

# Tor: Network vs. Web vs. Web Graph

The **Tor network** is an overlay network composed of thousands We scraped the Tor Web for five months, obtaining **three snapshots** of the Tor Web which constitute the **widest exploration** of routers, called relays. It is an **anonymity** network, representof the Tor Web in the literature. To capture persistent features, we defined a **CORE** Graph, which is induced by **stable edges**. ing one of the main infrastructures of the **DarkWeb**.

The Tor network allows:

- accessing the Internet anonymously;
- running anonymous Web servers, known as hidden services.





Hidden services:

- can only be accessed using a Tor-enabled browser:
- are identified by a non-mnemonic 16-character string;
- have a specific .onion domain.

Web resources hosted on Tor constitute the **Tor Web**. Tor Web pages may be connected through **hyperlinks**, thus defining a Web graph.



We study the properties of the Tor Web Graph aggregated by Hidden Service.



# Previously @CRANIC (IAC-CNR):

- [BCGL17] Massimo Bernaschi, Alessandro Celestini, Stefano Guarino, and Flavio Lombardi, Exploring and Analyzing the Tor Hidden Services Graph, ACM Trans. Web **11** (2017), no. 4, 24:1–24:26.
- Alessandro Celestini and Stefano Guarino, Design, implementation and test of a [CG17] flexible tor-oriented web mining toolkit, Proceedings of the 7th International Conference on Web Intelligence, Mining and Semantics (Amantea, Italy), WIMS '17, ACM, 2017, pp. 19:1-19:10.

# **Spiders like Onions: on the Network of Tor Hidden Services**

Alessandro Celestini Stefano Guarino

Institute for Applied Computing, National Research Council of Italy

# Data, Graphs, Persistence



# **Graph Structure**

Web graphs are typically scale-free (both in- and out-degrees are power-laws) and have a bow-tie structure. In all Tor Hidden Services Graphs (SNP1, SNP2, SNP3, CORE), the out-degree is not a power-law, it has heavier bulk and tail. The Graphs have a small average distance and a pronounced clustering/transitivity. No simple graph model fully explains this structure.

Erdos-Renyi		Paths length	Clustering/Transitivity	Degree distribution
$C \sim \langle \deg \rangle /  V  \\ \langle dist \rangle \sim \ln( V )$		8 6 10	8 6 10	8 6 10
Scale-Free				
$\lambda > 3 : \langle dist \rangle \sim \ln( V )$		3	3	3
$2 < \lambda < 3 : \langle dist \rangle \sim \ln \ln( V )$	graph	$\langle dist \rangle = \ln  V  - \ln \ln  V $	$C$ $T_1$ $T_2$ $\left  2 \langle \deg_{in} \rangle \right/$	V  deg <sub>in</sub> deg <sub>out</sub>
Small-World $C \gg \langle \deg \rangle /  V $ $\langle dist \rangle \sim \ln( V )$	SNP1 SNP2 SNP3 CORE	3.793169.4592.2474.9602810.1382.3163.664559.7672.2793.982918.9442.190	0.009430.006170.00390.00080.004070.007790.00160.00030.004920.002530.001970.00060.008730.005350.003560.0009	77 $\lambda = -2.88$ heavy-tail52 $\lambda = -3.32$ heavy-tail78 $\lambda = -3.35$ heavy-tail83 $\lambda = -2.7$ heavy-tail

The Graph is dominated by few **out-hubs** (mostly hidden directories) that link to most of the other hidden services. This also emerges in its highly asymmetrical bow-tie structure, which is always single-ended except for the CORE Graph.

- 83% to 95% of the services are **sinks**:
- Top **out-hubs** are 50× greater than top in-hubs;
- The 10 top out-hubs link to 80% to 95% of the entire graph.



# Flavio Lombardi

hidden service that do not host (ii) onion urls only accessing Tor-specific messaging services

There is **no** evident **correlation** among centrality metrics in the Tor Hidden Services Graph. Yet, **important** nodes have either:

## Large Out-degree & Betweenness





Ar	าล	ly:	Ζ

- pages).

# Studying the topology of the hidden services graph:

# pages contacted (by response type) 1821842 277813 197128 141205 2339718 471519 262403 324552 765876 393018 105406 67115





DISCONNECTED

In all snapshots, the LSCC is small, the OUT set is very large, all other

• The **CORE** is similar, but the IN is  $\approx 0.1\%$  and the TENDRILS are  $\approx 1\%$ .







# **Centralities and Central Hidden Services**

ion	Topic
derdj5ziov3ic7	Directory
l32teyptf4tvi	Directory
wxqmjaes2bae5	Directory

### Large In-degree & PageRank



### After **removing** the 10 top out-hubs from the CORE:

• C grew but  $\langle dist \rangle$  remained almost unvaried; • the **PageRank** can be used to identify **markets**, while the Betweenness becomes less informative.

# Conclusions

### zing Tor over time:

Exploring all Tor is virtually impossible (with a crawler). • Tor services are **highly volatile**, but there is a significant **stable** core, with slightly different properties. • **Connections** are also volatile (possibly due to **volatile** 

The Tor Web graph is (ultra?) small, only partially governed by **preferential attachment**, with clustering **not** large enough for being a (WS) small-world. The out-hubs determine the topology and may hide other features. Central services are mostly link directories and markets.

### Find more on http://www.cranic.it