

The Morphology of the Web is Changing!



Martin Lopatka

mlopatka@mozilla.com

Data Scientist/Applied Statistician



DATA NATIVES

moz://a

2018-11-22

└ The Morphology of the Web is

- Research Engineering and Data Science at Mozilla!
- The world's largest shared public resource... The Web.
- you may have heard of it.. it's called Firefox.
- "the study of form and structure without consideration of function"
- carrying out research on the nature, structure, and technology on the Modern web!



Disclaimer

Martin Lopatka provides this contribution to the Data Natives 2018 conference in a personal capacity. The views expressed are his own and do not necessarily represent the views of Mozilla Corporation or the Mozilla Foundation.



DATA NATIVES

moz://a

2018-11-22

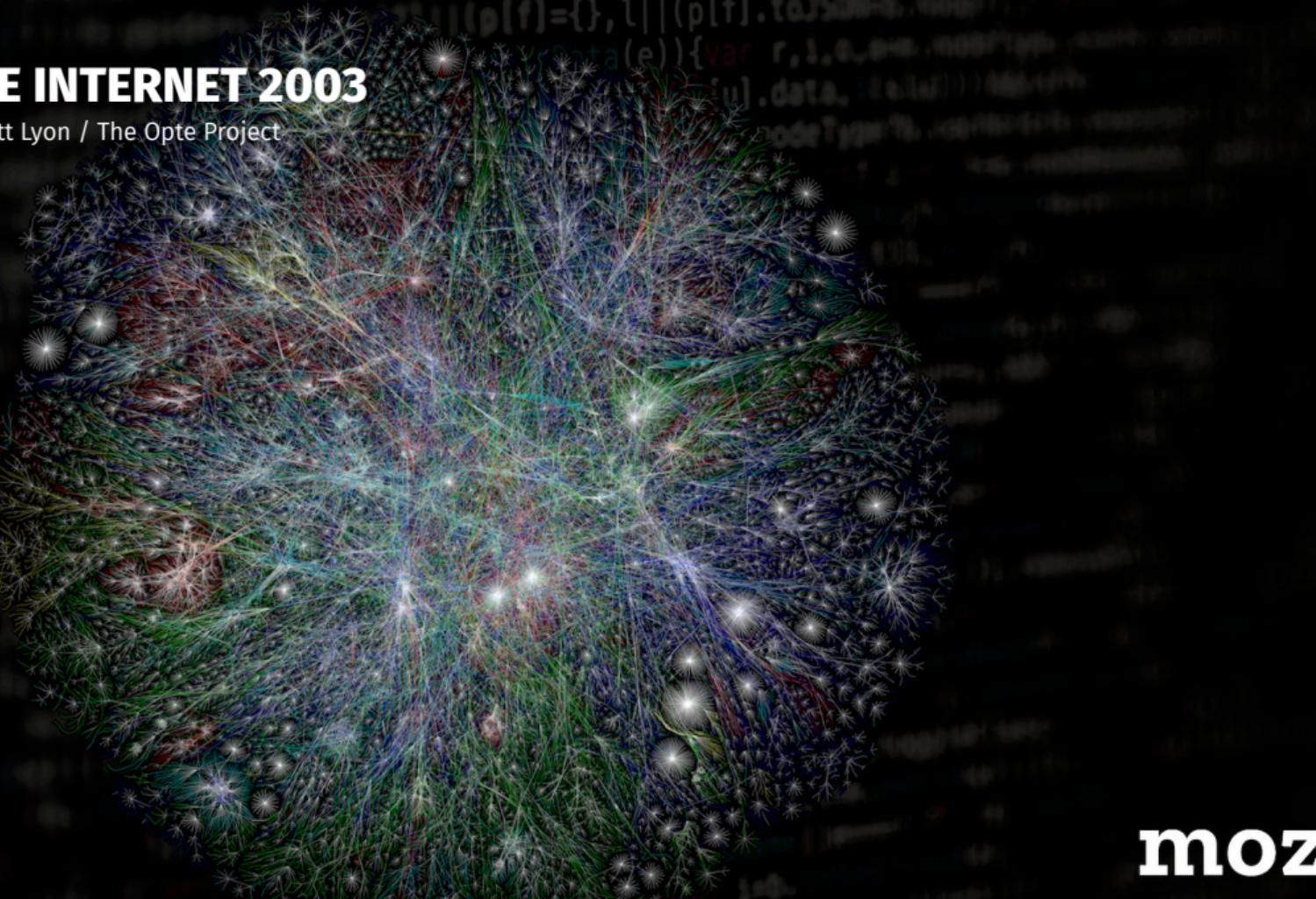
└ Disclaimer

Disclaimer! GO fast!

Martin Lopatka provides this contribution to the Data Natives 2018 conference in a personal capacity. The views expressed are his own and do not necessarily represent the views of Mozilla Corporation or the Mozilla Foundation.

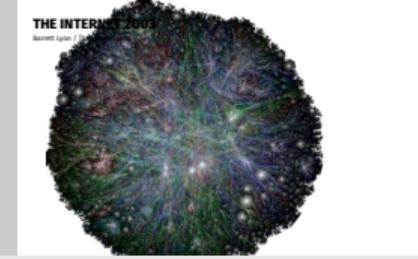
THE INTERNET 2003

Barrett Lyon / The Opte Project



moz://a

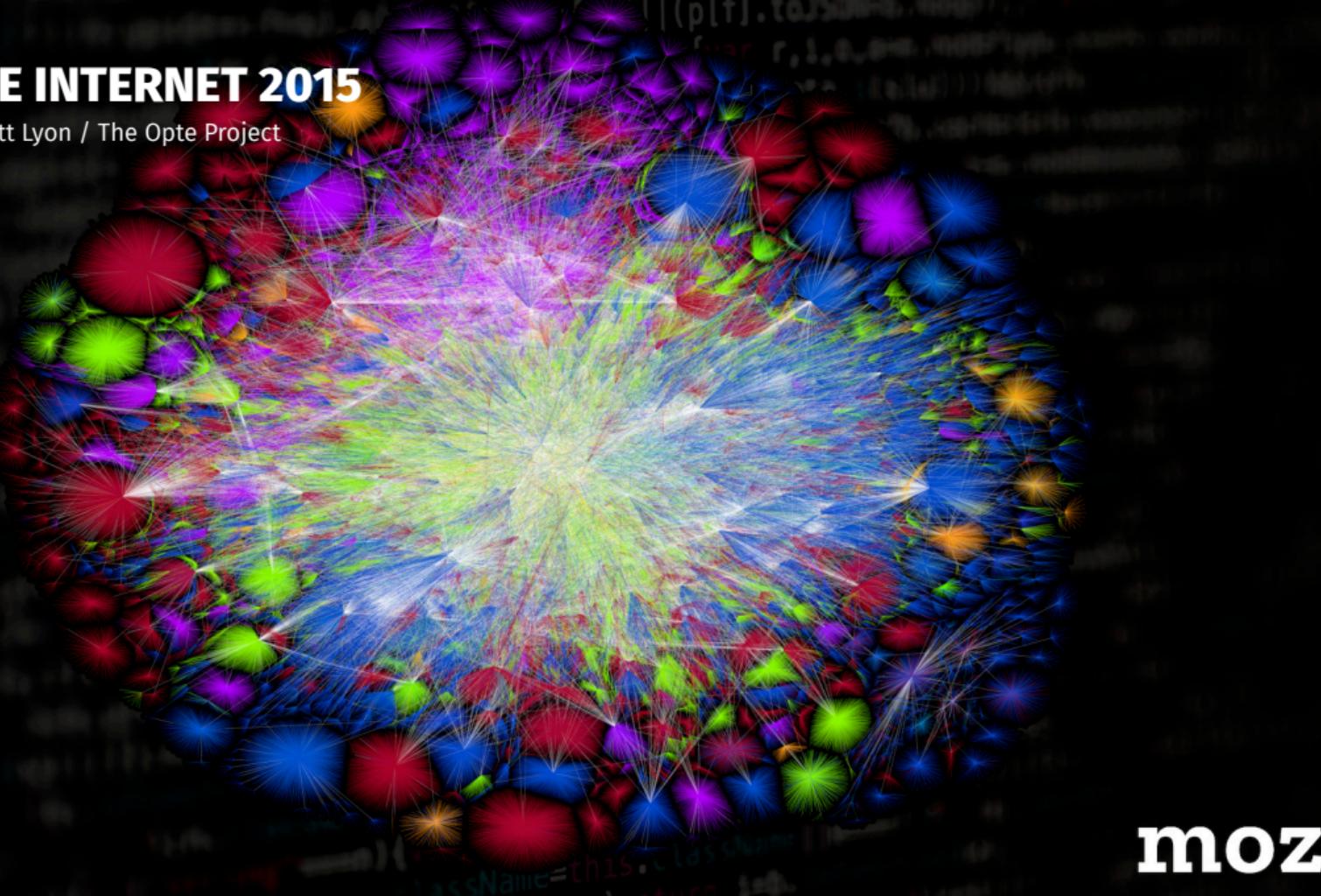
2018-11-22



- Paint a picture of the web's morphology in 2003
- Wikipedia young < 100K articles vs 6 million
- Reddit did not exist yet, Google had just released Page-rank (household name).
- Facebook won't launch until next year.
- Op-tee project visualisation made by a massive crawl fo the web,
- 30 million TLDs registered.
- links -> content creators. outbound links to other content. less log in. No SEO.

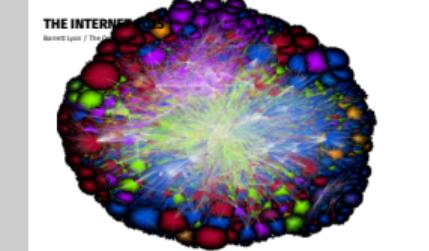
THE INTERNET 2015

Barrett Lyon / The Opte Project



moz://a

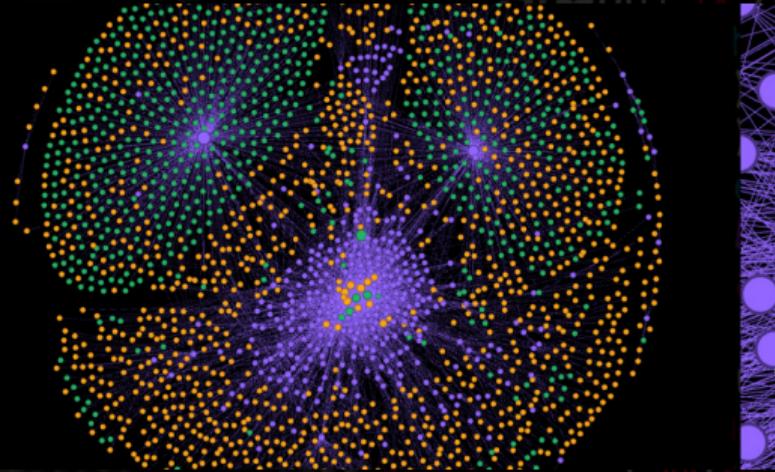
2018-11-22



The centrality of advertisers

Are trackers the new backbone of the Web?

<https://medium.com/firefox-context-graph/are-trackers-the-new-backbone-of-the-web-fb800435da15>



moz://a

2018-11-22

└ The centrality of advertisers

- investigation of these hyperlinked nodes to explore their nature
- we begin exploring the role of highly connected, hub nodes
- resources accessed in a third-party context.
Community detection was performed using Louvain Modularity
- fewer hyperlinks between pages belonging to separate tracking/advertising communities
- This is the first indication of a advertisement



How did we get to an ecosystem of silos?

- ▶ 40% of total Web browsing page views can be attributed to only 65 top level domains (TLD)
- ▶ Five (TLD+1) sites (Google, Facebook, Amazon, Yahoo, and Reddit) make up 22% of all traffic¹
- ▶ Of the 22,310,889 domains processed, 52.63% (11,742,112) were found to serve advertisements²

¹<https://medium.com/firefox-context-graph/are-trackers-the-new-backbone-of-the-web-fb800435da15>
²<https://commeica.com/2018/06/27/web-ad-prevalence/>



2018-11-22

└ How did we get to an ecosystem of

- We've already made some scary discoveries about the diversity
- "Comme ci comme ca" using the May 2018 batch of common crawl
- Easylist ruleset ->f Adblock adblocker
- Morphology is changing, important to discuss traffic and therefore content access

Social platforms were designed to facilitate; they became attention brokers, they are designed to help companies reach users as these sort of middlemen where they know a ton about you and thats the service they are providing to the advertisers.

Renee DiResta "The Internet's Original Sin"

27-Oct-2018 12:45



2018-11-22

- Mozfest, Renee DiResta: abuses of social media platforms
- platform purpose built for advertisement.
- *infrastructure for advertisement* content and social platforms have evolved to "maximise sustained engagement on site"
- Example: you consume your news, *on* facebook (60% US)

Social platforms were designed to facilitate; they became attention brokers, they are designed to help companies reach users as these sort of middlemen where they know a ton about you and thats the service they are providing to the advertisers.

Renee DiResta "The Internet's Original Sin"
27-Oct-2018 12:45

The highest trafficked pages on the Web³ are search engines, Social media platforms, and Commerce platforms

1. Google.com
2. Youtube.com
3. Facebook.com
4. Baidu.com
5. Wikipedia.org *
6. Qq.com
7. Taobao.com
8. Yahoo.com
9. Tmall.com
10. Amazon.com

2018-11-22

- highest trafficked pages on the web -> business model where engaged time on site is directly beneficial in a business sense.
- Web commerce platforms, advertisers, search engines, and Social media platforms.

³<https://www.alexa.com/topsites>; accessed 20-Nov-2018 11:17

The new ecology of outbound links

Dynamic pages feature changing content, showing different text, images and videos depending on who's visiting and when. This dynamic content includes different hyperlinks.

- ▶ retargeting
- ▶ advertisement
- ▶ realtime bidding
- ▶ social linking



DATA NATIVES

moz://a

2018-11-22

Dynamic pages feature changing content, showing different text, images and videos depending on who's visiting and when. This dynamic content includes different hyperlinks.

└ The new ecology of outbound links

- Retargeting, advertiser to tag you (with a cookie) and redirect as a specific media clicker, this requires a link through an advertiser's intermediate page
- Advertisement content, usually dynamically generated for the individual or profile of the page visitor
- Real Time bidding, the links surfaced may vary from one page visitor to another, based on estimates of *your* profile and value to advertisers, meaning crawlers will be served a fundamentally different



DATA NATIVES

moz://a

The problem is that the web is no longer built upon the simple premise of a collection of small static HTML and image files served up with a simple tag structure and readily parsed with a few lines of code. Today's web is richly dynamic, multimedia and increasingly broken into walled gardens and device-specific parallel webs.

Kalev Leetaru "Are Web Archives Failing The Modern Web: Video, Social Media, Dynamic Pages and The Mobile We"

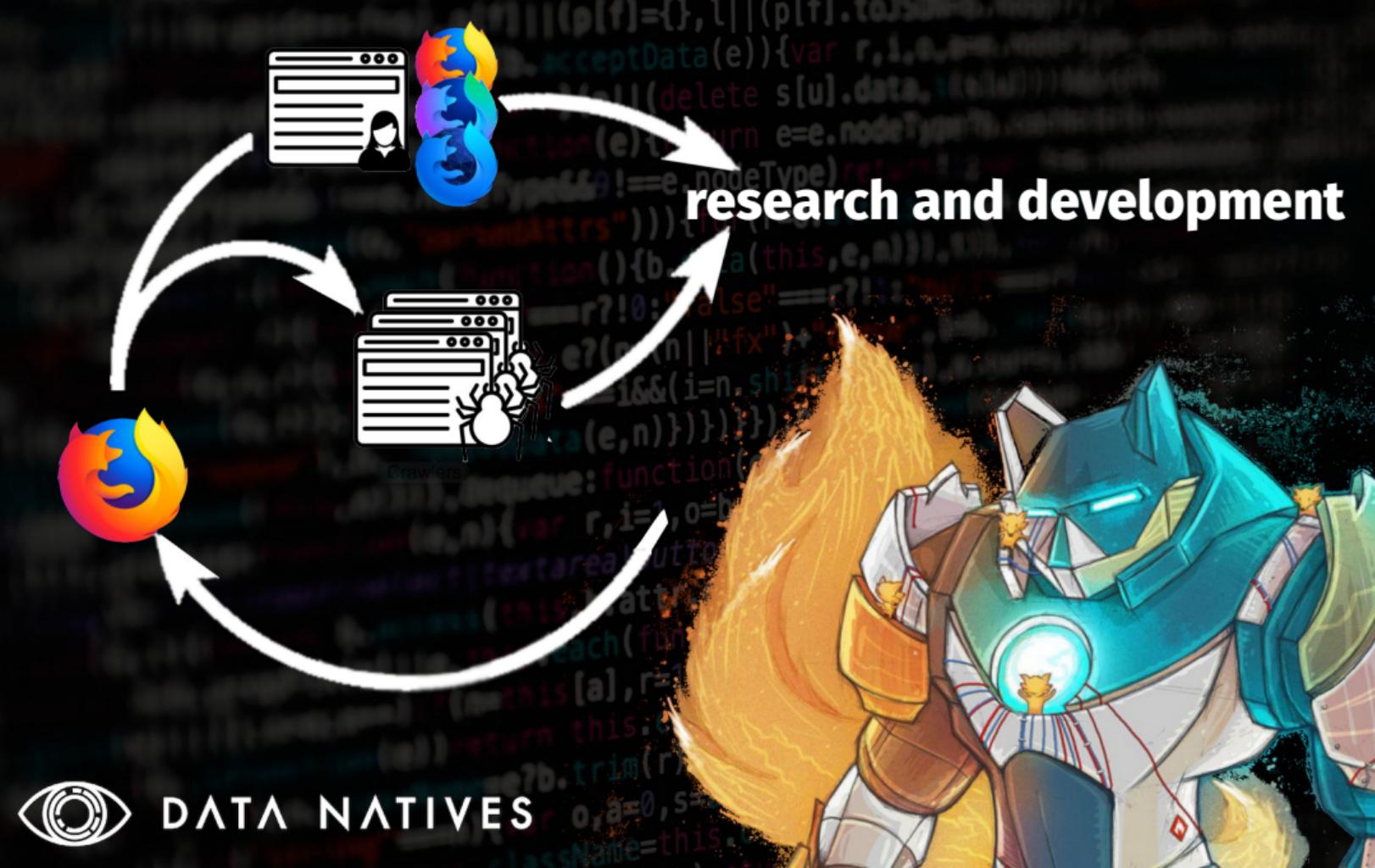
24-Feb-2017 22:41

2018-11-22

The problem is that the web is no longer built upon the simple premise of a collection of small static HTML and image files served up with a simple tag structure and readily parsed with a few lines of code. Today's web is richly dynamic, multimedia and increasingly broken into walled gardens and device-specific parallel webs.

Kalev Leetaru "Are Web Archives Failing The Modern Web: Video, Social Media, Dynamic Pages and The Mobile We"
24-Feb-2017 22:41

- from Kalev Letaru in a 2017 Forbes article discussing the role of the internet archive and common crawl initiatives in helping us make sense of the Web!



2018-11-22



- leverage crawler technology for it's robust technical measurement features.
- also want to user traffic patterns to ensure our feature development is relevant
- Crawlers, have a limited View - the open web makes up less and less of the Web.
- how do we study the nature of the web as people actually experience it in a Privacy respectful way that also takes advantage of scalable technology like Web crawlers?

Make Firefox better with pioneer



Firefox Pioneer
by Mozilla

<https://medium.com/firefox-context-graph/make-firefox-better-with-pioneer-10c82d0f9301>



DATA NATIVES

moz://a

2018-11-22

└ Make Firefox better with pioneer

- We ask permission,
- compliment the shortcomings of various data collection strategies
- The Firefox Pioneer program is an opt-in data collection
- real human readable consent policy



Firefox Pioneer
by Mozilla

<https://medium.com/firefox-context-graph/make-firefox-better-with-pioneer-10c82d0f9301>

Overscripted!

[https://github.com.mozilla/
overscripted/](https://github.com.mozilla/overscripted/)



DATA NATIVES

moz://a

2018-11-22

[https://github.com.mozilla/
overscripted/](https://github.com.mozilla/overscripted/)

- selective Javascript execution stack for tracking related activity
- 70Gb and touches over 2 million web pages seeded from the Alexa top 10K
- aspiring or established data scientist, come get your hands dirty
- data science can also work in a collaborative open source manner

Acknowledgements

- ▶ Sarah Bird
- ▶ Ruizhi You
- ▶ Victor Ng
- ▶ Louis Belleville
- ▶ David Zeber
- ▶ Calvin Luo
- ▶ Fredrik Wöllsen
- ▶ Zejun Yu
- ▶ Jason Thomas
- ▶ Vivian Jin
- ▶ Steven Englehardt
- ▶ Tyler Rubenuik
- ▶ Kyle Kung
- ▶ Alex McCallum



DATA NATIVES

moz://a

2018-11-22

└ Acknowledgements

- my team
- interns
- ops, secEng
- community



DATA NATIVES

moz://a

<https://github.com/mlopatka>

2018-11-22

- These slides are available on my GitHub

<https://github.com/mlopatka>