

# PolarStar: Expanding the Scalability Horizon of Diameter-3 Networks

Kartik Lakhotia  
Intel Labs  
USA  
kartik.lakhotia@intel.com

Laura Monroe  
Los Alamos National Laboratory  
USA  
lmonroe@lanl.gov

Kelly Isham  
Colgate University  
USA  
kisham@colgate.edu

Maciej Besta  
ETH Zürich  
Switzerland  
maciej.best@inf.ethz.ch

Nils Blach  
ETH Zürich  
Switzerland  
nils.blach@inf.ethz.ch

Torsten Hoefler  
ETH Zürich  
Switzerland  
htor@inf.ethz.ch

Fabrizio Petrini  
Intel Labs  
USA  
fabrizio.petrini@intel.com

## ABSTRACT

In this paper, we present PolarStar, a novel family of diameter-3 network topologies derived from the star product of two low-diameter factor graphs. The proposed PolarStar construction gives the largest known diameter-3 network topologies for almost all radixes. When compared to state-of-the-art diameter-3 networks, PolarStar achieves 31% geometric mean increase in scale over Bundlefly, 91% over Dragonfly, and 690% over 3-D HyperX.

PolarStar has many other desirable properties including a modular layout, large bisection, high resilience to link failures and a large number of feasible sizes for every radix. Our evaluation shows that it exhibits comparable or better performance than other diameter-3 networks under various traffic patterns.

## 1 INTRODUCTION

### 1.1 Motivation

The growth of datacenters and supercomputers is driving the need for extremely large-scale systems, requiring tens or hundreds of thousands of processing nodes. For example, Frontier, the most powerful system in the Top 500 list [35], has 9,408 CPUs and 37,632 GPUs [9], and the second-ranked supercomputer, Fugaku, has 158,976 processing nodes [14]. These systems were once the exclusive capability of scientific research and national labs, but are now the norm in commercial and industrial data centers, fueled by social media and large-scale simulation of industrial, financial and entertainment processes.

High performance networks are the backbone of these systems, and a switch may be considered to be a network building block. Given the number of links on a switch, the question is: what is the largest system that can be built having the smallest diameter? The size of the system is clearly important, as that determines the peak

compute performance, and the diameter is also important, as it affects communication latency and injection bandwidth per switch.

Low-diameter networks, i.e. diameter-2 [22] and diameter-3 [24] topologies, are of great interest, providing low-latency and cost-effective high-bandwidth communication infrastructure. Each packet consumes bandwidth on only a few links, limiting negative effects of tail latency and improving overall system performance [13].

The emergence of high-radix optical IO modules with high shoreline density has increased interest in *scalable low-diameter networks* [12, 23, 31, 38]. Co-packaging of such modules with compute nodes on the same chip greatly enhances the bandwidth available per node. Low-diameter networks are then required to efficiently utilize bandwidth on co-packaged chips. Since each router is integrated with a compute node, scalability of a co-packaged system is identified with the order of the graph defining the network topology.

### 1.2 Current Approaches and Limitations

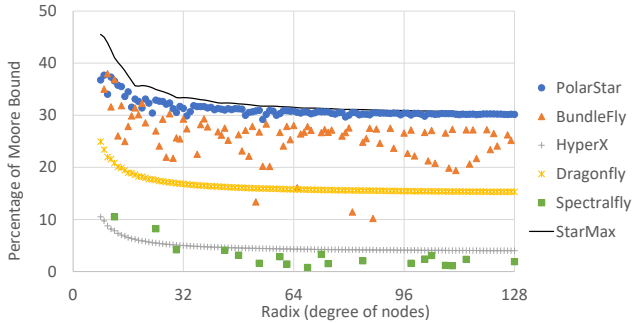
The formal problem of designing a network  $G$  with the largest order (number of nodes) for a given diameter  $D$  and router degree  $d$  is captured by the degree-diameter problem from [32]. The order of  $G$  is bounded above by the Moore bound [19]. Networks with good Moore-bound efficiency (proximity to Moore bound) are not only highly scalable, but also cost-effective and power-efficient as they can realize a system of given size with relatively lower radix switches and fewer cables. Unfortunately, the largest known graphs for  $D > 2$  and  $d > 2$  are much smaller in size than the Moore bound.

Some diameter-2 networks such as Polarfly [25] and Slimfly [7] do come close to the Moore bound. However, the scale of these networks is limited by the small diameter-2 Moore bound ( $d^2 + 1$  for radix  $d$ ). Thus, diameter-2 networks are not suitable for massive-scale datacenters and HPC systems with current or foreseeable technology, because they can only span a few thousands of nodes.

On the other hand, diameter-3 networks have a high enough Moore bound to address scalability requirements ( $d^3 - d^2 + d + 1$  for radix  $d$ ). Dragonfly [24] and HyperX [2] are popular diameter-3 topologies deployed in large systems [9, 26, 27]. However, these particular topologies exhibit poor Moore-bound efficiency, as shown in

The submitted manuscript has been authored by employee(s) of Triad National Security, LLC, under contract with the U.S. Department of Energy. Accordingly, the U.S. Government retains an irrevocable, nonexclusive, royalty-free license to publish, translate, reproduce, use, or dispose of the published form of the work and to authorize others to do the same for U.S. Government purposes. This paper has been assigned the LANL identification number LA-UR-23-20414.

Figure1, which drives up the network cost and power consumption.



**Figure 1: Scalability of direct diameter-3 topologies with respect to the Moore bound. For Bundlefly and PolarStar, the largest construction for each radix is shown. StarMax denotes an upper bound on the largest graphs theoretically achievable with star product – the mathematical construct used in state-of-the-art Bundlefly and PolarStar networks. Asymptotic Moore-bound efficiency is 16% for Dragonfly and less than 5% for 3-D HyperX. For Spectralfly, which is not a fixed diameter topology, we only compare design points with diameter  $\leq 3$  and largest scale for a given radix (if it exists).**

Recently, Lei et al. introduced Bundlefly [29], a diameter-3 network based on a star product of two graphs. Bundlefly has a modular design with support for bundling inter-module links into multi-core fibres. This reduces cabling complexity and cost. However, Bundlefly does not exist for several radices in the range [8, 128], and its Moore-bound efficiency varies significantly, as shown in Figure 1.

The mathematical problem of establishing the largest diameter-3 graphs is an open problem. There is a big gap between the best-known diameter-3 graphs and the Moore bound for diameter 3 [32]. In this paper, we extend the orders of the largest known diameter-3 graphs and design a network based on these.

### 1.3 Contributions

We propose a new family of network topologies called PolarStar that extends PolarFly [25] to large diameter-3 networks using a mathematical construct called the star product.

Degree	Best known Order in [32]	Moore-bound Efficiency	PolarStar Order	Moore-bound Efficiency
18	1,620	29.3%	1,830	33.3%
19	1,638	25.1%	2,128	32.6%
20	1,958	25.7%	2,394	31.4%

**Table 1: For radices 18-20, PolarStar surpasses the previously largest known diameter-3 graphs listed on Combinatorics Wiki degree-diameter table [32].**

- A novel graph construction is presented here for the first time, giving the *largest known graphs of diameter-3* for degrees 18–20 as per the Combinatorics Wiki leaderboard [32] (Table 1), superseding former records for the first time since 2010 on this open problem in an active field of mathematics. (The Wiki shows graphs only up to degree 20, so this construction likely gives "best" graphs at higher degrees as well).

- PolarStar is derived from these graphs, providing the *largest known direct networks* of diameter-3 for almost all radices. It achieves 31%, 91% and 672% geometric mean increase in scale over Bundlefly, Dragonfly and HyperX, respectively.
- We show that PolarStar reaches near-optimal scalability for diameter-3 star product graphs and further optimizations on star product are unlikely to provide notable benefits.
- Two alternative star product constructions supporting different networking requirements are discussed in detail. We also list other constructions with their respective desirable properties.
- PolarStar extends several networking benefits of PolarFly including a modular layout amenable to bundling of links into multi-core fibers, and a large bisection cut.
- PolarStar has a large design space: it exists for every radix in [8, 128] and has multiple configurations for each radix.

## 2 BACKGROUND

### 2.1 Network Model

An interconnection network can be modeled as an undirected graph  $G(V, E)$ :  $V(G)$  is the set of switching nodes, or vertices,  $|V(G)|$  is the order of  $G$ , and  $E$  is the set of links, or edges. In co-packaged networks, each node functions as both a router and a compute endpoint. Each node has  $d$  links to other nodes where  $d$  is the *network radix*, or *degree*. The maximum length of shortest paths between any node pair is the *diameter*  $D$ . In this paper, we consider networks of diameter 3.

### 2.2 Moore Bound and Low-Diameter Networks

The Moore bound [19] is an upper bound on the number of nodes  $N$  that a network with degree  $d$  and diameter  $D$  may have. This bound is given by

$$N \leq 1 + d \cdot \sum_{i=0}^{D-1} (d-1)^i$$

For diameter-3 networks, the Moore bound is  $N \leq d^3 - d^2 + d + 1$ .

The only graphs with  $D \geq 2$  and  $d \geq 2$  that actually meet the Moore bound are the cycles, the Hoffman-Singleton graph, the Petersen graph [5, 11, 19] and a hypothetical diameter-2 degree-57 graph [10]. These graphs are not suited for large-scale network design: the degree-2 cycles with low diameter are very small, and the others have only one design point each.

Few graphs even come close to the Moore bound. The latest leaderboard of degree-diameter problem from [32] shows that the best known diameter-3 graphs, with the most relevant degrees, reach only 25–30% of the Moore bound. The PolarStar construction proposed here is larger than the best graphs for degrees 18–20, the highest degrees in the leaderboard, as shown in Table 1.

### 2.3 Network Properties

Throughout the paper, we discuss desirable properties of network topologies and evaluate various networks on their basis. We analyze Slimfly [7], PolarFly [25], Dragonfly [24], HyperX [2], MegaFly [16, 37], SpectralFly [40] and widely used three stage Fat-trees. A summary of the evaluation is given in Table 2. Only *PolarStar* fully

supports all desired properties, while simultaneously achieving highest Moore-bound efficiency for diameter-3 networks.

Topology	Direct	Scalability	Stable Design-space	$D \leq 3$	Bundlability
Fat tree	✘	✘	▣	✘	▣
PolarFly	▣	✘	▣	▣	▣
Slimfly	▣	✘	▣	▣	▣
HyperX	▣	▣	▣	▣	▣
Dragonfly	▣	▣	▣	▣	✘
Bundlefly	▣	▣	▣	▣	▣
MegaFly	✘	▣	▣	▣	✘
Spectralfly	▣	▣	▣	▣	✘
<b>PolarStar</b>	▣	▣	▣	▣	▣

Table 2: Feasibility. “▣”: full support, “▣”: partial support, “✘”: no support,  $D$ : Network Diameter

**Directness:** Every switch in a direct network is attached to one or more endpoints. In contrast, indirect networks also have some switching nodes that are not attached to any endpoint. If co-packaged modules are used, indirect networks such as Fat tree and MegaFly require fabricating two types of chips, which increases their cost. Further, a switch-only chip in these topologies requires twice the number of ports than a co-packaged chip with an endpoint.

**Scalability:** A network’s scale depends on its diameter and Moore-bound efficiency (proximity to Moore bound for the given degree and diameter). Diameter-2 networks such as PolarFly [25] asymptotically approach the Moore bound but are limited in scale as the bound itself is small. Three-stage Fat-trees scale similarly to diameter-2 networks. Diameter-3 networks can scale to hundreds of thousands of nodes with currently available switches. However, HyperX and diameter-3 SpectralFly have poor Moore-bound efficiency, resulting in higher cost, compared to a more efficient topology such as the PolarStar presented in this paper.

**Stable Design-Space:** A desirable topology will provide feasible configurations for all radixes, and should scale smoothly with the radix. This allows network construction with a wide range of switches. However, the diameter-2 Slimfly [7] topology has few feasible radixes and Bundlefly’s [29] Moore bound efficiency fluctuates significantly with the radix. Although SpectralFly may be constructed for any radix  $p + 1$  with  $p$  an odd prime, SpectralFly of diameter-3 exists for very few radixes, as shown in Figure 1.

**Low-diameter:** Low-latency remote accesses are a necessity for performance scalability in Global Address Space (GAS) programing models. Networks with small diameter provide the desired communication latency and can sustain high ingestion bandwidth per switch. Diameter = 3 is preferred as it can provide both performance and scalability at low cost.

**Bundlability:** It is the property of a network that allows bundling of multiple (global) links into fewer multi-core fibers (MCFs, [3, 29]). Bundlable networks such as PolarFly, Bundlefly and PolarStar have multiple links between adjacent modules (logical groups of nodes) that can be packed together in an MCF. This significantly reduces cabling cost and complexity in large networks. However, Dragonfly and MegaFly are not amenable to bundling because each pair

of node groups in these topologies is connected by a single link. While multiple links can be used to connect a pair of groups, it significantly reduces the scalability of these topologies.

### 3 APPROACH

We use a graph theoretical formulation called *star product* to construct the scalable PolarStar topology (see Section 4). The inputs to star product are two *factor graphs*, and the output is a larger graph whose order is the product of the order of factors. By using the Erdős-Rényi polarity graphs with a novel construction as the factor graphs, we design a new family of diameter-3 networks *larger than any previously designed*, reaching  $\approx 30\%$  of the Moore bound. This is quite good: PolarStar is 31% geometric mean larger than state-of-the-art Bundlefly [29] for radixes in [8, 128].

Bermond, Delorme and Farhi defined the star product  $G * G'$  [6], over two graphs, which we call the *structure graph*  $G$  and the *supernode*  $G'$ , and also gave properties on  $G$  and  $G'$  that give a star product with minimal or no increase in the diameter over that of  $G$ .

We give new properties related to, but distinct from those in [6], and show that results on diameter from [6] also apply to graphs having our new properties. By careful choice of  $G$  and  $G'$ , we construct a star product graph of the desired diameter 3. To get the largest possible network, we then maximize the sizes of  $G$  and  $G'$ .

The *Erdős-Rényi* (ER) family of polarity graphs, introduced by Erdős and Rényi [15] and independently by Brown [8], has the property required for  $G$ , so we use this for a structure graph. To our knowledge, this family of graphs is larger than any other known family of diameter-2 graphs, and hence is the best  $G$  candidate that applies to a range of degrees.

In particular, since this family of graphs asymptotically approaches the Moore bound, any improvement to the choice for  $G$  would be small for the radixes of interest. As seen in PolarFly [25], ER graphs also exhibit a modular layout, high bisection bandwidth and other desirable networking properties.

For the supernode  $G'$ , the required properties impose an upper bound on the order of  $G'$ . We construct a new candidate graph for  $G'$ , called *Inductive-Quad*, and use that as our supernode. Inductive-Quad attains the upper bound on order, so is a better  $G'$  candidate than any existing construction, none of which meet the bound.

We also discuss other choices for the supernode  $G'$ , including Paley graphs and complete graphs. While the star products with these choices of  $G'$  do not achieve the same scale as those using Inductive-Quad, they provide flexibility.

## 4 CONSTRUCTION OF THE STAR PRODUCT

### 4.1 The Star Product Graph $G * G'$

Let  $G$  and  $G'$  be two graphs, and let each edge of  $G$  be ordered in some arbitrary way. For each ordered edge  $(x, y)$  in  $G$ , define a bijection  $f_{(x,y)}$  on the vertices of  $G'$ .

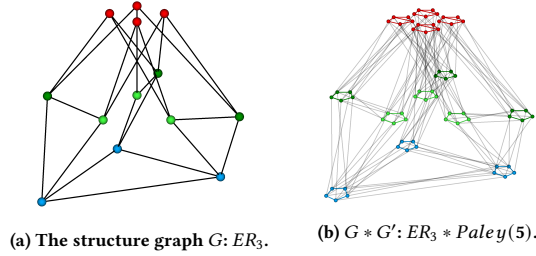
The *star product*  $G_* = G * G'$  is defined as follows [6]:

- The vertex set of  $G_*$  is the Cartesian product of  $G$  and  $G'$ :  $V(G_*) = V(G) \times V(G') = \{(x, x') \mid x \in G, x' \in G'\}$ .
- An edge  $e = ((x, x'), (y, y'))$  exists between vertices  $(x, x')$  and  $(y, y')$  if and only if either
  - (1)  $x = y$  and  $(x', y') \in E(G')$ , or

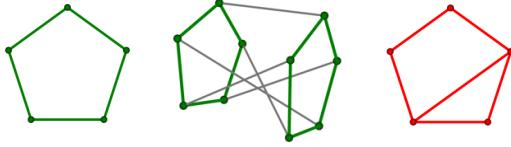
$$(2) (x, y) \in E(G) \text{ and } y' = f_{(x,y)}(x').$$

The construction of an example star product of an Erdős-Rényi polarity graph  $G$  (see Section 6.1) and a Paley graph  $G'$  (see Section 6.2.2) is shown in Figure 2.

Condition (1) on the edges inserts copies of  $G'$  in place of individual nodes of  $G$ . We call these nodes *supernodes*. A pair of these are illustrated in Figure 2c. Condition (2) joins these supernodes to each other, where each edge  $(x, y) \in G$  is replaced by a set of edges between all vertices in the supernodes  $x$  and  $y$ , as defined by  $f_{(x,y)}$ . Sample connectivity for this example is shown in Figure 2c.



(a) The structure graph  $G: ER_3$ . (b)  $G * G': ER_3 * Paley(5)$ .



(c) The supernode  $G': Paley(5)$ . We show here a non-self-loop supernode (in green), the connections between two non-loop supernodes (in green), based on the bijections defined in Section 6.2.2.

Figure 2: Construction of the star product  $ER_3 * Paley(5)$  in 2b. The star product takes the form of the structure graph shown in 2a with each node in 2a replaced by a supernode joined together as in 2c. The degree of the star product is the sum of the degrees of the structure graph and the supernode. The self-loop supernode (in red) has additional edges due to the structure graph self-loop, and degree 1 more than the other supernodes, but this does not change the overall degree, as self-loop nodes in the structure graph  $ER_q$  have degree one less than non-self-loop nodes, as is seen in 2a.

These connections follow the structure of  $G$ , as shown in Figure 2a and Figure 2b. We call  $G$  the *structure graph*. If there are no self-loops in  $G$ , all supernodes are just instances of the graph  $G'$ . However, if there are self-loops in  $G$ , the supernodes representing the nodes with self-loops will be instances of the graph  $G'$  augmented with additional edges. This can be seen in the red supernodes in Figure 2b that correspond to self-adjacent nodes in  $ER_3$ , and have an extra edge due to the self-loop.

## 4.2 Properties of the Star Product

The star product  $G_* = G * G'$  has the following properties:

- (1) The number of vertices is  $|V(G_*)| = |V(G)||V(G')|$ .
- (2) If the maximum degrees in  $G$  and  $G'$  are  $d$  and  $d'$ , respectively, the maximum degree of  $G_*$  is given by  $d_* \leq d + d'$ .
- (3) If the diameters of  $G$  and  $G'$  are  $D$  and  $D'$ , respectively, then the diameter of  $G_*$  is given by  $D_* \leq D + D'$ .

Thus, the star product may be used to construct a large diameter-constrained graph from two smaller graphs.

## 5 LOW-DIAMETER STAR PRODUCTS

The properties listed in Section 4.2 give an upper bound on degrees and diameters of star product constructions. In this section, we show properties on factor graphs that produce a star product of large size and diameter at most 3.

In [6], a set of properties  $P$  and  $P^*$  on graphs  $G$  and  $G'$  was given that guarantees star products  $G * G'$  of large size and small diameter. These graph properties are connected to the bijections  $f: V(G') \rightarrow V(G')$  defined on vertices of  $G'$ .

We present a novel and different set of properties  $R$  and  $R^*$  that also enable low-diameter star product graphs. These properties are similar in spirit to those in [6], but are not the same, and are similarly intertwined with the bijections  $f$  defined on  $V(G')$ .

Our properties also permit the construction of a  $G'$  larger than those found in [6], for a given degree. In fact, we develop a novel construction that meets the upper bound on graphs having the property  $R^*$ . For simplicity of notation, for a function  $f$  defined on vertices  $V$ , we use  $f(E)$  to denote the edges obtained by applying  $f$  to both vertices of each edge in  $E$ .

### 5.1 Properties for Low-Diameter Star Products

The first property applies to the structure graph  $G$ .

**PROPERTY R.** A graph  $G$  of diameter  $D$  has Property R if any vertex pair  $x, y \in V(G)$  can be joined by a path of length  $D$ .

Note that in the definition of Property R, self-loops (if they exist in  $G$ ) are permissible as part of the length- $D$  path.

**COROLLARY 5.1.** In a graph  $G$  having Property R, there exists a path of length  $D + 1$  between any pair of vertices  $x, y \in V(G)$ .

**PROOF.** Consider any neighbor  $z$  of  $y$ . By Property R, there exists a path  $p_D$  of length  $D$  from  $x$  to  $z$ . Appending edge  $(z, y)$  to  $p_D$  gives a path of length  $D + 1$  between  $x$  and  $y$ .  $\square$

Corollary 5.1 highlights the path diversity in  $G$ . We will later see that this diversity enables reachability between all vertex pairs in the star product within  $D + 1$  hops.

The next properties,  $R^*$  and  $R_1$ , apply to the supernode  $G'$ .

**PROPERTY  $R^*$ .** A graph  $G'$  satisfies Property  $R^*$  if there is a bijection  $f$ , with  $f^2$  being the identity on  $G'$  (i.e.,  $f$  is an involution), so that the set of edges

$$E(G') \cup f(E(G')) \cup \{(x', f(x')) \mid x' \in V(G')\}$$

is the entire set of edges in the complete graph on  $V(G')$ .

**COROLLARY 5.2.** A graph  $G'$  satisfies Property  $R^*$  if and only if for any  $x' \in V(G')$ ,

$$\begin{aligned} V(G') &= \{x'\} \cup \{f(x')\} \cup f(N(x')) \cup N(f(x')) \\ &= \{f(x')\} \cup \{x'\} \cup f(N(f(x'))) \cup N(x'). \end{aligned}$$

**PROOF.** We prove that a graph  $G'$  has Property  $R^*$  only if its vertex set  $V(G')$  satisfies the equation given in this corollary. The other direction can be proven similarly.

Consider the graph with edges  $E(G') \cup f(E(G')) \cup \{(x', f(x')) \mid x' \in V(G')\}$ . The neighbors of a vertex  $x'$  in this graph are given by

$\{f(x')\} \cup N(x') \cup f(N(f(x')))$ . By Property  $R^*$ , this is a complete graph on  $V(G')$ . Hence,

$$V(G') = \{x', f(x')\} \cup N(x') \cup f(N(f(x')))$$

For each  $x' \in V(G')$ ,  $f(x')$  is also in  $V(G')$ . Since  $f$  is an involution on  $V(G')$ , by substituting  $y' = f(x')$  in the above, we get that for all  $y' \in V(G')$ ,

$$\begin{aligned} V(G') &= \{y', f(y')\} \cup N(f^{-1}(y')) \cup f(N(y')) \\ \implies V(G') &= \{y', f(y')\} \cup N(f(y')) \cup f(N(y')) \end{aligned} \quad (1)$$

Clearly, the edges  $E(G') \cup f(E(G')) \cup \{(x', f(x')) \mid x' \in V(G')\}$  give a complete graph over  $V(G')$  only if  $V(G') = \{x', f(x')\} \cup f(N(x')) \cup N(f(x'))$ , for all vertices  $x'$ .

The second equality is derived from the first by substituting  $f(x')$  for  $x'$ .  $\square$

Intuitively, Property  $R^*$  and Corollary 5.2 illustrate the ways of covering vertices in the supernode  $G'$ . Starting from any  $x' \in V(G')$ , all vertices are reached by hopping to (a) neighbors of  $x'$  and then their images by  $f$ , or (b)  $f(x')$  and the neighbors of  $f(x')$ .

The following property  $R_1$  is the same as Property  $P_i$  in [6], but is stated here for the special case  $i = 1$  which is all that is needed for diameter 3 star product graphs.

**PROPERTY  $R_1$ .** [6] *A graph  $G'$  has Property  $R_1$  if there is a bijection  $f$ , with  $f^2$  an automorphism of  $G'$ , so that the set of edges*

$$E(G') \cup f(E(G'))$$

*is the entire set of edges in the complete graph on  $V(G')$ .*

## 5.2 Constructing Diameter-3 Graphs with the $R$ Properties

In this section, we show that if the structure graph  $G$  has diameter 2, the star product  $G * G'$  has diameter at most 3:

- If the structure graph  $G$  has Property  $R$  and the supernode  $G'$  has Property  $R^*$  (Theorem 5.3), or
- If the supernode  $G'$  has Property  $R_1$  (Theorem 5.4).

**THEOREM 5.3.** *Let  $G$  and  $G'$  be graphs that satisfy Property  $R$  and  $R^*$ , respectively. If diameter of  $G$  is  $D$ , the star product  $G_* = G * G'$  is a graph with diameter at most  $D + 1$ .*

**PROOF.** Consider any arbitrary vertices  $(x, x'), (y, y') \in V(G_*)$ . By Property  $R$ , there is a path of length  $D$  connecting  $x$  and  $y$  in  $G$ , given by  $p_D = (x, \dots, z, y)$ . Clearly,  $z$  is at distance  $D - 1$  from  $x$  in  $G$ , and  $(z, z' = f^{D-1}(x'))$  is at a distance  $D - 1$  from  $(x, x')$  in  $G_*$ . By Corollary 5.2,

$$y' \in \{z'\} \cup \{f(z')\} \cup f(N(z')) \cup N(f(z')).$$

We proceed by cases. Case 1) is illustrated in Figure 3b, and Case 2) in Figures 3c and 3d, for  $D = 2$  and  $z' = f(x')$ .

(1)  $y' = z' = f^{D-1}(x')$  – Since  $f$  is an involution,  $y' = f^2(y') = f^{D+1}(x')$ . By Corollary 5.1, there is a  $D + 1$  length path between  $x$  and  $y$  in  $G$ , so there is a  $D + 1$  length path between  $(x, x')$  and  $(y, f^{D+1}(x'))$  in  $G_*$ .

(2)  $y' \in \{f(z')\} \cup f(N(z')) \cup N(f(z'))$  – Vertices in  $(z, N(z'))$  and  $(y, f(z'))$  are adjacent to  $(z, z')$ , so vertices in  $(y, N(f(z')))$  and

$(y, f(N(z')))$  are 2 hops away from  $(z, z')$ . So  $(y, y')$  is at distance at most 2 from  $(z, z')$  and  $D + 1$  from  $(x, x')$ .  $\square$

The proof for Theorem 5.3 also highlights the 3–hop paths between vertex pairs in the star product, as shown in Figure 3.

We restate the below theorem from [6] for the property  $R_1$ . The proof is the same as that found in [6].

**THEOREM 5.4.** [6] *Let  $G$  be a graph of diameter  $D \geq 2$ , and let  $G'$  have property  $R_1$ . Define  $f_{(x,y)}(x') = f(x')$  for every arc  $(x, y)$  of an arbitrary orientation of the edges of  $G$ . Then  $G * G'$  has diameter at most  $D + 1$ .*

## 6 GOOD FACTOR GRAPHS

We choose structure graphs  $G$  and supernodes  $G'$  so that the star product  $G * G'$  has the largest possible order for its degree.

- A structure graph  $G$  having Property  $R$  and diameter 2 has a number of vertices limited in theory only by the Moore bound for its degree. For the structure graphs  $G$  in this paper, we use the Erdős-Rényi polarity graphs, as they asymptotically approach the Moore bound quite rapidly.
- On the other hand, a supernode  $G'$  of degree  $d'$  having either Property  $R^*$  or  $R_1$  has order  $\leq 2d' + 2$ . (Properties in [6] give the same bound.) We present here a new construction, Inductive-Quad, that is the first to our knowledge to attain this bound.

Thus, we asymptotically reach the maximum order for star product  $G * G'$ . We also discuss other supernodes  $G'$  with desirable features.

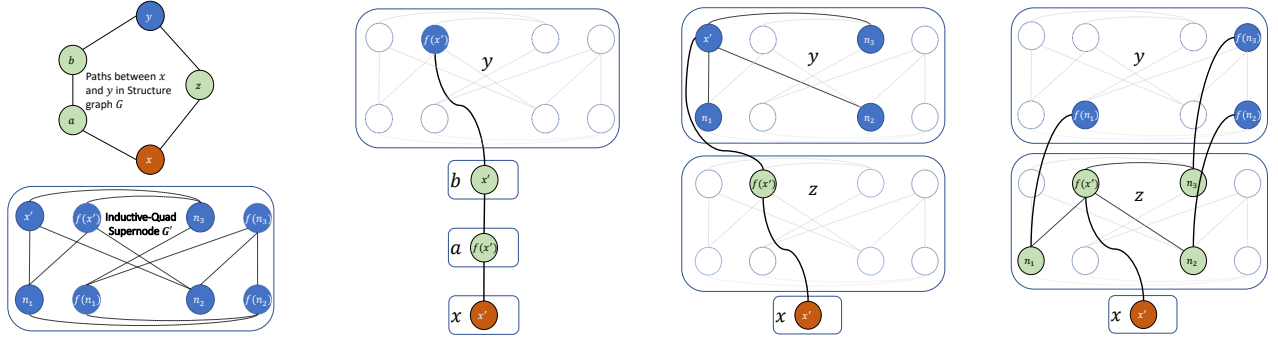
### 6.1 The Erdős-Rényi Structure Graph $G$

The Erdős-Rényi (ER) or Brown family of polarity graphs [8, 15] is based on finite projective geometry, where adjacency is defined by orthogonality. These graphs were used for the PolarFly network, due to their many advantages [25].

An ER graph may be generated for any degree  $q + 1$ , where  $q$  is a prime power. The order of such an ER graph is  $q^2 + q + 1$ . Thus, ER graphs asymptotically reach the Moore bound for diameter-2 graphs and are larger than other known diameter-2 graphs for almost all degrees, as shown in Figure 4. Thus, the ER graph family is a good candidate for the structure graph.

The vertices of  $ER_q$  are vectors  $(x, y, z)$ , with  $x, y, z \in \mathbb{F}_q$ , the finite field of order  $q$ . Vertices  $v$  and  $w$  are adjacent if their dot product  $v \cdot w = 0$ , with addition and multiplication as in  $\mathbb{F}_q$ . Since adjacency is defined by orthogonality of two vectors, all multiples of any two vectors retain the same adjacency relationship. Thus, we move into projective space and consider only the left-normalized form of each vector (so the leftmost non-zero entry of each vector is 1). The ER graph has these left-normalized vectors as the vertices and edges between all orthogonal vector pairs. Note that the arithmetic over finite field  $\mathbb{F}_q$  is used to compute orthogonality. See [33] for details of the arithmetic over  $\mathbb{F}_q$  and [25] for ER graph construction.

$ER_q$  is a diameter-2 graph. This may intuitively be seen by considering perpendicularity in Euclidean space. Each pair of distinct vectors  $v_0$  and  $v_1$  is orthogonal to a common  $w = v_0 \times v_1$ , the cross product of  $v_0$  and  $v_1$ . The 2-hop path from  $v_0$  to  $v_1$  is then given by  $(v_0, w, v_1)$ . The intuition is similar in the case of finite geometry.



(a) Structure graph  $G$  (on top) with Property R. In  $G$ , only the length-2 and 3 paths between  $x$  and  $y$  are shown. Supernode  $G'$  on bottom. (b) Path between  $(x, x')$  and  $(y, f(x'))$ , passing through supernodes corresponding to  $a$  and  $b$  in the structure graph. (c) Paths between  $(x, x')$  and all vertices  $(y, \{x'\} \cup N(x'))$ , passing through the supernode corresponding to  $z$  in the structure graph. (d) Paths between  $(x, x')$  and all vertices  $(y, f(N(f(x'))))$ , passing through the supernode corresponding to  $z$  in the structure graph.

Figure 3: A graphical illustration of the diameter-3 property of the star product on factor graphs  $G$  and  $G'$  with Properties R and  $R^*$ . As an example, we use an arbitrary diameter-2 graph  $G$  having Property R as the structure graph, and the Inductive-Quad  $G'$  from Section 6.2.1, which has Property  $R^*$ , as the supernode. From Property R, there exists a 2-hop and a 3-hop path between any vertices in  $G$ . These paths for two vertices  $x$  and  $y$  are shown in Figure 3a –  $(x, z, y)$  and  $(x, a, b, c)$ . Figures 3b, 3c, and 3d show 3-hop paths in the star product  $G * G'$ , from an arbitrary vertex  $(x, x')$  in supernode  $x$  to all vertices  $(y, y')$  in supernode  $y$ , where  $y' \in \{f(x')\} \cup \{x'\} \cup N(x') \cup f(N(f(x')))$ . This set satisfies corollary 5.2 and hence, is the entire set of vertices in the supernode.

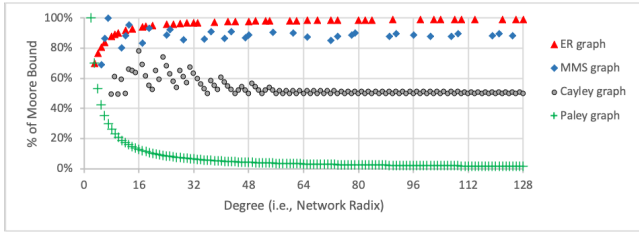


Figure 4: Moore-bound comparison for some known families of diameter-2 graphs: the ER graph, the McKay-Miller-Širáň graphs [34], the best Cayley graphs [1], and the Paley graph. The ER graph has Property R and may be used as a structure graph with  $R^*$  or  $R_1$  graphs. This family asymptotically reaches the Moore bound quickly, so a larger structure graph to be discovered later would only marginally increase the size of the star product.

## 6.2 The Supernodes $G'$

For the supernode  $G'$ , we construct a new family of graphs, called Inductive-Quad. This family, where it exists, meets the upper bound on the order of a supernode. We also explore the Paley graph as a  $G'$  candidate, which is only slightly smaller, and has symmetricity, useful in network design. We discuss here only the Inductive-Quad and Paley graphs, as they give the largest star-products for almost all radices. However, other topologies may also be of interest. For instance, complete graphs provide densely connected regions of locality, and Cayley graphs are highly symmetric [34]. The BDF graph is a graph family designed in [6] specifically for large star products. We list some choices for the supernode in Table 3.

Supernodes	Order	Permitted $d'$	Symmetric	$R^*$	$R_1$
Inductive-Quad	$2d' + 2$	0 or 3 (mod 4)	N	Y	N
Paley	$2d' + 1$	1 (mod 4)	Y	N	Y
BDF [6]	$2d'$	all	N	Y	N
Cayley [34]	$2d' + \delta, \delta \in \{0, \pm 1\}$	$2d' + \delta$ a prime power	Y	N	Y
Complete	$d' + 1$	all	Y	Y	Y

Table 3: Comparison of parameters of degree  $d'$  supernodes.

6.2.1 *Inductive-Quad Graphs (Property  $R^*$ )*. Consider a graph  $G'$  of maximum degree  $d'$  having Property  $R^*$ . By Corollary 5.2, the scale of  $G'$  has upper bound  $|V(G')| \leq 2d' + 2$ . We devise a new construction for  $G'$  that reaches this upper bound.

To get started with this construction, we introduce the following lemma, which describes how a graph of scale  $2d' + 2$  can have Property  $R^*$ .

LEMMA 6.1. *Let  $G'$  be a graph of degree  $d'$  and order  $|V(G)| = 2d' + 2$ .  $G'$  satisfies Property  $R^*$  iff for any pair of vertices  $x, y \in V(G')$  such that  $y \neq x$  and  $y \neq f(x)$ , either*

- (1) Condition  $C_0 \rightarrow x, f(x) \in N(y)$ , or
- (2) Condition  $C_1 \rightarrow x, f(x) \in N(f(y))$ , or
- (3) Condition  $C_2 \rightarrow y, f(y) \in N(x)$ , or
- (4) Condition  $C_3 \rightarrow y, f(y) \in N(f(x))$

For a vertex  $x$ , let  $X_0, X_1, X_2$  and  $X_3$  be the sets of all the vertices  $y$  that satisfy the conditions  $C_0, C_1, C_2$  and  $C_3$ , respectively, and  $X_n$  be the vertices that satisfy neither. Clearly,

$$V(G') = \{x\} \cup \{f(x)\} \cup X_0 \cup X_1 \cup X_2 \cup X_3 \cup X_n$$

Note that  $X_1 = f(X_0), X_2 = f(X_2)$  and  $X_3 = f(X_3)$ . Neighbors of  $x$  and  $f(x)$  are given by  $N(x) = X_0 \cup X_2$  and  $N(f(x)) = X_0 \cup X_3$ , respectively. Consider the vertex set from Property  $R^*$

$$\begin{aligned} V^* &= \{x\} \cup \{f(x)\} \cup f(N(x)) \cup N(f(x)) \\ &= \{x\} \cup \{f(x)\} \cup X_1 \cup X_2 \cup X_0 \cup X_3 \end{aligned}$$

A graph satisfies  $R^*$  if  $V^* = V(G')$ , which is true if and only if  $X_n$  is an empty set.

Intuitively, Lemma 6.1 says that if any graph  $G'$  with  $2d' + 2$  vertices satisfies  $R^*$  for an involution  $f$ , then for any two distinct vertex pairs  $(x, f(x)), (y, f(y)) \in V(G')$ , either both  $x$  and  $f(x)$  connect to a vertex in the other pair or vice-versa.

Here, the involution  $f$  is a construction device: the vertices come in pairs  $(x, f(x))$ , and the edges are drawn to satisfy Property  $R^*$ , and meet the upper bound on  $R^*$  graphs.



**COROLLARY 6.2.** Consider a graph  $G'$  of degree  $d'$  and scale  $V(G) = 2d' + 2$ . For any vertex  $x \in V(G')$ , let  $X$  and  $X_f$  denote the set of pairs  $(y, f(y)) \in N(x)$  and  $(w, f(w)) \in N(f(x))$ , respectively. If  $G'$  satisfies  $R^*$ , then  $X$  and  $X_f$  are disjoint sets of same cardinality.

**PROOF.** Clearly,  $f(X) = X$ . Let  $X' = N(x) \setminus X$  be the neighbors of  $x$  not in  $X$ . Since  $|V(G')| = 2d' + 2$ , from Lemma 6.1,  $X'$  are also neighbors of  $f(x)$ . By Corollary 5.2,  $f(N(x))$  and  $N(f(x))$  should be disjoint sets of same cardinality  $d'$ . Since  $X'$  are common neighbors of both  $x$  and  $f(x)$ , this is only feasible if  $X$  and  $X_f$  are disjoint and  $|X| = |X_f|$ .  $\square$

**PROPOSITION 6.3.** For every non-negative integer  $n$ , there exists a graph  $G'_d$  of degree  $d' = 4n$  or  $d' = 4n + 3$  that satisfies  $R^*$  and has  $2d' + 2$  vertices.

**Construction**  $\rightarrow$  We develop an inductive construction starting from the graphs  $G'_0$  and  $G'_3$  shown in Figure 5a. The graph  $G'_0$  is simply two vertices  $\{x', f(x')\}$  with no edges. It satisfies Property  $R^*$  as adding the edge  $(x', f(x'))$  gives a complete graph.  $G'_3$  has 8 vertices  $\{x', f(x'), y', f(y'), z', f(z'), w', f(w')\}$ .  $f(y'), f(z')$  and  $f(w')$  are adjacent to  $\{z', f(z')\}$ ,  $\{w', f(w')\}$  and  $\{y', f(y')\}$ , respectively. Both  $\{x', f(x')\}$  are adjacent to  $y', z'$  and  $w'$ . Clearly, for any vertex  $v' \in V(G'_3)$ ,

$$V(G'_3) = \{v'\} \cup \{f(v')\} \cup N(f(v')) \cup f(N(v')).$$

Hence, by Corollary 5.2,  $G'_3$  satisfies Property  $R^*$ .

Now we show how to construct a graph of degree  $d' + 4$  with Property  $R^*$ , from a graph of degree  $d'$  (Figure 5). Assume that we have a graph  $G'_d$  of degree  $d'$  that satisfies  $R^*$  and has  $2d' + 2$  vertices. As shown in Figure 5b,  $V(G'_d)$  can be partitioned into two disjoint sets of  $d + 1$  vertices each –  $A$  and  $f(A)$ . To construct  $G'_{d'+4}$ , we add 8 vertices  $\{x', f(x'), y', f(y'), z', f(z'), w', f(w')\}$  with the subgraph  $G'_3$  induced between them. Next, we add edges

- (1) between  $\{x', f(x'), z', f(z')\}$  and all vertices in  $A$ , and
- (2) between  $\{y', f(y'), w', f(w')\}$  and all vertices in  $f(A)$ .

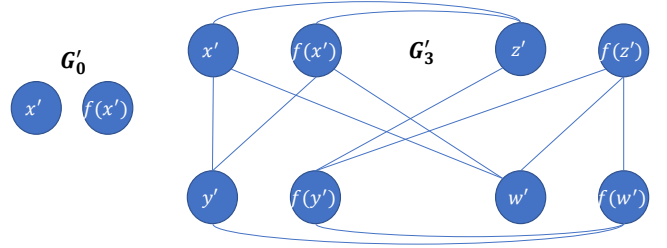
Clearly,  $G'_{d'+4}$  has  $2d' + 10$  vertices and degree  $d' + 4$  and using Lemma 6.1, we can easily verify that it satisfies  $R^*$ .

**PROPOSITION 6.4.** A graph  $G'_d$  of degree  $d'$  that satisfies  $R^*$  and has  $2d' + 2$  vertices, can only exist if  $d' = 4n$  or  $d' = 4n + 3$  for some non-negative integer  $n$ .

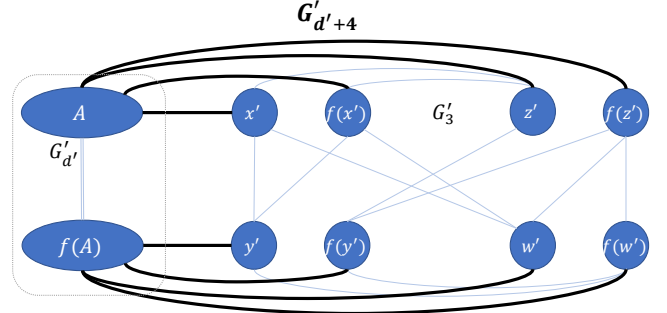
**PROOF.** Proposition 6.3 shows existence of such graphs for  $d = 4n$  and  $d' = 4n + 3$ . Assume that there is a graph  $G'_d$  of degree  $d' = 4n + 1$  or  $4n + 2$  and  $2d' + 2$  vertices that satisfies  $R^*$ . Create a directed graph  $G_D$  consisting of a unique vertex  $X = (x', f(x'))$  for each unordered pair  $\{x', f(x')\} \in V(G'_d)$ . Clearly,  $|V(G_D)| = d' + 1$ . For any pair of vertices  $\{X, Y\} \in V(G_D)$ , insert an oriented arc  $(X \rightarrow Y)$  if either  $y'$  or  $f(y')$  is adjacent to both  $x'$  and  $f(x')$ .

From Corollary 5.2,  $\{x'\}, \{f(x')\}, f(N(x'))$  and  $N(f(x'))$  must be disjoint for all  $x'$  if  $|V(G'_d)| = 2d' + 2$ . Hence,  $G_D$  cannot have self-loops, or bidirectional edges. From Corollary 6.2, the in-degree of every vertex in  $V(G_D)$  should be even and hence, total number of arcs should be even.

From Lemma 6.1, for any pair of vertices  $\{X, Y\} \in V(G_D)$ , either  $(X \rightarrow Y)$  or  $(Y \rightarrow X)$  arc exists. Therefore, total number of arcs is given by  $\binom{d'+1}{2} = \frac{d'(d'+1)}{2}$  arcs which is an odd number if  $d' =$



(a) Base Inductive-Quad graphs of degree  $d' = 0$  and  $d' = 3$ .



(b) The construction of the Inductive-Quad graph of degree  $d' + 4$  from PolarStar graphs of degrees  $d'$  and 3.

**Figure 5: Inductive construction of Inductive-Quad topology with embedded bijection  $f$  that satisfies Property  $R^*$ .**

$4n + 1$  or  $d' = 4n + 2$ . This is a contradiction and hence,  $G'_d$  cannot exist.  $\square$

**6.2.2 Paley Graphs (Property  $R_1$ ).** Let  $q = 4e + 1$  be a prime power. The Paley graph is a well-known graph having degree  $d' = 2e$  and  $q = 2d' + 1$  vertices [17]. The vertex set is the set of elements in  $\mathbb{F}_q$  and there is an edge  $(x, y)$  in the graph if  $y - x$  is a square in  $\mathbb{F}_q$ .

Bermond, Delorme, and Farhi [6] show the Paley graph satisfies their Property  $P_1$ , which is equivalent to Property  $R_1$ , using the following bijection  $f$ : Let  $\zeta$  be a primitive root of  $\mathbb{F}_q$ , i.e. an element for which the sequence  $\zeta, \zeta^2, \zeta^3, \dots$  covers all of  $\mathbb{F}_q$  except 0. For each arc  $(x, y)$  in the Paley graph, we define  $f_{(x,y)}(\alpha) = \zeta\alpha$ . Using this  $f$ , Paley graphs satisfy Property  $P_1$ , and hence Property  $R_1$ .

**PROPOSITION 6.5.** [6] Let  $q = 4e + 1$  be a prime power. The Paley graph satisfies Property  $R_1$ .

The Paley graphs only exist when  $q = 4e + 1$ . They do not achieve the maximum scale of  $2d' + 2$  vertices.

## 7 DESIGN SPACE OF POLARSTAR

We evaluate the scale of network achievable by PolarStar and compare it against the existing diameter-3 topologies.

### 7.1 Theoretical Analysis

Recall that the degree of star product  $G * G'$  is  $\deg(G) + \deg(G')$ , and the order is  $|V(G)| \cdot |V(G')|$ . Our structure graph is an  $ER_q$ , which has degree  $d = q + 1$  and order  $q^2 + q + 1$ , where  $q$  is a prime power. If we use a Inductive-Quad supernode of degree  $d'$ , we get

a PolarStar  $G_*$  of degree  $d_* = d + d'$  and order

$$|V(G_*)| = (q^2 + q + 1)(2d_* - 2q)$$

The order is maximized for

$$\begin{aligned} \arg \max_q V(G_*) &= \frac{(d_* - 1) + \sqrt{(d_* - 1)(d_* - 2)}}{3} \\ &\approx \frac{d_* - 1 + d_* + \frac{1}{2}}{3} \approx \frac{2d_*}{3} \end{aligned} \quad (2)$$

Substituting this value of  $q$ , we get that the maximum order of PolarStar for a given degree  $d_*$  is

$$\max_{\text{Inductive-Quad}} |V(G_*)| \approx \frac{8d_*^3 + 12d_*^2 + 18d_*}{27}. \quad (3)$$

Similarly, if we use Paley graphs for supernodes,  $\max_{\text{Paley}} |V(G_*)| \approx \frac{8d_*^3}{27}$ . Thus, PolarStar asymptotically reaches  $\frac{8}{27}$ th of Moore bound for diameter-3 topologies.

In practice, the degree distribution among factor graphs is constrained – (a)  $q$  must be a prime power in an ER graph of degree  $q + 1$ , and (b) Inductive-Quad and Paley graphs exist for a subset of integer radices, as shown in Table 3. Therefore, we evaluate all feasible combinations of  $d$  and  $d' = d_* - d$  for both Inductive-Quad and Paley supernodes, and select the combination with maximum order for degree  $d_*$  PolarStar.

## 7.2 Scalability in Practice

Figure 1 compares the scalability of PolarStar and other direct diameter-3 topologies, in terms of their Moore-bound efficiency. Clearly, PolarStar *exceeds the scalability* of all known diameter-3 topologies. Compared to HyperX [2] and Dragonfly [24], it achieves 672% and 91% geometric mean increase in the order, respectively.

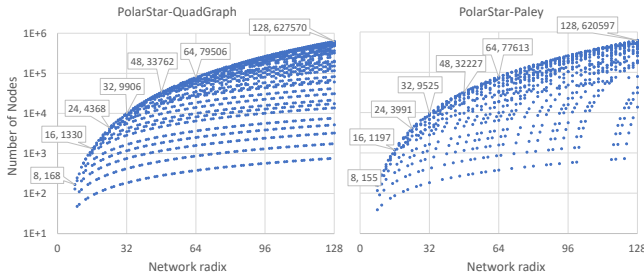


Figure 6: Feasible combinations of network radix and order in PolarStar .

Unlike the state-of-the-art Bundlefly [29], PolarStar offers a feasible construction for every network radix and more *stable scaling* with respect to the Moore bound. Ignoring outliers in Bundlefly, PolarStar is 22% geometric mean larger. If Bundlefly outliers are also considered, PolarStar is 31% geometric mean larger. This increase in scale results from PolarStar’s use of more scalable structure graphs and supernodes in its star product, with larger number of achievable degree distributions for better optimization of scale.

PolarStar also approaches the optimal scale for a star-product diameter-3 network. This is because: (a) the Erdős-Rényi topology of structure graph asymptotically reaches the diameter-2 Moore bound, and (b) the Inductive-Quad supernode topology reaches the

optimal order for graphs satisfying Property R\*. In other words, PolarStar is about as large as it gets for diameter-3 star product graphs. In order to improve upon it, one would need to either

- (1) Design a larger diameter-2 graph than the ER graphs, which would provide little improvement given the proximity of ER graphs to the Moore bound, or
- (2) Develop new mathematical properties that give a star product with diameter-3.

Both of these problems seem difficult and require sophisticated mathematical breakthroughs.

Except for  $k = 23, 50, 56, 80$ , the largest PolarStar order for degrees  $k \in [8, 128]$  is constructed with the Inductive-Quad supernode. It has higher order and more feasible radices than Paley graphs.

Besides the largest construction, PolarStar also offers a wide range of network orders for each radix, as shown in Figure 6. This diversity of feasible designs is enabled by varying (a) the degree distribution between external and supernode graphs, and (b) the choice of supernode graph.

## 8 LAYOUT

For physical deployment, a modular network topology is desirable, comprised of smaller identical subgraphs that can be implemented as blades or racks in a system. Further, if the adjacent modules share multiple links between them, they can be bundled into Multicore Fibers (MCFs) to reduce the number of cables in the system [3, 29]. Since PolarStar is a star product of two graphs, it exhibits a *hierarchically modular structure* as shown in Figure 7.

We also show that PolarStar *supports bundling* of inter-module links. For analysis, we assume a degree  $d_*$  PolarStar with  $ER_q$  structure graph of degree  $q + 1$ , and Inductive-Quad supernode of degree  $d' = d_* - (q + 1)$ . PolarStar with Paley supernode exhibits similar properties. We exploit the PolarFly layout for  $ER_q$  proposed in [25].

The *supernode* topology is the smallest building block of  $2d_* - 2q$  nodes in PolarStar, and is replicated  $q^2 + q + 1$  times in its topology (once per node of  $ER_q$ ). There are  $2(d_* - q)$  links between each pair of adjacent supernodes that can be bundled together (Figure 7b), resulting in  $q(q + 1)^2$  inter-module MCFs (non-self-loop edges in  $ER_q$  [25]). Thus, from Equation (2), bundling inter-supernode links can reduce the global cables by a factor of  $\frac{2d_*}{3}$ .

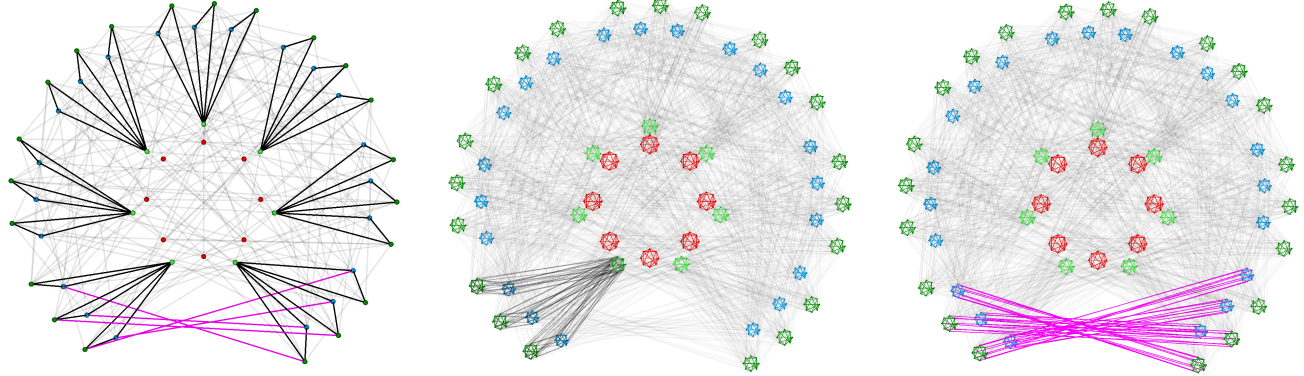
The next level in modular hierarchy is the *clusters of supernodes*. As shown in Figure 7a, the  $ER_q$  structure graph can be divided into  $q + 1$  clusters, classified as follows:

- *Quadric* cluster : single cluster of all  $q + 1$  quadrics (supernodes incident with a self-loop in  $ER_q$ ). There are no links between these supernodes.
- *Non-quadric* clusters :  $q$  clusters, each having  $q$  non-quadric supernodes in  $ER_q$ . These  $ER_q$  clusters contain  $\frac{3q-1}{2}$  edges and  $\frac{q-1}{2}$  edge disjoint triangles. In PolarStar, these clusters will contain  $\frac{3q-1}{2}$  bundles of inter-supernode links.

Multiple bundles of links connect each pair of supernode clusters:

- The Quadric cluster shares exactly  $q + 1$  bundles of inter-supernode links ( $2(q + 1)(d_* - q)$  total links) with every non-quadric cluster.
- Each pair of non-quadric clusters share exactly  $q - 2$  bundles of inter-supernode links (i.e.  $2(q - 2)(d_* - q)$  individual links).





(a) Modular layout for  $ER_7$  graph [25]. Each group of 3 triangles with a common node is a non-quadric cluster. All red quadric nodes form another cluster. Magenta edges connect two clusters.

(b) Layout of  $PolarStar_{11}$  with an  $ER_7$  structure graph. Each  $ER_7$  node becomes an  $IQ_3$  supernode in  $PolarStar$ . The highlighted bundles of links and incident supernodes are a supernode cluster.

(c) Each pair of supernode clusters in  $PolarStar$  are connected by multiple link bundles (in magenta). Each bundle of links corresponds to a single inter-cluster link of  $ER_7$  shown in Figure 7a.

Figure 7: Hierarchical Modular Layout for  $PolarStar$  derived from a layout for  $ER$  structure graphs used in the  $PolarFly$  network [25]. Adjacent supernodes are connected by a bundle of links and adjacent supernode clusters are connected by multiple such bundles.

The hierarchical modular structure of  $PolarStar$  is highly suitable for bundling. Consider the maximum order  $PolarStar$  for degree  $d_s$  with  $q \approx 2d_s/3$  (Equation (2)). Hypothetically, if it is feasible to bundle all links between a pair of supernode clusters into a single MCF, then  $PolarStar$  will only have these inter-module cables:

- (1) Approximately  $\frac{2d_s^2}{9}$  MCFs, each connecting a pair of supernode clusters.
- (2) Approximately  $\frac{2d_s^2}{3}$  MCFs, each connecting a pair of supernodes.

To put it into perspective, a 64-radix  $PolarStar$  with 79,506 nodes will have only 910 and 2,730 inter-cluster and inter-supernode MCFs, respectively, after bundling.

## 9 PERFORMANCE EVALUATION

Network	Parameters	# Routers	Network Radix	# Endpoints
$PolarStar$ with Inductive-Quad (PS-IQ)	$d=12, d'=3, p=5$	1,064	15	5,320
$PolarStar$ with Paley (PS-Pal)	$d=9, d'=6, p=5$	993	15	4,965
Bundlefly (BF)	$d=11, d'=4, p=5$	882	15	4,410
3-D HyperX (HX)	$S=10, L=3, p=9$	1,000	27	9,000
Dragonfly (DF)	$a=12, h=6, p=6$	876	17	5,256
Spectrally (SF)	$\rho=23, q=13, p=8$	1,092	24	8,736
Megaflly (MF)	$\rho=8, a=16, p=8$	1,040	16	4,160
3-level Fat Tree (FT)	$n=3, p=18$	972	36	5,832

Table 4: Topology configurations used in simulations.

### 9.1 Topologies

We compare  $PolarStar$  with Bundlefly [29], Megaflly [16, 37] and Spectrally [40] as state-of-the-art diameter-3 networks, 3-D HyperX [2] and Dragonfly [24] as popular diameter-3 networks in practice, and 3-level Fat trees [30] as the most widely used network. Networks such as torus, hypercube or Flattened Butterfly have been shown to have lower performance than these baselines [7, 24]. We also explored the Galaxyfly family of flexible low-diameter topologies [28]. A diameter-3 Galaxyfly is isomorphic to a Dragonfly,

which is included in the comparison. The configurations of topologies used are shown in Table 4. For direct diameter-3 networks, number of endpoints per router ( $p$ ) is 1/3 of the network radix. In Megaflly and Fat-tree, half and one-third of the routers, respectively, have endpoints on half of their ports.

### 9.2 Routing

We use the following well-known routing schemes to analyze the performance of  $PolarStar$  and other networks:

- *Minimal Static Routing (MIN)*: Every packet between a source and destination is routed along a pre-determined shortest path.
- *Multiple Minpath Routing (M\_MIN)*: It uses multiple minimal paths between source and destination, if they exist. At each hop, ties are broken on the basis of the local output buffer occupancy.
- *Load-balancing Adaptive Routing (UGAL)*: Valiant routing uniformly distributes the network load by misrouting each packet to a randomly selected intermediate router and from there, routing it to the destination. In our UGAL implementation, Valiant misrouting is employed when local output buffers on shortest path(s) have more than 25% occupancy. For misrouting, we sample 4 feasible intermediate routers at random and predict overall path latency via these, as a product of the corresponding local output buffer occupancy and estimated path lengths. The router with smallest latency estimate is used as the intermediate.

For DF and MF, MIN and M\_MIN are identical. DF only has a single shortest path between any router pair. MF only has min-path diversity between routers within the same supernode (group), which we already account for in MIN.

### 9.3 Simulation Setup

We analyze network performance using the cycle-accurate BookSim simulator [20]. Simulation parameters (latency, bandwidth) are

normalized to the values of a single link. We use packets of size 4 flits each and input-queued routers with 128 flit buffers per port and 4 virtual channels (VCs). Dragonfly and Megafly respectively use 2 and 1 VCs for minpath routing, and 3 and 2 VCs for adaptive routing. A warm-up phase precedes all simulations, for the network to reach steady-state before measurements.

To analyze network performance, we use synthetic traffic patterns that represent crucial applications. Synthetic patterns are widely used to compare network topologies [2, 7, 24, 25, 29].

- (1) *Uniform random traffic* – the destination for each packet is selected uniformly at random (represents graph processing, sparse linear algebra solvers, and adaptive mesh refinement methods [7, 25, 41]).
- (2) *Random permutation traffic* – a fixed permutation mapping  $\tau$  of source to destination routers is chosen uniformly at random. All endpoints on a router  $R_s$  transmit only to corresponding endpoints on router  $\tau(R_s)$ . This pattern also emulates permutation traffic under a co-packaged setting where compute node is integrated with the router. Permutation traffic is commonly seen in FFT, physics simulations and collectives [7]. The random permutation pattern represents the traffic generated by these applications when process IDs are randomly assigned to the nodes.
- (3) *Bit Shuffle traffic* - address of the destination is obtained by shifting the source address bits to the left by 1 ( $d_i = s_{(i-1) \bmod b}$ ). This pattern is common in FFT and sorting algorithms [4].
- (4) *Bit Reverse traffic* - address of the destination is obtained by reversing the bit order of source address ( $d_i = s_{b-i-1}$ ). This pattern occurs in Cooley-Tukey FFT, binary search and dynamic tree data structures [18, 36, 39].

Bit Shuffle and Bit Reverse traffic use  $2^b$  endpoints, where  $2^b$  is the largest power of two no more than the total endpoints. The endpoints on a router have contiguous addresses, and in hierarchical topologies (PolarStar, Bundlefly, Dragonfly, Megafly, Fat tree), endpoints in each supernode/sub-tree are also contiguously addressed. In such topologies, almost all endpoints in any supernode communicate with only two other supernodes under Bit Shuffle.

## 9.4 Results

Figure 8 shows a comparison of PolarStar performance against the baseline topologies, for different routing schemes and traffic patterns. The labels follow the scheme < topology > - < routing > and the load is normalized to the peak injection bandwidth.

Overall, PS-Pal and PS-IQ perform well for most of the patterns. With multiple minpaths, PS-Pal and PS-IQ sustain more than 75% of full injection bandwidth (load) on uniform traffic. With adaptive UGAL routing, they sustain between 0.4 to 0.6 of the full load on various traffic patterns. Their performance is comparable to BF which is also a star-product topology, and significantly better than DF with adaptive routing. At small load, latencies of all diameter-3 topologies with minpath routing, are comparable. However, for PS-IQ, PS-Pal, BF and HX, the use of multiple minimum paths significantly improves the maximum sustained load.

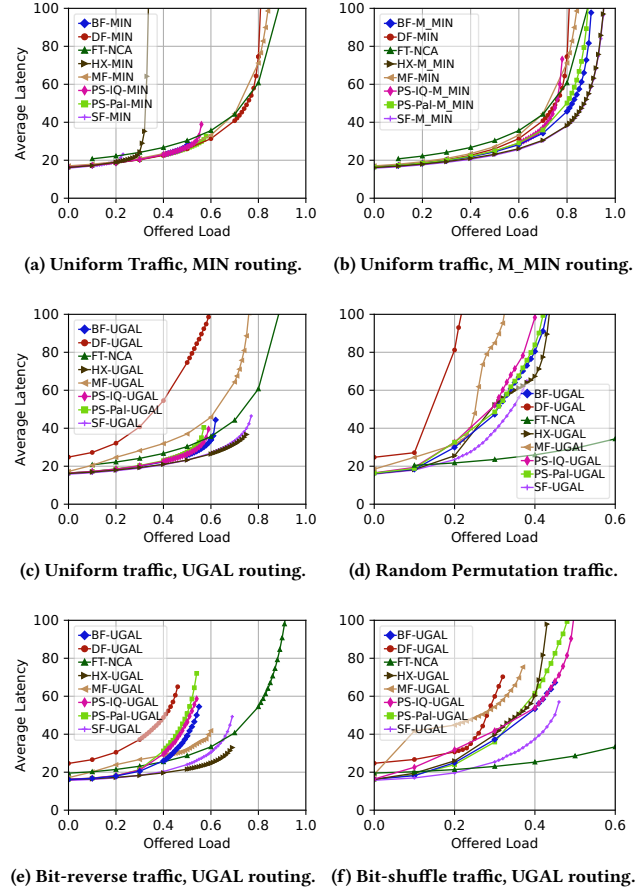


Figure 8: Performance comparison of PolarStar and various topologies. We only measure latency till the highest injection rate where simulation is stable, beyond which the network is saturated and average latency increases with simulation time. For Random Permutation Traffic and Shuffle traffic, FT exceeds average 100 units latency at approximately 0.9 load – some FT data points are omitted to improve display clarity for performance of diameter-3 networks.

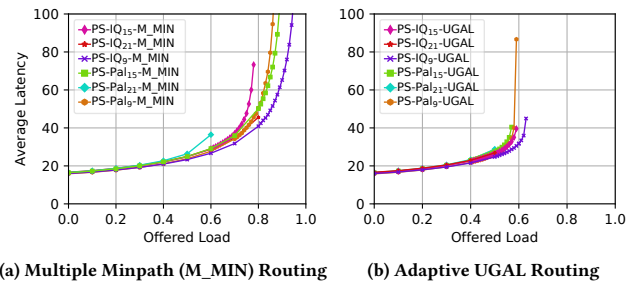


Figure 9: Performance comparison of PolarStar of different sizes

For Random permutation and Bit Shuffle, PS-\* perform better than MF and DF. Bit Shuffle performance of MF and DF is poor because they only have one link between each pair of supernodes [16], as opposed to BF and PS-\* which have a multiple links between the supernodes. This pattern thus highlights the benefits of star-product

topologies over DF and MF. Comparatively, MF performs better on Bit Reverse pattern which has more balanced load distribution.

For most traffic patterns, HX and SF sustain the highest injection rate. This is because of high degree of symmetry, and importantly, the high link density in these topologies. However, HX and SF trade off scaling efficiency (Figure 1 – 6.7× and 12.8× geometric mean lower than PolarStar, respectively), resulting in higher construction cost relative to PolarStar.

*PolarStar Size Comparison* → Figure 9 illustrates the performance of PolarStar for radixes 9 (Paley: 189 routers, IQ: 248 routers), 15 (Paley: 993 routers, IQ: 1064 routers) and 21 (Paley: 2457 routers, IQ: 2928 routers). We keep the ratio of endpoints per router to network radix is 1:3. Both PS-IQ and PS-Pal exhibit consistent performance across different sizes for both the routing schemes. There is a slight variation in saturation bandwidth under M\_MIN routing. This is likely due to the different ratios of supernode and structure graph degrees, which can cause variations in the number of shortest paths available between endpoints.

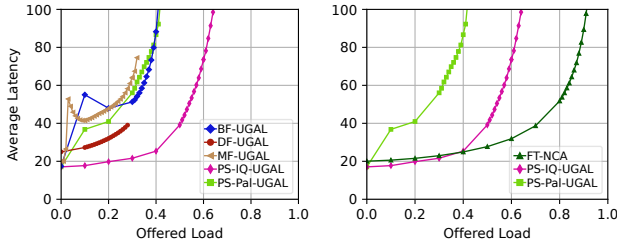


Figure 10: Performance of different topologies under adversarial traffic.

## 9.5 Adversarial Traffic

For hierarchical topologies – PS-\*, BF, DF and MF, we implement an adversarial pattern where all endpoints in a supernode transmit to endpoints in only one other supernode, resulting in congestion on global (inter-supernode) links. To stress the global links, for every source and destination pair, we enforce the longest possible minpath (3-hops in PS-\*, BF, DF and MF, and 4-hops in FT), and in PS-\* and BF, also the maximum number of global hops (3 in PS-IQ, 2 in PS-Paley and BF), as shown in Figure 3. This pattern has been used as a worst-case scenario for BF, DF and MF [16, 24, 29]. Similar to the random permutation pattern in sec.9.3, all endpoints on a router transmit to only one destination router. This may not be the worst-case pattern for PS-\* because the true worst-case may depend on the routing algorithm and the minimum bisection (which is NP-hard to compute for arbitrary graphs).

Figure 10 shows the network performance under this traffic pattern. DF and MF saturate at the lowest bandwidth as they only have a single link between supernode pairs. Comparatively, BF and PS-\* perform better because they have a bundle of multiple links between every pair of supernodes. PS-IQ performance is superior to DF, MF, BF and PS-Pal because of the relatively larger proportion of global links (table 4). Fat tree has the highest saturation, as expected.

## 10 STRUCTURAL ANALYSIS

### 10.1 Bisection Analysis

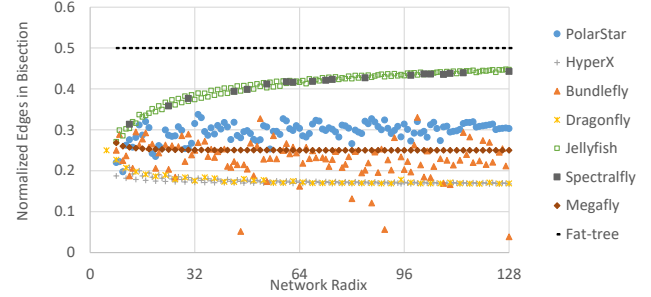


Figure 11: Fraction of links crossing the minimum bisection estimated by METIS [21]. PolarStar, Bundlefly and Spectralfly use their largest feasible (diameter-3) constructions for each radix. Jellyfish has the same radix and scale as PolarStar. Fat-tree and Megafly bisection is normalized by the network links incident with routers that have attached endpoints.

Figure 11 shows the minimum bisection of different topologies for network radix in range [8, 128]. The minimum bisection is estimated using METIS [21] for PolarStar, Spectralfly, Megafly, Bundlefly, Jellyfish and Dragonfly. Among the direct networks, Jellyfish has the highest fraction of links in bisection due to the random connectivity between vertices, although its diameter is more than 3. Spectralfly uses Ramanujan graphs that optimize the expansion properties. Hence, it has a large bisection (comparable to Jellyfish), but it has very few feasible diameter-3 constructions. Among the other diameter-3 topologies, PolarStar has the *highest* proportion of links crossing the bisection, with an average of 29.6% across all radixes. Comparatively, Bundlefly, Dragonfly, HyperX and even the indirect Megafly only have 22.9%, 17.8%, 17.4% and 25.5% links in the bisection cut, respectively. The improved bisection cut can be attributed to the near-optimal

- (1) expansion of ER topology of the structure graph [25], and
- (2) radix distribution across supernode and structure graphs due to plethora of choices for supernode radix.

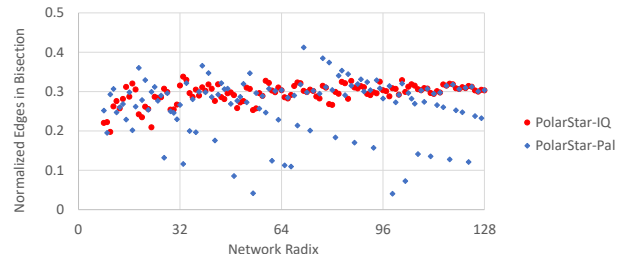


Figure 12: Minimum bisection of PolarStar with Inductive-Quad and Paley supernodes, approximated by METIS [21].

Figure 12 shows the size of bisection cut of PolarStar as a function of radix and supernode topologies. PolarStar with Inductive-Quad and Paley supernodes have an average 29.5% and 26.6% edges in the bisection cut, respectively. The former also offers a more stable bisection across a range of radixes. This is because Inductive-Quad has more feasible radixes and allows better distribution of radix



between structure graph and supernode. Comparatively, the limited choice of radices for Paley graphs may result in a PolarStar with large supernodes and small structure graphs. Such a network will have small bisection because most of the links are concentrated within dense supernode subgraphs.

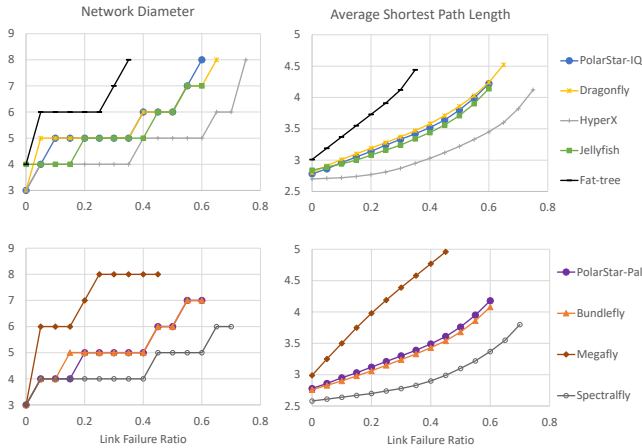


Figure 13: Network Diameter and Average Path Length as a function of random link failures. For Fat-tree and Megafly, we only consider the distance between nodes that have endpoints. For clarity, topologies are split across two graphs.

## 10.2 Fault Tolerance

To estimate fault tolerance, we simulate 100 random link failure scenarios until disconnection, for the networks given in Table 4. We randomly select a simulation with median disconnection ratio, and report the variation in network diameter and average shortest path length in Figure 13.

PolarStar, Bundlefly and Jellyfish have similar resilience with a 60% disconnection ratio. Dragonfly has a higher 65% disconnection ratio. However, at low failure ratios, Dragonfly’s diameter and average shortest path length increases more rapidly. This is likely because if a global link fails, traffic between corresponding Dragonfly groups is routed via another group. HyperX and Spectralfly are the most resilient of diameter-3 topologies due to higher connection density, although they suffer from poor scalability (Figure 1).

All of the direct topologies evaluated in Figure 13 have a much higher disconnection ratio than the indirect topologies Fat-tree and Megafly. Similar to Dragonfly, Megafly has only one global link between each pair of groups. Hence, its diameter increases to 6 with just 5% failed links, and its average shortest path length increases more rapidly than the Fat-tree.

## 11 CONCLUSION

We presented PolarStar – a novel diameter-3 network topology based on star product of factor graphs. PolarStar exhibits state-of-the-art scalability for diameter-3 networks with 91% and 31% geometric mean increase in scale over Dragonfly and state-of-the-art Bundlefly, respectively. PolarStar has several desirable properties including a modular layout, a large design-space, high bisection

bandwidth, and is amenable to bundling of global links into cost-effective multi-core fibers.

## REFERENCES

- [1] Marcel Abas. 2017. Large Networks of Diameter Two Based on Cayley Graphs. In *Computer Science On-line Conference*. Springer, 225–233.
- [2] Jung Ho Ahn, Nathan Binkert, Al Davis, Moray McLaren, and Robert S Schreiber. 2009. HyperX: topology, routing, and packaging of efficient large-scale networks. In *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*. 1–11.
- [3] Yoshinari Awaji, Kunimasa Saitoh, and Shoichiro Matsuo. 2013. *Optical Fiber Telecommunications VIB: Chapter 13. Transmission Systems Using Multicore Fibers*. Elsevier Inc. Chapters.
- [4] Jun Ho Bahn and Nader Bagherzadeh. 2008. A generic traffic model for on-chip interconnection networks. *Network on Chip Architectures* (2008), 22.
- [5] Eiichi Bannai and Tatsuro Ito. 1973. On finite Moore graphs. *Journal of the Faculty of Science, the University of Tokyo. Sect. 1 A, Mathematics* 20 (1973), 191–208.
- [6] J.C. Bermond, C. Delorme, and G. Farhi. 1982. Large graphs with given degree and diameter III. *Ann. of Discrete Math.* 13 (1982), 23–32.
- [7] Maciej Besta and Torsten Hoefler. 2014. Slim Fly: A cost effective low-diameter network topology. In *SC’14: proceedings of the international conference for high performance computing, networking, storage and analysis*. IEEE, 348–359.
- [8] W. G. Brown. 1966. On Graphs that do not Contain a Thomsen Graph. *Can. Math. Bull.* 9, 3 (1966), 281–285. <https://doi.org/10.4153/CMB-1966-036-2>
- [9] Charles Q Choi. 2022. The Beating Heart of the World’s First Exascale Supercomputer. <https://spectrum.ieee.org/frontier-exascale-supercomputer>.
- [10] C. Dalfó. 2019. A survey on the missing Moore graph. *Linear Algebra Appl.* (2019).
- [11] R. M. Damerell. 1973. On Moore graphs. *Proc. Camb. Phil. Soc.* 74 (1973), 227–236.
- [12] Stuart Daudlin, Anthony Rizzo, Nathan C Abrams, Sunwoo Lee, Devesh Khilwani, Vaishnavi Murthy, James Robinson, Terence Collier, Alyosha Molnar, and Keren Bergman. 2021. 3D-Integrated Multichip Module Transceiver for Terabit-Scale DWDM Interconnects. In *Optical Fiber Communications, OFC 2021*.
- [13] Jeffrey Dean and Luiz André Barroso. 2013. The Tail at Scale. *Commun. ACM* 56 (2013), 74–80. <http://cacm.acm.org/magazines/2013/2/160173-the-tail-at-scale/fulltext>
- [14] Jack Dongarra. 2020. *Report on the Fujitsu Fugaku System*. Technical Report ICL-UT-20-06. University of Tennessee, Knoxville.
- [15] Paul Erdős and Alfred Rényi. 1962. On a problem in the theory of graphs. *Publ. Math. Inst. Hungar. Acad. Sci.* 7A (1962), 623–641.
- [16] Mario Flajslik, Eric Borch, and Mike A Parker. 2018. Megafly: A topology for exascale systems. In *High Performance Computing (ISC High Performance 2018)*. Springer, 289–310.
- [17] Chris Godsil and Gordon F Royle. 2013. *Algebraic graph theory*. Springer Science & Business Media.
- [18] Bernard Gold and Charles M Rader. 1969. Digital processing of signals. *Digital processing of signals* (1969).
- [19] A. J. Hoffman and R. R. Singleton. 1960. On Moore Graphs with Diameters 2 and 3. *IBM Journal of Research and Development* 4, 5 (1960), 497–504. <https://doi.org/10.1147/rd.45.0497>
- [20] Nan Jiang, Daniel U Becker, George Michelogiannakis, James Balfour, Brian Towles, David E Shaw, John Kim, and William J Dally. 2013. A detailed and flexible cycle-accurate network-on-chip simulator. In *2013 IEEE international symposium on performance analysis of systems and software (ISPASS)*. IEEE.
- [21] George Karypis and Vipin Kumar. 2009. MeTis: Unstructured Graph Partitioning and Sparse Matrix Ordering System, Version 4.0. <http://www.cs.umn.edu/~metis>.
- [22] G. Kathareios, C. Minkenber, B. Prasadari, G. Rodriguez, and Torsten Hoefler. 2015. Cost-Effective Diameter-Two Topologies: Analysis and Evaluation. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (Austin, TX, USA). ACM.
- [23] Gordon Keeler. [n. d.]. ERI Programs Panel - Phase II Overview. DARPA ERI Summit 2019.
- [24] John Kim, William J. Dally, Steve Scott, and Dennis Abts. 2008. Technology-Driven, Highly-Scalable Dragonfly Topology. In *Proceedings of the 35th ISCA*. IEEE Computer Society, Washington, DC, USA.
- [25] Kartik Lakhota, Maciej Besta, Laura Monroe, Kelly Isham, Patrick Iff, Torsten Hoefler, and Fabrizio Petrini. 2022. PolarFly: A Cost-Effective and Flexible Low-Diameter Topology. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–15.
- [26] Kartik Lakhota, Fabrizio Petrini, Rajgopal Kannan, and Viktor Prasanna. 2021. Accelerating Allreduce with in-network reduction on Intel PIUMA. *IEEE Micro* 42, 2 (2021), 44–52.
- [27] Kartik Lakhota, Fabrizio Petrini, Rajgopal Kannan, and Viktor Prasanna. 2021. In-network reductions on multi-dimensional HyperX. In *2021 IEEE Symposium on High-Performance Interconnects (HOTI)*. IEEE, 1–8.

- [28] Fei Lei et al. 2016. Galaxyfly: A novel family of flexible-radix low-diameter topologies for large-scales interconnection networks. In *Proceedings of the 2016 International Conference on Supercomputing*. 1–12.
- [29] Fei Lei, Dezun Dong, Xiangke Liao, and José Duato. 2020. Bundlefly: A low-diameter topology for multicore fiber. In *Proceedings of the 34th ACM International Conference on Supercomputing*.
- [30] Charles E. Leiserson. 1985. Fat-trees: universal networks for hardware-efficient supercomputing. *IEEE Trans. Comput.* 34, 10 (Oct. 1985), 892–901.
- [31] Lightmatter. [n. d.]. <https://lightmatter.co/>.
- [32] E. Loz et al. 2010. *The Degree-Diameter Problem for General Graphs*. [http://www.combinatoricswiki.org/wiki/The\\_Degree\\_Diameter\\_Problem\\_for\\_General\\_Graphs](http://www.combinatoricswiki.org/wiki/The_Degree_Diameter_Problem_for_General_Graphs).
- [33] Robert J. McEliece. 1987. *Finite Fields for Computer Scientists and Engineers*. Springer, Boston, MA.
- [34] Brendan D. McKay, Mirka Miller, and Jozef Širáň. 1998. A note on Large Graphs of Diameter Two and Given Maximum Degree. *Journal of Combinatorial Theory, Series B* (1998).
- [35] Hans Meuer, Erich Strohmaier, Jack Dongarra, Horst Simon, and Martin Meuer. 2021. Top 500: The List. <https://top500.org/lists/top500/>.
- [36] Hamid Sarbazi-Azad, Mohamed Ould-Khaoua, and Lewis M. Mackenzie. 2001. Communication delay in hypercubes in the presence of bit-reversal traffic. *Parallel Comput.* 27, 13 (2001), 1801–1816.
- [37] Alexander Shpiner, Zachy Haramaty, Saar Eliad, Vladimir Zdornov, Barak Gafni, and Eitan Zahavi. 2017. Dragonfly+: Low cost topology for scaling datacenters. In *2017 IEEE 3rd International Workshop on High-Performance Interconnection Networks in the Exascale and Big-Data Era (HiPINEB)*. IEEE, 1–8.
- [38] Mark Wade, Erik Anderson, Shaha Ardalán, Pavan Bhargava, Sidney Buchbinder, Michael Davenport, John Fini, Haiwei Lu, Chen Li, and Roy Meade. 2020. TeraPHY: a Chiplet Technology for Low-Power, High-Bandwidth In-Package optical I/O. *IEEE Micro* 40, 2 (2020), 63–71.
- [39] Robert Wilber. 1989. Lower bounds for accessing binary search trees with rotations. *SIAM journal on Computing* 18, 1 (1989), 56–67.
- [40] Stephen Young, Sinan Aksoy, Jesun Firoz, Roberto Gioiosa, Tobias Hagge, Mark Kempton, Juan Escobedo, and Mark Raugas. 2022. SpectralFly: Ramanujan Graphs as Flexible and Efficient Interconnection Networks. In *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 1040–1050.
- [41] Xin Yuan, Santosh Mahapatra, Wickus Nienaber, Scott Pakin, and Michael Lang. 2013. A new routing scheme for Jellyfish and its performance with HPC workloads. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. 1–11.