

MARS

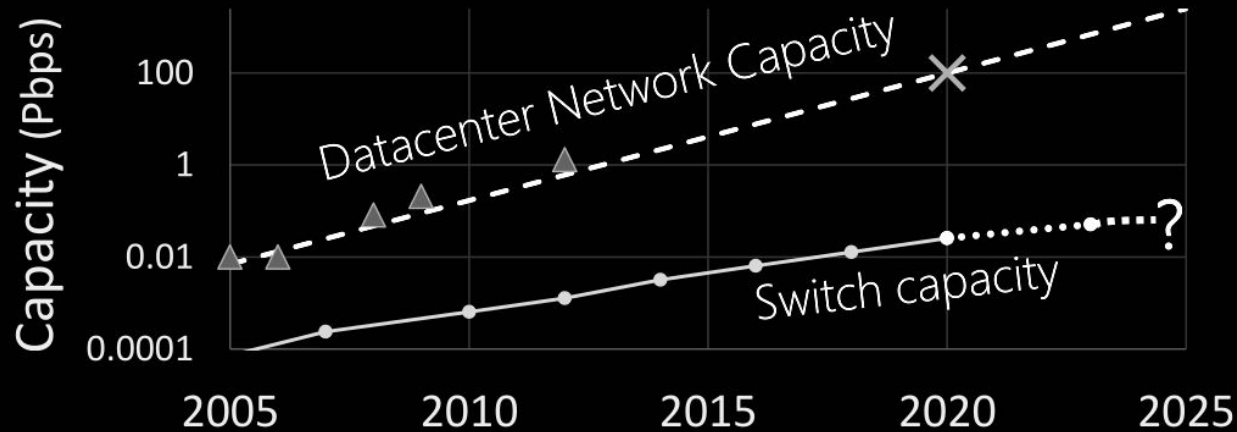
Near-Optimal Throughput with Shallow Buffers in Reconfigurable Datacenter Networks

Vamsi Addanki, Chen Avin, Stefan Schmid



SIGMETRICS 2023

Network Demand vs Capacity **Mismatch**



Traditional Approach Static Data Center Topologies

A Scalable, Commodity Data Center Network Architecture

Mohammad Al-Fares
maalfares@cs.ucsd.edu

Alexander Loukissas
aloukiss@cs.ucsd.edu

Amin Vahdat
vahdat@cs.ucsd.edu

Department of Computer Science and Engineering
University of California, San Diego
La Jolla, CA 92093-0404

ABSTRACT

Today's data centers may contain tens of thousands of computers with significant aggregate bandwidth requirements. The network architecture typically consists of a tree of routing and switching elements with progressively more specialized and expensive equipment moving up the network hierarchy. Unfortunately, even when deploying the highest-end IP switches/routers, resulting topologies may only support 50% of the aggregate bandwidth available at the edge of the network, while still incurring tremendous cost. Non-uniform bandwidth among data center nodes complicates application design and performance. In this paper, we show how to leverage largely commodity Ethernet switches to support the full aggregate bandwidth of clusters consisting of tens of thousands of devices. Similar to how clusters of commodity computers have largely replaced more specialized SMPs and MPPs, we argue that appropriately architected and interconnected commodity switches may deliver more performance at less cost than available from today's higher-end solutions. Our approach requires no modifications to the end host network interface, operating system, or applications, critically, it is fully backward compatible with Ethernet, IP, and TCP.

Categories and Subject Descriptors

C.2.1 [Network Architecture and Design]: Network topology; C.2.2 [Network Protocols]: Routing protocols

General Terms

Design, Performance, Management, Reliability

Keywords

Data center topology, equal-cost routing

1. INTRODUCTION

Growing expense with clusters of commodity PCs has enabled a number of innovations in the design of computation servers and patterns of storage in a cost-efficient manner. Clusters consisting of tens of thousands of PCs are not unheard of in the largest

institutions and thousand-node clusters are increasingly common in universities, research labs, and companies. Important applications classes include scientific computing, financial analysis, data analysis and warehousing, and large-scale network services.

Today, the principle bottleneck in large-scale clusters is often inter-node communication bandwidth. Many applications must exchange information with remote nodes to proceed with their local computation. For example, MapReduce [12] must perform significant data shuffling to transport the output of its map phase before proceeding with its reduce phase. Applications running on cluster-based systems [18, 28, 33, 20] often require remote-node access before proceeding with their IO operations. A proxy to a web search engine often requires parallel communication with every node in the cluster housing the inverted index to return the most relevant results [7]. Even between logically distinct clusters, there are often significant communication requirements, e.g., when updating the inverted index for individual clusters performing search from the site responsible for building the index. Internet services increasingly employ server oriented architectures [13], where the retrieval of a single web page can require coordination and communication with literally hundreds of individual end-services ranging from hardware and communication protocols, such as InfiniBand [2] or Myrinet [6]. While these solutions can scale to clusters of thousands of nodes with high bandwidths, they do not leverage commodity parts (and are hence more expensive and are not naturally compatible with TCP/IP applications). The second choice leverages commodity Ethernet switches and routers to interconnect cluster machines. This approach supports a familiar management infrastructure along with unmodified applications, operating systems, and hardware. Unfortunately, aggregate bandwidths are limited poorly with cluster size, and achieving the highest levels of bandwidth incur non-linear cost increases with cluster size.

For compatibility and cost reasons, most cluster communication systems follow the second approach. However, communication bandwidth in large clusters may become oversubscribed by a significant factor depending on the communication patterns. That is, two nodes connected to the same physical switch may be able to communicate at full bandwidth (e.g., 10Gbps) but moving between switches, potentially across multiple levels in a hierarchy, may limit available bandwidth severely. Addressing these bottlenecks requires non-commodity solutions, e.g., large 10Gbps switches and routers. Further, typical single path routing along lines of interconnected switches means that overall cluster bandwidth is limited by the bandwidth available at the root of the communication hierarchy.

VL2: A Scalable and Flexible Data Center Network

Albert Greenberg
Srinivas Kandula
David A. Maltz

James R. Hamilton
Changhoon Kim
Parveen Patel

Naveendu Jain
Shravanthi Lahoti
Sudipta Sinha Gupta

Microsoft Research

Abstract

To be agile and cost effective, data centers should allow dynamic resource allocation across large server pools. In particular, the data center network should enable any server to be assigned to any service. To meet these goals, we present VL2, a practical network architecture that scales to support huge data centers with uniform high capacity between servers, performance isolation between services, and Ethernet layer semantics. VL2 uses (1) flat addressing to allow service instances to be placed anywhere in the network, (2) Valiant Load Balancing to spread traffic uniformly across network paths, and (3) end-system based address resolution to scale to large server pools without introducing complexity to the network control plane. VL2's design is driven by detailed measurements of traffic and fault data from a large open network technology provider. VL2's implementation leverages proven network technologies, already available at low cost in high-speed hardware implementations, to build a scalable and reliable network architecture. As a result, VL2 networks can be deployed today, and we have built a working prototype. We evaluate the merits of the VL2 design using measurement, analysis, and experiments. Our VL2 prototype shuffles 1.7 Tbit of data among 75 servers in 300 seconds—sustaining a rate that is 94% of the maximum possible.

Categories and Subject Descriptors: C.1 [Computer-Communication Networks]: Network Architecture and Design
General Terms: Design, Performance, Reliability
Keywords: Data center network, commoditization

1. INTRODUCTION

Cloud services are driving the creation of data centers that hold tens to hundreds of thousands of servers and that concurrently support a large number of distinct services (e.g., search, email, music, video, computations, and utility computing). The motivations for building such shared data centers are both economic and technical: to leverage the economies of scale available to bulk deployments and to benefit from the ability to dynamically reallocate servers among services as workload changes and equipment fails. The latter also allows large—upwards of a million per month—for a 100,000 server data center—with the servers themselves comprising the largest cost component. To be profitable, these data centers must achieve high utilization, and key to this is the property of agility—the capacity to assign any server to any service.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear the notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/publisher.
SIGCOMM '09, August 17–21, 2009, Barcelona, Spain.
Copyright 2009 ACM 978-1-60559-500-0/09...\$10.00.

Agility promises improved risk management and cost savings. Without agility, each service must pre-allocate enough servers to meet difficult to predict demand spikes, or risk failure at the brink of success. With agility, the data center operator can meet the fluctuating demands of individual services from a large shared server pool, resulting in higher server utilization and lower costs.

Unfortunately, the design for today's data center network prevents agility in several ways. First, existing architectures do not provide enough capacity between the servers they interconnect. Conventional architectures rely on the tree-like network configurations built from end host hardware. Due to the cost of the equipment, the capacity between different branches of the tree is typically oversubscribed by factors of 1 or more, with paths through the highest levels of the tree oversubscribed by factors of 10 to 1,000. This limits communication between servers to the point that it fragments the server pool—congestion and computation hot spots are prevalent even when spare capacity is available elsewhere. Second, while data centers host multiple services, the network does little to prevent a traffic flood in one service from affecting the other services around it—when one service experiences a traffic flood, it is common for all those sharing the same network sub-tree to suffer collateral damage. Third, the routing design in conventional networks achieves scale-by assigning servers topologically significant IP addresses and dividing servers among VLANs. Such fragmentation of address space limits the utility of virtual machines, which cannot migrate out of their original VLAN while keeping the same IP address. Further, the fragmentation of address space creates an enormous configuration burden when servers must be managed among services, and the human involvement typically required in these reconfigurations limits the speed of deployment.

To overcome these limitations in today's design and achieve agility, we arrange for the network to implement a familiar and compact design that gives each server the address that the client assigned to it, and only those servers, are connected by a single non-interfering Ethernet switch—a *Virtual Layer 2*—and maintain the illusion even as the size of each service varies from a few servers to thousands. Realizing this vision concretely translates into building a network that meets the following three objectives:

- **Uniform high capacity:** The maximum rate of a server-to-server traffic flow should be limited only by the available capacity on the network interface cards of the sending and receiving servers, and assigning servers to a service should be independent of network topology.
- **Performance isolation:** Traffic of one service should not be affected by the traffic of any other service, just as if each service was connected by a separate physical network interface card.
- **Layer 2 semantics:** Just as if the servers were on a LAN—where any IP address can be connected to any part of an Ethernet switch due to flat addressing—data center management software should be able to easily assign any server to any service and configure

Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network

Arjun Singh, Junyi Ong, Amit Agarwal, Glen Anderson, Ashby Armstrong, Roy Bannan, Seb Boving, Gautrav Desai, Bob Felderman, Paulie Germano, Anand Kanagala, Jeff Provoost, Jason Simmons, Elichi Tenda, Jim Wanderer, Urs Hölzle, Stephen Stuart, and Amin Vahdat
Google Inc.
jupiter_sigcomm@google.com

ABSTRACT

We present our approach for overcoming the cost, operational complexity, and limited scale endemic to datacenter networks a decade ago. Three themes unify the five generations of datacenter networks detailed in this paper. First, multi-stage Clos topologies built from commodity switch silicon can support cost-effective deployment of building-scale networks. Second, much of the general, but complex, decentralized network routing and management protocols supporting arbitrary deployment scenarios were overkill for single-operator, pre-planned datacenter networks. We built a centralized control mechanism based on a global configuration pushed to all datacenter switches. Third, modular hardware design coupled with simple, robust software allowed our design to also support inter-cluster and wide-area networks. Our datacenter networks run at dozens of sites across the planet, scaling in capacity by 100x over ten years to more than 21Pbps of inaction bandwidth.

Clos Concepts

• **Networks** – Data center networks;

Keywords

Datacenter Networks; Clos topology; Merchant Silicon; Centralized control and management

1. INTRODUCTION

Datacenter networks are critical to delivering web services, modern storage infrastructure, and as a key en-

abler for cloud computing. Bandwidth demands in the datacenter are doubling every 12–15 months (Figure 1), even faster than the wide area Internet. A number of recent trends drive this growth. Dataset sizes are continuing to explode with more photo/video content, logs, and the proliferation of Internet-connected sensors. As a result, network-intensive data processing pipelines must operate over ever-larger datasets. Next, Web services can deliver higher quality results by accessing more data on the critical path of individual requests. Finally, constellations of co-resident applications often share substantial data with one another in the same cluster, consider index generation, web search, and serving ads.

Ten years ago, we found the cost and operational complexity associated with traditional datacenter network architectures to be prohibitive. Maximum network scale was limited by the cost and capacity of the highest end switches available at any point in time [23]. Those switches were expensive marvels, typically recycled from products targeting wide area deployments. WAN switches were differentiated with hardware support/download for a range of protocols (e.g., IP multi-protocol) or by pushing the envelope of chip memory (e.g., Internet-scale routing tables, off chip DRAM for deep buffers, etc.). Network control and management protocols targeted autonomous individual switches rather than pre-configured and largely static datacenter fabrics. Most of those features were not useful for datacenters: increased cost, complexity, delayed time to market, and made network management more difficult.

Datacenter switches were also built as complex devices targeting the highest levels of availability. In a WAN Internet deployment, using a single switch/router can have substantial impact on applications. Because WAN links are so expensive, it makes sense to invest in high availability. However, more plentiful and cheaper datacenter bandwidth makes it prudent to trade cost for somewhat reduced interconnect capacity. Finally, switches operating in a multi-tenant WAN environment with arbitrary and hosts require support for many protocols to ensure interoperability. In single-operator data-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear the notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/publisher.
SIGCOMM '09, August 17–21, 2009, London, United Kingdom
© 2009 Copyright held by the owner(s).
ACM ISBN 978-1-60559-500-0/09...\$10.00
DOI: <http://dx.doi.org/10.1145/1628508.1787506>

Traditional Approach vs. Scalable Data Center Topologies

A Scalable, Commodity Data Center Network Architecture

Mohammad Al-Fares
malfares@cs.ucsd.edu

Alexander Loukissas
aloukiss@cs.ucsd.edu

Amin Vahdat
vahdat@cs.ucsd.edu

Jellyfish: Networking Data Centers, Randomly

Ankit Singla*, Chi-Yao Hong[†], Lucian Popa[‡], P. Brighten Godfrey[§]
[†]University of Illinois at Urbana-Champaign
[‡]University of California, Berkeley

1 Introduction

Data centers today form the backbone of cloud operations. A well provisioned data center network is important to ensure that services do not face bandwidth bottlenecks to utilize servers from across each other, and to gain more freedom in scaling by adding capacity by having to take placement of workloads to where bandwidth is available [16]. As a result, a significant body of work has tackled the problem of building high network capacity interconnects [1, 10–13, 26, 28, 29].

One crucial problem that has been ignored in these designs is that of incremental expansion of the network, i.e., adding servers and network capacity incrementally to the data center. This may be motivated by growth of the user base, which requires more servers, or by the deployment of more bandwidth-intensive applications. Such expansion can be made feasible by either planned over-provisioning of space and power, or by an upgrade of the server base. The latter enables some of the old servers to be replaced by a larger number of more powerful and, at the same time, less power-consuming new servers.

Industry experience indicates that incremental expansion is an important problem. Consider the growth of Facebook's data center server population from ~30,000 in September 2009 to more than 60,000 by June 2010 [24]. While Facebook has added new data center facilities too, much of this growth is more incremental in existing facilities ("adding capacity on a daily basis" [23]). For instance, Facebook announced that it will double the size of its facility at Prineville, Oregon by early 2012 [9]. Industry experts have also identified incremental build-out as a useful strategy to reduce capex up-front [20].

Do current high-bandwidth data center network proposals allow incremental growth? Consider the fat-tree proposal [1] as an illustrative example. The entire structure is completely determined by the port-count of all the switches available. This is limiting in at least two ways. First, it makes the design space very coarse: full-bisection bandwidths can only be built as sizes $4k, 6k, 8k, 12k, 16k$, and $64k$ corresponding to the commonly

available port counts of 24, 32, 48, and 64. Second, even if (for example) 50-port switches were available, the smallest incremental upgrade from the 48-port switch design would be 4,000 servers. Moreover, this "incremental" growth would require replacing all the 48-port switches to utilize servers from across each other, and to gain more freedom in scaling by adding capacity by having to take placement of workloads to where bandwidth is available [16]. As a result, a significant body of work has tackled the problem of building high network capacity interconnects [1, 10–13, 26, 28, 29].

One crucial problem that has been ignored in these designs is that of incremental expansion of the network, i.e., adding servers and network capacity incrementally to the data center. This may be motivated by growth of the user base, which requires more servers, or by the deployment of more bandwidth-intensive applications. Such expansion can be made feasible by either planned over-provisioning of space and power, or by an upgrade of the server base. The latter enables some of the old servers to be replaced by a larger number of more powerful and, at the same time, less power-consuming new servers.

Industry experience indicates that incremental expansion is an important problem. Consider the growth of Facebook's data center server population from ~30,000 in September 2009 to more than 60,000 by June 2010 [24]. While Facebook has added new data center facilities too, much of this growth is more incremental in existing facilities ("adding capacity on a daily basis" [23]). For instance, Facebook announced that it will double the size of its facility at Prineville, Oregon by early 2012 [9]. Industry experts have also identified incremental build-out as a useful strategy to reduce capex up-front [20].

Do current high-bandwidth data center network proposals allow incremental growth? Consider the fat-tree proposal [1] as an illustrative example. The entire structure is completely determined by the port-count of all the switches available. This is limiting in at least two ways. First, it makes the design space very coarse: full-bisection bandwidths can only be built as sizes $4k, 6k, 8k, 12k, 16k$, and $64k$ corresponding to the commonly

available port counts of 24, 32, 48, and 64. Second, even if (for example) 50-port switches were available, the smallest incremental upgrade from the 48-port switch design would be 4,000 servers. Moreover, this "incremental" growth would require replacing all the 48-port switches to utilize servers from across each other, and to gain more freedom in scaling by adding capacity by having to take placement of workloads to where bandwidth is available [16]. As a result, a significant body of work has tackled the problem of building high network capacity interconnects [1, 10–13, 26, 28, 29].

*Funding: a co-mentor decided the first among the first two authors.

VL2: A Scalable and Flexible Data Center Network

Albert Greenberg
Srikanth Kandula
David A. Maltz

James R. Hamilton
Changhoon Kim
Parveen Patel

Naveendu Jain
Parantap Lahiri
Sudipta Sengupta

Microsoft Research

Slim Fly: A Cost Effective Low-Diameter Network Topology

Maciej Besta
ETH Zurich
maciej.best@inf.ethz.ch

Torsten Hofer
ETH Zurich
hoer@inf.ethz.ch

1. INTRODUCTION

Interconnection networks play an important role in today's large-scale computing systems. The importance of the network grows with ever increasing per-node (multi-core) performance and memory bandwidth. Large networks with tens of thousands of nodes are deployed in warehouse-sized HPC and data centers [8]. Key properties of such networks are determined by their *topologies*: the arrangement of nodes and cables.

Several metrics have to be taken into account while designing an efficient topology. First, high bandwidth is indispensable as many applications perform all-to-all communication [38]. Second, networks can account for as much as 33% of the total system cost [27] and 50% of overall system energy consumption [2] and thus they should be cost and power efficient. Third, low end-to-end latency is important for many applications, e.g., in high frequency trading. Finally, topologies should be resilient to link failures.

In this paper we show that *lowering network diameter* not only reduces the latency but also the cost of a network and the amount of energy it consumes while maintaining high bisection bandwidths. Lowering the diameter of a network has two effects. First, it reduces energy consumption as each packet traverses a smaller number of SerDes. Another consequence is that packets visit fewer switches and router buffers and will thus be less likely to contend with other packets flowing through the network. This enables us to reduce the number of costly routers and connections while maintaining high bisection bandwidth.

The well-known fat tree topology [30] is an example of a network that provides high bisection bandwidth. Slim Fly, our new packet has to traverse many connections as it first has to move up the tree to reach a core router and only then go down to its

destination. Other topologies, such as Dragonfly [28], reduce the diameter to three, but their structure also limits bandwidth and, as we will show, has a negative effect on resiliency.

In this work, we propose a new topology, called Slim Fly, which further reduces the diameter and thus cost, energy consumption, and the latency of the network while maintaining high bandwidth and resiliency. Slim Fly is based on graphs with low-diameter for a given router radix and is, in this sense, approaching the optimal diameter for a given router radix. Figure 1 motivates Slim Fly by comparing the average number of network hops for random uniform traffic using minimal path routing on different network topologies.

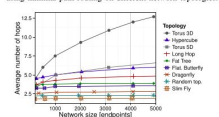


Fig. 1. Comparison of the average number of hops (uniform traffic) in Slim Fly and other networks. Topologies are in balanced or close to balanced configuration (explained in Section III). *allowing for 33% of the total system cost* [27] and 50% of overall system energy consumption [2] and thus they should be cost and power efficient. Third, low end-to-end latency is important for many applications, e.g., in high frequency trading. Finally, topologies should be resilient to link failures.

- We design and analyze a new class of cost effective low-diameter network topologies called Slim Flies.
- We discuss and evaluate different deadlock-free minimal and adaptive routing strategies and compare them to existing topologies and approaches.
- We show that, in contrast to the first intuition, Slim Fly, using fewer cables and routers, is more tolerant towards link failures than comparable Dragonflies.
- We show a physical layout for a datacenter or an HPC center network and a detailed cost and energy model.

*Numbers for random topologies are updated from values obtained using the Bonnet simulator in the lower ones substituted with analytical formulas.

Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network

Xpander: Towards Optimal-Performance Datacenters

Asaf Valadarsky*
asaf.valadarsky@mail.huji.ac.il

Gal Shaafat*
gal.shaafat@mail.huji.ac.il

Michael Dinitz*
mdinitz@cs.huji.ac.il

Michael Schapira*
schapira@huji.ac.il

Abstract

Despite extensive efforts to meet ever-growing demands, today's datacenters often exhibit far from optimal performance in terms of network utilization, resiliency to failures, cost efficiency, incremental expandability, and more. Consequently, many novel architectures for high performance datacenters have been proposed. We show that the benefits of state-of-the-art proposals are, in fact, derived from the fact that they are (implicitly) utilizing "expander graphs" (aka expanders) in their network topologies, thus unifying a unifying theme of these proposals. We observe, however, that these proposals are not optimal with respect to performance, do not scale, or suffer from seemingly insurmountable deployment challenges. We leverage these insights to present Xpander, a novel datacenter architecture that achieves near-optimal performance and provides a tangible alternative to existing datacenter designs. Xpander's design turns ideas from the rich graph-theoretic literature on constructing optimal expanders into an operational reality. We evaluate Xpander via theoretical analysis, extensive simulations, experiments with a network emulator, and an implementation on an SPOC-capable network switch. Our results demonstrate that Xpander significantly outperforms both traditional and proposed datacenter designs. We discuss challenges to real-world deployment and explain how these can be resolved.

1 Introduction

The rapid growth of Internet services is placing tremendous demands on datacenters. Yet, as evidenced by the extensive research on improving datacenters' performance [2, 6, 15, 16, 23, 25, 26], today's datacenters often exhibit far from optimal performance in terms of network utilization, resiliency to failures, cost efficiency, amenability to incremental growth, and beyond.

1.1 The Secret to High Performance

We show that state-of-the-art proposals for next-generation datacenters, e.g., low-diameter networks such as Slim Fly [8] or random networks like Jellyfish [8], have an an-

the network topology and exploiting the diversity of short paths afforded by expanders for efficient delivery of data traffic. Thus, our first contribution is shedding light on the underlying reason for the empirically good performance of previously proposed datacenter architectures, by showing that these proposals are specific points in a much larger design space of "expander datacenters". We observe, however, that these points are either not sufficiently close to optimal performance-wise, are inherently not scalable, or face significant deployment and maintenance challenges (e.g., in terms of unpredictability and wiring complexity).

We argue that the quest for high-performance datacenter designs is inextricably intertwined with the rich body of research in mathematics and computer science on building good expanders. We seek a point in this design space that provides *near-optimal performance* guarantees while providing a *practical alternative* for today's datacenters (in terms of cabling, physical layout, backwards compatibility with today's protocols, and more). We present Xpander, a novel expander-datacenter architecture carefully engineered to achieve both these desiderata.

Importantly, utilizing expanders as network topologies has been proposed in a large variety of contexts, ranging from parallel computing and high-performance computing [6, 9, 15, 16, 23] to optical networks [4] and peer-to-peer networks [10, 39]. Our main contributions are examining the performance and operational implications of utilizing expanders in the *datacenter networking context*, and the optimal design points in this specific domain (namely, Xpander). Indeed, despite the large body of research on expanders, many aspects of using expanders as datacenter networks (e.g., throughput-related performance measures, specific routing and congestion control protocols, deployment costs, incremental growth, etc.) remain little understood. We next elaborate on expanders, expander datacenters, and Xpander.

1.2 Why Expanders?

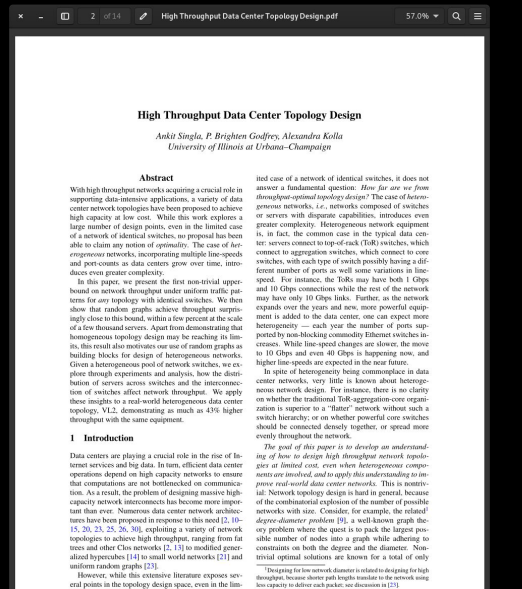
Intuitively, in an expander graph the total capacity from any set of nodes S to the rest of the network is large with respect to the size of S . We present a formal definition of this property in the next section.

Traditional Approach: Static Datacenter Topologies

Which topology has better throughput?

Traditional Approach: Static Datacenter Topologies

What's wrong with that?

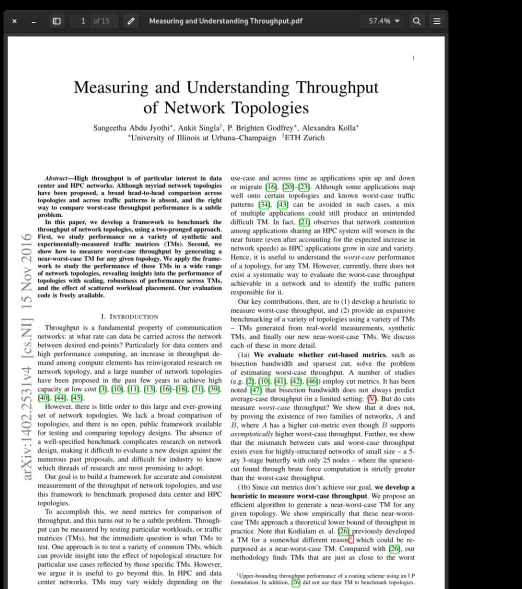


High Throughput Data Center Topology Design

Ankit Singla, P. Brighton Godfrey, Alexandra Kolla
University of Illinois at Urbana-Champaign

Abstract
With high throughput networks acquiring a crucial role in supporting data-intensive applications, a variety of data center network topologies have been proposed to achieve high capacity at low cost. While this work explores a large number of design points, even in the limited case of a network of identical switches, no proposal has been able to claim any notion of optimality. The case of heterogeneous networks, incorporating multiple line-speeds and port counts at data centers grow over time, introduces even greater complexity. In this paper, we present the first non-trivial approach on networks throughput under uniform traffic patterns for any topology with identical switches. We then show that random graphs achieve throughput surprisingly close to the bound, within a few percent of the scale of a few thousand servers. Apart from demonstrating that homogeneous topology design may be reaching its limits, this result also motivates use of random graphs as building blocks for design of heterogeneous networks. Given a heterogeneous pool of network switches, we explore through experiments and analysis, how the distribution of servers across switches and the interconnection of switches affect network throughput. We apply these insights to a real-world heterogeneous data center topology, VL2, demonstrating as much as 43% higher throughput with the same equipment.

1 Introduction
The goal of this paper is to develop an understanding of how to design high throughput network topologies or limited cases, even when heterogeneous components are involved, and to apply this understanding to improve real-world data center networks. This is a motivating Network topology design is hard in general, because of the combinatorial explosion of the number of possible networks with size. Consider, for example, the related degree-diameter problem [9], a well-known graph theoretic problem where the quest is to pack the largest possible number of nodes into a graph while adhering to constraints on the degree and the diameter. Standard optimal solutions are known to exist only for very small values of d and D .
Designing low network diameter is related to designing high throughput because wider paths translate to network nodes less capacity to deliver each packet to destination [21].



Measuring and Understanding Throughput of Network Topologies

Sangeetha Abha Jochi¹, Ankit Singla¹, P. Brighton Godfrey¹, Alexandra Kolla¹
¹University of Illinois at Urbana-Champaign ^{†ETH Zurich}

Abstract—High throughput is of particular interest in data centers and HPC networks. Although myriad network topologies have been proposed, a broad head-to-head comparison across topologies and across traffic patterns is absent, and the right way to compare worst-case throughput performance is unclear. In this paper, we develop a framework to benchmark the throughput of network topologies using a two-pronged approach. First, we study performance on a variety of synthetic and experimentally-motivated traffic matrices (TMs). Second, we show how to measure worst-case throughput by generating and measuring TMs for any given network. We apply the framework to study the performance of three TMs in a wide range of network topologies, relating insights into the performance of topologies with scaling, robustness of performance to TM, and the effect of network workload placement. Our evaluation code is freely available.

1. INTRODUCTION
Throughput is a fundamental property of communication networks, at what size can data be carried across the network between desired end-points? Particularly for data centers and high performance computing, an increase in throughput demands among computer clusters has reignited research on network topology, and a large number of network topologies have been proposed in the past few years to achieve high capacity [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]. However, there is little order to this large and ever-growing sea of network topologies. We lack a broad comparison of topologies, and there is no open, public framework available for testing and comparing network topologies. The absence of a well-specified benchmark complicates research on network design, making it difficult to design and compare networks. Our goal is to build a framework for accurate and consistent measurement of the throughput of network topologies, and this framework to benchmark proposed data center and HPC topologies.

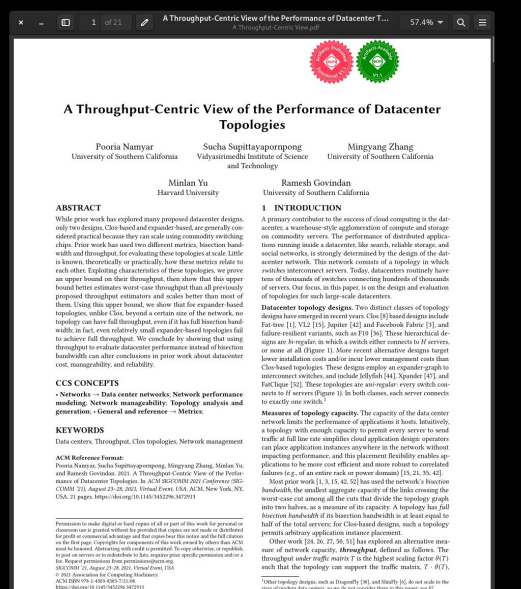
To accomplish this, we need metrics for comparison of throughput, and this turns out to be a subtle problem. What does it mean to compare two particular topologies, or traffic matrices (TMs), in the immediate question is throughput? One approach is to use a variety of common TMs, which can provide insight into the effect of topological structure for particular use cases reflected by these specific TMs. However, we argue it is useful to go beyond this. In HPC and data center networks, TMs may vary widely depending on the

use-case and across time as applications spin up and down or migrate [13, 19, 23]. Although some applications may still use certain topologies and known worst-case traffic patterns [13, 19, 23] can be avoided in such cases, a mix of multiple applications could still produce an unintended difficult TM. In fact, [23] observed that network contention among applications sharing an HPC system will worsen in the near future (even after accounting for the expected increase in network speeds) as HPC applications grow in size and variety. Hence, it is useful to understand the worst-case throughput of a topology, for any TM. However, currently, there does not exist a systematic way to evaluate the worst-case throughput achievable in a network and to identify the traffic pattern responsible for it.

Our key contributions, then, are to (1) develop a heuristic to measure worst case throughput, and (2) provide an extensive benchmarking of a variety of topologies using a variety of TMs — TMs generated from real-world measurements, synthetic TMs, and finally our new worst-work-case TMs. We discuss each of these in more detail.

(a) **We evaluate whether cost-based metrics**, such as bisection bandwidth and spanned cut, solve the problem of estimating worst-case throughput. A number of studies (e.g., [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]) that bisection bandwidths do not always predict worst-case throughput in a limited setting [24]. But do cost-based metrics solve worst-case throughput? We show that it does not, by proving the existence of two families of networks, A and B , where A has a higher cut metric even though B supports a consistently higher worst-case throughput. Further, we show that the mismatch between cut and worst-case throughput exists even in highly-interconnected networks of small size: a 3-ary 3-stage butterfly with only 25 nodes — where the spectral bound would have been fourfold in capacity is strictly greater than the worst-case throughput.

[†]Upon browsing through performance of a certain switch, such as 1P switches. In addition, [25] did not use their TM to benchmark topologies.



A Throughput-Centric View of the Performance of Datacenter Topologies

Poorita Namyar¹, Sudeha Saptharishyapong¹, Mingyang Zhang¹
¹University of Southern California ¹Vidyasasthina Institute of Science and Technology

Milhan Yu¹, Ramesh Govindan¹
¹University of Southern California

ABSTRACT
While prior work has explored many proposed datacenter designs, only two designs, Clos-based and expander-based, are generally considered practical because they can scale using commodity switching chips. Prior work has used two different metrics, bisection bandwidth and throughput, for evaluating these topologies at scale. However, bisection bandwidth, like flow metrics, tends to misbehave in network settings. Today, datacenters routinely have tens of thousands of switches connecting hundreds of thousands of servers. Our focus, on this paper, is on the design and evaluation of topologies for such large-scale datacenters.
Datacenter topology designs. Two distinct classes of topology designs have emerged in recent years. Clos [3] based designs include Fat Tree [1], VL2 [13], Spine [14] and Facebook Fabric [1], and failure-resilient variants, such as Fat Tree [3]. These hierarchical designs are regular in which a switch either connects to servers, or more at all (Figure 1). More recent alternative designs target lower installation costs and/or lower power management costs than Clos-based topologies. These designs employ an expander graph to interconnect switches, and include Jellyfish [4], Xpander [4], and FatSpine [2]. These topologies use an expander, every switch connects to k servers (Figure 1). In both cases, each server connects to exactly one switch.
Measure of topology capacity. The capacity of the data center network limits the performance of application traffic. Intuitively, a topology with enough capacity to permit every server to send traffic at full line rate simplifies cloud application design, liberates cloud client application designers to experiment with the network without impacting performance, and this placement flexibility enables application-specific design choices. Capacity is measured in terms of applications (e.g., of an entire rack or power domain) [13, 21, 35, 42]. Clos-based topologies [1, 13, 14, 15, 43] use fixed network bisection bandwidths, the smaller aggregate capacity of the links connecting the switch cuts cut across all the links that divide the topology graph into two halves, as a measure of its capacity. A topology has full bisection bandwidth if its bisection bandwidth is at least equal to half of the total servers; for Clos-based designs, such a topology permits arbitrary uniform inter-server placement.

Other work [13, 26, 27, 35, 51] has explored an alternative measure of network capacity, **throughput**, defined as follows. The throughput under traffic matrix T is the highest total traffic (in T) such that the topology can support the traffic matrix T . (Fig. 1)
Other topology designs, such as Jellyfish [4] and FatSpine [2], do not divide the space of datacenters into two network families in this paper. (Fig. 1)

Permission to make digital or hard copies of all part of this work for personal or classroom use is granted by ACM, provided that the fee code and the copyright notice are preserved in this version of the work. For all other uses, permission should be sought from ACM. Copyright 2019 ACM 978-1-4503-6931-4/19/0000-0000...\$15.00. DOI: <https://doi.org/10.1145/3328078>
ACM 978-1-4503-6931-4/19/0000-0000...\$15.00
ACM 978-1-4503-6931-4/19/0000-0000...\$15.00

Reconfigurable Datacenter Networks

- Generalization of the design space: *Topology can change over time*
- Static networks are a special case

Reconfigurable Datacenter Networks

- Generalization of the design space: *Topology can change over time*
- Static networks are a special case



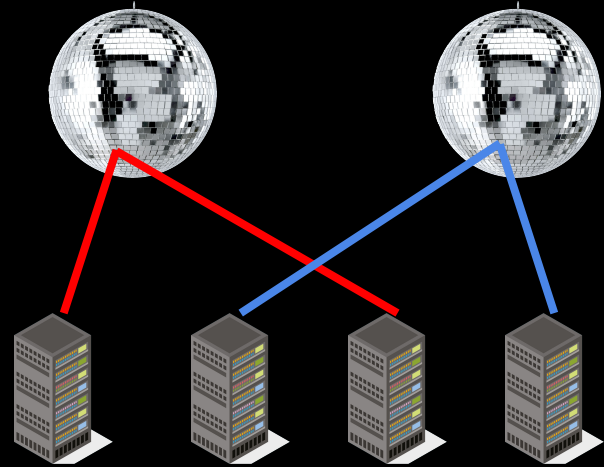
Sirius [Sigcomm 2020]

Reconfigurable Datacenter Networks

- Generalization of the design space: *Topology can change over time*
- Static networks are a special case



Sirius [Sigcomm 2020]



ProjecToR [Sigcomm 2016]

Reconfigurable Datacenter Networks



MA

Reconfigurable Datacenter Networks

Which topology has better throughput?

Reconfigurable Datacenter Networks

Which topology has better throughput?

Main focus of this paper.
Stay tuned!



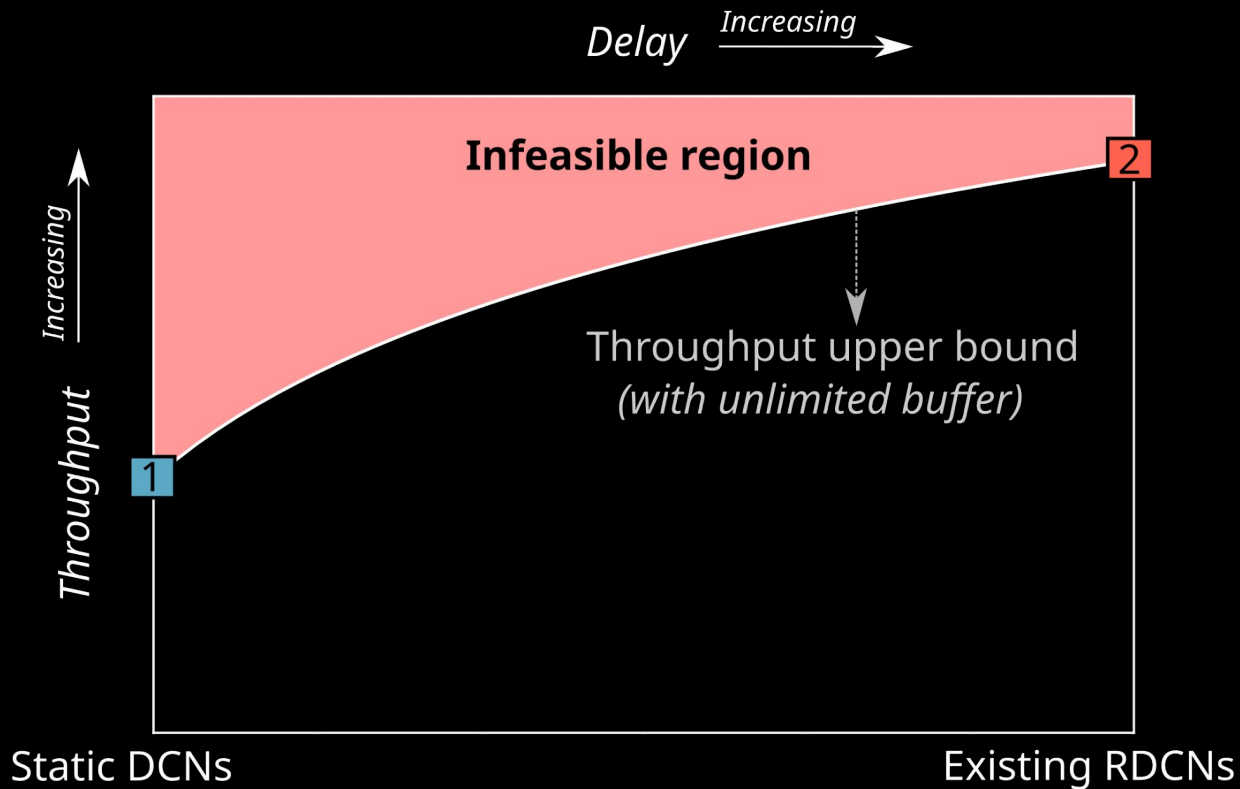
Static DCNs

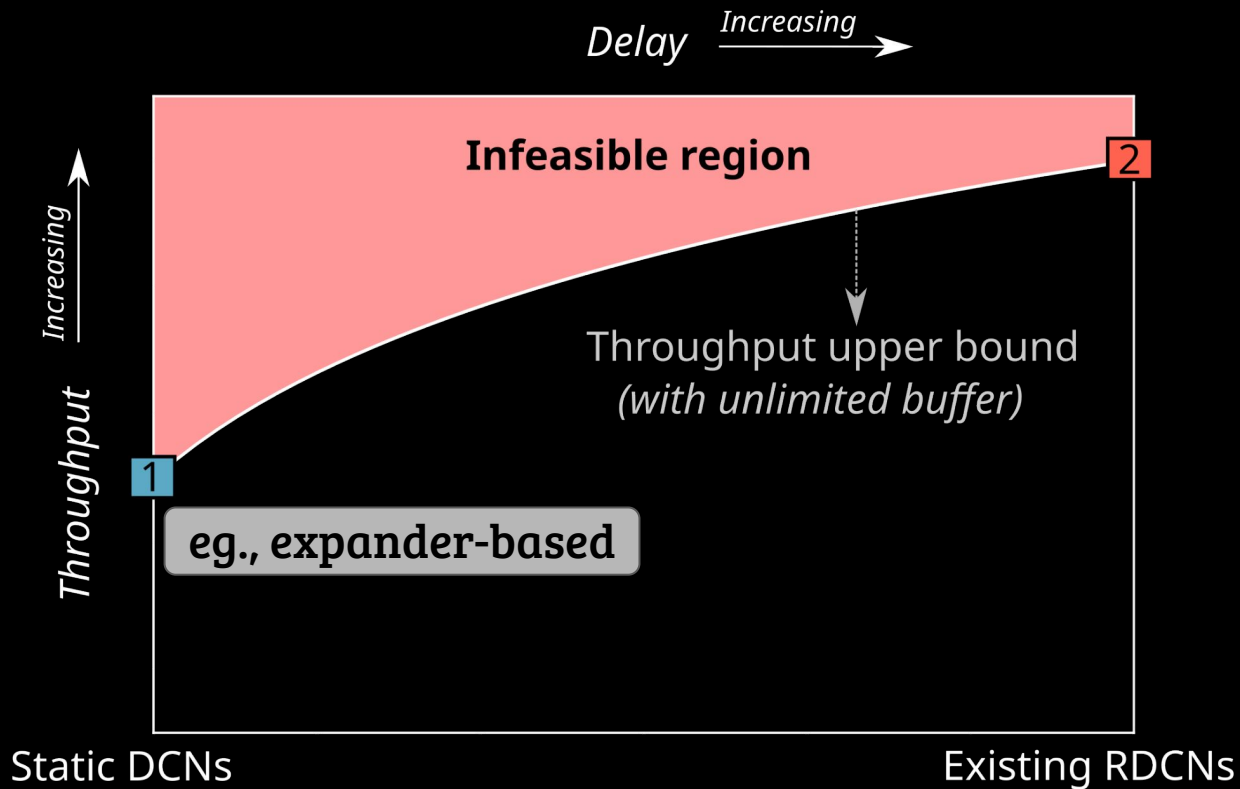
Existing RDCNs

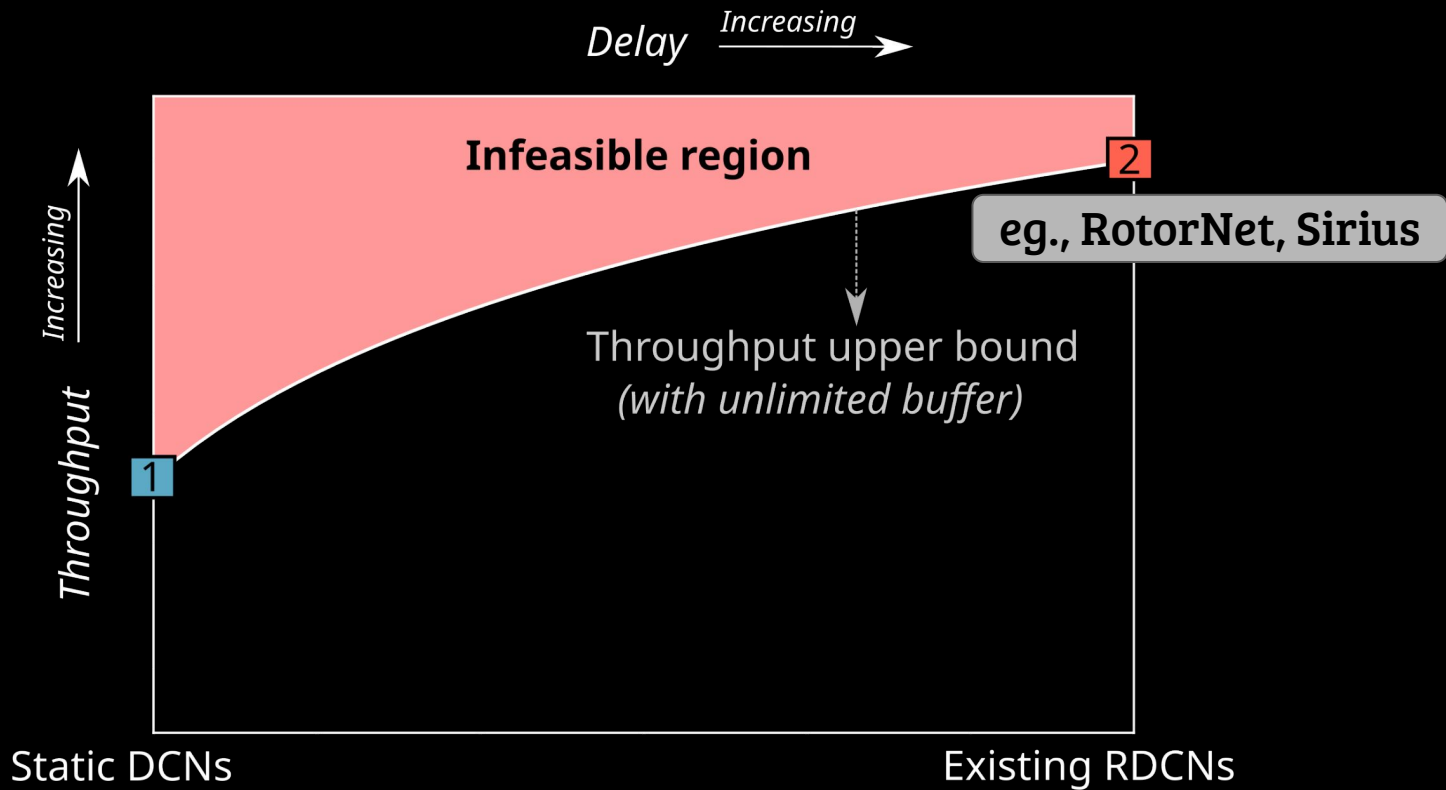
Topology Over Time



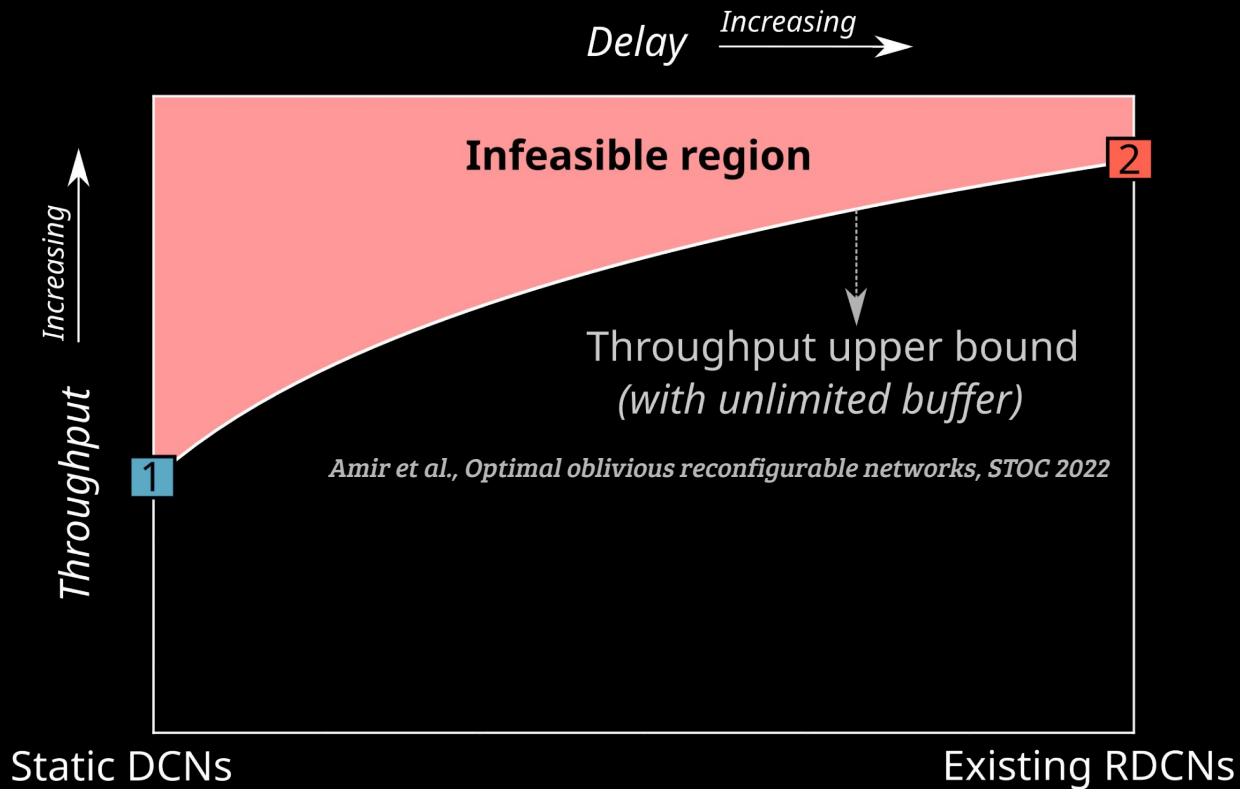




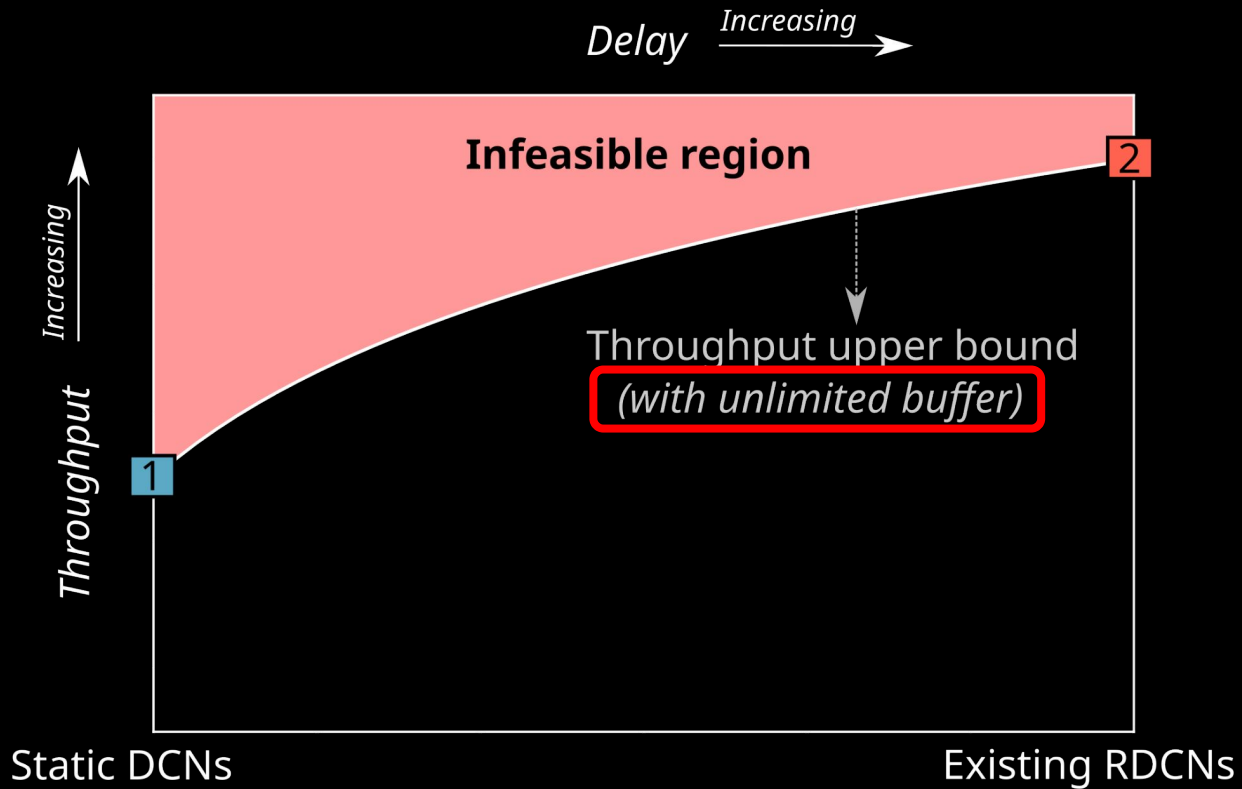


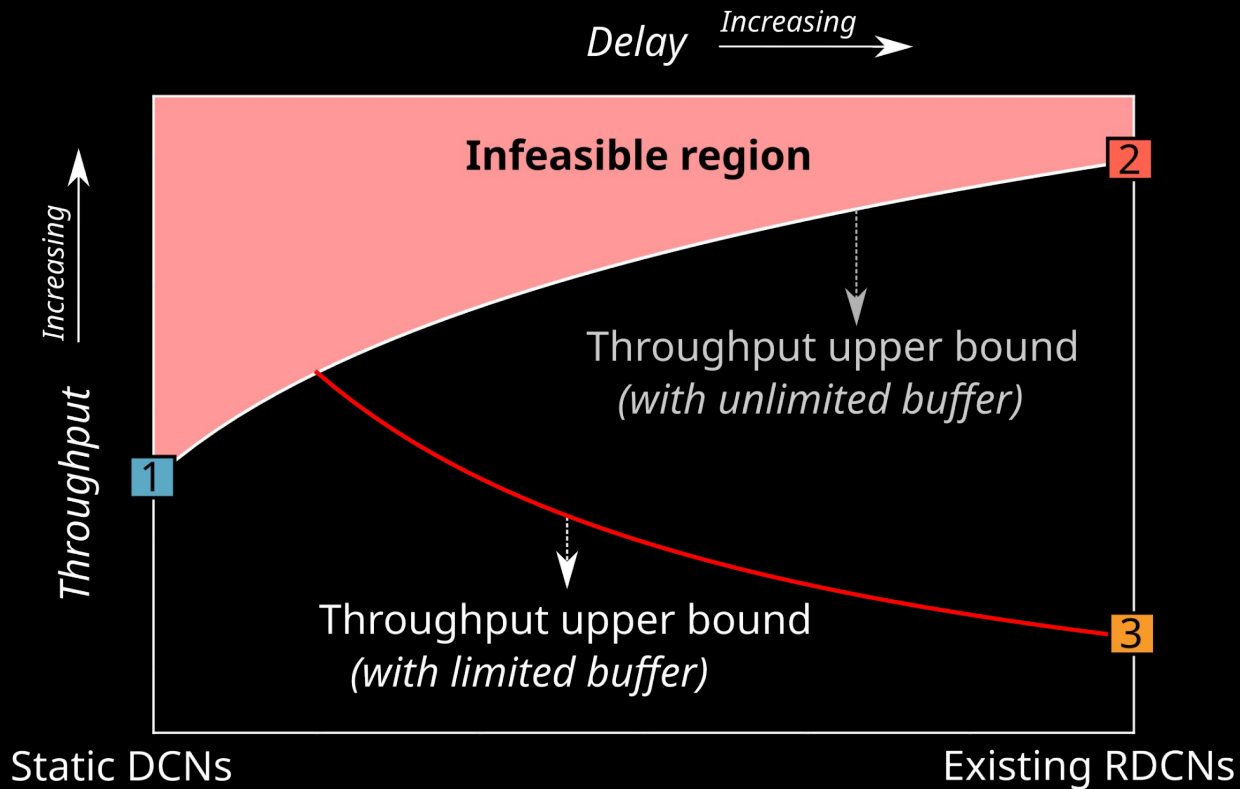


Topology Over Time →

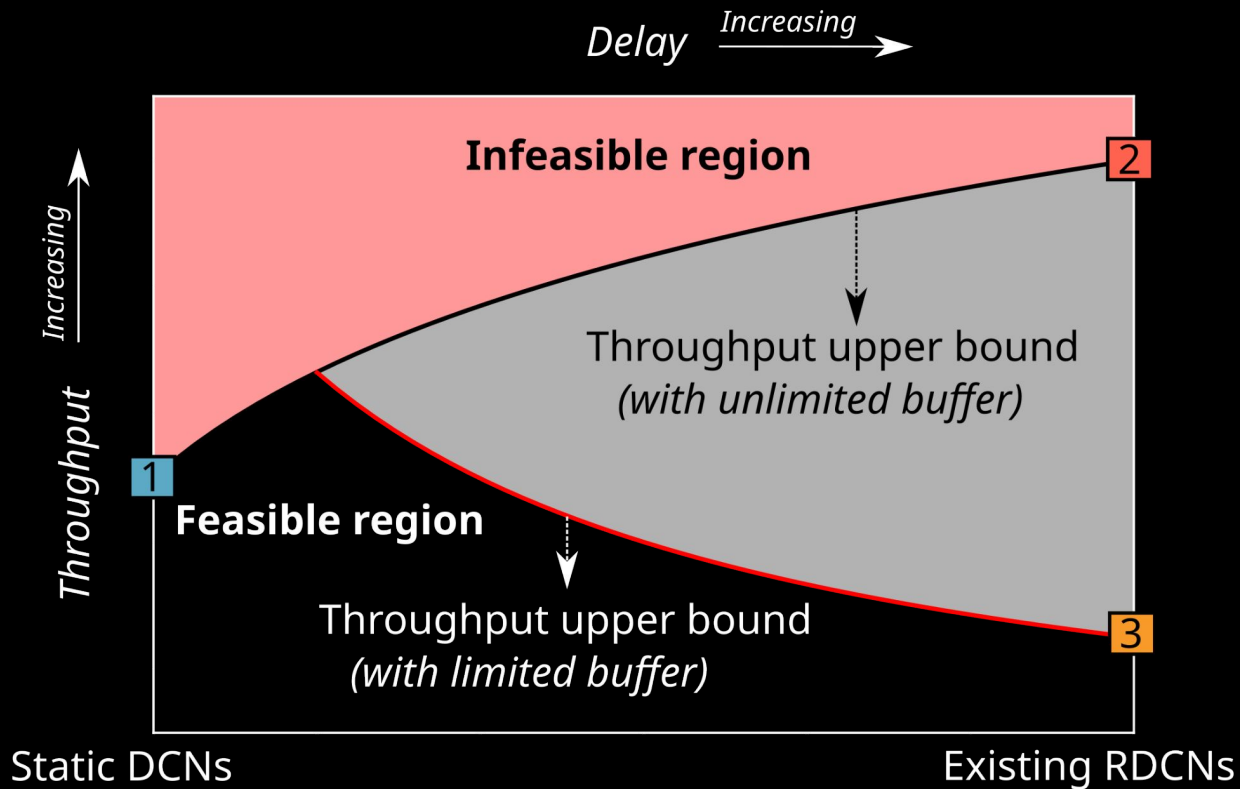


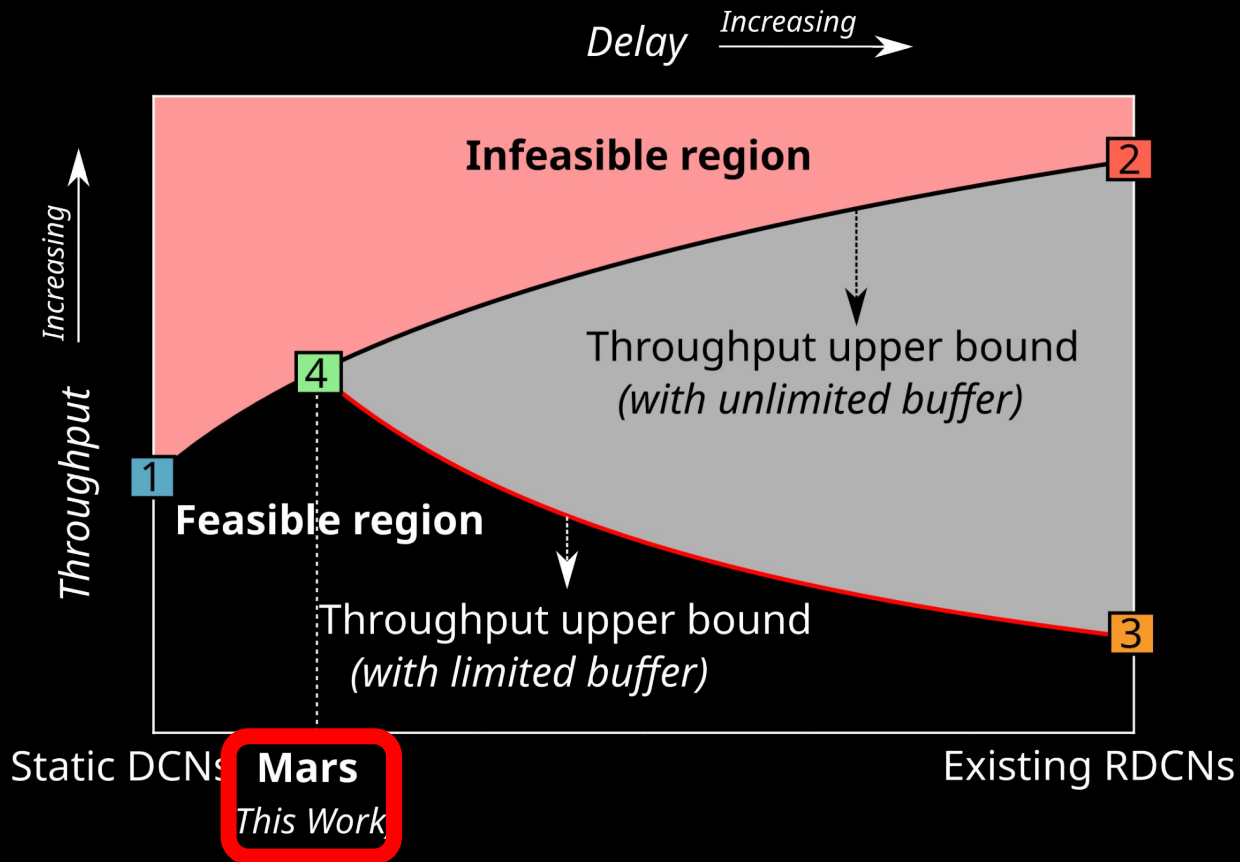
Topology Over Time →

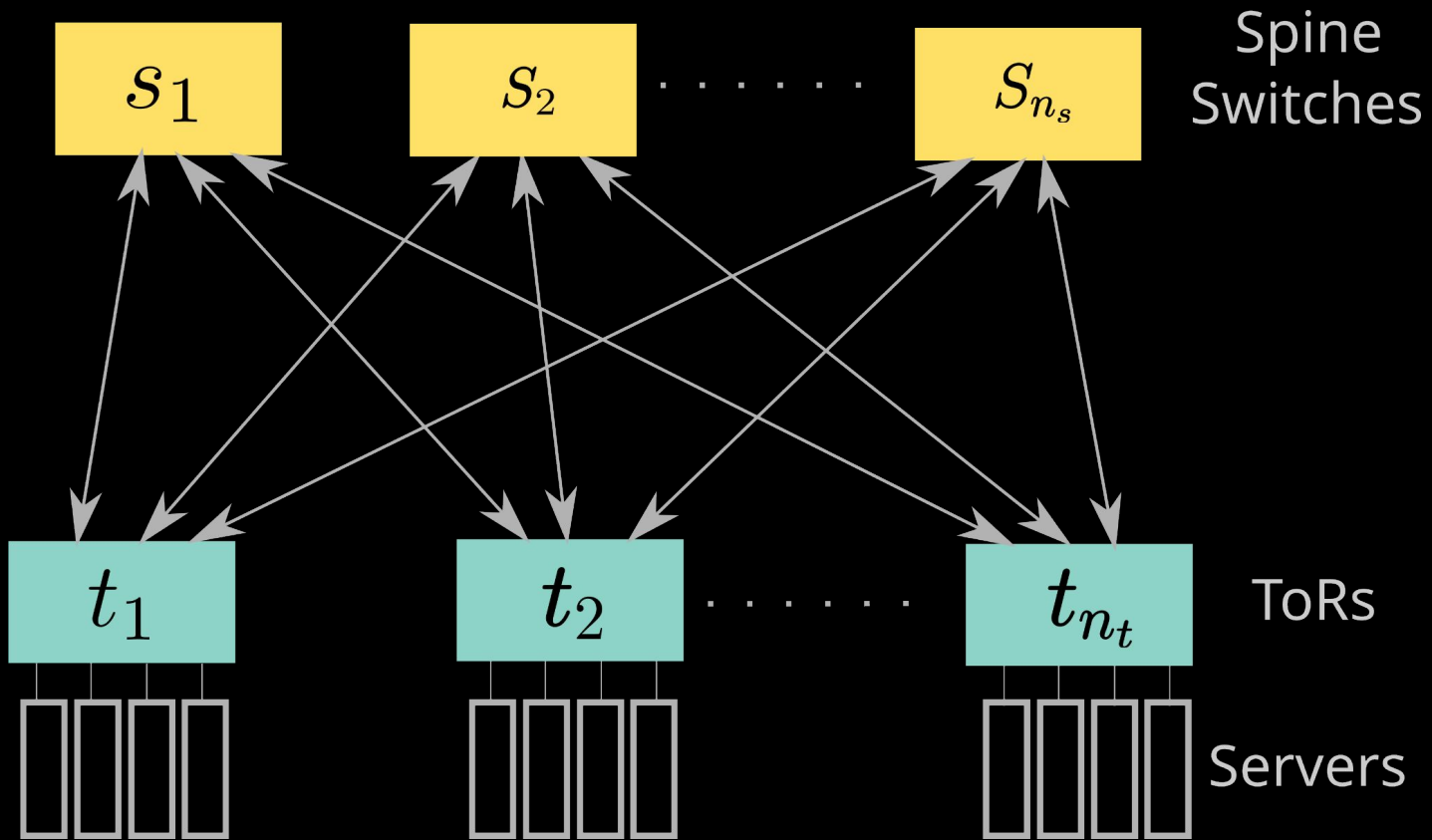




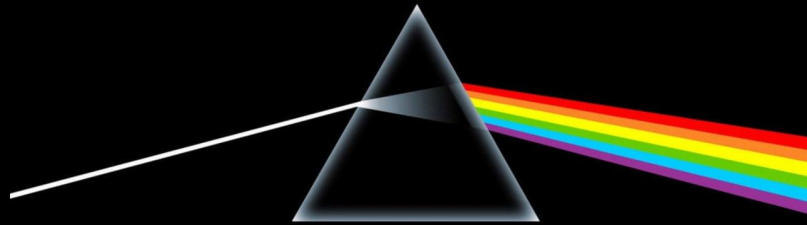
Topology Over Time



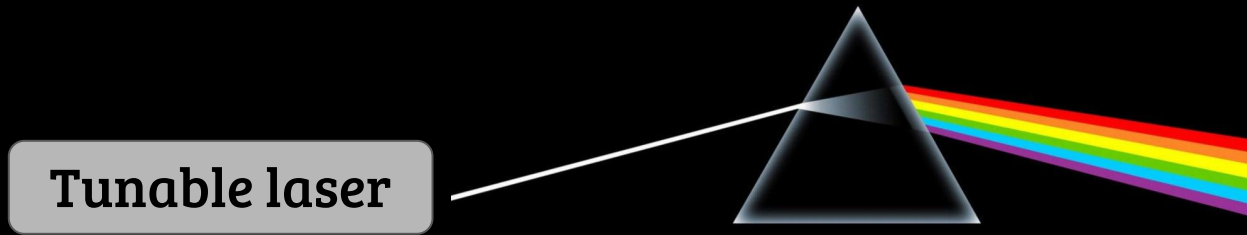




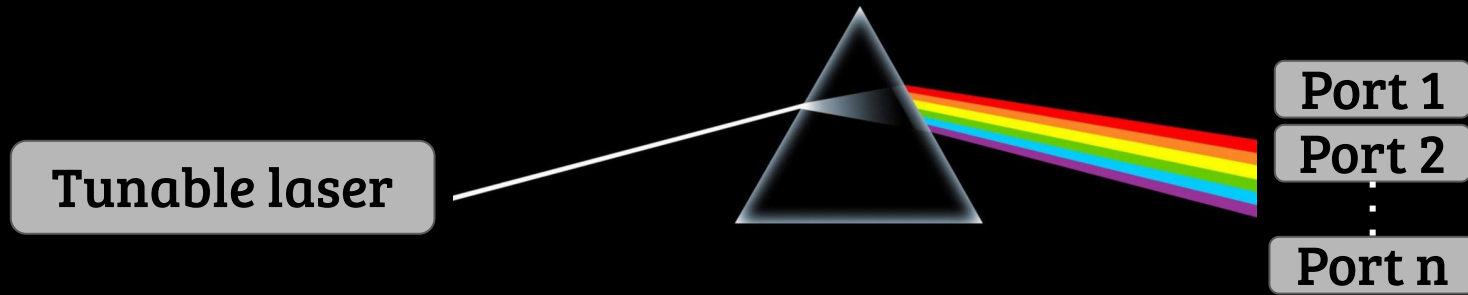
Emerging Technologies - Optical Networks on the Rise



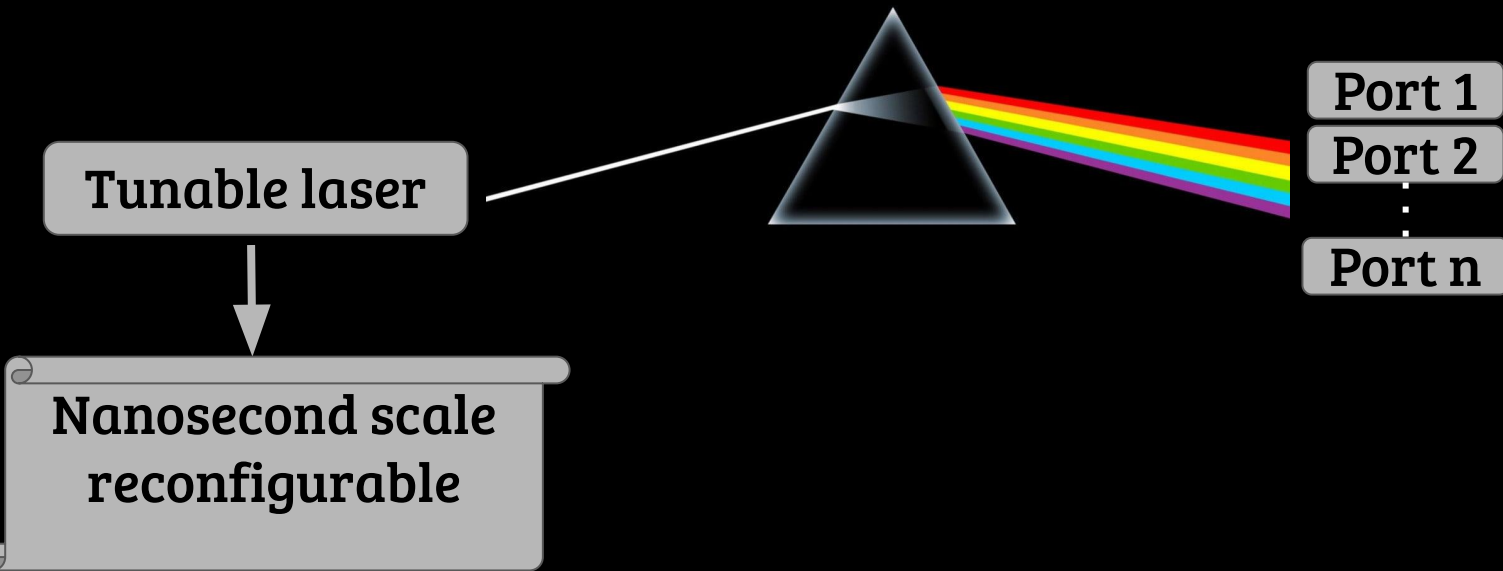
Emerging Technologies - Optical Networks on the Rise



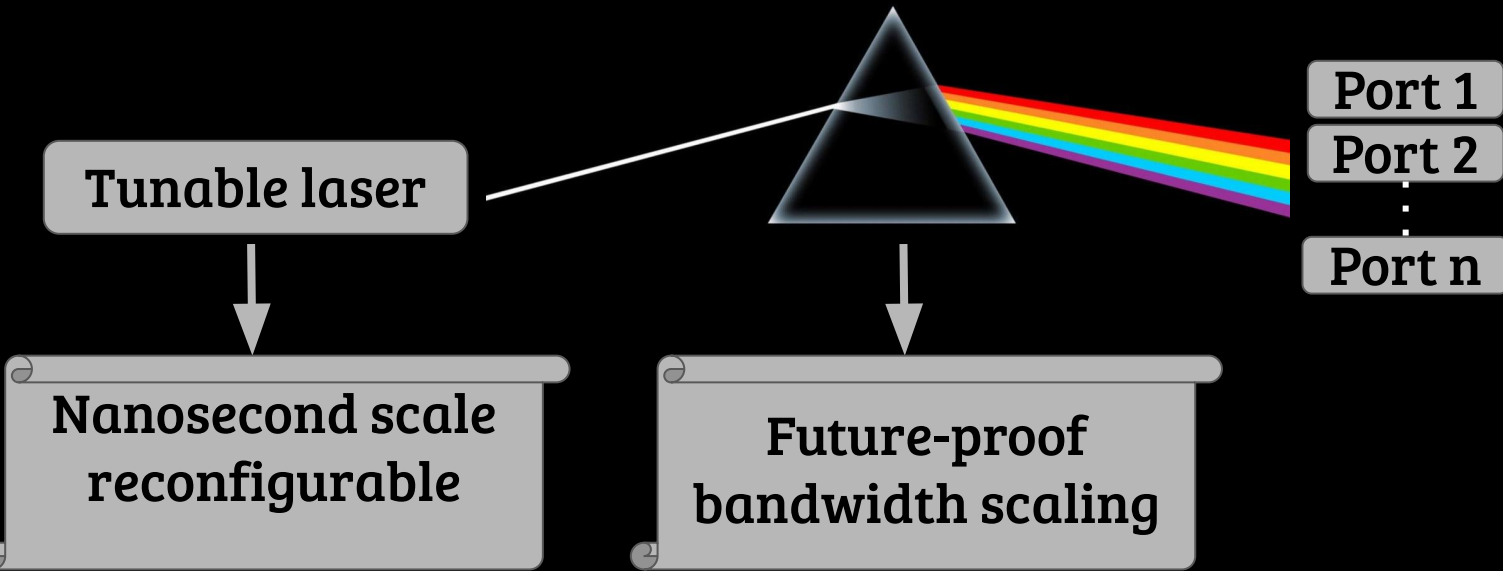
Emerging Technologies - Optical Networks on the Rise



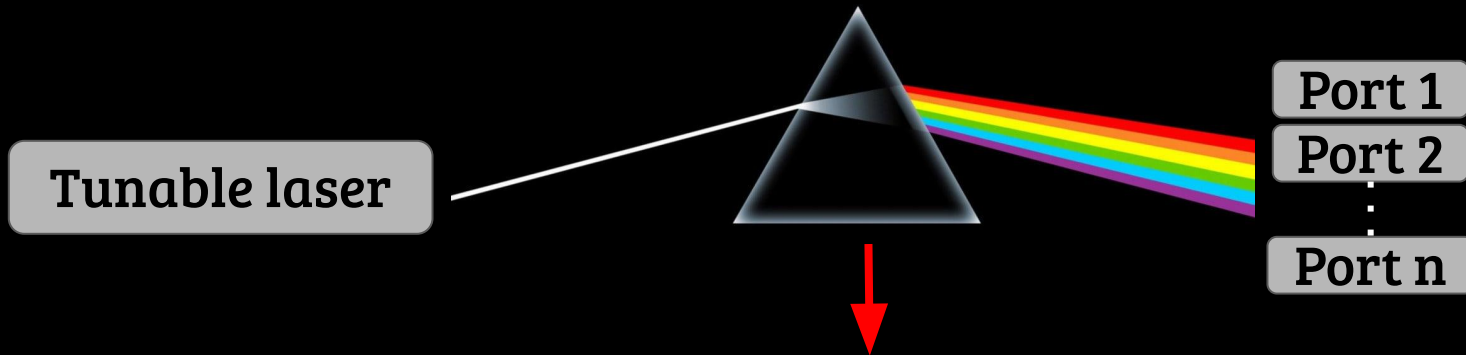
Emerging Technologies - Optical Networks on the Rise



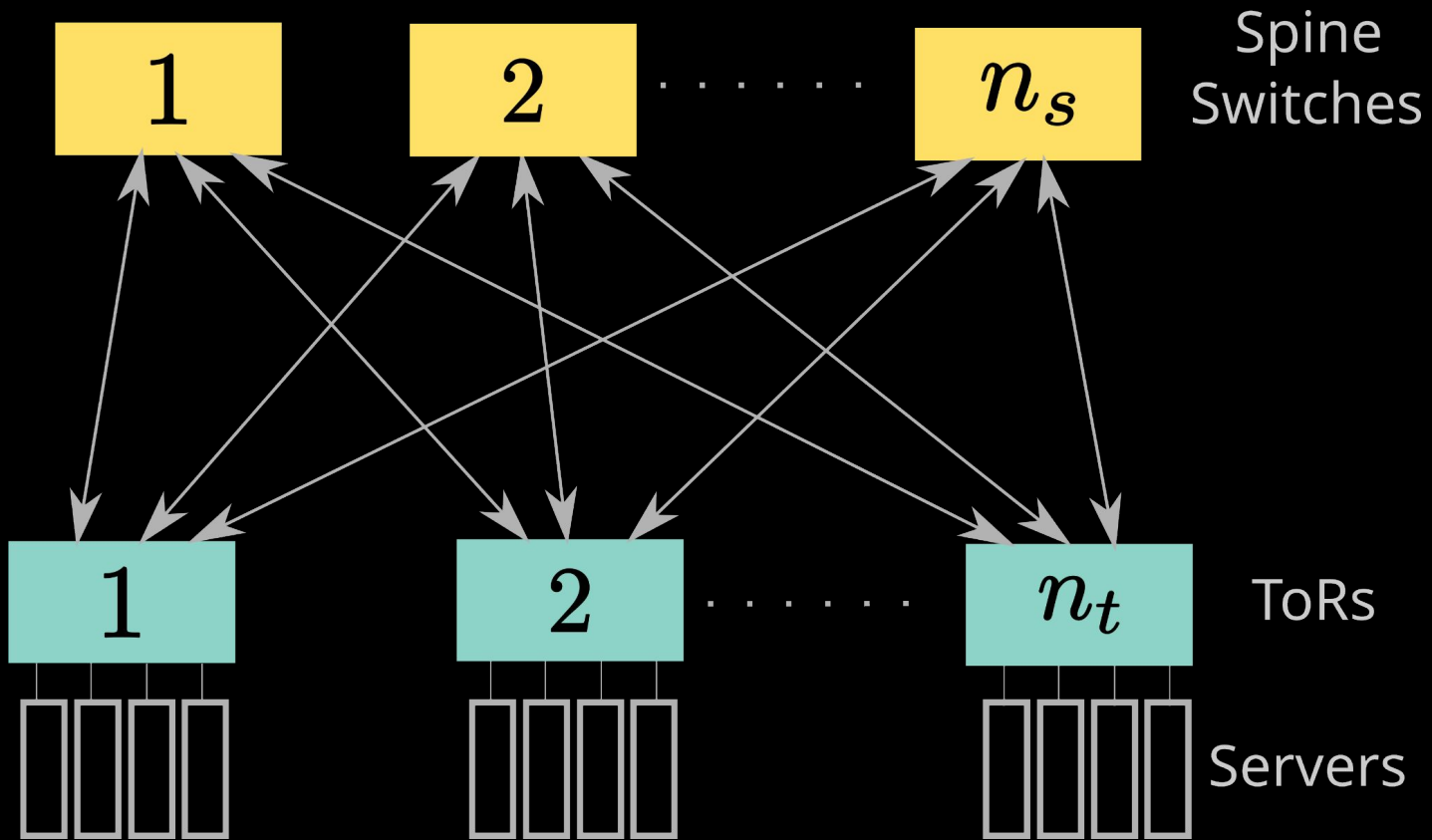
Emerging Technologies - Optical Networks on the Rise

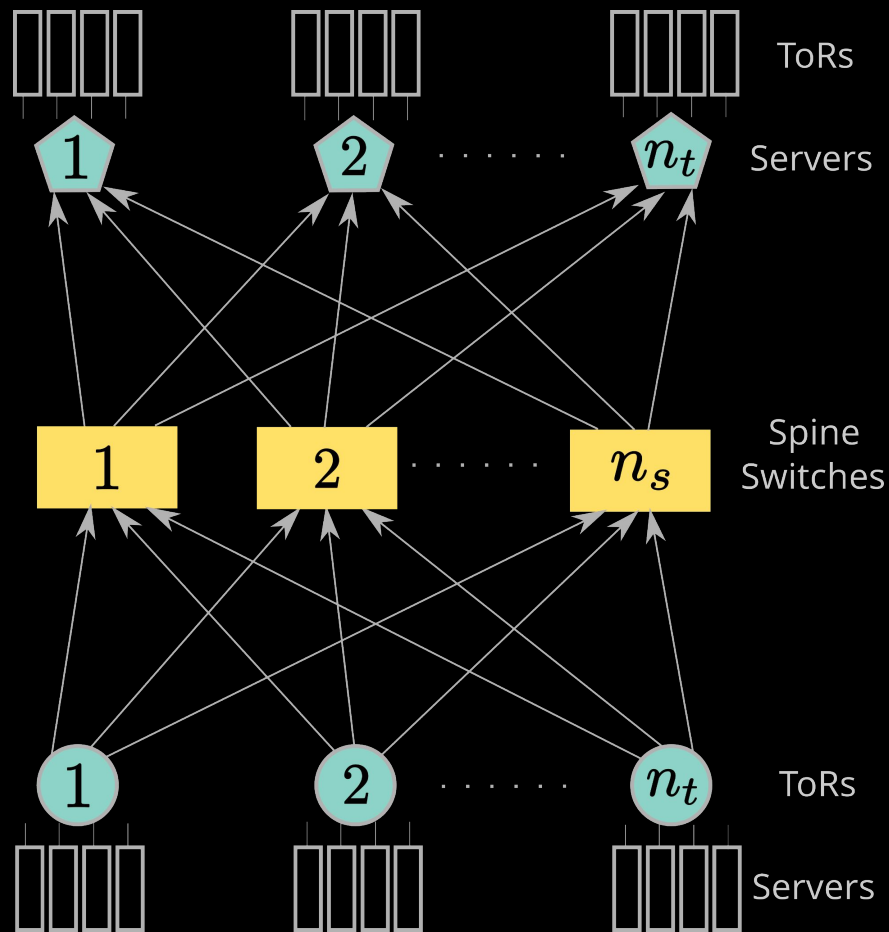
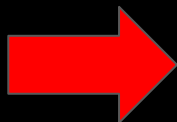
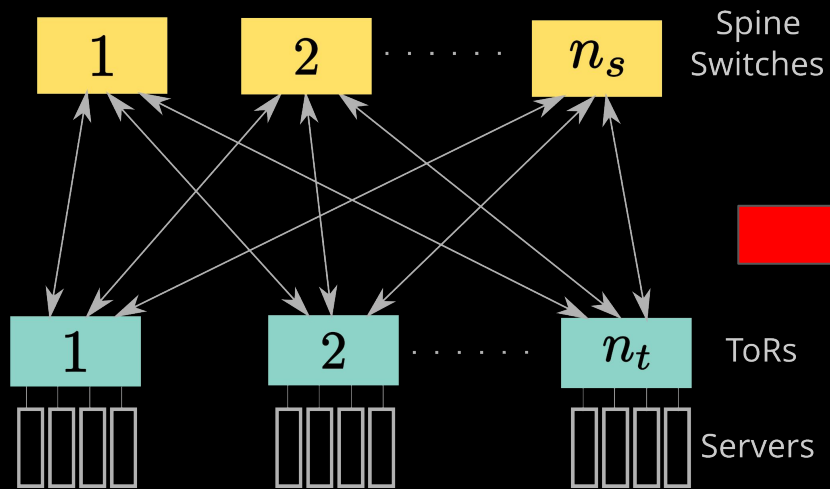


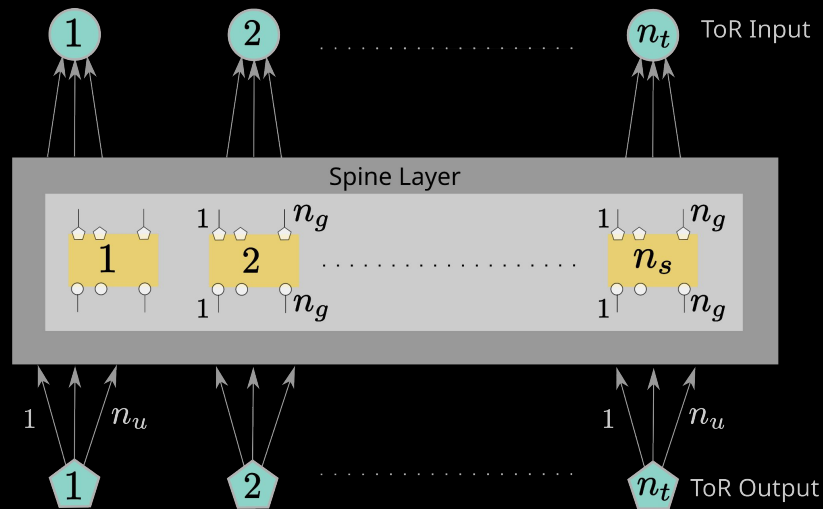
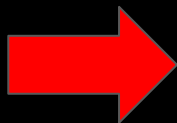
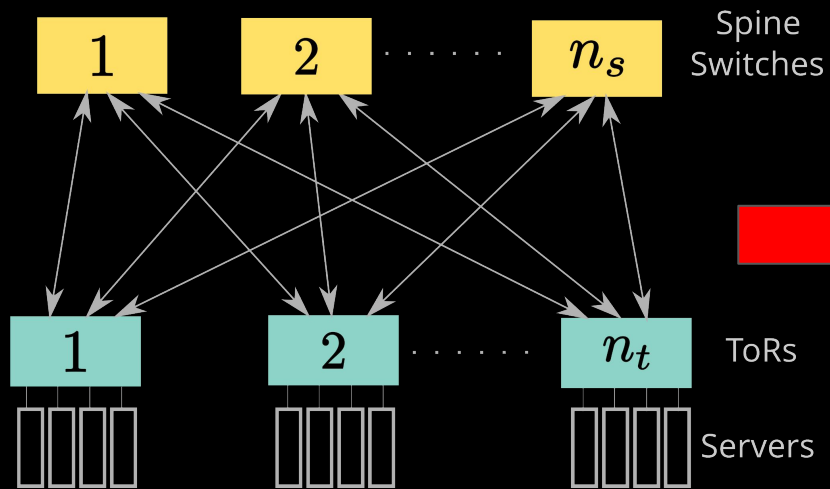
Emerging Technologies - Optical Networks on the Rise

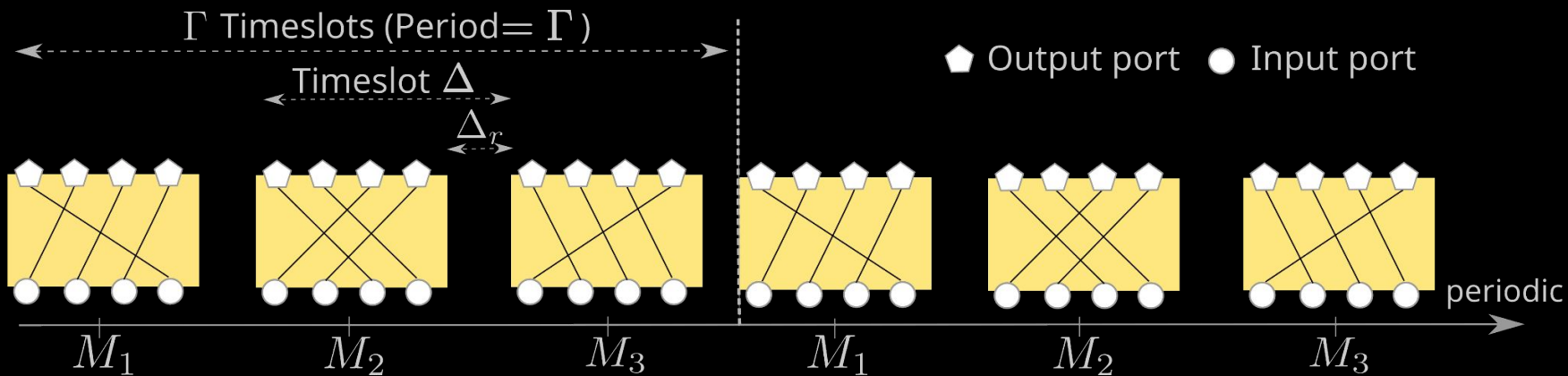


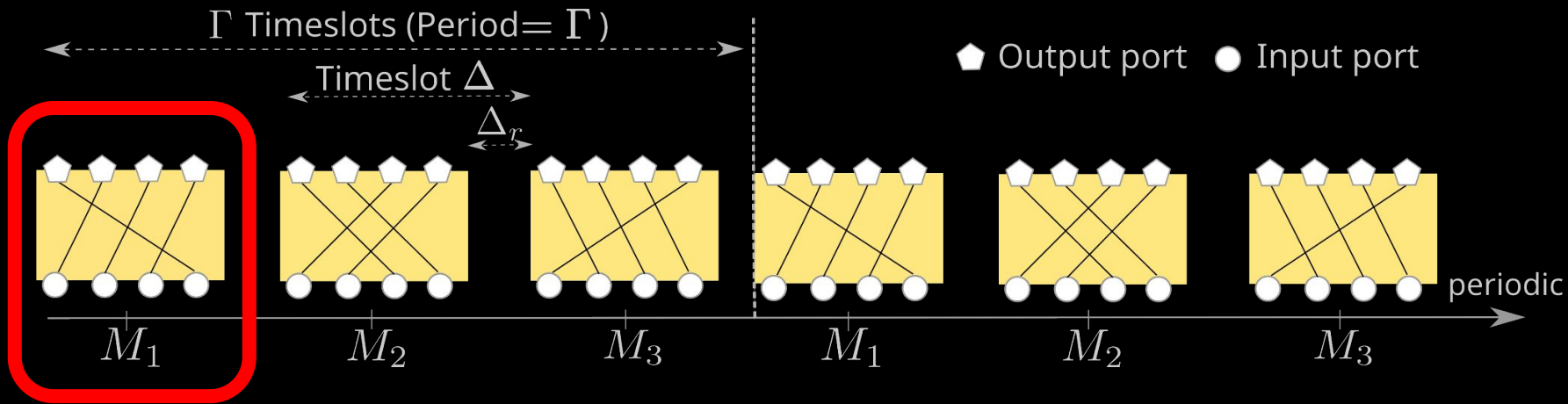
There is a Catch: Bufferless & Circuit Switched

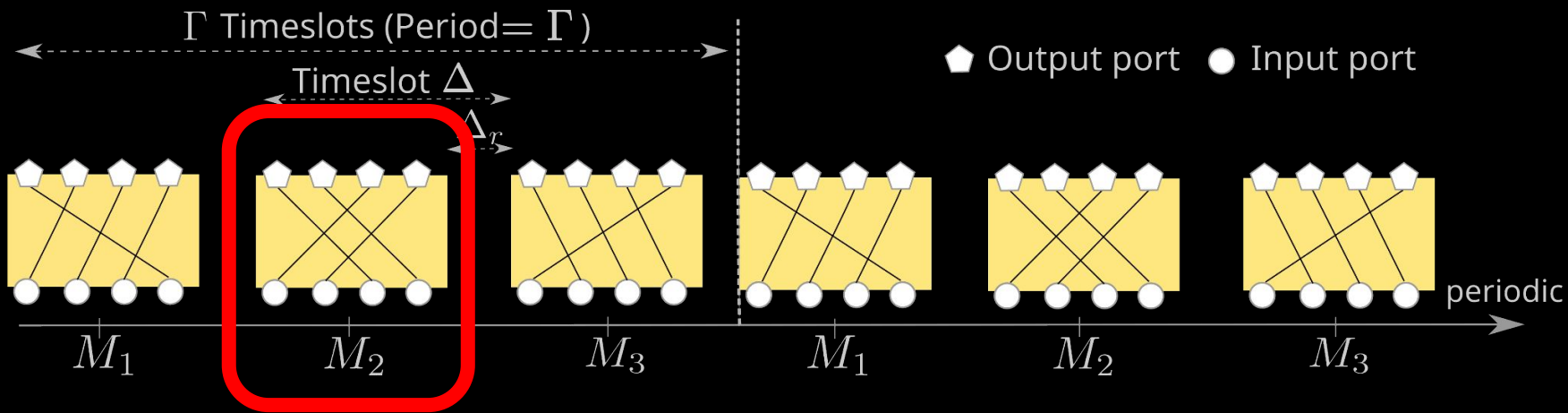


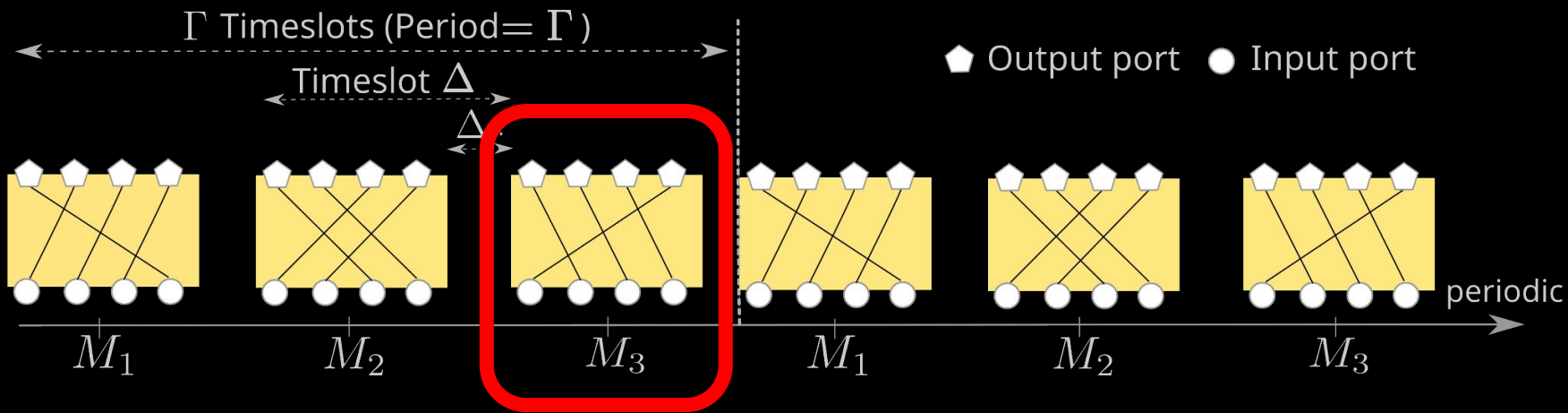


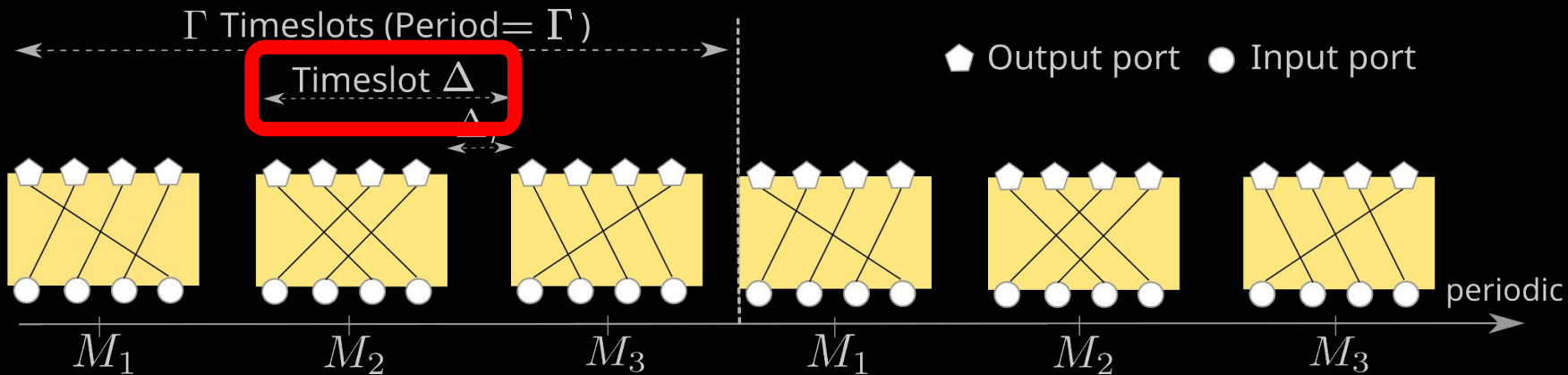


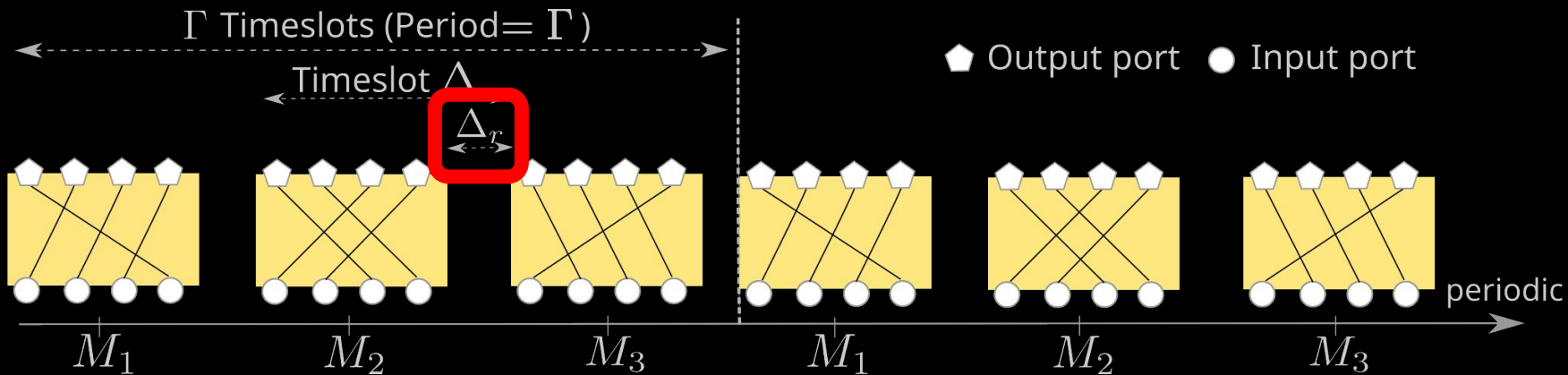


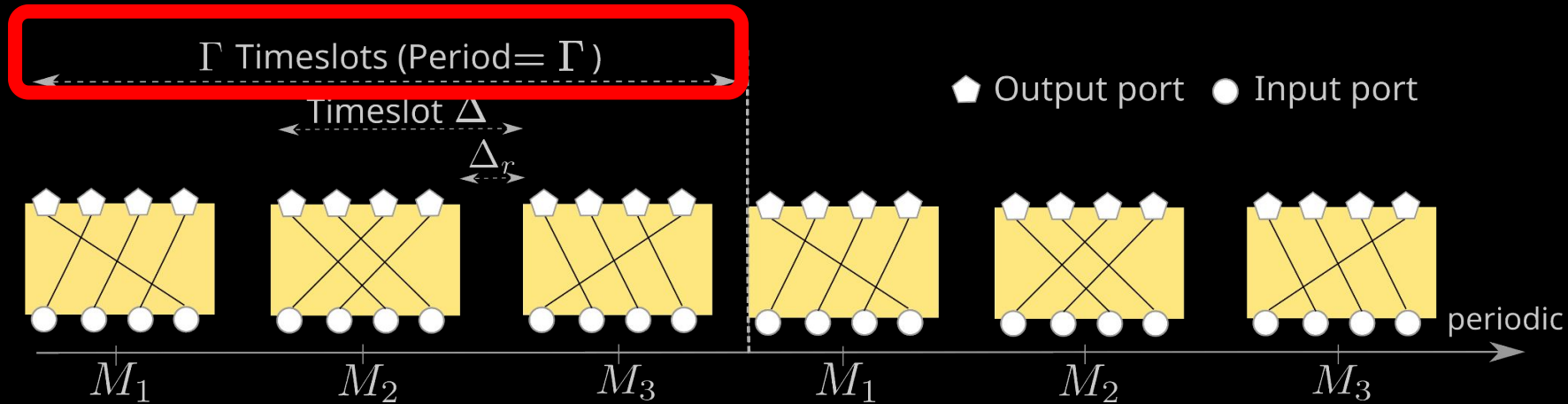


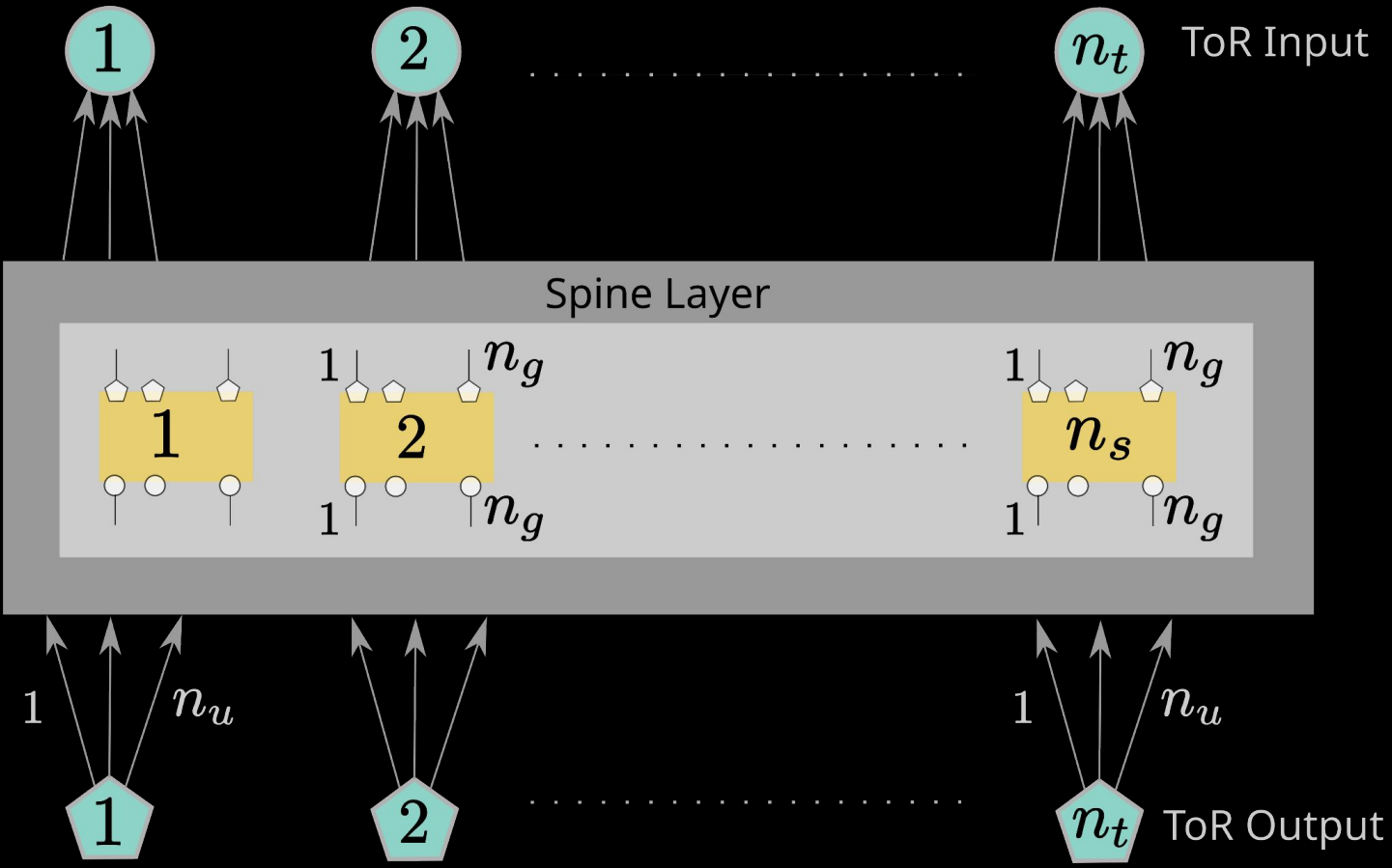




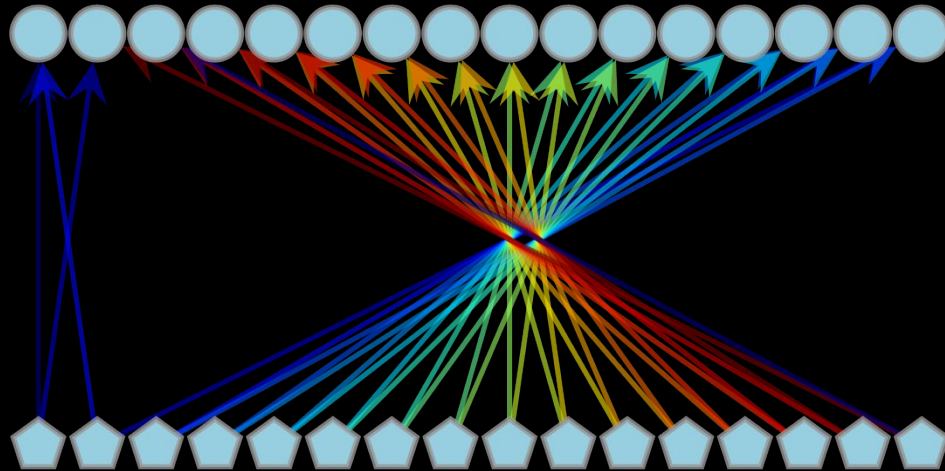




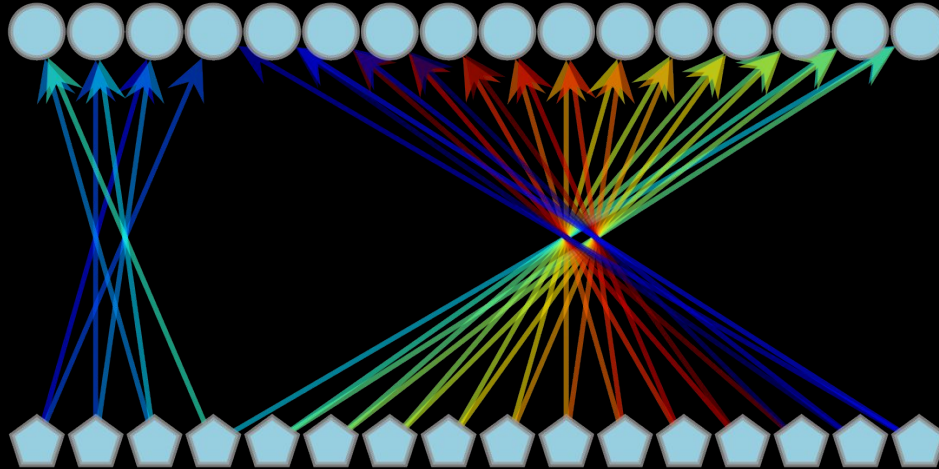




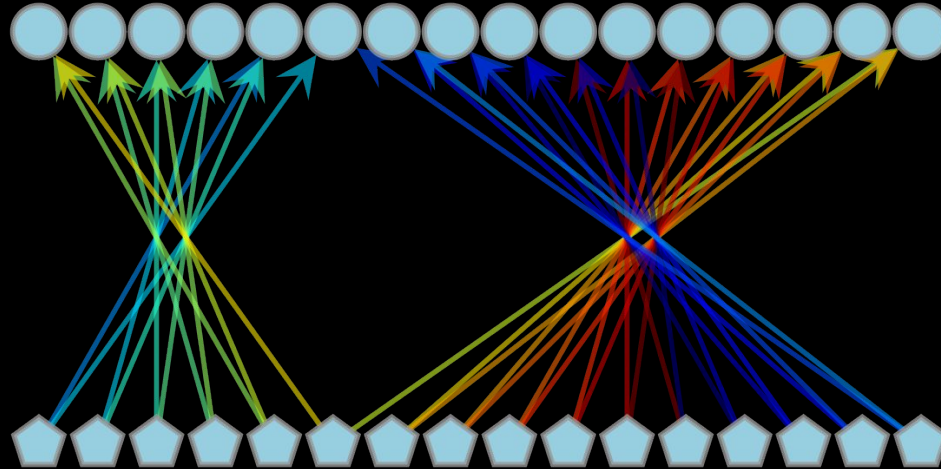
Timeslot 1



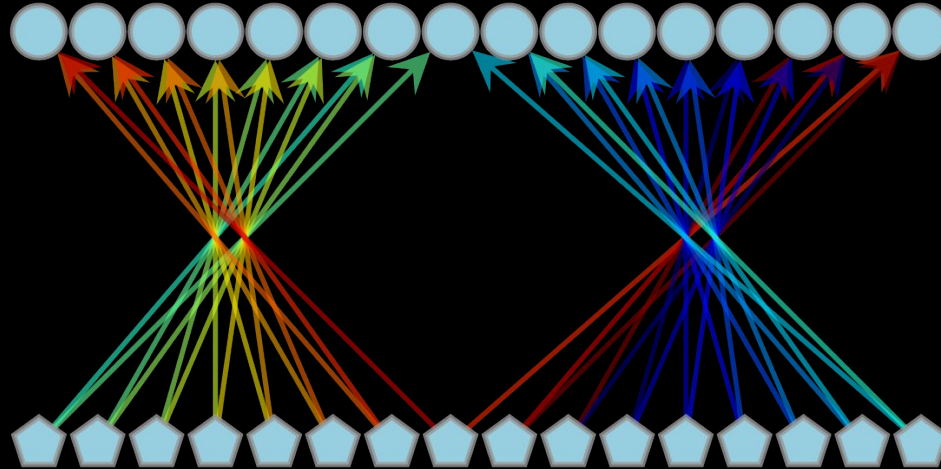
Timeslot 2



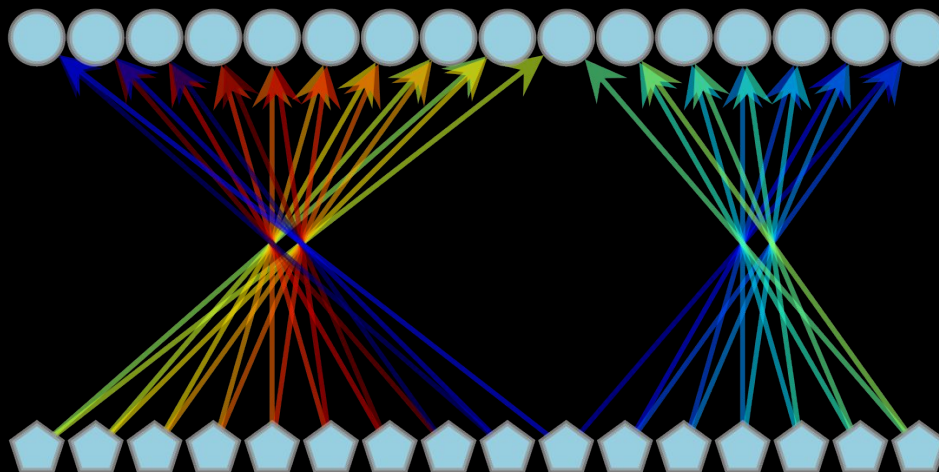
Timeslot 3



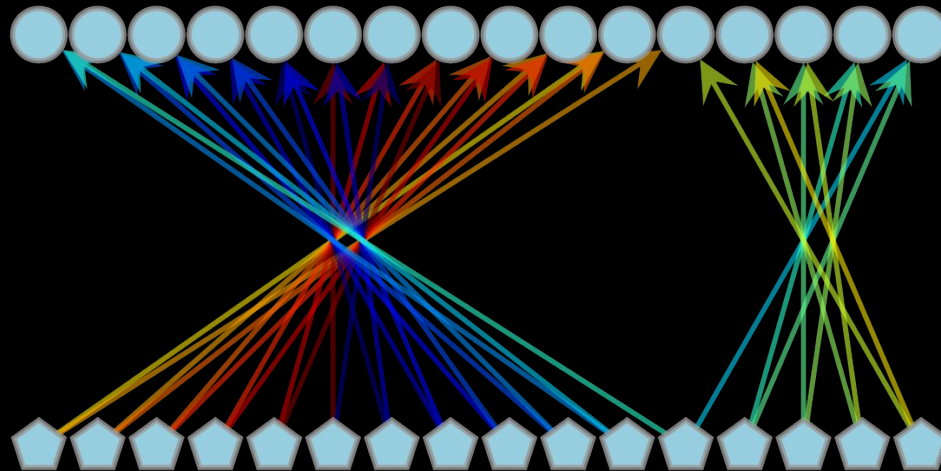
Timeslot 4



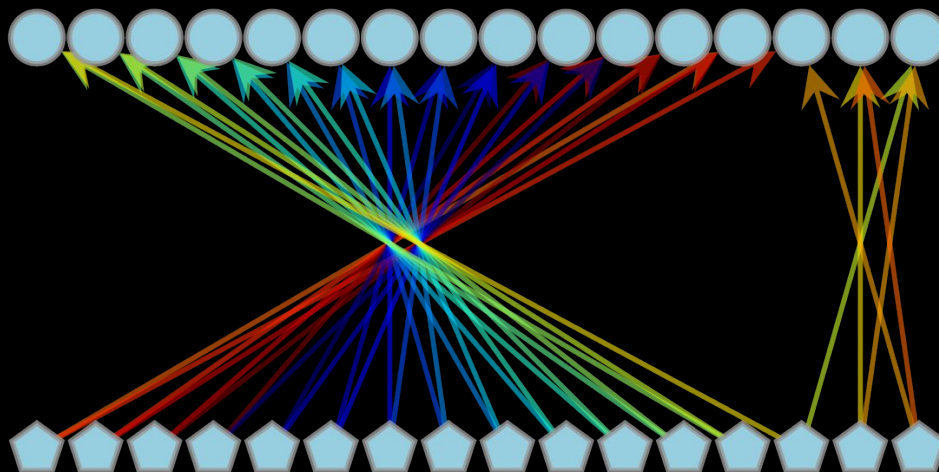
Timeslot 5



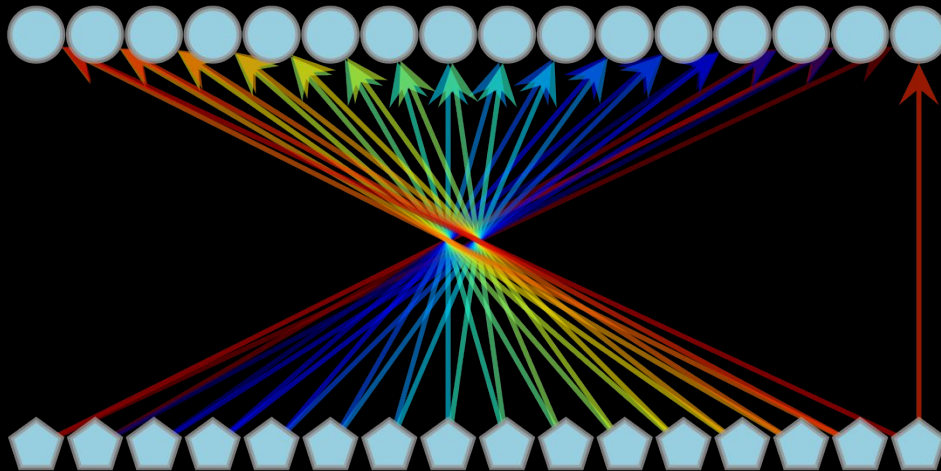
Timeslot 6



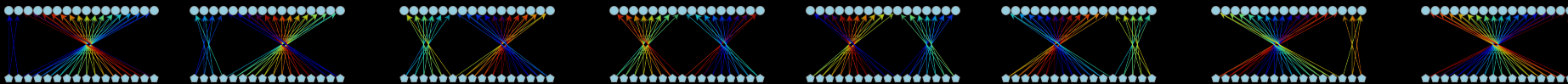
Timeslot 7



Timeslot 8



Periodic Graph



Optimal Topology

Input: Demand Matrix \mathcal{M}

Available buffer size

Output: Periodic graph

Maximize: Throughput

Throughput

	1	2	3	4	5	6	7	8
1								
2								
3								
4								
5								
6								
7								
8								

Demand Matrix

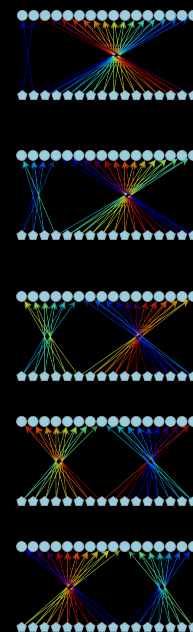
\mathcal{M}

\times

$\theta(\mathcal{M})$



Periodic Graph



t

Throughput

*Highest scaling factor such
that the scaled demand θ
(\mathcal{M}) is feasible in the
periodic graph*

Throughput of the Periodic Graph

Input: Periodic Graph \mathcal{G}

Demand Matrix \mathcal{M}

Objective: Maximize $\theta(\mathcal{M})$

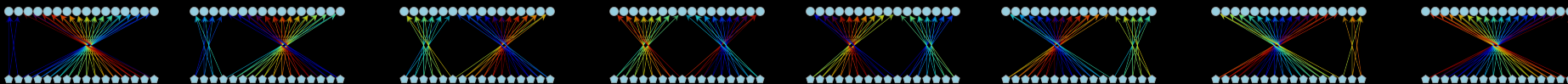
Output: $\theta(\mathcal{M})$ and a feasible *flow**

**subject to conservation, demand and capacity constraints*

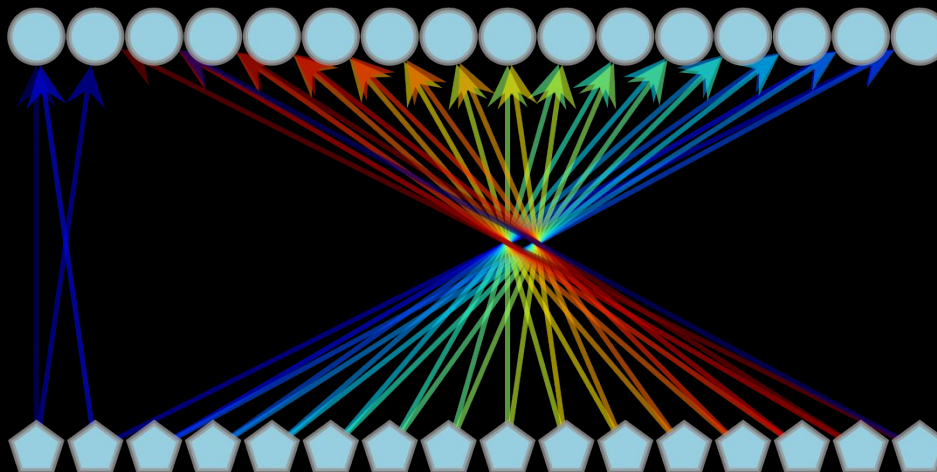
Theorem 1: Periodic Graph \longleftrightarrow Static Graph

- The periodic graph has the **same throughput** as that of a static graph it emulates - *Static Emulated Graph*

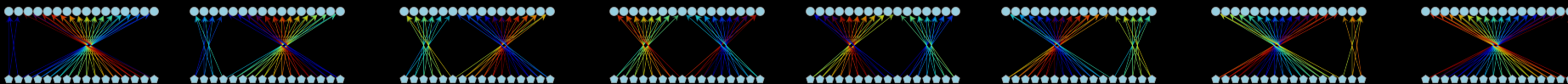
Periodic Graph



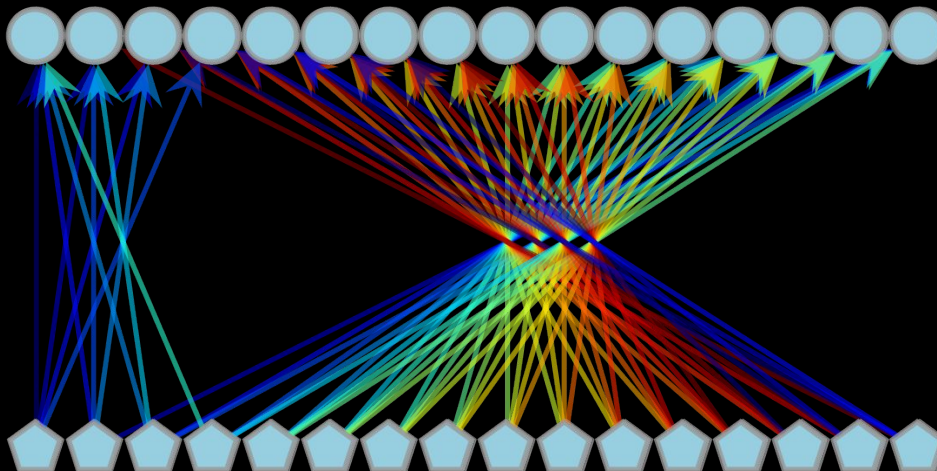
Static Emulated Graph



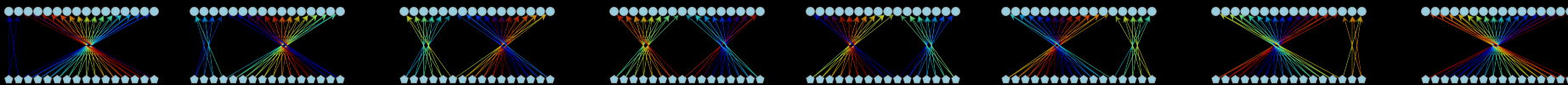
Periodic Graph



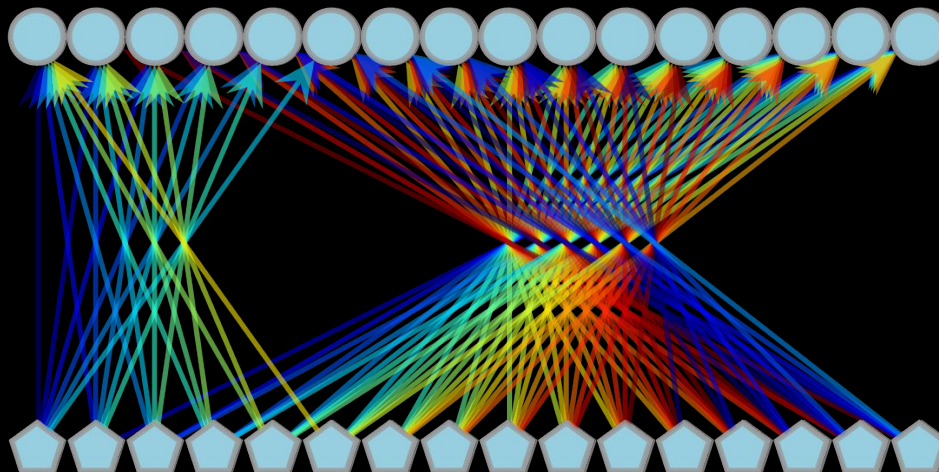
Static Emulated Graph



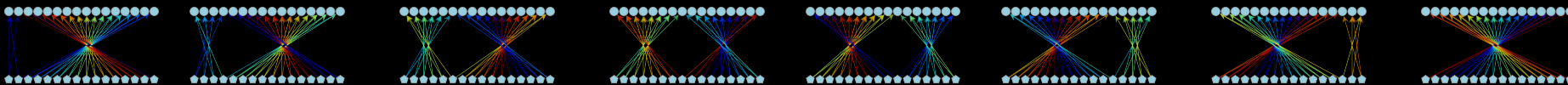
Periodic Graph



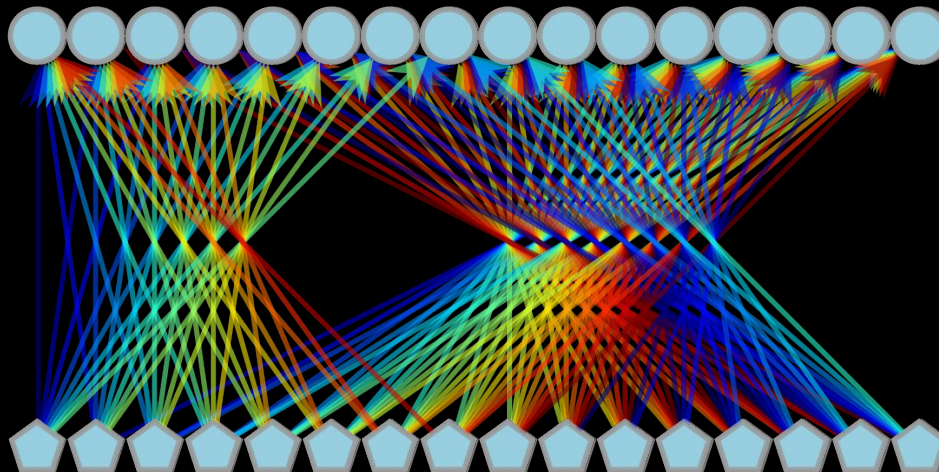
Static Emulated Graph



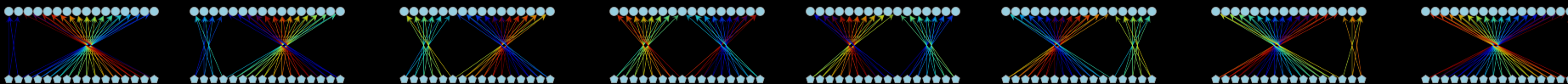
Periodic Graph



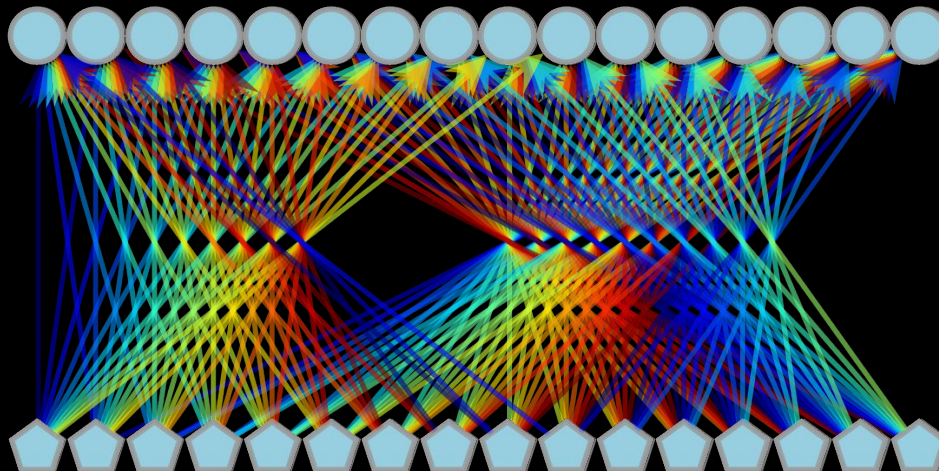
Static Emulated Graph



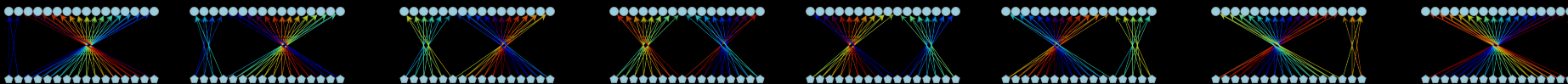
Periodic Graph



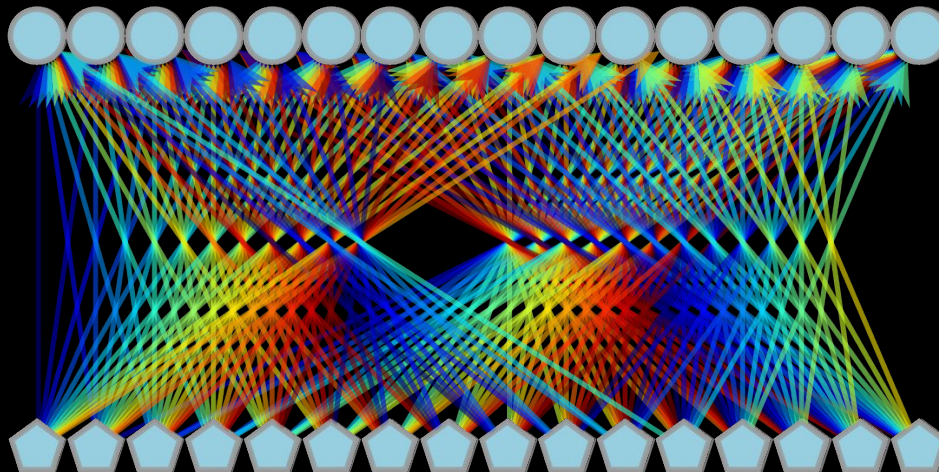
Static Emulated Graph



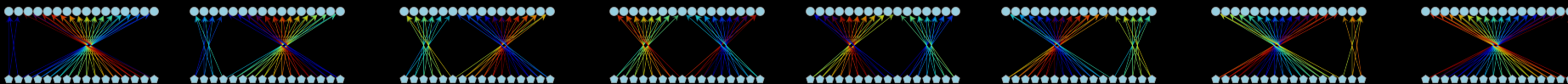
Periodic Graph



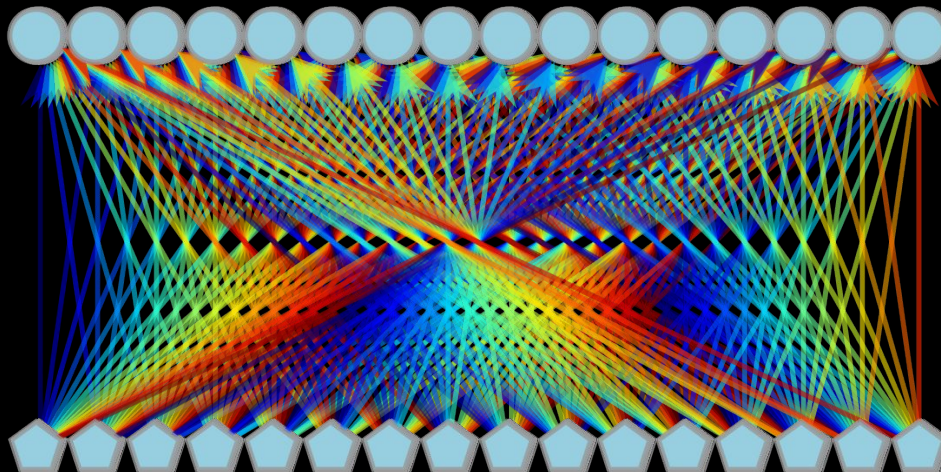
Static Emulated Graph



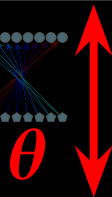
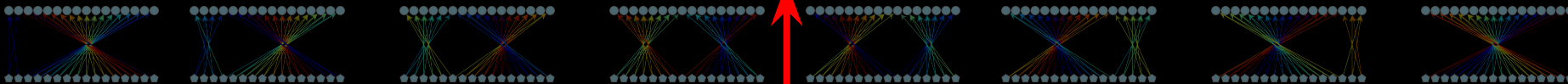
Periodic Graph



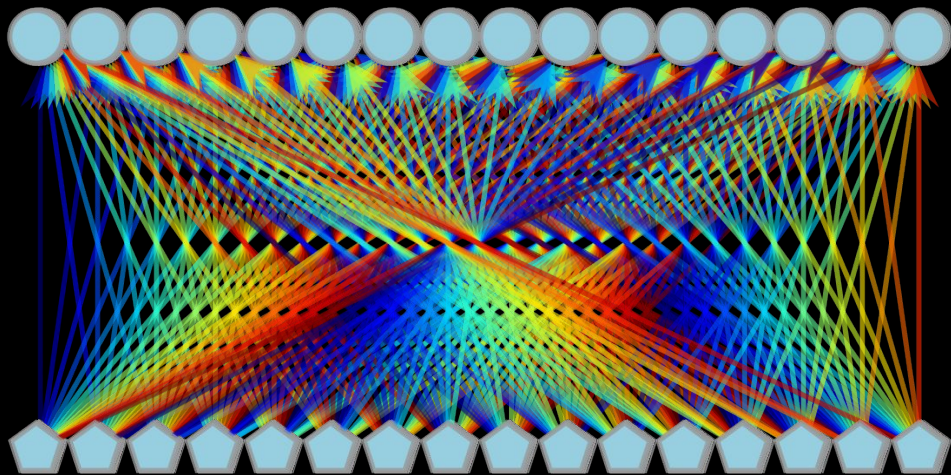
Static Emulated Graph



Periodic Graph



Static Emulated Graph



Throughput of the Periodic Graph

Input: Periodic Graph \mathcal{G}

Demand Matrix \mathcal{M}

Objective: Maximize $\theta(\mathcal{M})$

Output: $\theta(\mathcal{M})$ and a feasible *flow**

**subject to conservation, demand and capacity constraints*

Throughput of the Periodic Graph

Input: ~~Periodic Graph~~ Static Emulated Graph G

Demand Matrix \mathcal{M}

Objective: Maximize $\theta(\mathcal{M})$

Output: $\theta(\mathcal{M})$ and a feasible *flow**

**subject to conservation, demand and capacity constraints*



Uni-regular
(Static DCNs)

Complete Graph
(Existing RDCNs)

Degree $\xrightarrow{\text{Increasing}}$

Emulated Graph (Topology Over Time)

Theorem 2: Throughput Upper Bound

- Throughput is a function of *Average Route Length*

$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \cdot \text{ARL}(\mathcal{M}, F)}$$

Theorem 2: Throughput Upper Bound

- Throughput is a function of *Average Route Length*

$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \cdot \text{ARL}(\mathcal{M}, F)}$$

Theorem 2: Throughput Upper Bound

- Throughput is a function of *Average Route Length*

$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \cdot \text{ARL}(\mathcal{M}, F)}$$

Theorem 2: Throughput Upper Bound

- Throughput is a function of *Average Route Length*

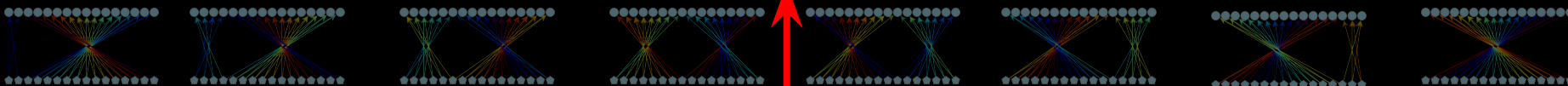
$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \text{ ARL}(\mathcal{M}, F)}$$

Theorem 2: Throughput Upper Bound

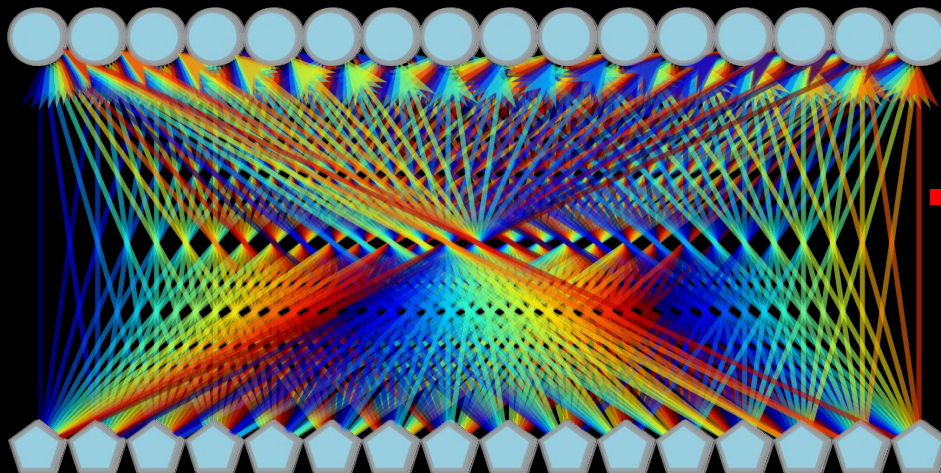
- Throughput is a function of *Average Route Length*

$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \cdot \text{ARL}(\mathcal{M}, F)}$$

Periodic Graph

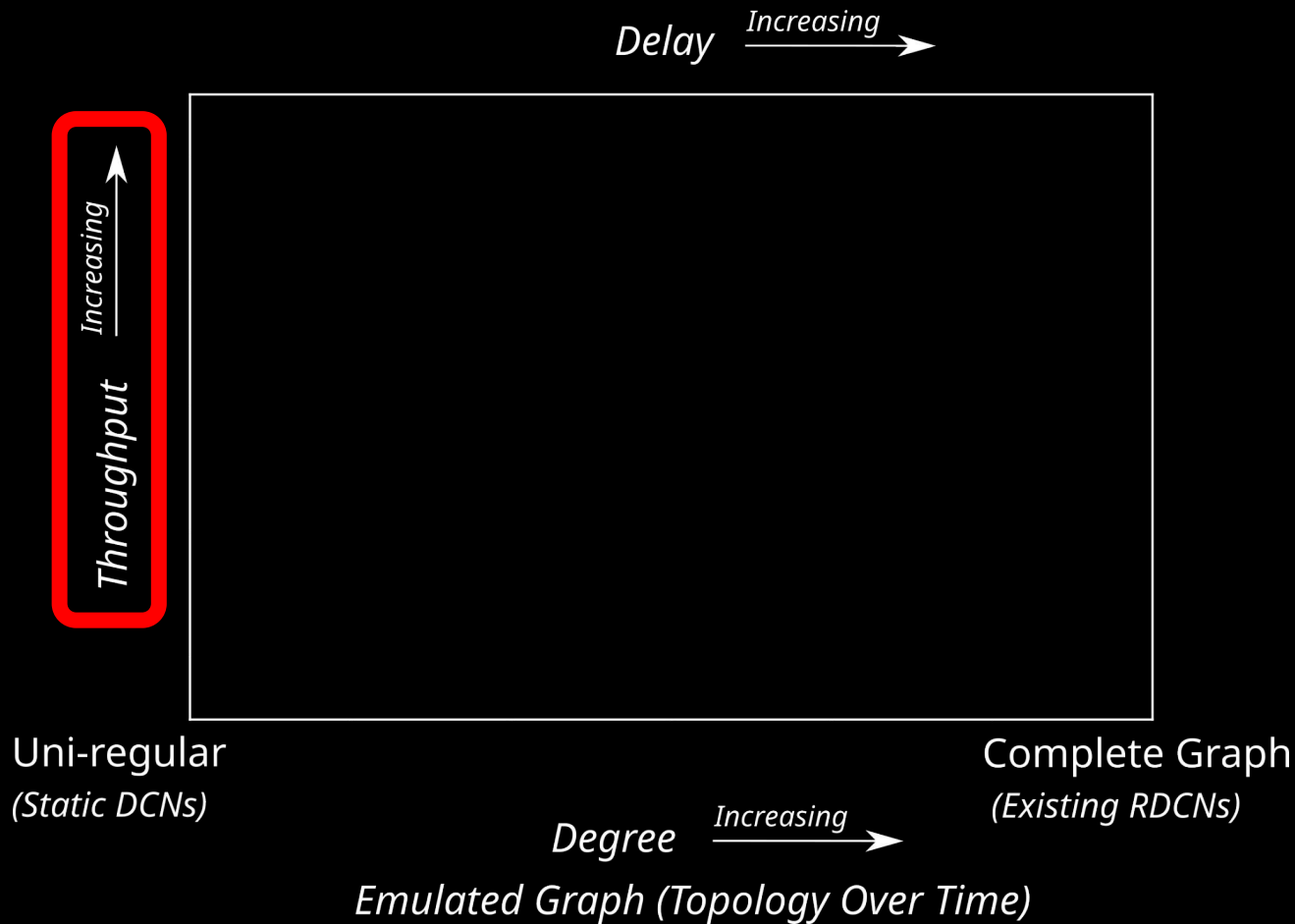


Static Emulated Graph



$$\theta(\mathcal{M}, F) \leq \frac{\hat{C}}{M \cdot \text{ARL}(\mathcal{M}, F)}$$

Capacity
Demand x ARL



Delay

Input: Periodic Graph \mathcal{G}
Demand Matrix \mathcal{M}
Throughput θ and a feasible flow \mathcal{F}

Output: Minimum worst-case delay

Theorem 3: Delay

- Delay bound is a function of:
 - Degree d of the emulated graph
 - Duration of the period $\Gamma \cdot \Delta$
 - Throughput θ

$$\begin{aligned} L_{max} &\geq \text{ARD}(\mathcal{M}, \mathcal{F}) = \text{ARL}(\mathcal{M}, \mathcal{F}) \cdot \Gamma \cdot \Delta \\ &\geq \Omega \left(\frac{d \cdot \Delta}{n_u \cdot \theta(\mathcal{M}, \mathcal{F})} \right) \end{aligned}$$

Theorem 3: Delay

- Delay bound is a function of:
 - Degree d of the emulated graph
 - Duration of the period $\Gamma \cdot \Delta$
 - Throughput θ

$$\boxed{L_{max}} \geq \text{ARD}(\mathcal{M}, \mathcal{F}) = \text{ARL}(\mathcal{M}, \mathcal{F}) \cdot \Gamma \cdot \Delta$$
$$\geq \Omega \left(\frac{d \cdot \Delta}{n_u \cdot \theta(\mathcal{M}, \mathcal{F})} \right)$$

Theorem 3: Delay

- Delay bound is a function of:
 - Degree d of the emulated graph
 - Duration of the period $\Gamma \cdot \Delta$
 - Throughput θ

$$\begin{aligned} L_{max} &\geq \text{ARD}(\mathcal{M}, \mathcal{F}) = \text{ARL}(\mathcal{M}, \mathcal{F}) \cdot \Gamma \cdot \Delta \\ &\geq \Omega \left(\frac{d \cdot \Delta}{n_u \cdot \theta(\mathcal{M}, \mathcal{F})} \right) \end{aligned}$$

Theorem 3: Delay

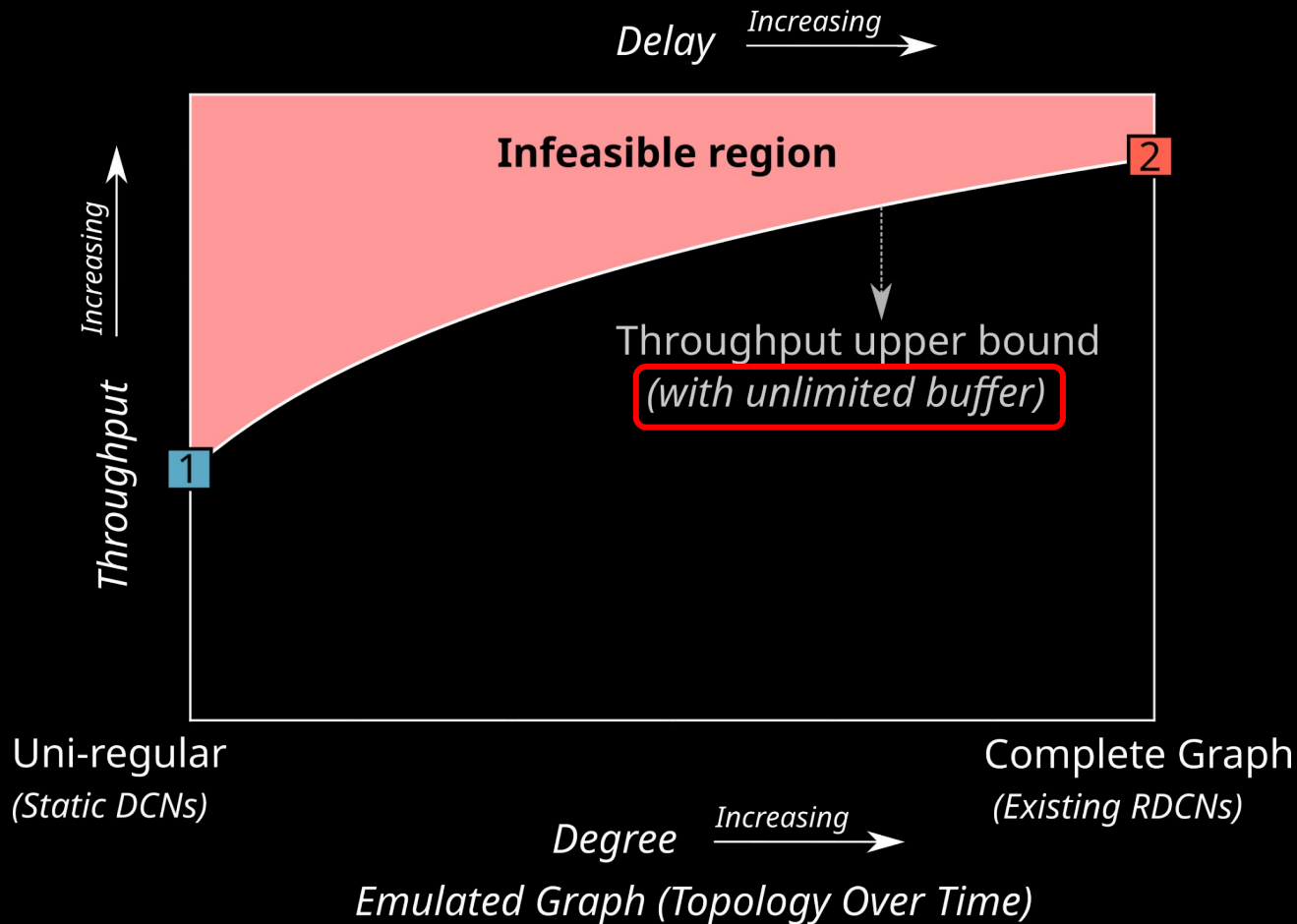
- Delay bound is a function of:
 - Degree d of the emulated graph
 - Duration of the period $\Gamma \cdot \Delta$
 - Throughput θ

$$\begin{aligned} L_{max} &\geq \text{ARD}(\mathcal{M}, \mathcal{F}) = \text{ARL}(\mathcal{M}, \mathcal{F}) \cdot \Gamma \cdot \Delta \\ &\geq \Omega \left(\frac{d \Delta}{n_u \cdot \theta(\mathcal{M}, \mathcal{F})} \right) \end{aligned}$$

Theorem 3: Delay

- Delay bound is a function of:
 - Degree d of the emulated graph
 - Duration of the period $\Gamma \cdot \Delta$
 - Throughput θ

$$\begin{aligned} L_{max} &\geq \text{ARD}(\mathcal{M}, \mathcal{F}) = \text{ARL}(\mathcal{M}, \mathcal{F}) \cdot \Gamma \cdot \Delta \\ &\geq \Omega \left(\frac{d \cdot \Delta}{n_u \cdot \theta(\mathcal{M}, \mathcal{F})} \right) \end{aligned}$$



Buffer Requirements

Input: Periodic Graph \mathcal{G}

Demand Matrix \mathcal{M}

Throughput θ and a feasible flow \mathcal{F}

Output: Minimum required buffer to achieve throughput θ

Theorem 4: Buffer Requirements

- The required buffer is at least the *throughput delay product*

$$\hat{B} \geq (\theta(\mathcal{M}, \mathcal{F}) \cdot M) \cdot \text{ARD}(\mathcal{M}, \mathcal{F})$$

Theorem 4: Buffer Requirements

- The required buffer is at least the *throughput delay product*

$$\hat{B} \geq (\theta(\mathcal{M}, \mathcal{F}) \cdot M) \cdot \text{ARD}(\mathcal{M}, \mathcal{F})$$

Theorem 4: Buffer Requirements

- The required buffer is at least the *throughput delay product*

$$\hat{B} \geq (\theta(\mathcal{M}, \mathcal{F}) \cdot M) \cdot \text{ARD}(\mathcal{M}, \mathcal{F})$$

Theorem 4: Buffer Requirements

- The required buffer is at least the *throughput delay product*

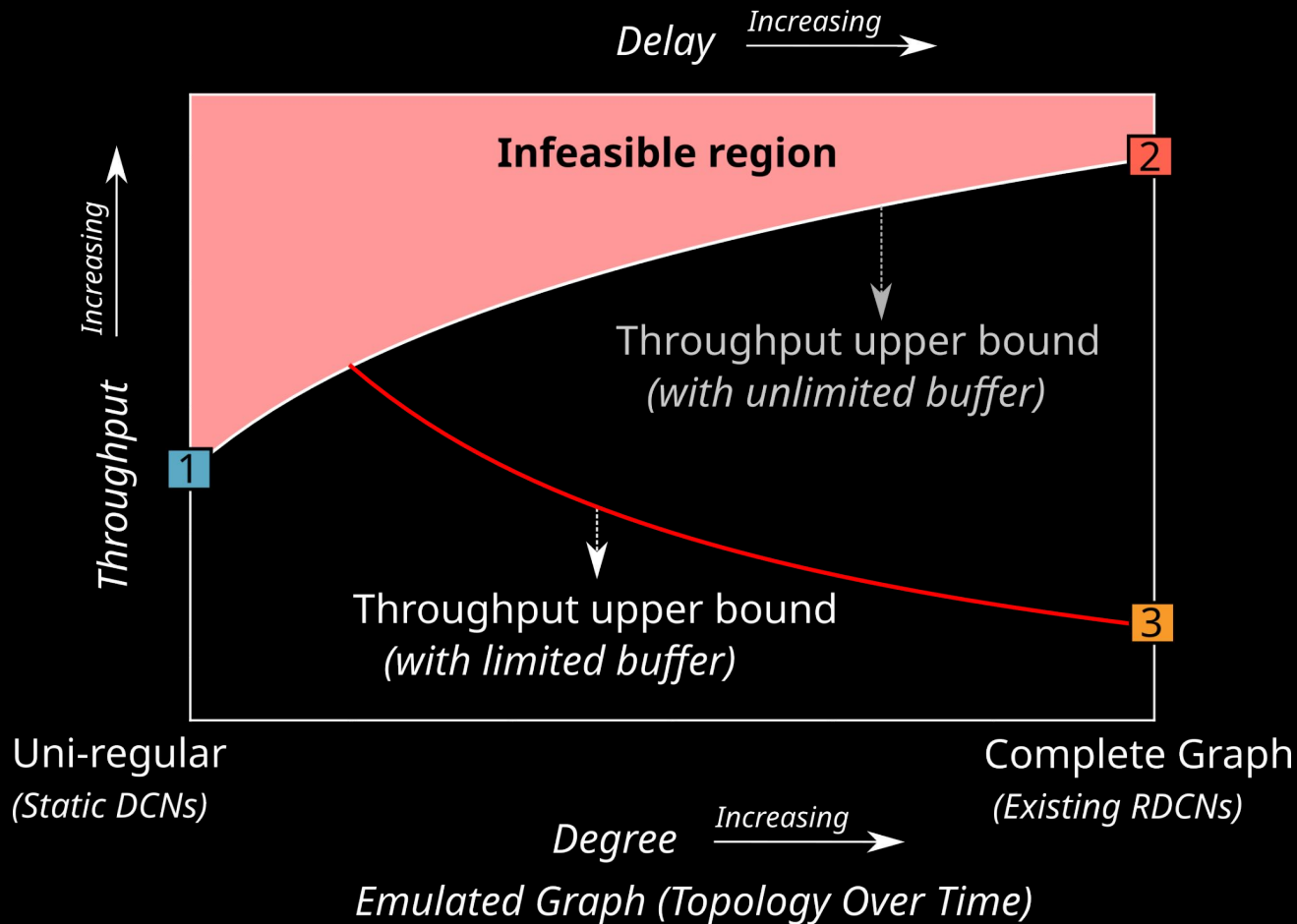
$$\hat{B} \geq (\theta(\mathcal{M}, \mathcal{F}) \cdot M) \cdot \text{ARD}(\mathcal{M}, \mathcal{F})$$

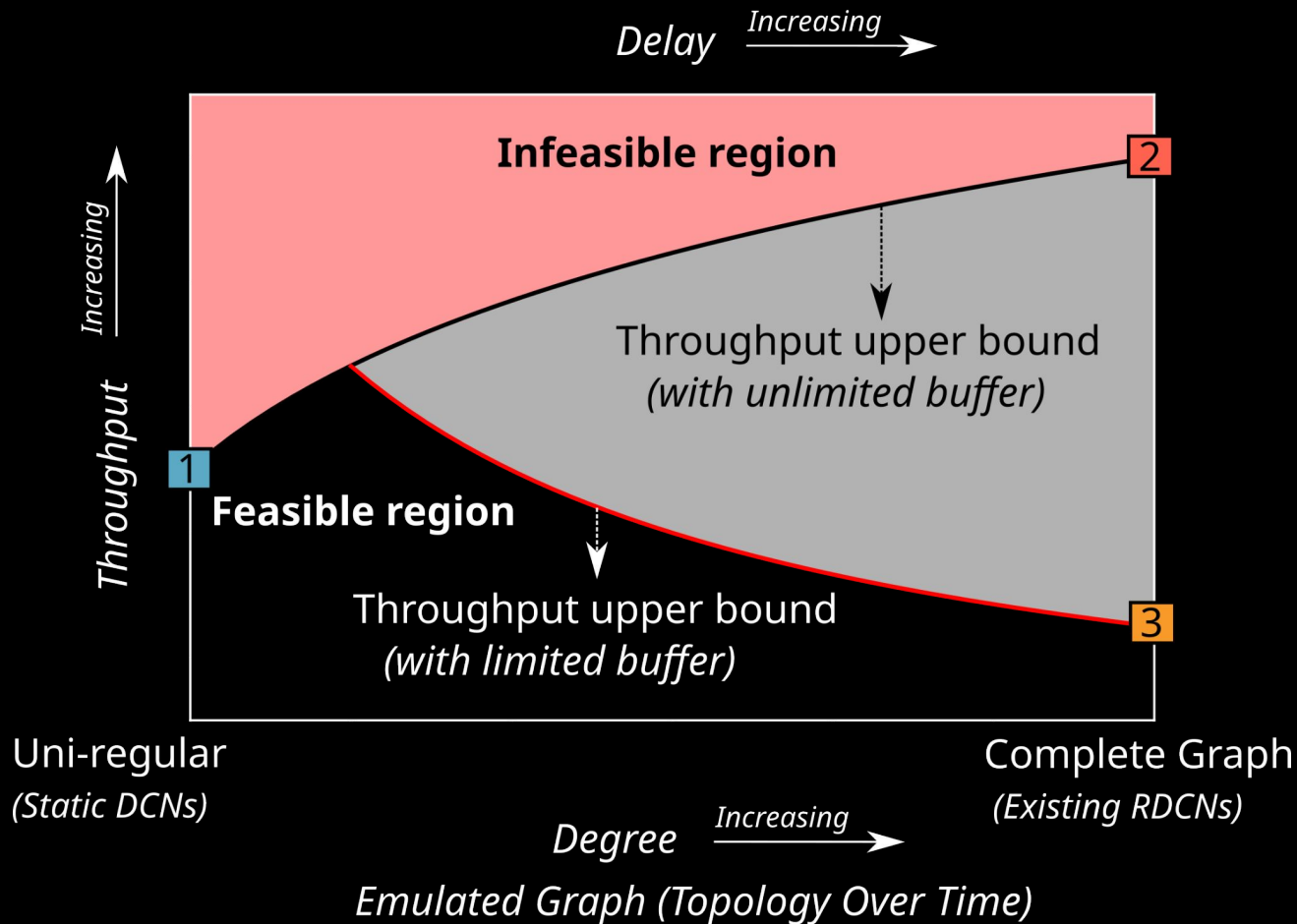
Theorem 4: Buffer Requirements

- The required buffer is at least the *throughput delay product*

$$\hat{B} \geq (\theta(\mathcal{M}, \mathcal{F}) \cdot M) \cdot \text{ARD}(\mathcal{M}, \mathcal{F})$$

Buffer > Bandwidth x Delay





Goals

- **Maximize Throughput**
- **Minimize Latency**
- **Minimize Buffer Requirements**

Optimal Topology with Buffer Constraints

Input: ~~Periodic Graph~~ \mathcal{G}

Demand Matrix \mathcal{M}

Available buffer size B at each node

Output: Degree d of the emulated graph \longrightarrow Periodic graph

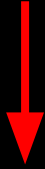
Maximize: Throughput θ

Optimal Topology with Buffer Constraints

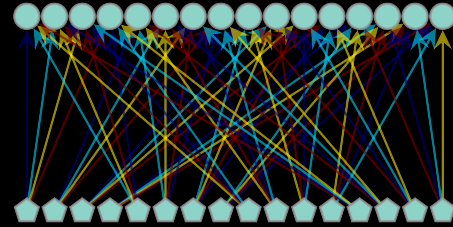
Output: $d = B / (c \cdot \Delta)$

Optimal Topology with Buffer Constraints

Output: $d = B / (c \cdot \Delta)$

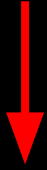


d-regular directed deBruijn graph



Optimal Topology with Buffer Constraints

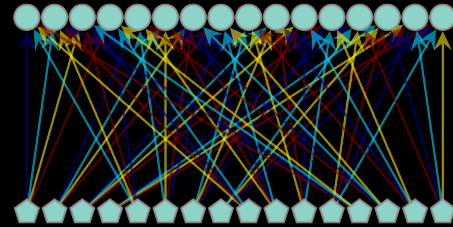
Output: $d = B / (c \cdot \Delta)$



d-regular directed deBruijn graph



Decomposition to d matchings



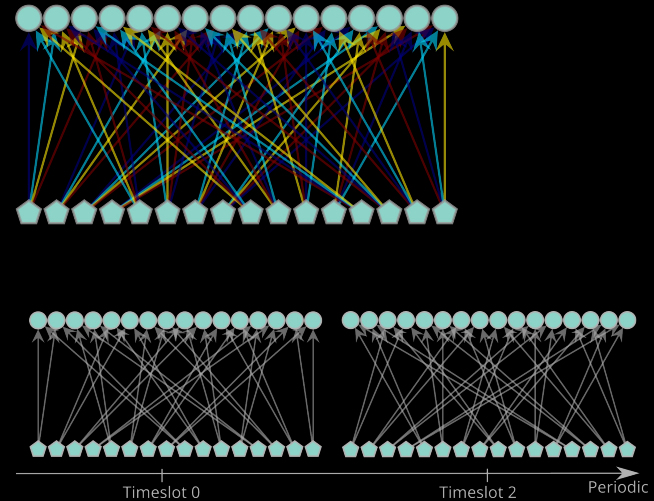
Optimal Topology with Buffer Constraints

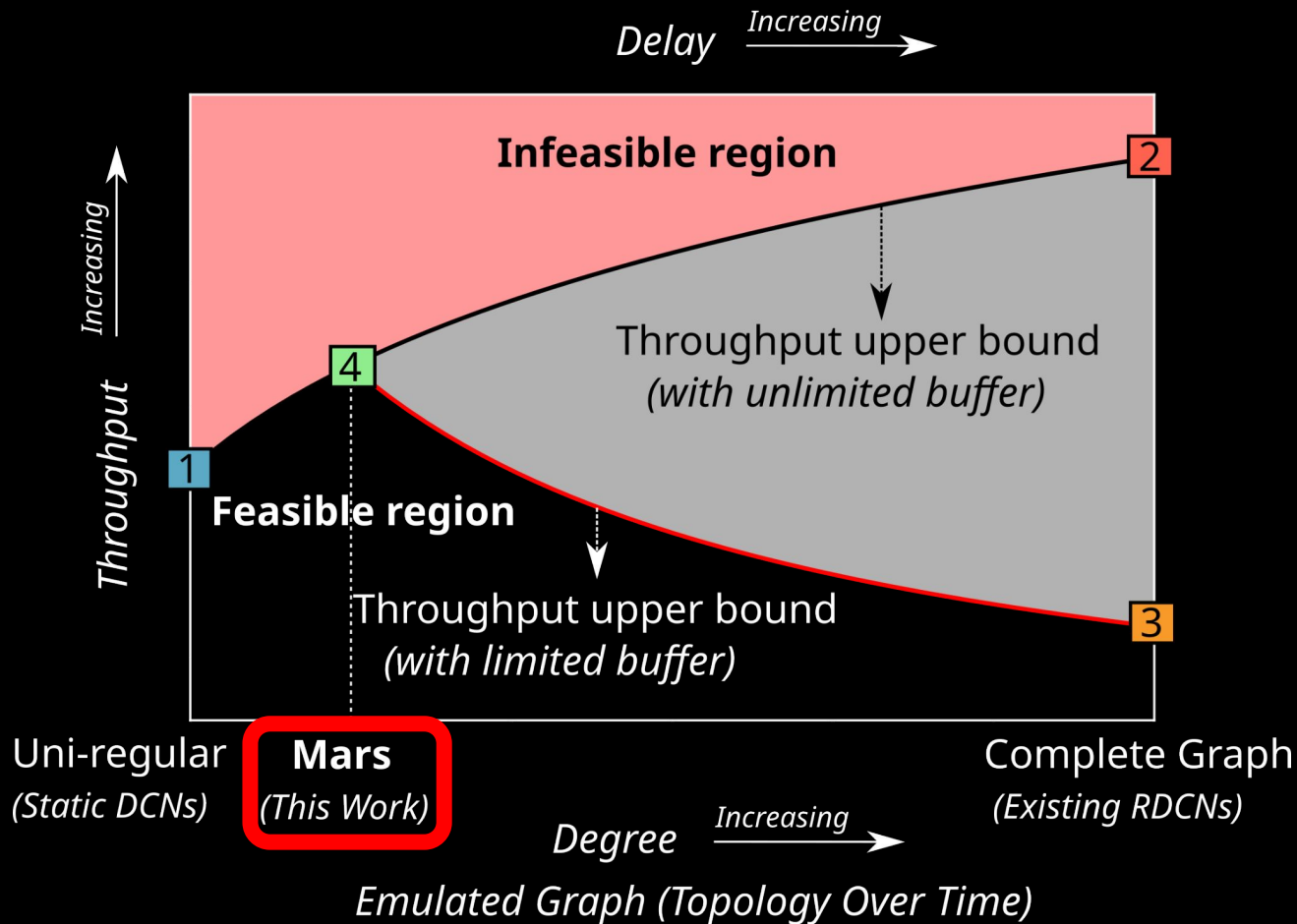
Output: $d = B / (c \cdot \Delta)$

d-regular directed deBruijn graph

Decomposition to d matchings

Periodic graph





Optimal Topology Implications and Future Outlook

Output: $d = B / (c \cdot \Delta)$

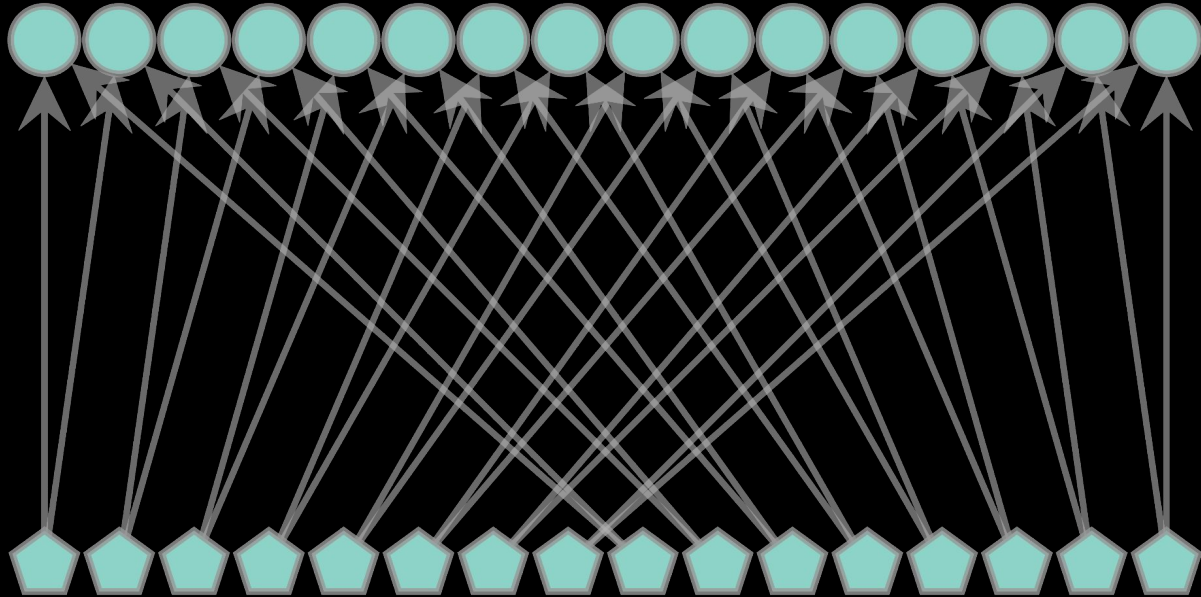
If the reconfiguration technology (Δ) remains same:

- Buffer sizes (B) \uparrow must keep up with the increase in capacity (c) \uparrow

If Buffer sizes (B) do not keep up:

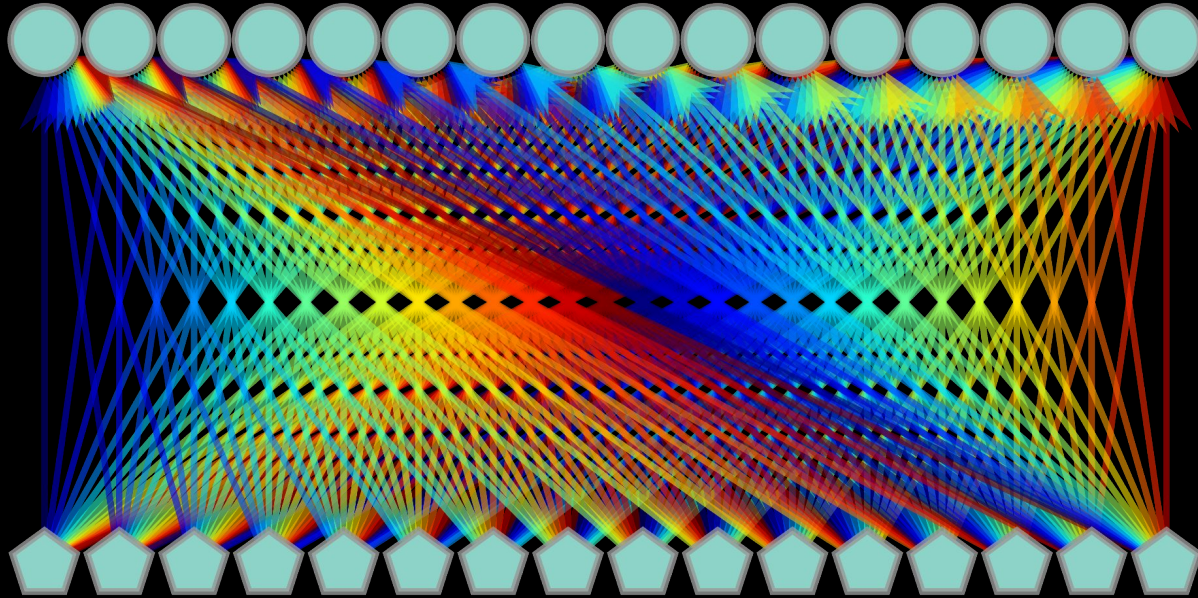
- Increase in capacity (c) \uparrow must be accompanied by decrease in reconfiguration times (Δ) \downarrow
- If not, reducing the degree (d) \downarrow of the emulated graph is inevitable to optimize throughput \rightarrow eventually reaching the case of static topologies.

Static DCNs (*uni-regular*)



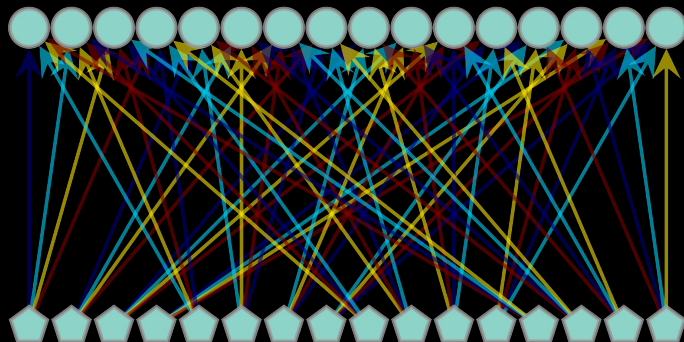
Low throughput but **low delay and buffer requirements**

Existing RDCN designs (*Emulating a complete graph*)

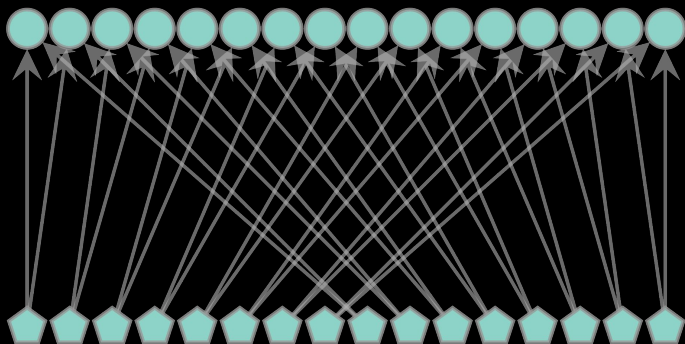


High throughput **but high delay and buffer requirements**

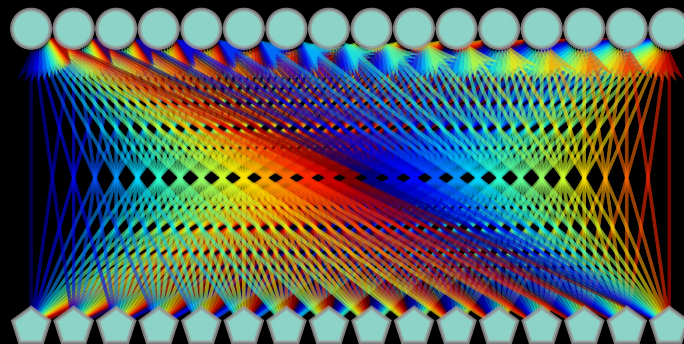
MARS



Near-optimal throughput within the available buffer



Static DCN: Low Throughput



Existing RDCN: High Delay and buffer

MARS

Thank you

Vamsi Addanki

vamsi@inet.tu-berlin.de

 @Vamsi_DT

MARS