v})$ is the angle between $\vec{u}$ and $\vec{v}$. This yields a sketching scheme for estimating the cosine similarity measure between two vectors, as well as a simple alternative to minwise independent permutations for estimating set similarity.

2. A collection of distributions on $n$ points in a metric space, with distance between distributions measured by the Earth Mover Distance (**EMD**), (a popular distance measure in graphics and vision). Our hash functions map distributions to points in the metric space such that, for distributions $P$ and $Q$,

$$\mathbf{EMD}(P, Q) \leq \mathbf{E}_{h \in \mathcal{F}}[d(h(P), h(Q))]$$
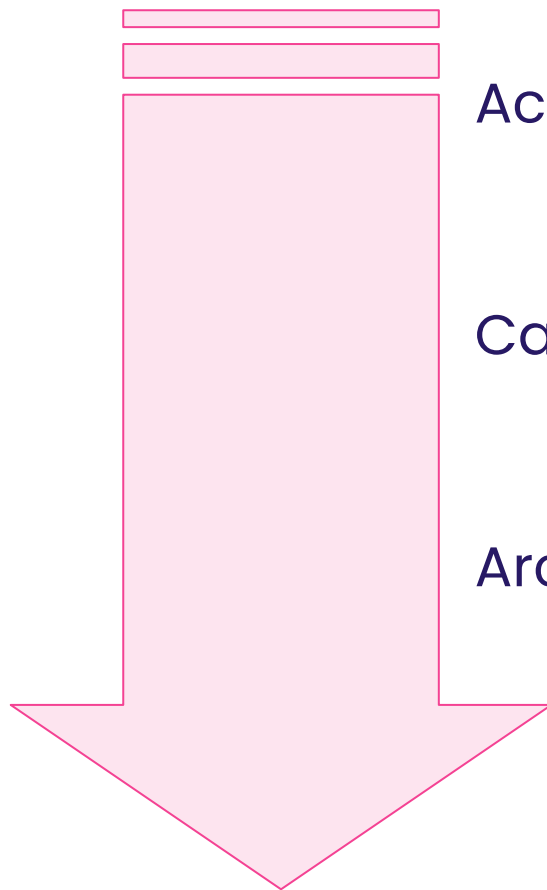$$\leq O(\log n \log \log n) \cdot \mathbf{EMD}(P, Q).$$

SimHash (Moses S. Charikar)

1 Collecting data

2 Preparing data

3 Analyzing data

4 Storing data

Accessible information

Calculatable information

Archived information

oncrawl

1 Collecting data

2 Preparing data

3 Analyzing data

4 Storing data

If you think of a datamart as a store of bottled water – cleansed and packaged and structured for easy consumption – the data lake is a large body of water in a more natural state. **The contents of the data lake stream in from a source** to fill the lake, and various users of the lake can come to **examine**, **dive in**, or **take samples**.

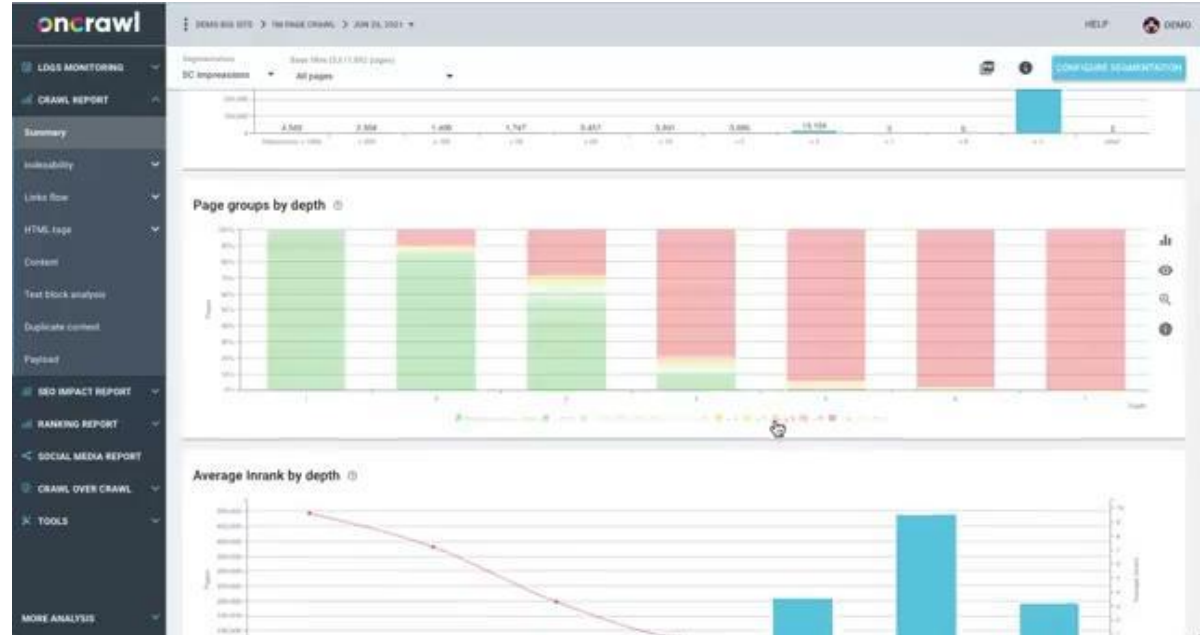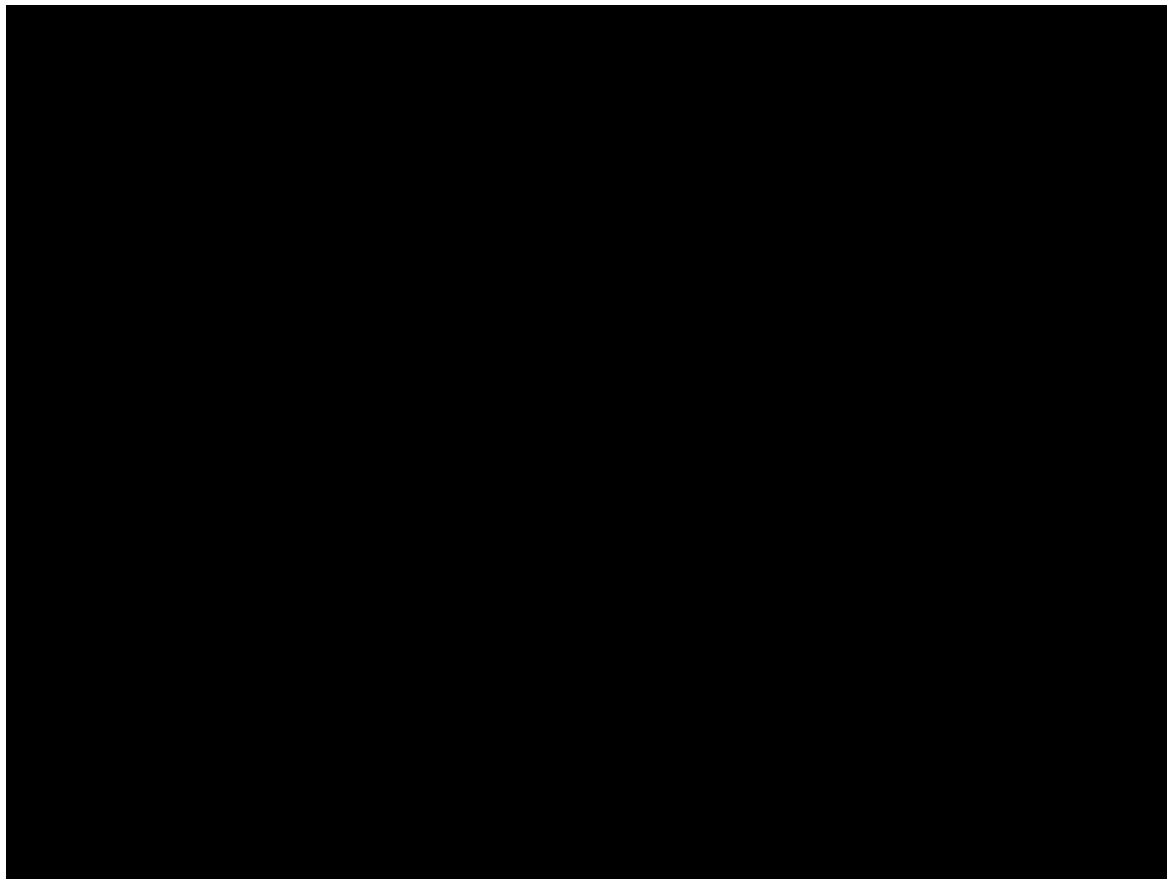*—James Dixon, then Chief Technology Officer, Pentaho*

□ △ ○ ✕

oncrawl

1. Collecting data
2. Preparing data
3. Analyzing data
4. Storing data
5. Retrieving data

1 Collecting data

2 Preparing data

3 Analyzing data

4 Storing data

5 Retrieving data

oncrawl

1. Collecting data

2. Preparing data

3. Analyzing data

4. Storing data

5. Retrieving data

We use one of the best systems – ElasticSearch – to store, index and retrieve your data, for **over 1400 native metrics** as well as additional connectors, ingested or scraped data, etc. Our goal is **fast and reliable retrieval**.

—*Vincent Terrasi, Data SEO Educator & Product Manager, Oncrawl*

oncrawl

**1** Collecting data

**2** Preparing data

**3** Analyzing data

**4** Storing data

**5** Retrieving data

**6** Interpreting data

**1** Collecting data

**2** Preparing data

**3** Analyzing data

**4** Storing data

**5** Retrieving data

**6** Interpreting data

**Page groups by depth**

**By number of GA sessions**
green = pages with most sessions

**Page groups by depth**

**By profit margin**
darker = greater margin

**Pages in structure: ratio of ranking pages**

85.70%   99.29%   87.04%   68.34%   58.64%   47.81%   3.97%   14.23%

**By publication date**
lefthand columns = most recent

**Crawl behavior breakdown**

Bots hits    Pages crawled    Pages newly crawled    Crawl frequency per day

**By page template margin**
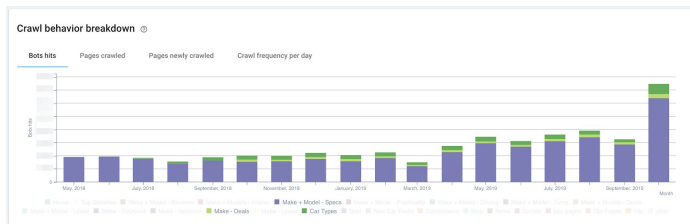each color = different template

oncrawl

1. Collecting data
2. Preparing data
3. Analyzing data
4. Storing data
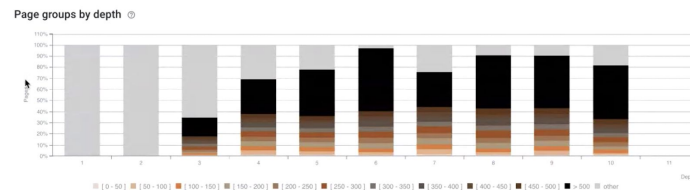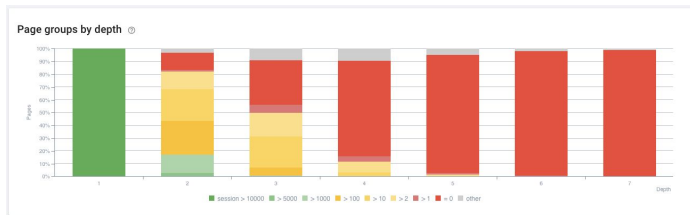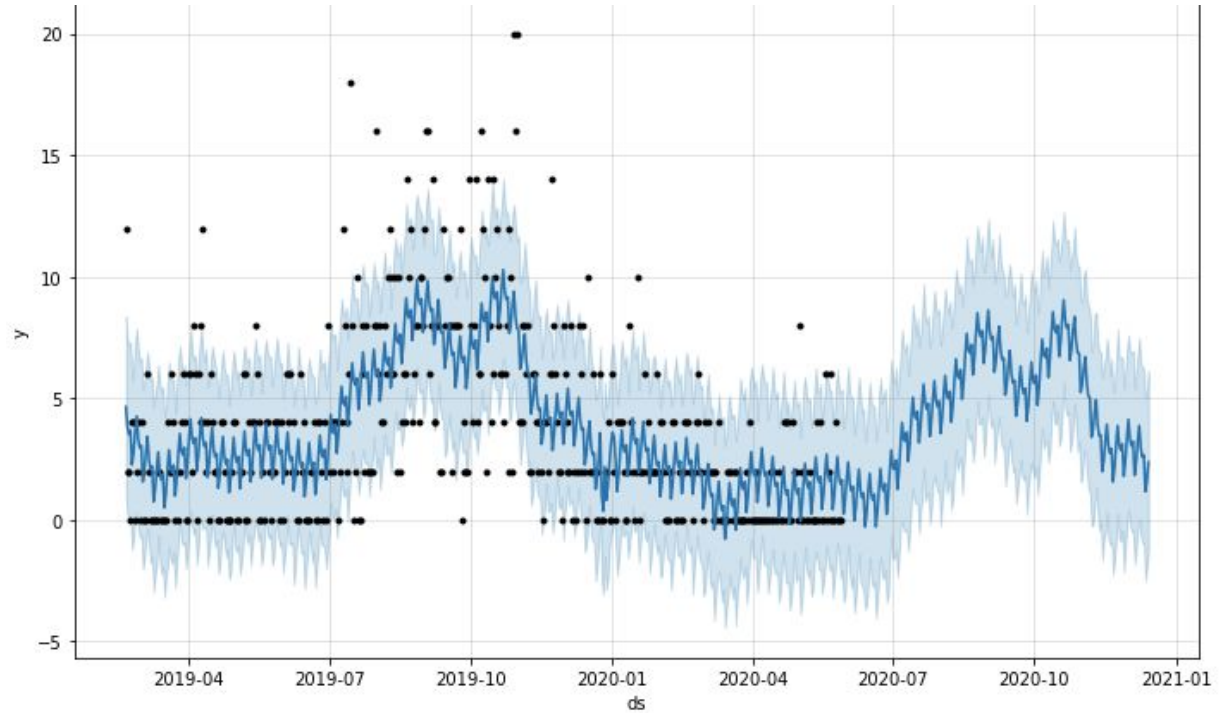5. Retrieving data
6. Interpreting data

# What does it mean to think **big**?

# Get the right support for SEO

Use data to illustrate:

- **Low risk**: the project is not complex or hard to predict

- **Low cost**: the technical solution envisioned is reasonable

- **Easy to measure/monitor**: quantitative data is available to measure the actions and the results

- **Impactful**: there is a business benefit to be gained from the success of this project

oncrawl

I could pitch this project as the first step to **improve the ROI of all of our future SEO projects**. At the **C-level** when I present my ideas, I can tell that there's a lot of confidence. I'd still like to bring the data to the table, because you can't make good decisions without the data to back it up. But when I say, 'here's what I'd like to do next,' the response now is mostly, 'that sounds like **a good direction** to move in.'

*— Murat Yatağan, VP of Growth, Brainly*

oncrawl

# Define strategy adapted to your budget

# Set up a sustainable crawl strategy

- **Infrequent massive crawls**: Obtain and maintain a view of the full site

- **Regular smaller crawls**: Monitor and improve
  - Top levels of the site
  - Limited to a certain template
  - Limited by parameters
  - Lists of important pages
  - Lists of pages with known problems
  - Limited to areas where changes have occurred

# Set up a sustainable crawl strategy

- **Representative crawls**: Can you generalize? What are the limitations?

- **Setting up small(er) crawls**: How do you guide a crawler? Why?

  - URL lists

  - Links that are "followed" by the crawler

  - Subdomains

  - Robots.txt rules

# Where are the **big** wins?

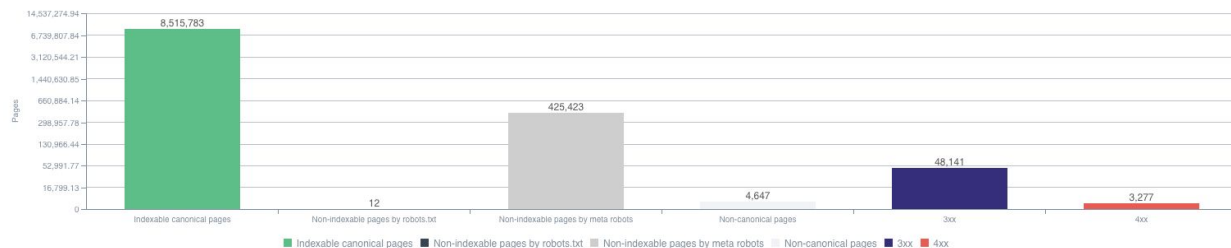"I'd probably start by looking at **indexability** and the **canonical strategy**."

— *Mickaël Serantes, Technical SEO Strategist, Oncrawl*

"Particularly with ecommerce sites with parameters, I'd focus on **indexability** and **crawl budget**."
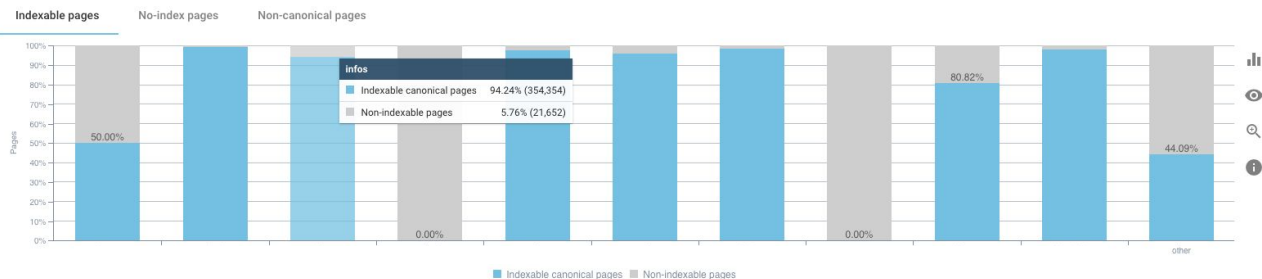
— *Adeline Lecellier, SEO Strategist, Oncrawl*

**oncrawl**

# Indexing & Strategy



**Pages by state of indexation**

| | 8,515,783 | 12 | 425,423 | 4,647 | 48,141 | 3,277 |
|---|---|---|---|---|---|---|

Indexable canonical pages · Non-indexable pages by robots.txt · Non-indexable pages by meta robots · Non-canonical pages · 3xx · 4xx

**Indexability breakdown**

Indexable pages | No-index pages | Non-canonical pages

infos
| | |
|---|---|
| ■ Indexable canonical pages | 94.24% (354,354) |
| ■ Non-indexable pages | 5.76% (21,652) |

50.00% · 0.00% · 0.00% · 80.82% · 44.09%
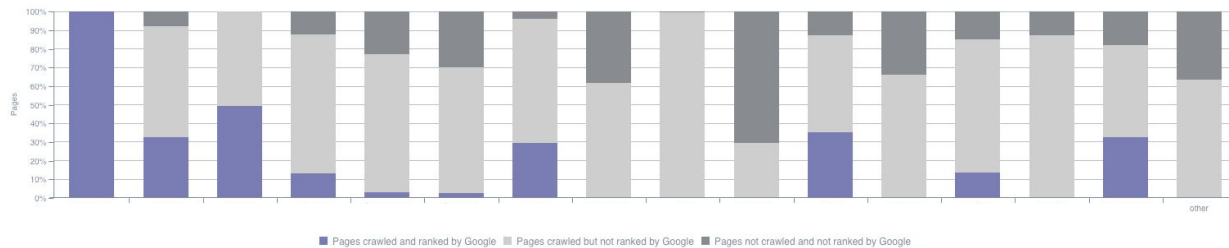
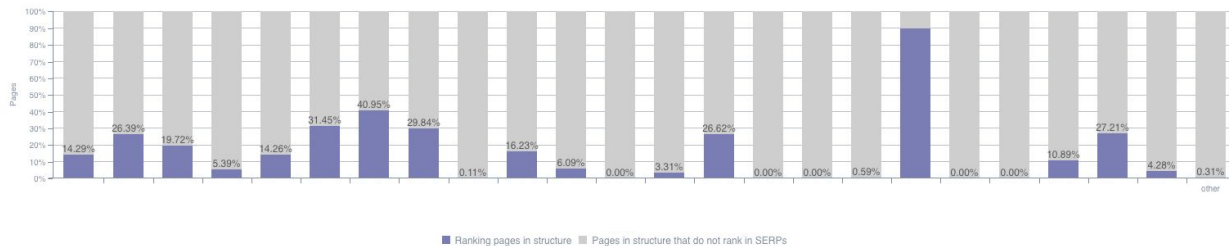■ Indexable canonical pages   ■ Non-indexable pages

# Indexing & Ranking



Pages crawled by Google: ratio of ranking pages



Pages in structure: ratio of ranking pages

# Indexing & Crawl budget
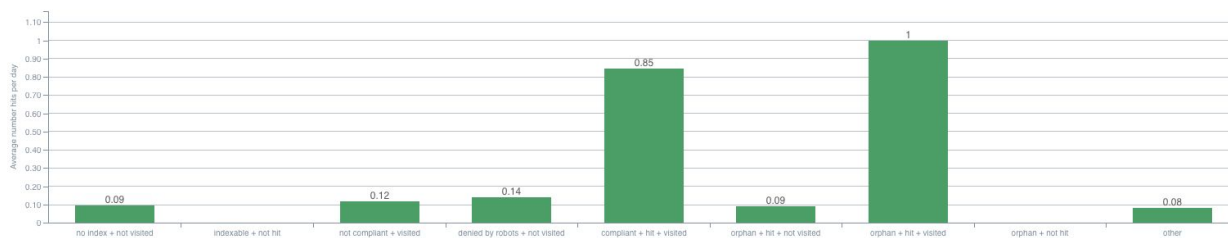
## Unique pages crawled ⓘ

| | no index + not visited | indexable + not hit | not compliant + visited | denied by robots + not visited | compliant + hit + visited | orphan + hit + not visited | orphan + hit + visited | orphan + not hit | other |
|---|---|---|---|---|---|---|---|---|---|
| | 22.25% | 0.00% | 26.30% | 0.00% | | | | 0.00% | 7.82% |

■ Known pages crawled by Google  ■ Known pages not crawled by Google

## Crawl frequency per day ⓘ

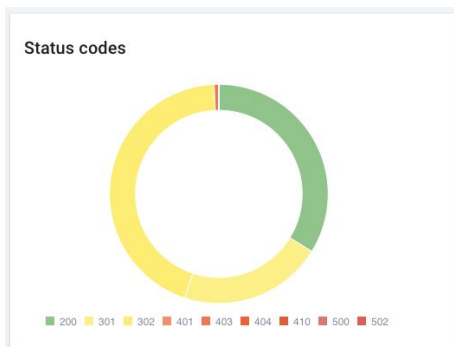| no index + not visited | indexable + not hit | not compliant + visited | denied by robots + not visited | compliant + hit + visited | orphan + hit + not visited | orphan + hit + visited | orphan + not hit | other |
|---|---|---|---|---|---|---|---|---|
| 0.09 | | 0.12 | 0.14 | 0.85 | 0.09 | 1 | | 0.08 |

"It's important for big sites to look at **overall health**, particularly **status codes**. Often if they've been using less enterprise-level tools, they have **never seen a full view** of their site health before."

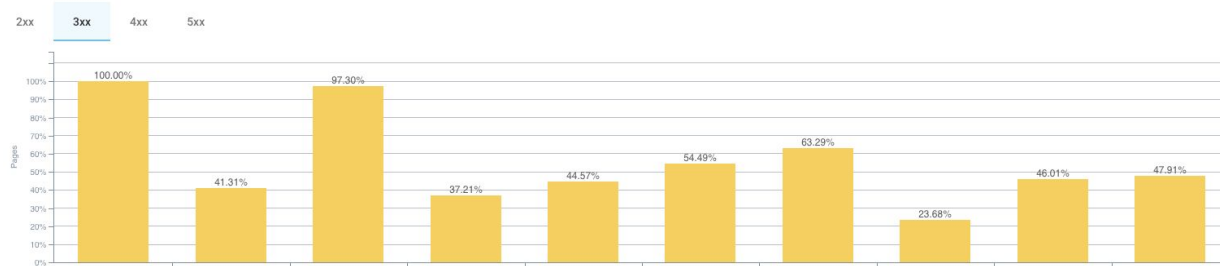*— Janaina Barreto-Romero, Customer Success Manager, Oncrawl*

"A good first focus is to check **sanity**. This works well with a sample of pages that is **representative of page templates**."

*— Jérôme Salomon, Technical SEO Strategist, Oncrawl*
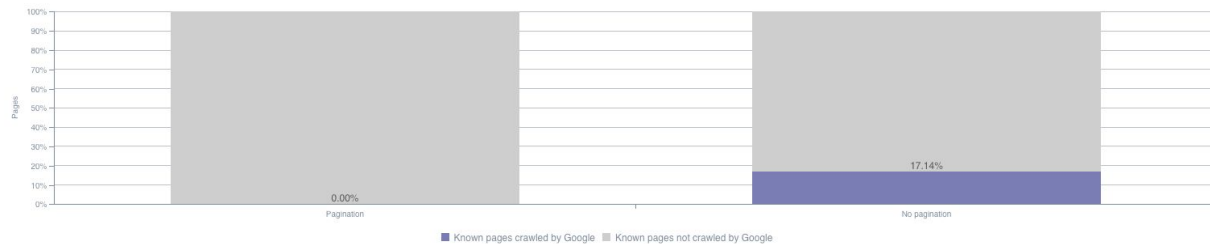
oncrawl

# Sanity checks & 3xx redirects

# Pagination

## Unique pages crawled ⓘ



Pagination: 0.00%
No pagination: 17.14%

■ Known pages crawled by Google  ■ Known pages not crawled by Google

## Page groups by depth ⓘ

ⓘ **Show only first 50 depths**
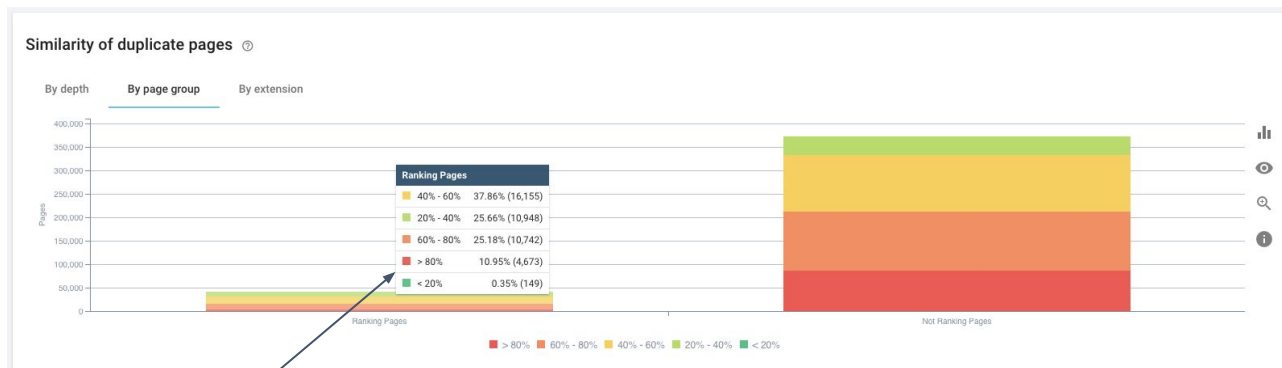


Depth

■ Pagination  ■ other

# Duplicate content



pages in these clusters declare a canonical URL...
BUT there's **more than one canonical URL** per cluster
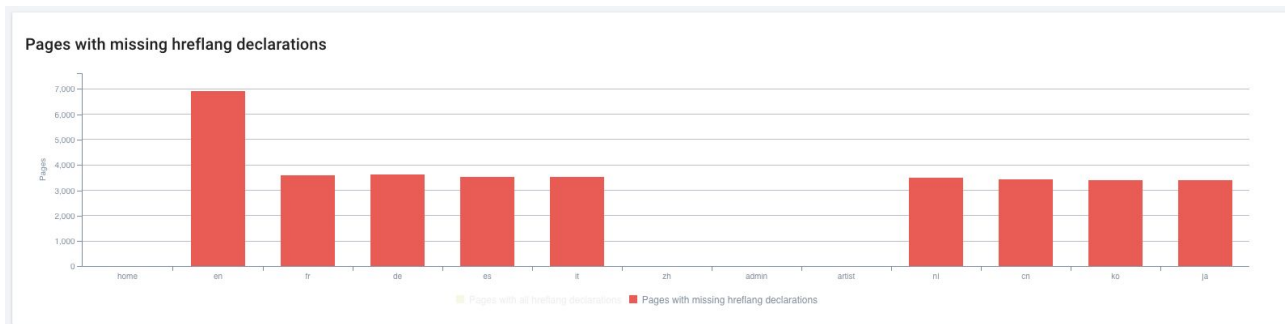
# Duplicate content



the pages that rank
that are also in a cluster with >80% similarity

# Hreflang



pages in these clusters use hreflang declarations to explain the similarities between the pages

# Hreflang



Pages with missing hreflang declarations



22,003
Non-indexable pages declared as hreflang

N/A

"It's also important to **ingest additional information**, like exports from Ahrefs or SEMrush."

— *Janaina Barreto-Romero, Customer Success Manager, Oncrawl*

"What's important? Sanity, indexing, and logs. And you should also be **ingesting conversion data** if you really want to improve SEO on a big site."

— *Julien Picot, Product Manager, Oncrawl*

**oncrawl**

Big sites are
**different**.

oncrawl

- Insist on what you need
- Use the right tools
- Focus
- Segment
- Prioritize
- "Quick wins"

Grab your controller,
it's time for Q&A!

oncrawl