



# Hypergraph Analytics of Domain Name System Relationships

Cliff A. Joslyn<sup>1</sup>(✉), Sinan Aksoy<sup>2</sup>, Dustin Arendt<sup>2</sup>, Jesun Firoz<sup>1</sup>, Louis Jenkins<sup>3</sup>, Brenda Praggastis<sup>1</sup>, Emilie Purvine<sup>1</sup>, and Marcin Zalewski<sup>4</sup>

<sup>1</sup> Pacific Northwest National Laboratory, Seattle, WA, USA  
cliff.joslyn@pnnl.gov

<sup>2</sup> Pacific Northwest National Laboratory, Richland, WA, USA

<sup>3</sup> University of Rochester, Rochester, NY, USA

<sup>4</sup> NVIDIA, Santa Clara, CA, USA

**Abstract.** We report on the use of novel mathematical methods in hypergraph analytics over a large quantity of DNS data. Hypergraphs generalize graphs, as used in network science, to better model complex multiway relations in cyber data. Specifically, casting DNS data from Georgia Tech’s ActiveDNS repository as hypergraphs allows us to fully represent the interactions between *collections* of domains and IP addresses. To facilitate large-scale analytics, we fielded an analytical pipeline of two capabilities: HyperNetX (HNX) is a Python package for the exploration and visualization of hypergraphs; while on the backend, the Chapel HyperGraph Library (CHGL) is a library for high performance hypergraph analytics written in the exascale programming language Chapel. CHGL was used to process gigascale DNS data, performing compute-intensive calculations for data reduction and segmentation. Identified portions are then sent to HNX for both exploratory analysis and knowledge discovery targeting known tactics, techniques, and procedures.

**Keywords:** Hypergraphs · DNS · High performance computing · Chapel

## 1 Introduction

Many problems in data analytics involve rich interactions amongst multiple entities, for which graph representations are commonly used. High order (high dimensional) interactions abound in cyber and social networks, and can only be represented in graphs as highly inefficiently coded, “reified” labeled subgraphs. Lacking multi-dimensional relations, it is hard to address questions of “community interaction” in graphs: how is a collection of entities  $A$  connected to another collection  $B$  through chains of other communities?; where does a particular community stand in relation to other communities in its neighborhood?

*Hypergraphs* [4] are generalizations of graphs which allow edges to connect any number of vertices. Hypergraph methods are well known in discrete mathematics, and are closely related to important objects in data science such as bipartite graphs, set systems, partial orders, finite topologies, and especially graphs proper, which they directly generalize (every graph is a 2-uniform hypergraph). In HPC, hypergraph-partitioning methods help enable parallel matrix computations [8], and have applications in VLSI [13]. In the network science literature, researchers have devised several path and motif-based hypergraph data analytics (albeit fewer than their graph counterparts), such as in clustering coefficients [15] and centrality metrics [9].

Complex data commonly analyzed using network science methods, and especially including cyber data, often contain multi-way interactions. But while they thus present naturally as hypergraphs, still hypergraph treatments are very unusual compared to graph representations of the same data. This is due at least to the greater mathematical, conceptual, and computational complexity of hypergraph methods, since as data objects, hypergraphs scale as  $O(2^n)$  in the number of vertices  $n$ , as opposed to  $O(n^2)$  for graphs. In the face of this, complex data are typically collapsed or are simplified to graphs to ease analysis.

We are accepting the challenge of the complexity of hypergraphs in order to gain the formal clarity and support for analysis of complex data they provide. A substantial high-performance computing (HPC) component is thus necessary, despite hypergraph analytics not receiving much attention in the software engineering community at large, and the HPC community in particular. We thus pursue a two-fold approach to developing our methods:

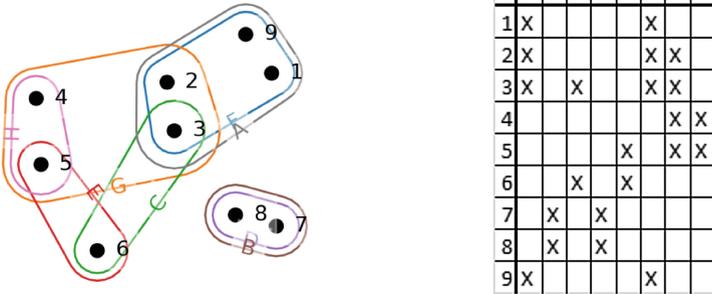
1. The **Chapel Hypergraph Library** (CHGL, <https://github.com/pnml/chgl>) [12] is a library for hypergraph computation in the emerging Chapel programming language [6, 7], for HPC hypergraph processing, large scale analysis, and data segmentation.
2. We explore single hypergraphs or collections of hypergraphs using **HyperNetX** (HNX, <https://github.com/pnml/HyperNetX>), a Python library for exploratory data analytics and visualization of hypergraphs.

In our work, CHGL and HNX are two stages of an analytical pipeline: CHGL provides a highly abstract interface for implementation of HPC hypergraph algorithms over large data, identifying segments and subsets which can then be passed to HNX for more detailed analysis.

In this paper we first introduce the foundations of hypergraph mathematics and hypernetwork science in the context of our CHGL and HNX capabilities. We then describe the DNS data set, selections of the ActiveDNS data sets from the Georgia Institute of Technology [1]. We then describe CHGL, before going on to describe the results of our demonstration analyses. These include both basic global statistics like degree and edge size distributions, as well as exploratory discovery of small components involving motif mining and computation of simple hypergraph metrics to discover outliers.

## 2 Hypergraph Analytics

An undirected **hypergraph** is a pair  $\mathcal{H} = \langle V, \mathcal{E} \rangle$  with  $V$  a finite, non-empty set of **vertices**, and  $\mathcal{E}$  a non-empty multiset of **hyperedges**  $e \in \mathcal{E}$  (or just “edges” when clear), where  $\forall e \in \mathcal{E}, e \subseteq V$ . Hypergraphs can be represented in many forms, two of which are shown in Fig. 1 for a small example  $\mathcal{H}$  with  $V = \{1, 2, \dots, 9\}$ , representing  $|V| = 9$  IP addresses.<sup>1</sup> On the left is an Euler diagram showing each of eight hyperedges  $A, B, \dots, H$ , representing domains, as a “lasso” around its vertices. On the right is a  $V \times \mathcal{E}$  incidence matrix  $I$ , where a non-null  $\langle v, e \rangle \in I$  cell indicates that  $v \in e$  for some  $v \in V, e \in \mathcal{E}$ .



**Fig. 1.** (Left) An Euler diagram of an example hypergraph  $\mathcal{H}$ . (Right) Its incidence matrix  $I$ .

We call each hyperedge  $e \in \mathcal{E}$  an  $s$ -edge where  $s = |e|$ . Thus all graphs are hypergraphs, in that all graph edges are 2-edges, for example  $H = \{4, 5\}$ , saying that the domain  $H$  has two IPs 4 and 5. But  $F = \{1, 2, 3, 9\}$  is a 4-edge, with domain  $F$  having those four IPs. Where each column of the incidence matrix of a graph has exactly two cells, those of hypergraphs are unrestricted.

Our research group is pursuing hypergraph analytics as an analog to graph analytics [14]. While our development is consistent with others in the literature [9, 16], our notation and concepts are somewhat distinct. For a more comprehensive development see [2].

We say that two edges  $e, f \in \mathcal{E}$  are  $s$ -adjacent if  $|e \cap f| \geq s$  for  $s \geq 1$ . An  $s$ -star is a set of edges  $\mathcal{F} \subseteq \mathcal{E}$  sharing exactly a common intersection  $f \subseteq V$ , with  $|f| \geq s$ , so that  $\forall e_i, e_j \in \mathcal{F}$  we have  $e_i \cap e_j = f$ . An  $s$ -path is a sequence of edges  $\langle e_0, e_1, \dots, e_n \rangle$  such that each  $e_{i-1}, e_i$  are  $s$ -adjacent for  $1 \leq i \leq n$ ; and an  $s$ -component is a maximal collection of edges any pair of which is connected by an  $s$ -path. The  $s$ -diameter of an  $s$ -component is the length of its longest shortest  $s$ -path. Comparing again to graphs, graph paths are all 1-paths, and graph

<sup>1</sup>  $\mathcal{H}$  can also be represented as a bipartite graph on the disjoint union  $V \sqcup \mathcal{E}$ , with each component a distinct part.

components all 1-components. Our example has two 1-components (shown obviously), but also four 2-components (listed edge-wise)  $\{A, F, G, H\}$ ,  $\{B, D\}$ ,  $\{C\}$  and  $\{E\}$ . Its 3- and 4-components are each single edges of size larger than 3 or 4 (respectively), and it has no 5 or higher components.

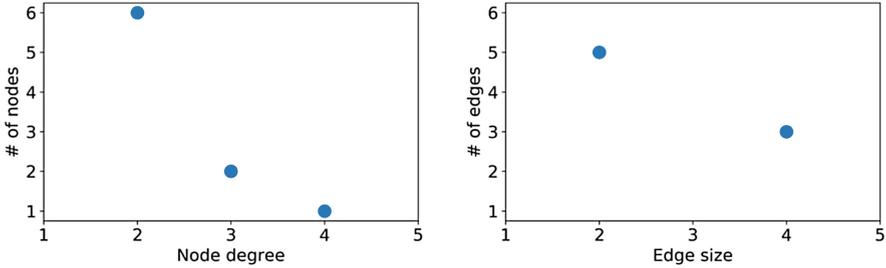
Given a hypergraph  $\mathcal{H}$ , it is possible to construct smaller representations which capture important properties:

- Note that in our example, the edges  $A = F$  and  $B = D$ , and the vertices  $1 = 9$  and  $7 = 8$ , are equivalent, represented as duplicate columns and rows in  $I$  respectively. **Collapsing** is the process of combining these and replacing them with a representative, while also possibly maintaining a multiplicity count to be used for a weighting. The edges  $\mathcal{E}$  are hereby transformed from a multiset to a set.
- Additionally, note that after collapsing, the smaller 1-component becomes an **isolated singleton**, effectively a collection of non-interacting vertices, or a diagonal block in  $I$ . These are especially common in DNS data. Pre-collapse, an isolated singleton would indicate the normal, *uninteresting* behavior in DNS where a single IP is associated with a single domain, and *vice versa*. But post-collapse, they indicate a collection of IPs and domains which are universally associated only with themselves, effectively forming a set of domain and IP aliases. In this work, these are counted and pruned, but in the future they could be attended to with respect to their multiplicities.
- Finally, note that  $H \subset G$  is an **included** edge. Non-included edges are called **toplexes**, and not only is the collection of toplexes much smaller than  $\mathcal{H}$ , but it is sufficient to derive some hypergraph information, for example  $s$ -components.

Table 1 shows some important statistics for our example, first for the initial hypergraph, then after collapsing, and finally after removing isolated singletons from the collapsed hypergraph. For hypergraph data, a vastly high or low aspect ratio can indicate difficulty in analysis. Note that as reductions commence, the number of vertices, edges, and cells reduces, while density increases. Finally, Fig. 2 shows the distribution of node degree (# edges per node) and edge size.

**Table 1.** Basic hypergraph statistics for our example.

	Initial	Collapsed	Non-singleton components
$ V $	9	7	6
$ E $	8	6	5
Aspect ratio	1.125	1.167	1.200
# Cells	23	14	13
Density	0.319	0.333	0.433



**Fig. 2.** (Left) Distribution of node degree (# edges per node) in our example. (Right) Distribution of edge size  $s$ .

In our pipeline the segmentation steps of collapsing, removing isolated singletons, and computing  $s$ -components are all performed using CHGL, as are node degree, edge size, and  $s$ -component size distributions. Subsequent exploration of the structures found within the components themselves, e.g., identification of stars and computation of diameters, are done via HNX. HNX builds on the popular library NetworkX [10], which offers metrics and algorithms for the analysis of graph data. Euler diagram visualizations that appear in this paper are provided directly by the HNX package.

### 3 Hypergraph Representations of DNS Data

The Domain Name System (DNS) provides a decentralized service to translate from the domain names that humans keep track of (e.g., [www.google.com](http://www.google.com)) to IP addresses that computers require to communicate. Perhaps somewhat counter-intuitively, DNS data present naturally as a hypergraph, in being a many-many relationship between domains and IPs. While typically this relationship is one-to-one, with each domain uniquely identifying a single IP address and *vice versa*, there are a number of circumstances which can violate this:

- Some domains have aliases so that multiple domains (e.g., misspellings) resolve to the same IP address.
- There are large hosting services where one IP serves multiple websites.
- Some domains are used so frequently that they must be duplicated across hosts and therefore map to multiple IPs.
- IP addresses are randomly reassigned within some small IP block so the same domain may map to multiple IP addresses when queried over time.

In order to explore large volumes of DNS mappings we used ActiveDNS (ADNS), a data set maintained by the Astrolavos Lab at Georgia Institute of Technology (<https://activednsproject.org>). ADNS submits daily DNS lookups for popular zones (e.g., .com, .net, .org) and lists of domain names. The data is stored in Avro format (<https://avro.apache.org>) which provides structured

records for each DNS lookup in a compressed binary file. Each record contains information including: query date, lookup input (often a domain name), data returned by a DNS server (often a list of IP addresses), and IP addresses of the DNS servers that answered the query. DNS records are also typed according to different properties (recorded as the `qtype` field in ADNS) such as the format of the data and to indicate its intended use. This initial analysis accepted any reasonable pairing of domain name IP address and did not restrict to any particular `qtypes`. Future work will restrict to `qtype = 1`, which map hostnames to an IPv4 address of the host.

Our group acquired data from the time period April 24–May 29, 2018, and in this paper we focus on the single day of April 26, 2018. This day consists of 1,200 Avro files with each file containing on average 900 K records. There was some data cleaning necessary to remove records with empty lookup input or empty returned data. Additionally we removed any records in which the lookup input was an IP address or the returned data was a domain name. After cleaning, each file was reduced to approximately 180 K records.

We structured these DNS data as a hypergraph on a vertex set  $V$  of IPs and edge set  $\mathcal{E}$  of domains. Thus our hypergraphs  $\mathcal{H}$  coded each domain as a collection of its IPs. We show results of our analysis below in Sect. 5, including global statistics and the results of targeted exploration.

## 4 Chapel Hypergraph Library (CHGL)

The Chapel HyperGraph Library (CHGL) [12] is a prototype exascale library written in Chapel [6, 7] that brings generation, representation, and computation of hypergraphs to the world of high performance computing (HPC). Thanks to Chapel, CHGL provides scalability in both shared memory and distributed memory contexts.

In most cases, data underlying a hypergraph is more complex than CHGL’s internal representation of vertices and hyperedges as consecutive integers. In such situations, a hash table that maps user-defined generic properties to the consecutive identifiers of vertices and hyperedges is used for translation. The properties are embedded in the internal representation of the hyperedges and vertices, allowing  $O(1)$  bidirectional lookup as well as locality when iterating over the graph, shared-memory and distributed alike.

CHGL performs *segmentation*, or reduction, of the data in multiple highly-parallel phases. Segmentation reduces both the size of the graph to one that HNX can process in a reasonable amount of time and the computational workload on CHGL when computing metrics. Proper care is taken to ensure that references to the collapsed hyperedges and vertices are taken forward to the hyperedge or vertex that they collapsed into, and that all references to removed hyperedges and vertices are removed. This is performed in linear time and applies to both the graph and property map.

To prune away redundant entities, which is generally useful for computation, hyperedges and nodes are placed into equivalence classes through the process of

collapsing described in Sect. 2. All but one arbitrarily chosen representative is removed from the graph. Determining the equivalence class of a vertex or hyper-edge can be done by using a set or hash table, and can be performed in  $O(|V|)$  or  $O(|V| \log |V|)$  time, depending on the data structure used. In practice, the time complexity is often linear or quasilinear, but in the worst-case scenario when the hypergraph is fully connected, the time complexity is  $O(|V|^2)$  or  $O(|V|^2 \log |V|)$ .

Isolated singletons, as described in Sect. 2, tend to be uninteresting. After collapsing, these are pruned away in a straightforward manner.

We implemented computation of  $s$ -components using a parallel search method, where we iterate over edges in parallel, and every edge begins an independent search. The  $s$ -neighbors of an edge are marked with the component number originating from the initial edge. The component number is taken from a global atomic counter at the beginning of every parallel search. 1-Components are implemented by simply traversing the edges by following included vertices (edge to vertex to neighbor edge), but 2-components and higher require an implementation with set intersections to check the cardinalities of adjacencies. This implementation is well suited for a large number of small components because most components end up being searched by a single task. The best case scenario complexity of the parallel search algorithm is linear, and the worst is quadratic if the maximum number of component collisions occur. The average complexity in our case is close to linear since the DNS data has a large number of small components, and most components are handled by a single task.

Obtaining the vertex degree and edge cardinality distributions is simple and intuitive in CHGL, thanks to Chapel’s high-level abstractions. This particular operation is short enough that it can be presented in full in Fig. 3. We compute these both pre- and post-collapsing.

```

1 // Find largest degree of all vertices in the hypergraph
2 var N = max reduce [v in graph.getVertices()] graph.degree(v);
3 var degreeDist : [1..N] int;
4 forall v in graph.getVertices() with (+ reduce degreeDist) {
5   degreeDist[graph.degree(v)] += 1;
6 }

```

**Fig. 3.** Obtaining the vertex degree distribution in CHGL.

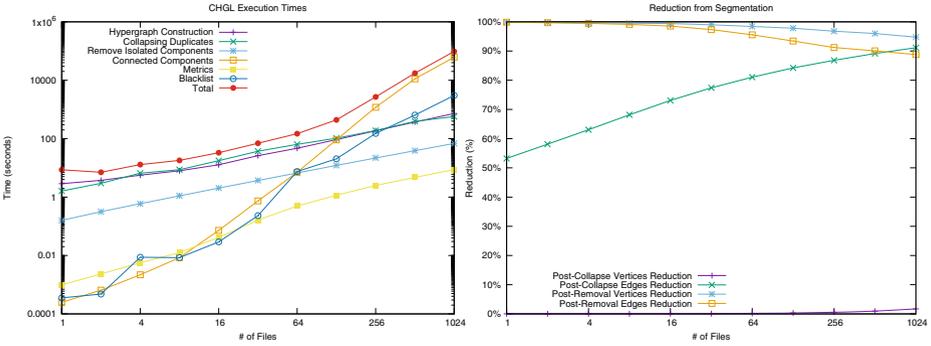
Finally, the  $s$ -component size distributions were computed, recording the number of nodes and edges in each  $s$ -component and how many  $s$ -components have each size. This allows us to understand how nodes and edges are distributed, e.g., is there one giant component and a few small components or are component sizes more uniformly distributed.

## 5 Computational Results

We ran experiments on one of the compute nodes of an Infiniband cluster, each equipped with a 20-core Intel Xeon processor and 132 GB memory. All cores

were involved in the experiments. CHGL v0.1.3 was compiled against Chapel pre-release version 1.18.0 with `--fast` flag to enable all compiler optimizations.

Execution times of the stages of the CHGL DNS processing pipeline are shown on the left side of Fig. 4. *s*-component computation dominates the execution time for 128 or more files. The *s*-components are reused when computing the *s*-component size distributions, leading to them taking significantly less time. Collapsing duplicates and removing isolated components scale linearly, as is expected for their time complexity. The hypergraph is constructed in about the same amount of time it takes to collapse it, showing that processing DNS data is mostly compute-bound.



**Fig. 4.** (Left) Execution times (log-log scale). (Right) Effectiveness of reduction from segmentation.

The purpose of segmentation is to reduce the size of the graph while also maintaining the data that is of interest. The right side of Fig. 4 shows the compression as a result of performing segmentation. Collapsing of duplicate edges results in the most compression, reducing the graph from 55% at one file to over 90% at 1024 files, which can be expected to improve further when more data is processed. Removing isolated components results in less compression as data size increases, likely due to the premature marking of components as isolated prior to having all of the data. Perhaps with larger amounts of data, there will be a convergence to a stable number of isolated components in the entirety of the DNS network. Note that there are very few duplicate IP addresses on smaller samples, but that may change as more data is processed; nonetheless, collapsing duplicate vertices may be unnecessary and can possibly save some time.

Above we reported on scaling of loading and compute time using CHGL on varying numbers of ActiveDNS files, from 1 to 1,024. Here we report on analysis of the hypergraph built from one full day, April 26, 2018, comprising 1,200 files. See Table 2 for basic count statistics.

The node degree and edge size distributions are shown in Figs. 6a and 6c. Except for the small increase around  $x = 10^2$  the node degree distribution looks

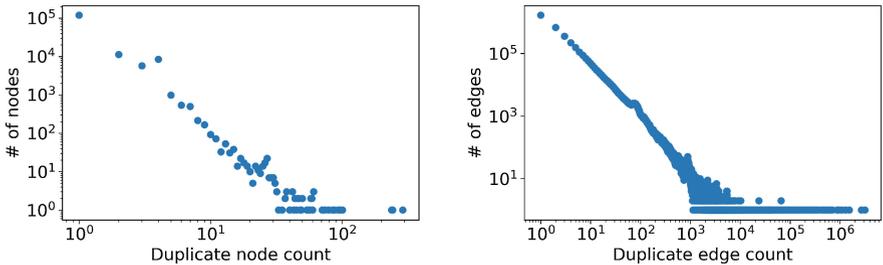
**Table 2.** Basic hypergraph statistics for ActiveDNS data for April 26, 2018.

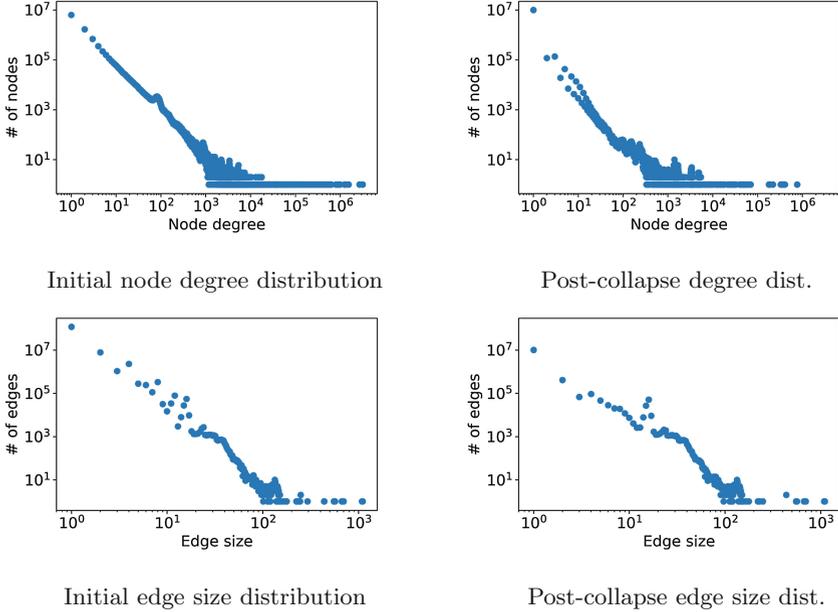
	Initial	Collapsed	Non-singleton components
$ V $	10.6M	10.3M	557K
$ E $	131.2M	11.0M	1.2M
Aspect ratio	0.081	0.941	0.460
# cells	157.4M	25.7M	15.9M
Density	1.14 E-7	2.26 E-7	2.35 E-5

like a power law or heavy tailed distribution typical in real-world graphs [3]. The degree distribution has a general decreasing tendency from  $x = 1$  to  $x = 70$ , it increases by roughly 1,000 through  $x = 80$ , and then returns to the downward trend. We do not know why this occurs, but it is possible that it could be an artifact of DNS server configuration practices. Edge size distribution also seems to be heavy-tailed although somewhat more noisy for low edge sizes than the degree distribution.

See the second column in Table 2 for the simple count statistics of the collapsed hypergraph. Notice that collapsing resulted in a much more square incidence matrix since only 2% of nodes were collapsed while 92% of edges were collapsed. The number of cells in the collapsed hypergraph incidence matrix is now reduced to 16% of the full hypergraph.

The distributions of node and edge duplicate counts are shown in Fig. 5. Notice that the distribution of duplicate edge counts has a similar shape as the node degree distribution of the original hypergraph with a slight increase around  $x = 10^2$ . After seeing this it is possible that the nodes which had degree around 70–80, where this increase occurs, were actually in many duplicate edges which are now collapsed. The node degree distribution for the collapsed hypergraph found in Fig. 6b further supports this hypothesis since the increase around  $x = 10^2$  in the node degree distribution is absent.

**Fig. 5.** Distribution of duplicate node counts (top) and edge counts (bottom).



**Fig. 6.** Node degree and edge size distributions, on a log-log scale, for April 26, 2018 DNS hypergraph. The  $x$  and  $y$  axes are the same across both node plots and across both edge plots to illustrate the changes through the collapsing procedure.

The edge size distribution post collapse is shown in Fig. 6d. This distribution is very similar to that of the original hypergraph, although it appears less noisy up through approximately  $x = 20$ . This is not surprising since there were not many duplicate nodes removed, so edges that remained likely stayed close to their original size.

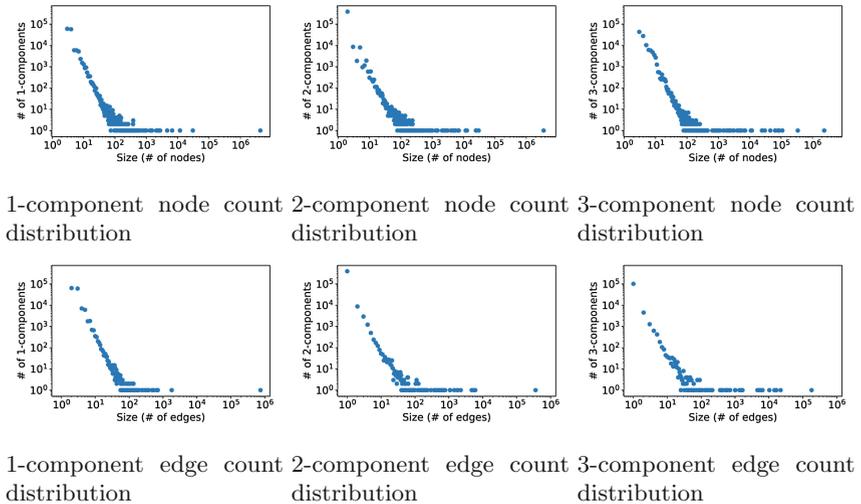
After collapsing duplicate nodes and edges we removed all 9,784,763 isolated singleton edges, or 89% of all remaining edges. The only differences between the collapsed hypergraph and the hypergraph after removal of isolated singleton components is the number of degree 1 nodes and the number of size 1 edges. Therefore, we omit the final node degree and edge size distributions since they are identical to the post-collapse distributions except for the points at  $x = 1$ .

Comparing the pre-collapse (left), post-collapse (right), and post-removal distributions (not pictured) in Fig. 6, we observe that hypergraph collapsing and removal significantly alters the shape of degree and edge size distributions. In addition to the qualitative differences apparent from the plots, these differences can also be quantified using the Kolmogorov-Smirnov (KS) distance metric, a normalized statistic between 0 and 1 in which larger values indicate greater degree distribution dissimilarity. In the case of the degree distributions (top row), KS distance suggests the pre-collapsing hypergraph differs significantly from the post-collapse and post-removal degree distributions, with KS values of 0.36 and

0.34, respectively. In the case of the edge-size distributions (bottom row), the most pronounced difference is between the pre-collapsing and post-removal edge size distribution, with a KS value of 0.60. Here, the large KS distance reflects the dramatic changes at the head of the distribution, where the number of 1-edges decreases from 118 million to 369 thousand.

## 6 Analytical Results

The next step toward finding interesting subgraphs within the single day of ActiveDNS data was to compute and explore  $s$ -components. CHGL computed  $s$ -components of the hypergraph post-collapse and post-removal of isolated singletons for  $s = 1, 2, 3$ . Before exploring these components themselves we report the distribution of component sizes (both node and edge counts) which are found in Fig. 7. As  $s$  increases the shapes of these distributions do not change much but they do tend to skew more toward smaller components and the distribution flattens slightly. This is required since every  $s$ -component is contained within some  $s'$ -component for  $s' < s$ : as  $s$  increases components can only decrease in size. These distributions also show that while there are some very large  $s$ -components the majority are very small. Additionally, we see that the notion of a “giant component” is much more prevalent in the set of 1-components than for  $s = 2$  or 3. Indeed, as  $s$  increases the largest component breaks up and the jump between the largest component and second largest becomes smaller.

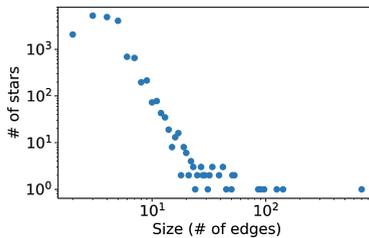


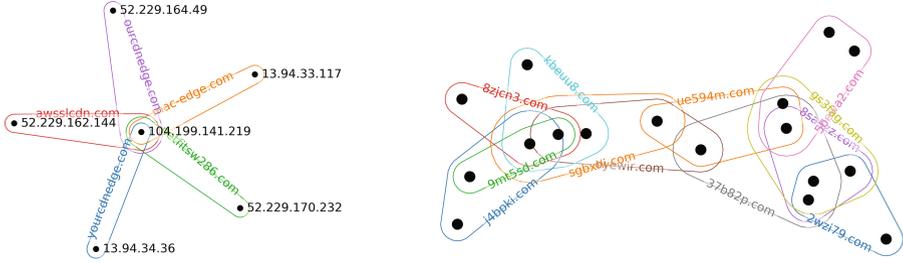
**Fig. 7.** Node and edge count distributions, on a log-log scale, for  $s$ -components within simplified April 26, 2018 DNS hypergraph. The  $x$  and  $y$  axes are the same across all three node count plots and across three edge count plots to illustrate the changes as  $s$  increases.

Once the hypergraphs were segmented into  $s$ -components by CHGL we proceeded to do exploratory analysis using HNX. In particular, we looked for:

- Occurrences of 1-stars within the 1-components, and
- $s$ -components with maximum  $s$ -diameter for  $s = 2, 3$ .

Recall that a 1-star is a small hypergraph in which all edges pairwise intersect in one node, and that one node is the same across all pairwise intersections. The simplest 1-star has all edges of size 2, see Fig. 9a for an example of this case. In our DNS use case a star is a collection of domains which all share exactly one IP address but each also have their own separate IP address(es). These are consistent with the behavior of content delivery networks (CDN), geographically distributed networks of servers with the goal of quickly and reliably serving up content to a variety of users, which could explain the existence of stars with a diverse set of IP addresses since a consideration for IP assignment is geographic location. Star motifs are also consistent with DNS sinkholes and domain hosting services.





**Fig. 9.** (a) A small star seen in the ActiveDNS data. (b) The 2-component with largest 2-diameter.

from Google Cloud whereas the leaves are from Microsoft Corporation. All five domains are registered through the hosting site GoDaddy.com.

To discover interesting 2-components (resp. 3-components) we calculated 2-diameters (resp. 3-diameters) of each of the components and look more closely at those with maximal diameter. In the case of the 2-components the maximum 2-diameter is 6 and there is only one 2-component with that 2-diameter, shown in Fig. 9b. The IP addresses in this component all belong to the IP range 103.86.122.0/24 and the domains are registered to GMO INTERNET, INC according to WHOIS records. Moreover, current DNS queries for most of these domains at a later date resolve to IPs in the range 103.86.123.0/24 and have a time to live of only 120 s. This pattern of quickly changing of IP address is consistent with the *fast flux* DNS technique which can be used by botnets to hide malicious content delivery sites and make networks of malware more difficult to discover [11].

The large diameter 3-components tell different stories. The maximum 3-diameter is 3 and there are four 3-components with this 3-diameter. One has only one topex with six sub-edges. Two others are fairly simple and, like the large 2-diameter 2-component, are somewhat chain-like tracing out a long path. The fourth is quite large with 70 nodes, 189 edges, and all IPs belonging to an IP range from Amazon Technologies Inc.

## 7 Conclusions and Future Work

While our research group has been developing hypergraph methods and mathematics over a moderate period, this paper reflects the first application of CHGL to cyber data.

The current approach is limited in a number of ways. First, ActiveDNS records data from DNS lookups on a daily basis (or perhaps multiple times per day), but it does not do continual monitoring. This discrete sampling may mean that the pipeline misses patterns that would normally be seen in a more continuous approach. Additionally, the current analysis is for a single day, and extending to multiple days in the current architecture will exacerbate issues with

memory bounds. This might be mitigated using a theory of dynamic hypergraphs (much like that of dynamic graphs) to understand the time-evolution of DNS or similar data.

Additionally, certain DNS relationships are ignored, such as recursive DNS records where one domain resolves not to an IP address but to another domain name. This would require more complicated mathematics than just hypergraphs, likely cell complexes or partial orders, which we have started to consider in our research but not yet in our analysis. We also ignore other pieces of metadata like the authority IP addresses (those servers which answered the DNS request).

Additional future work includes:

- We are extending our prior theoretical work [2,14] to a full consideration of the mathematical foundations of hypergraphs for data science, including spectral approaches and consideration of multiplicity weightings.
- A range of hypernetwork methods generalizing network science centrality, connectivity, clustering coefficients, etc. are available [2].
- Also central to our approach is the consideration of hypergraphs as multidimensional objects, and thus inherently available for topological applications, including homology measurement for identification of loops and potential gaps in the underlying data.
- CHGL is also under active development to include topology, homology measures, a proper graph library, and a distributed data model.
- Finally, application and data analysis continues, including DNS, additional cyber data beyond DNS, and additional application domains including computational biology and social hypernetworks.

**Acknowledgements.** This work was partially funded by a US Department of Energy Computational Science Graduate Fellowship (grant DE-SC0020347).

This work was also partially funded under the High Performance Data Analytics (HPDA) program at the Department of Energy’s Pacific Northwest National Laboratory. Pacific Northwest National Laboratory is operated by Battelle Memorial Institute under Contract DE-ACO6-76RL01830.

Special thanks to William Nickless for helpful conversations surrounding the DNS analysis and interpretation.

## References

1. Active DNS project. <https://activednsproject.org/>. Accessed 26 Nov 2019
2. Aksoy, S.G., Joslyn, C., Marrero, C.O., Praggastis, B., Purvine, E.: Hypernetwork science via high-order hypergraph walks. arXiv preprint [arXiv:1906.11295](https://arxiv.org/abs/1906.11295) (2019, Submitted)
3. Barabási, A.-L., Bonabeau, E.: Scale-free networks. *Sci. Am.* **288**(5), 60–69 (2003)
4. Berge, C., Minieka, E.: *Graphs and Hypergraphs*. North-Holland, Amsterdam (1973)
5. Guy Bruneau. DNS Sinkhole. <https://www.sans.org/reading-room/whitepapers/dns/dns-sinkhole-33523>

6. Chamberlain, B.L., Callahan, D., Zima, H.P.: Parallel programmability and the chapel language. *Int. J. High Perform. Comput. Appl.* **21**(3), 291–312 (2007)
7. Chamberlain, B.L., et al.: Chapel comes of age: Making scalable programming productive. Cray Users Group (2018)
8. Devine, K.D., Boman, E.G., Heaphy, R.T., Bisseling, R.H., Catalyurek, U.V.: Parallel hypergraph partitioning for scientific computing. In: *Proceedings 20th IEEE International Parallel & Distributed Processing Symposium*. IEEE (2006)
9. Estrada, E., Rodríguez-Velázquez, J.A.: Subgraph centrality and clustering in complex hyper-networks. *Phys. A* **364**, 581–594 (2006)
10. Hagberg, A.A., Schult, D.A., Swart, P.J.: Exploring network structure, dynamics, and function using networkx. In: Varoquaux, G., Vaught, T., Millman, J. (eds.) *Proceedings of the 7th Python in Science Conference, Pasadena, CA USA*, pp. 11–15 (2008)
11. Riden, J.: How fast-flux service networks work. <http://www.honeynet.org/node/132>. Accessed 26 Nov 2018
12. Jenkins, L.P., et al.: Chapel hypergraph library (CHGL). In: *2018 IEEE High Performance Extreme Computing Conference (HPEC 2018)* (2018)
13. Karypis, G., Kumar, V.: Multilevel k-way hypergraph partitioning. *VLSI Des.* **11**(3), 285–300 (2000)
14. Purvine, E., Aksoy, S., Joslyn, C., Nowak, K., Praggastis, B., Robinson, M.: A topological approach to representational data models. In: Yamamoto, S., Mori, H. (eds.) *HIMI 2018*. LNCS, vol. 10904, pp. 90–109. Springer, Cham (2018). [https://doi.org/10.1007/978-3-319-92043-6\\_8](https://doi.org/10.1007/978-3-319-92043-6_8)
15. Robins, G., Alexander, M.: Small worlds among interlocking directors: network structure and distance in bipartite graphs. *Comput. Math. Organ. Theory* **10**(1), 69–94 (2004)
16. Wang, J., Lee, T.T.: Paths and cycles of hypergraphs. *Sci. China, Ser. A Math.* **42**(1), 1–12 (1999)