

Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists

Oliver Gasser, Quirin Scheitle, Paweł Foremski, Qasim Lone, Maciej Korczyński, Stephen D. Strowes, Luuk Hendriks, Georg Carle

IMC 2018, Boston



Joint work





















THE EXPANSE

EXPANSE

of the IPv6 Address Space

Previous work

ПΠ

Previous work on IPv6 address space analysis

- Dhamdhere et al. (2012)
- Czyz et al. (2014)
- Plonka and Berger (2015, 2017)
- Ullrich et al. (2015)
- Gasser et al. (2016)
- Rohrer et al. (2016)
- Foremski et al. (2016)
- Murdock et al. (2017)
- Fiebig et al. (2017, 2018)
- Borgolte et al. (2018)

Open questions



1. How balanced are different hitlist sources?

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 5





- 1. How balanced are different hitlist sources?
- 2. Can we identify addressing schemes to find new addresses?

Open questions



- 1. How balanced are different hitlist sources?
- 2. Can we identify addressing schemes to find new addresses?
- 3. What is the influence of aliased prefixes on IPv6 hitlists?

Open questions



- 1. How balanced are different hitlist sources?
- 2. Can we identify addressing schemes to find new addresses?
- 3. What is the influence of aliased prefixes on IPv6 hitlists?
- 4. How does cross-protocol responsiveness in IPv6 differ from IPv4?



- 1. How balanced are different hitlist sources?
- 2. Can we identify addressing schemes to find new addresses?
- 3. What is the influence of aliased prefixes on IPv6 hitlists?
- 4. How does cross-protocol responsiveness in IPv6 differ from IPv4?
- 5. Is there a benefit of using more than one address learning tool?



1. How balanced are different hitlist sources?



Where can we learn potential IPv6 addresses?



Where can we learn potential IPv6 addresses?



ТШ

Hitlist sources

Where can we learn potential IPv6 addresses?



Figure 1: Cumulative runup of IPv6 addresses.

ТШ

7

Hitlist sources

Where can we learn potential IPv6 addresses?



Figure 1: Cumulative runup of IPv6 addresses.

Address distribution

- Many addresses from domainlists, CT, and scamper
- Rapid increase of scamper addresses due to CPE routers



How balanced are the addresses from different sources?



How balanced are the addresses from different sources?



Figure 2: AS distribution for hitlist sources.



8

How balanced are the addresses from different sources?



Figure 2: AS distribution for hitlist sources.

Autonomous System distribution

• Unbalanced (CT, domainlists) vs. balanced (RIPE Atlas)

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists



How much of the announced address space do we cover?

Excursion: Visualizing prefixes



Visualizing prefixes using Hilbert space-filling curves



Figure 3: IPv4

Excursion: Visualizing prefixes



Visualizing prefixes using Hilbert space-filling curves





Figure 3: IPv4



Figures by Ben Cartwright-Cox https://blog.benjojo.co.uk/post/scan-ping-the-internet-hilbert-curve

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 10



How much of the announced address space do we cover?



How much of the announced address space do we cover?



Figure 5: Number of addresses per prefix.



How much of the announced address space do we cover?



Figure 5: Number of addresses per prefix.

zesplot

- IPv6 prefix visualization tool
- Input: set of IPv6 prefixes
- Each plotted as rectangle
- Prefixes of same AS and size are plotted adjacently
- Color based on metric (e.g. number of addrs. in prefix)



How much of the announced address space do we cover?



Figure 5: Number of addresses per prefix.

BGP prefix distribution

zesplot

- IPv6 prefix visualization tool
- Input: set of IPv6 prefixes
- Each plotted as rectangle
- Prefixes of same AS and size are plotted adjacently
- Color based on metric (e.g. number of addrs. in prefix)

- Good coverage of BGP prefixes: 25.5 k of 51.2 k
- Some prefixes with many addresses

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 11



2. Can we identify common addressing schemes in our hitlist?

Entropy clustering



Understand addressing patterns in IPv6 hitlists



пп

Networks have different entropy fingerprints



пп

1. Fingerprint each network



Networks have different entropy fingerprints

- 1. Fingerprint each network
- 2. Feed to k-means clustering



Networks have different entropy fingerprints

- 1. Fingerprint each network
- 2. Feed to k-means clustering
- 3. Plot median fingerprints and cluster popularity

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 13

Entropy clustering



IPv6 interface identifiers (IIDs)



Figure 6: Hitlist addressing schemes for IIDs.

Entropy clustering



IPv6 interface identifiers (IIDs)



Figure 6: Hitlist addressing schemes for IIDs.

- The IPv6 networks we cover employ predictable IIDs
- Also visible: privacy extensions, modified EUI-64 (ff:fe)



Entropy clustering Full IPv6 fingerprints



Figure 7: Hitlist addressing schemes for full addresses.



Entropy clustering Full IPv6 fingerprints



Figure 7: Hitlist addressing schemes for full addresses.

- Just a handful of schemes on the Internet
- Addresses largely predictable

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 15



3. What is the influence of aliased prefixes on IPv6 hitlists?



Taxonomy:

- Alias: another address of the same host
- Aliased prefix: whole prefix bound to the same host
- Bias: some hosts overrepresented due to aliased prefixes



Taxonomy:

- Alias: another address of the same host
- Aliased prefix: whole prefix bound to the same host
- Bias: some hosts overrepresented due to aliased prefixes



Figure 8: Multi-level aliased prefix detection using pseudo-random probing.



Results



Results



Figure 9: All prefixes covered by hitlist.



Figure 10: Aliased prefixes.



Results



Figure 9: All prefixes covered by hitlist.



Figure 10: Aliased prefixes.

- Only 3.2 % of prefixes are aliased
- But 46.6 % of addresses are in aliased prefixes
- Validated using fingerprinting (iTTL, TCP opts, timestamps)



4. How does cross-protocol responsiveness in IPv6 differ from IPv4?

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 19

- If address responds on protocol X, how likely is it to respond on protocol Y?
- Goal: Identify relevant addresses for specific measurements





Figure 11: Cross-protocol responsiveness between services.





Figure 11: Cross-protocol responsiveness between services.

- If responsive to any of the probes \rightarrow at least 89% probability it will answer to ICMPv6





Figure 11: Cross-protocol responsiveness between services.

- If responsive to any of the probes \rightarrow at least 89% probability it will answer to ICMPv6 vs. 73% in IPv4





Figure 11: Cross-protocol responsiveness between services.

- If responsive to any of the probes \rightarrow at least 89% probability it will answer to ICMPv6 vs. 73% in IPv4
- Web protocols: QUIC \rightarrow HTTPS and HTTP, HTTPS \rightarrow HTTP; but not the other way around

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 21



5. Is there a benefit of using more than one address learning tool?



Techniques to learn new addresses

- Entropy/IP: Generate new addresses by leveraging entropy of seed addresses
 - Similar approach to grouping addresses based on their structure as shown earlier



Techniques to learn new addresses

- Entropy/IP: Generate new addresses by leveraging entropy of seed addresses
 - Similar approach to grouping addresses based on their structure as shown earlier
- 6Gen: Generate new addresses in dense address regions
 - If we see addresses
 - 2001:0db8:0407:8000::4
 - 2001:0db8:0407:8000::5
 - 2001:0db8:0407:8000::**8**
 - Likely other valid addresses
 - 2001:0db8:0407:8000::6
 - 2001:0db8:0407:8000::7



How well do Entropy/IP and 6Gen perform?

- Input: All previously found IPv6 addresses
- Responsiveness: 278 k (of 118 M) and 489 k (of 129 M)



How well do Entropy/IP and 6Gen perform?

- Input: All previously found IPv6 addresses
- Responsiveness: 278 k (of 118 M) and 489 k (of 129 M)
- Overlap of only 675 k generated addresses
- 10x higher response rate for overlapping addresses



How well do Entropy/IP and 6Gen perform?

- Input: All previously found IPv6 addresses
- Responsiveness: 278 k (of 118 M) and 489 k (of 129 M)
- Overlap of only 675 k generated addresses
- 10x higher response rate for overlapping addresses

Table 1: Top 5 responsive protocol combinations for Entropy/IP and 6Gen.

ICMPv6	TCP/80	TCP/443	UDP/53	UDP/443	Entropy/IP	6Gen
 Image: A start of the start of	×	×	×	×	41.1%	66.8%
1	✓	\checkmark	×	×	12.3 %	9.2%
×	×	×	\checkmark	×	23.1 %	7.3%
\checkmark	\checkmark	×	×	×	3.4 %	4.9%
\checkmark	\checkmark	1	×	1	6.1 %	3.2%

Different host populations

Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 24

Community contributions



Reproducibility

- We publish data, code, and analysis scripts
- DOI: 10.14459/2018mp1452739

Software and tools published on GitHub

- ZMapv6
- zesplot
- Entropy clustering
- New Entropy/IP generator
- Entropy/IP open-sourced (thanks to Akamai)

IPv6 Hitlist Service



A one-off analysis is all well and good, but what if I need an up-to-date IPv6 hitlist for my research starting next month?

IPv6 Hitlist Service



A one-off analysis is all well and good, but what if I need an up-to-date IPv6 hitlist for my research starting next month?

ipv6hitlist.github.io

IPv6 Hitlist Service



A one-off analysis is all well and good, but what if I need an up-to-date IPv6 hitlist for my research starting next month?

ipv6hitlist.github.io

- Daily IPv6 hitlists and aliased prefixes available for download
- Interactive zesplots
- Continuously updated graphs



Gasser et al. — Clusters in the Expanse: Understanding and Unbiasing IPv6 Hitlists 26



Repeating addressing schemes

Repeating addressing schemes

Aliased prefixes

Repeating addressing schemes

Aliased prefixes

Conditional responsiveness

Repeating addressing schemes

Aliased prefixes

Conditional responsiveness

Learning unknowns

Repeating addressing schemes

Aliased prefixes

ipv6hitlist_github_io

Conditional responsiveness

Learning unknowns

Oliver Gasser <gasser@net.in.tum.de>