

# Virtual silos on the web

Exploring website clusters using hyperlinks and tracking code

Jelmer Neeven

Martin Lopatka (Mozilla)

Maarten Marx (UvA)

Hosein Azarbonyad (UvA PhD)

A collaboration with Mozilla

June 25, 2017

# Introduction

There is an increasing trend of recording, analyzing, and even selling users information and browsing behavior by third-party websites to companies.

Hypothesis: relationship between structural organisation of the Web and advertisement revenue model.

# Goals

- 1 Detect trackers using the embedded scripts and placed cookies for roughly 5 million webpages.
- 2 Build a communication graph of webpages, trackers, and stakeholders.
- 3 Use community detection algorithms to find virtual silos on the web.
- 4 Analyse and compare the silos to gain insight into hypothesised relationship.

## Seed URLs of pre-existing communities

- Need a starting point to get a list of websites.
- Using different communities, maximum diversity.
- Do the silos correspond to the communities?
- URLs gathered from `reddit.com`.

# Crawling the pages

- Get domain of URL.
- Get all third-party scripts and cookies.
- Count internal/external links.
- *robots.txt* permissions.

## Crawler results

- Data from roughly 355,000 URLs.
- 83,657 domains.
- jQuery not in top 50.

## Crawler results

A few of the most frequently used script sources:

- `google-analytics.com/analytics.js`
- `platform.twitter.com/widgets.js`
- `securepubads.g.doubleclick.net/gampad/ads`
- `b.scorecardresearch.com/beacon.js`

# Communication graphs

Three graphs:

- Hyperlink
- Tracking
- Total

Detect communities with Louvain Modularity Optimisation algorithm.



# Communication graphs

<b>Graph:</b>	<b>Communities</b>	<b>Modularity</b>
Hyperlink	671	0.7
Tracking	4169	0.62
Total	435	0.64

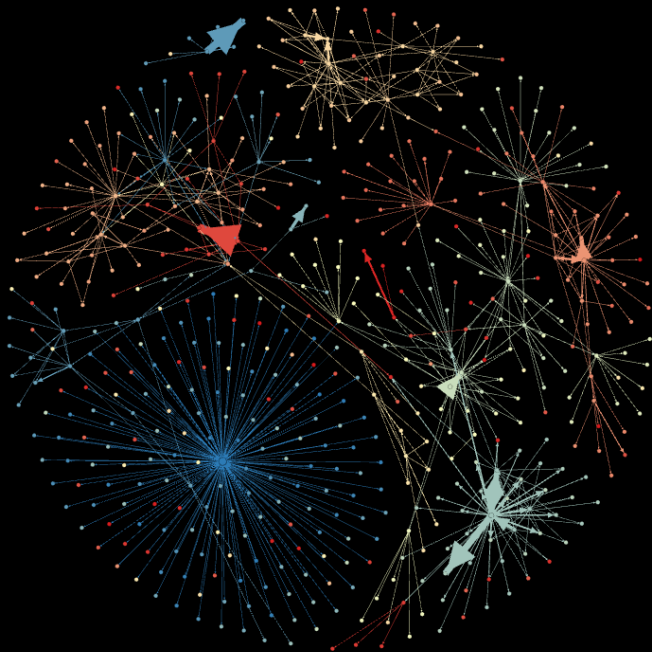
# Communication graphs

Calculate metrics of communities:

- size
- density
- degree
- connectivity

Calculated on induced graph so they can be compared meaningfully.







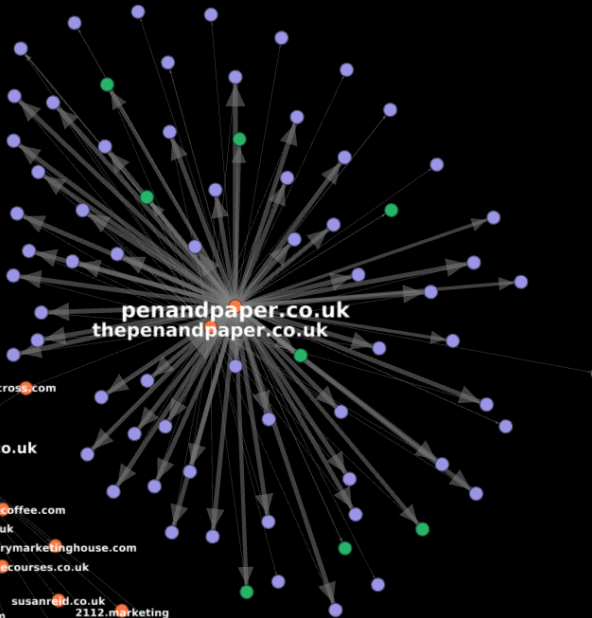
## Total silos

Manually analysed twenty silos:

- Twelve contain exactly the same domains.
- Six silos were (part of) a bigger silo, now separate because of cookies/scripts. However, can also be identified by Louvain.
- One lost a domain because it uses popular cookies/scripts.
- The following silo only exists within the total communities, which is remarkable:

clubmed.com  
lidacucina.co.uk  
envestry.com  
linn.co.uk  
sunseekerlondon.com  
eightraymusic.com  
luxurystudent.com  
itcluxurytravel.co.uk  
aji.co.uk  
whitehouse-cox.co.uk  
dfjewellery.co.uk  
berkeleyhomes.co.uk  
theluxurynetwork.co.uk  
hunter-design.co.uk  
call-systems.com  
differencecoffee.com  
redeyeevents.co.uk  
luxurymarketghouse.com  
tinini.com  
newmarketracescourses.co.uk  
inspiredlive.co.uk  
hl-cruises.com  
susanreid.co.uk  
2112.marketing  
moethennessy-selection.co.uk  
cornerstone.co.uk  
smokesilvertravel.com  
edgeretreats.com

penandpaper.co.uk  
thepenandpaper.co.uk

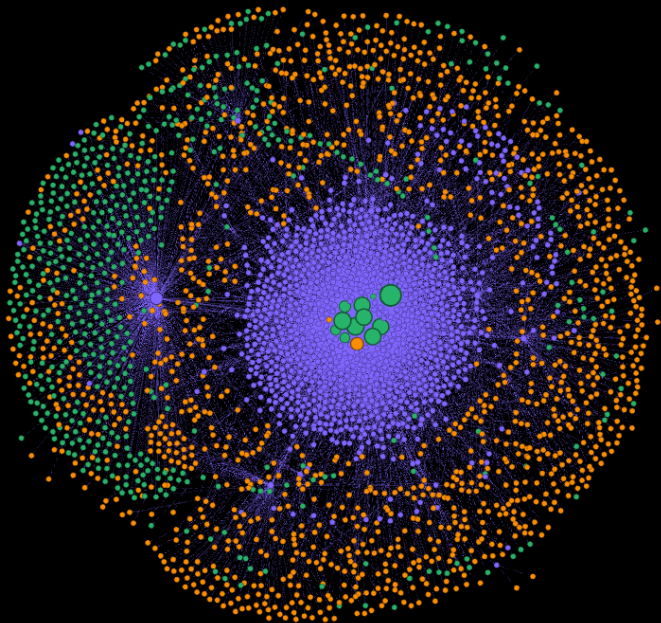


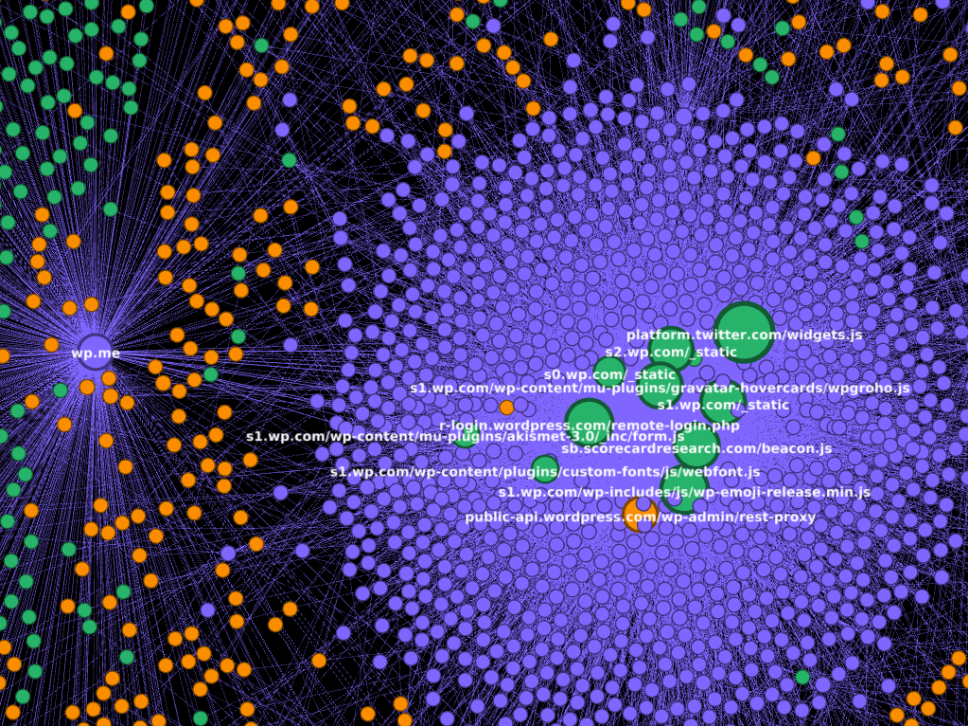
## Preliminary conclusion

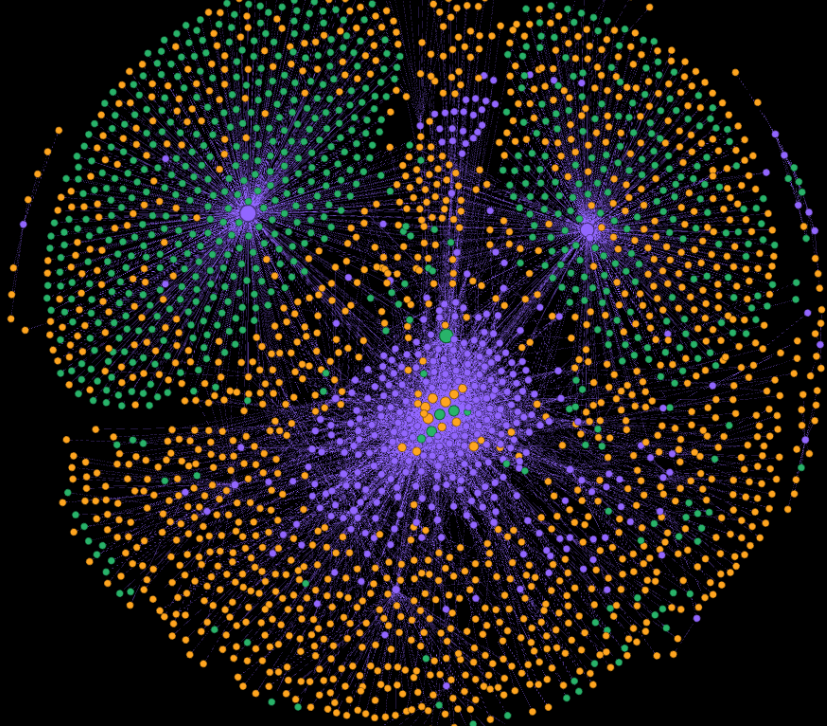
- Incorporating tracking data in communication graph does not yield many interesting results.
- In fact, sometimes sites are clustered together that are not linked in any way.
- Reduced connectivity, density, and average degree.

However, the bigger silos with a higher degree have not yet been considered.











s.pubmine.com/adj/114160/300/250

s.pubmine.com/adj/26942/300/250

aax-eu.amazon-adsystem.com/s/iu3

ads.pubmatic.com/AdServer/js/gshowad.js

us-u.openx.net/w/1.0/cm

gads.pubmatic.com/AdServer/AdCallAggregator

as-sec.casalemedia.com/cygnus

s.pubmine.com/head.js

impreump.basebanner.com/st

cmbestsrv.com/vpaid/units/13\_8\_9/creatives/creative\_js.js

static.criteo.net/js/ld/publishertag.js

cdn.teads.tv/media/format/v3/teads-format.min.js

sync.teads.tv/iframe

cas.criteo.com/delivery/ajs.php

sync.teads.tv/iframe/redirect

# Conclusion

- Although very labour-intensive, it is indeed possible to detect tracking networks using this approach.
- The tracking networks seem to be 'hidden' within larger networks.
- In general, it appears that the detected communities can be further split up.

## Future research

- Automated approach: recursive Louvain
- Tracking probabilities
- Different community detection algorithms

# End of presentation

Thank you for your attention!

## Appendix

Advertisement network:

- Inner edges: 66972
- Outer edges: 112843

Hyperlink connected:

- Total hyperlinks found: 827724
- Amount of connected domains sharing at least one piece of tracking code: 84234



