

Virtualized Congestion Control

Bryce Cronkite-Ratcliff, Aran Bergman, Shay Vargaftik,
Madhusudhan Ravi, Nick McKeown, Ittai Abraham, Isaac Keslassy



Data Center TCP (DCTCP)

Mohammad Alizadeh^{†‡}, Albert Patel[†], Balaji Prabhakar[†]

[†]Microsoft Research
{albert, dmaltz, padhye, balaji}@microsoft.com
{alizadeh, shyang, msharif, skatti, nmc}@stanford.edu

ABSTRACT

Cloud data centers host diverse applications, mixing workloads that require small predictable latency with others requiring sustained throughput. In this environment, today's state-of-the-art protocol falls short. We present measurements of a production cluster and reveal impairments that lead to long connection latencies, rooted in TCP's demands on the limited bandwidth available in data center switches. For example, bandwidth "background" flows build up queues at the switches, impacting the performance of latency sensitive "foreground" flows.

To address these problems, we propose DCTCP, a new TCP-like protocol for data center networks. DCTCP leverages Explicit Congestion Notification (ECN) in the network to provide rate feedback back to the end hosts. We evaluate DCTCP at 1 and 10 Gbps using commodity, shallow buffered switches. We find DCTCP delivers the same or better throughput than TCP, while using less buffer space. Unlike TCP, DCTCP also provides high tolerance and low latency for short flows. In handling flows derived from operational measurements, we found DCTCP allows the applications to handle 10X the current background traffic without impacting foreground traffic. Further, a 10X increase in background traffic does not cause any timeouts, thus largely eliminating incast problems.

Categories and Subject Descriptors: C.2.2 [Computer-Communication Networks]: Network Architecture and Design
General Terms: Measurement, Performance
Keywords: Data center network, ECN, TCP

1. INTRODUCTION

In recent years, data centers have transformed computing by large scale consolidation of enterprise IT into data centers and with the emergence of cloud computing services from Amazon, Microsoft and Google. A consistent theme in data center design has been to build highly available, high performance computing and storage infrastructure using low cost components [16]. A corresponding trend has also emerged in data center networks. In particular, low-cost switches are deployed at the top of the rack, providing up to 48 ports at 1Gbps, all for under \$2000 — roughly the price of one data center server.

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 this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGCOMM'10, August 30–September 3, 2010, New Delhi, India
Copyright 2010 ACM 978-1-4503-0201-2/10/08 ...\$10.00.

pFabric: Minimal Near-Optimal Datacenter Transport

Mohammad Alizadeh^{†‡},
Nick McKeown[†],
Shyamnath Gollakota[†]

[†]Stanford University
{alizadeh, shyang, msharif, skatti, nmc}@stanford.edu

ABSTRACT

In this paper we present pFabric, a minimalistic datacenter transport design that provides near theoretically optimal flow completion times even at the 99th percentile for short flows, minimizing average flow completion time for long flows. Moreover, pFabric delivers this performance with a very simple design that is based on a key conceptual insight: datacenter transports decouple flow scheduling from rate control. For flow scheduling, packets carry a single priority number set independent of the flow; switches have very small buffers and implement a simple priority-based scheduling/dropping mechanism. Rate control is also correspondingly simpler; flows start at line rate and back off only under high and persistent packet loss. We provide theoretical intuition and show via extensive simulations that the combination of these two simple mechanisms is sufficient for near-optimal performance.

Categories and Subject Descriptors: C.2.1 [Computer-Communication Networks]: Network Architecture and Design

General Terms: Design, Performance

Keywords: Datacenter network, Packet transport, Flow scheduling

1. INTRODUCTION

Datacenter workloads impose unique and stringent requirements on the transport fabric. Interactive soft real-time workloads, such as the ones seen in search, social networking, and retail, consist of a large number of small requests and responses across threads that are stitched together to perform a user-requested operation (e.g., delivering search results). These application workloads require low latency for each of the short request/response flows, and their perceived performance is dictated by how quickly responses (or a large fraction of) the requests are collected and delivered to the user. However, in currently deployed TCP-based datacenters, the latency for these short flows is poor — flow completion time (FCT) can be as high as tens of milliseconds while in the best case, flows could complete in 10-20 microseconds. The reason for this is that these flows often get queued up behind bursts of packets from other flows or co-existing workloads (such as backup, replication, mining, etc) which significantly increases their completion times.

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 this notice and the full citation on the first page. Copyrights for this work owned by others than ACM must be honored. Abstracting or republication of part of this work without the permission of the copyright owner is illegal. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Requests for permission should be addressed to permissions@acm.org.
SIGCOMM'13, August 12–16, 2013, Hong Kong, China
Copyright 2013 ACM 978-1-4503-2056-6/13/08 ...\$15.00.

TIMELY: RTT-based Congestion Control for the Datacenter

Radhika Mittal*(UC Berkeley), Vinh The Lam, Nandita Dukkipati, Emily Blem, Hassan Wassel, Monia Ghobadi*(Microsoft), Amin Vahdat, Yaogong Wang, David Wetherall, David Zetsche

Google, Inc.

ABSTRACT

Datacenter transports aim to deliver low latency messaging together with high throughput. We show that simple round-trip delay, measured as round-trip times at hosts, is an effective congestion signal without the need for switch feedback. First, we show that advances in NIC hardware have made RTT measurement possible with microsecond accuracy, and that these RTTs are sufficient to estimate switch queueing. Then we describe how TIMELY can adjust transmission rates using RTT gradients to keep packet latency low while delivering high bandwidth. We implement our design in host software running over NICs with OS-bypass capabilities. We show using experiments with up to hundreds of machines on a Clos network topology that it provides excellent performance: turning on TIMELY for OS-bypass messaging over a fabric with PFC lowers 99 percentile tail latency by 9X while maintaining near line-rate throughput. Our system also outperforms DCTCP running in an optimized kernel, reducing tail latency by 13X. To the best of our knowledge, TIMELY is the first delay-based congestion control protocol for use in the datacenter, and it achieves its results despite having an order of magnitude fewer RTT signals (due to NIC offload) than earlier delay-based schemes such as Vegas.

CCS Concepts

•Networks → Transport protocols;

Keywords

datacenter transport; delay-based congestion control; OS-bypass; RDMA

*Work done while at Google

Permission to make digital or hard copies of part or all 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 this 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/author(s).

SIGCOMM '15 August 17–21, 2015, London, United Kingdom

© 2015 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-3542-3/15/08

DOI: <http://dx.doi.org/10.1145/2785956.2787510>

1. INTRODUCTION

Datacenter networks run tightly-coupled computing tasks that must be responsive to users, e.g., thousands of backend computers may exchange information to serve a user request, and all of the transfers must complete quickly enough to let the complete response to be satisfied within 100 ms [24]. To meet these requirements, datacenter transports must simultaneously deliver high bandwidth (\gg Gbps) and utilization at low latency (\ll msec), even though these aspects of performance are at odds. Consistently low latency matters because even a small fraction of late operations can cause a ripple effect that degrades application performance [21]. As a result, datacenter transports must strictly bound latency and packet loss.

Since traditional loss-based transports do not meet these strict requirements, new datacenter transports [10, 18, 30, 35, 37, 47], take advantage of network support to signal the onset of congestion (e.g., DCTCP [35] and its successors use ECN), introduce flow abstractions to minimize completion latency, cede scheduling to a central controller, and more. However, in this work we take a step back in search of a simpler, immediately deployable design.

The crux of our search is the congestion signal. An ideal signal would have several properties. It would be fine-grained and timely to quickly inform senders about the extent of congestion. It would be discriminative enough to work in complex environments with multiple traffic classes. And, it would be easy to deploy.

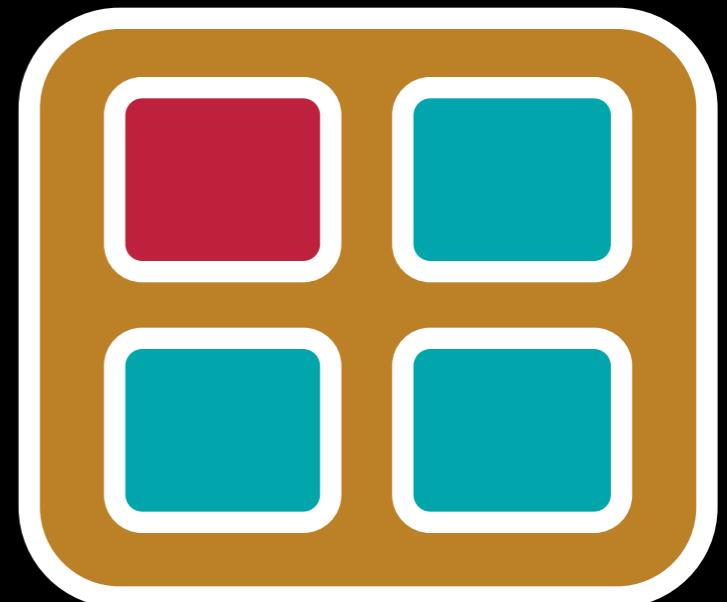
Surprisingly, we find that a well-known signal, properly adapted, can meet all of our goals: delay in the form of RTT measurements. RTT is a fine-grained measure of congestion that comes with every acknowledgment. It effectively supports multiple traffic classes by providing an inflated measure for lower-priority transfers that wait behind higher-priority ones. Further, it requires no support from network switches.

Delay has been explored in the wide-area Internet since at least TCP Vegas [16], and some modern TCP variants use delay estimates [44, 46]. But this use of delay has not been without problems. Delay-based schemes tend to compete poorly with more aggressive, loss-based schemes, and delay



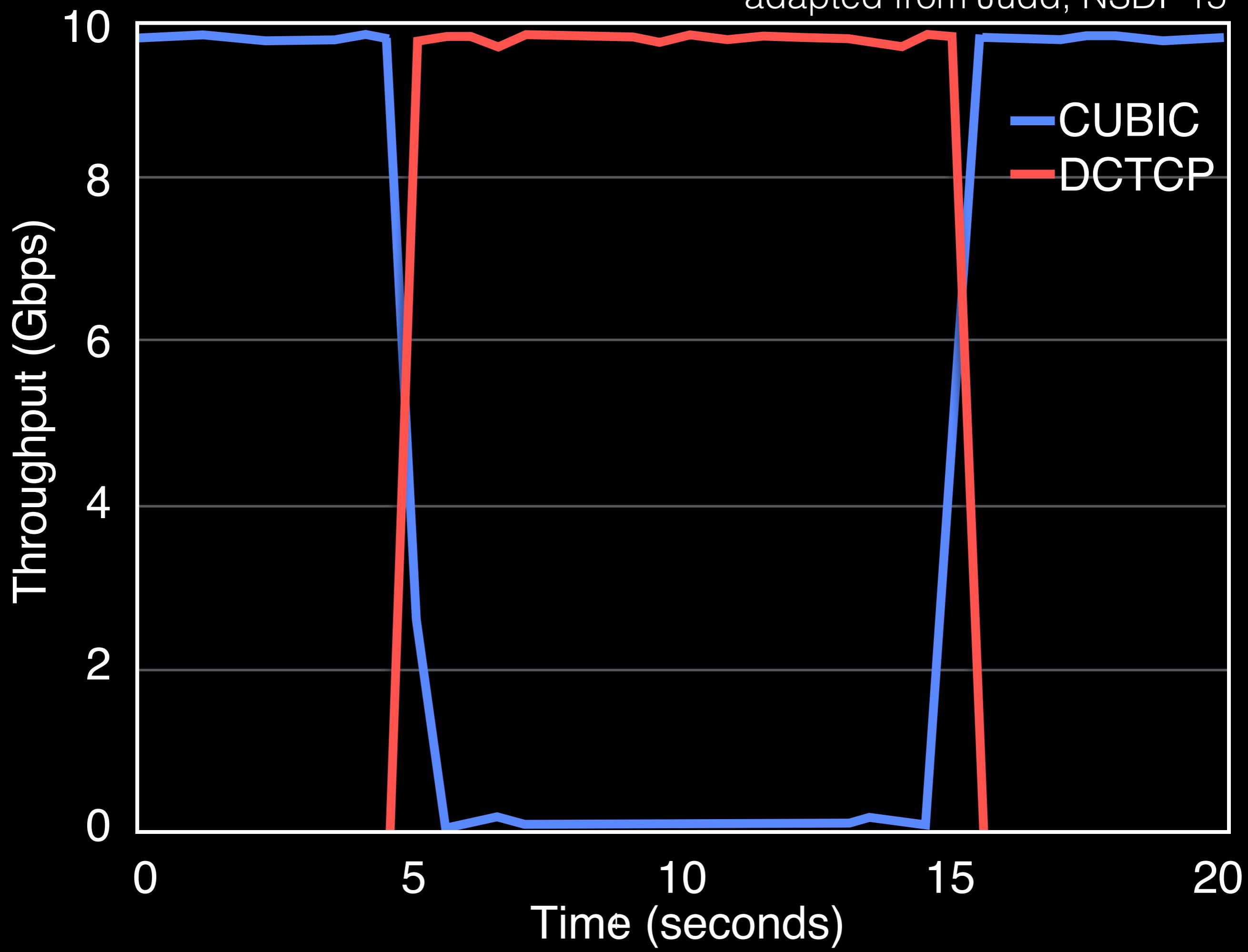
Multitenant Datacenters

Tenants choose Congestion Control



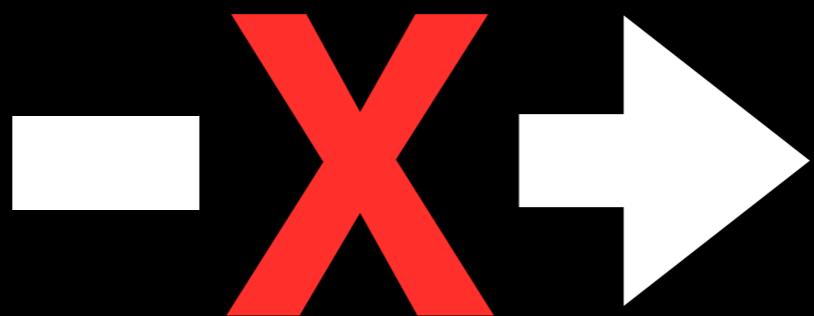
Enterprise Datacenters

Can't upgrade legacy applications

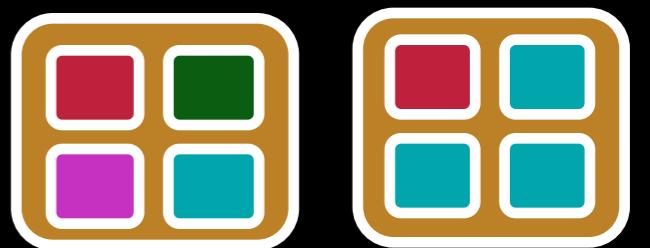


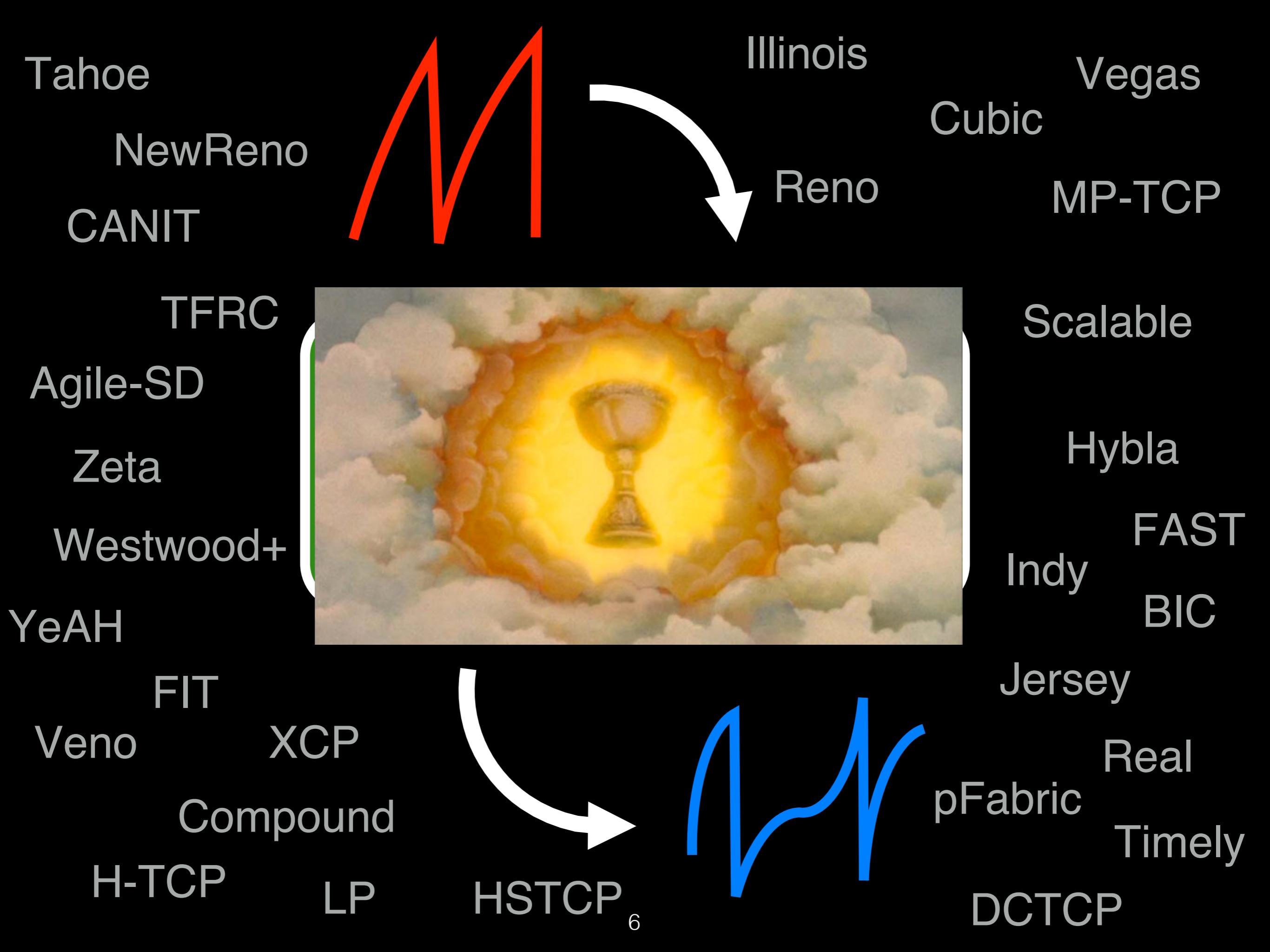
New Congestion
Control Algorithms

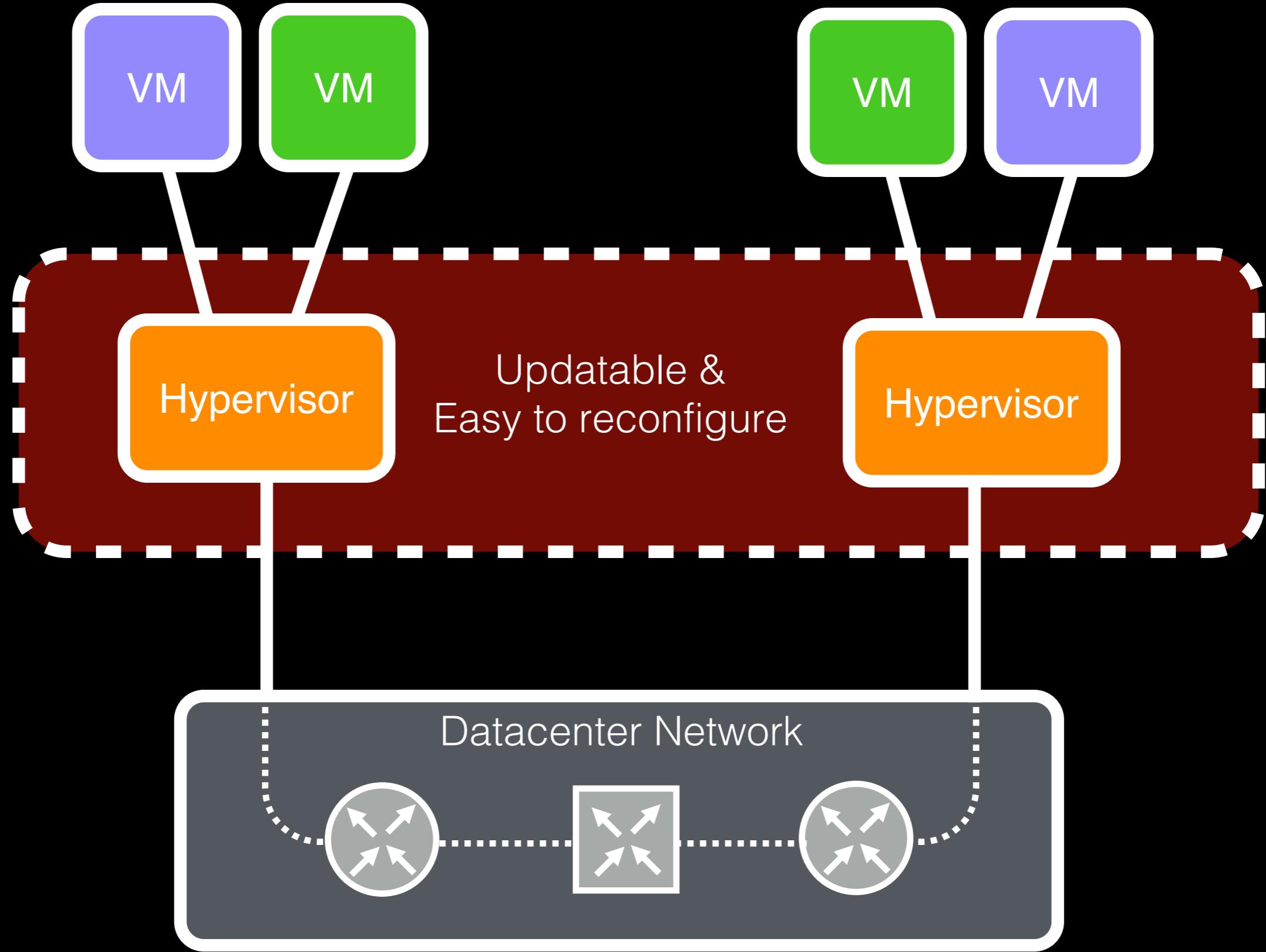
pFabric Timely
DCTCP

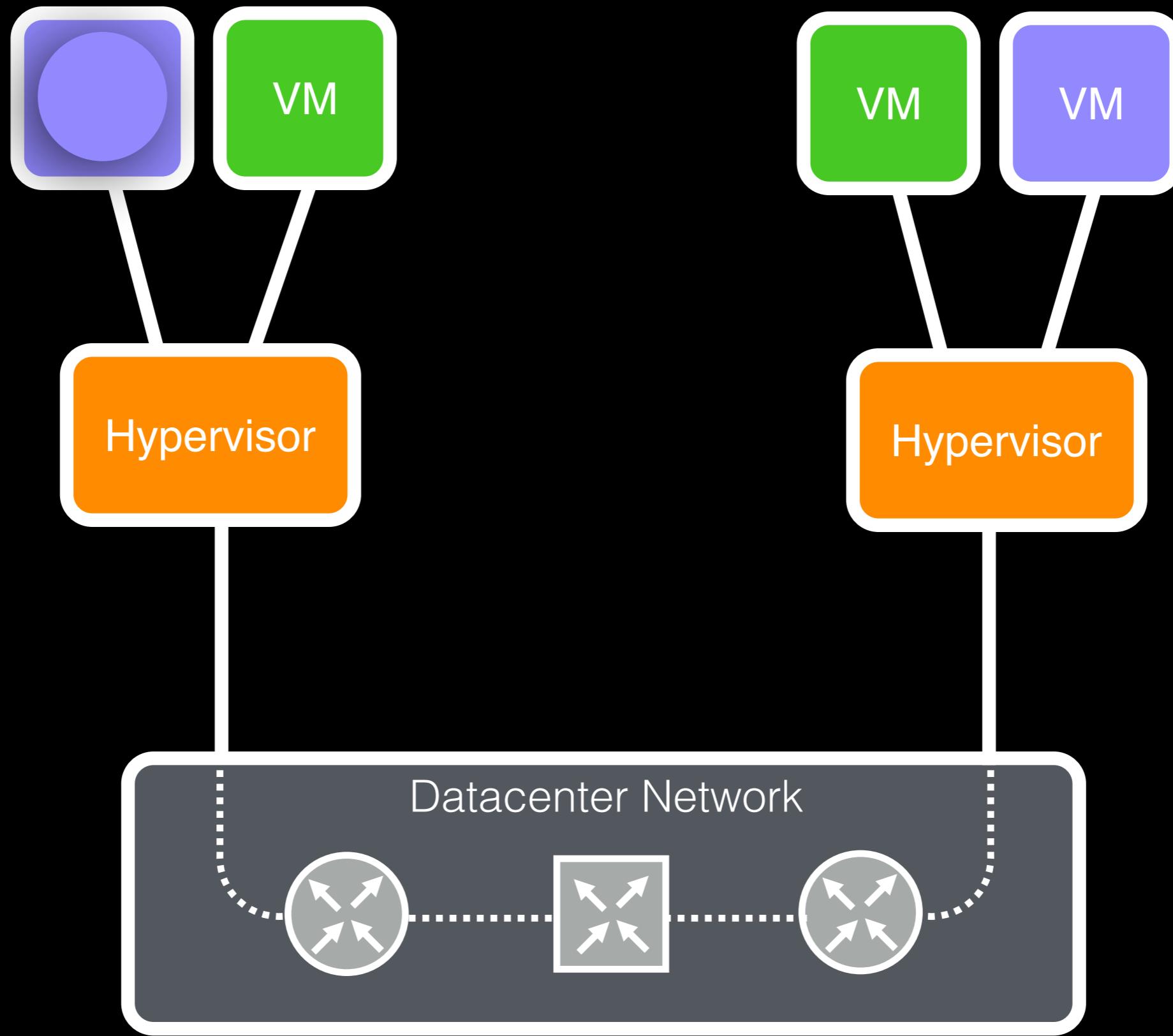


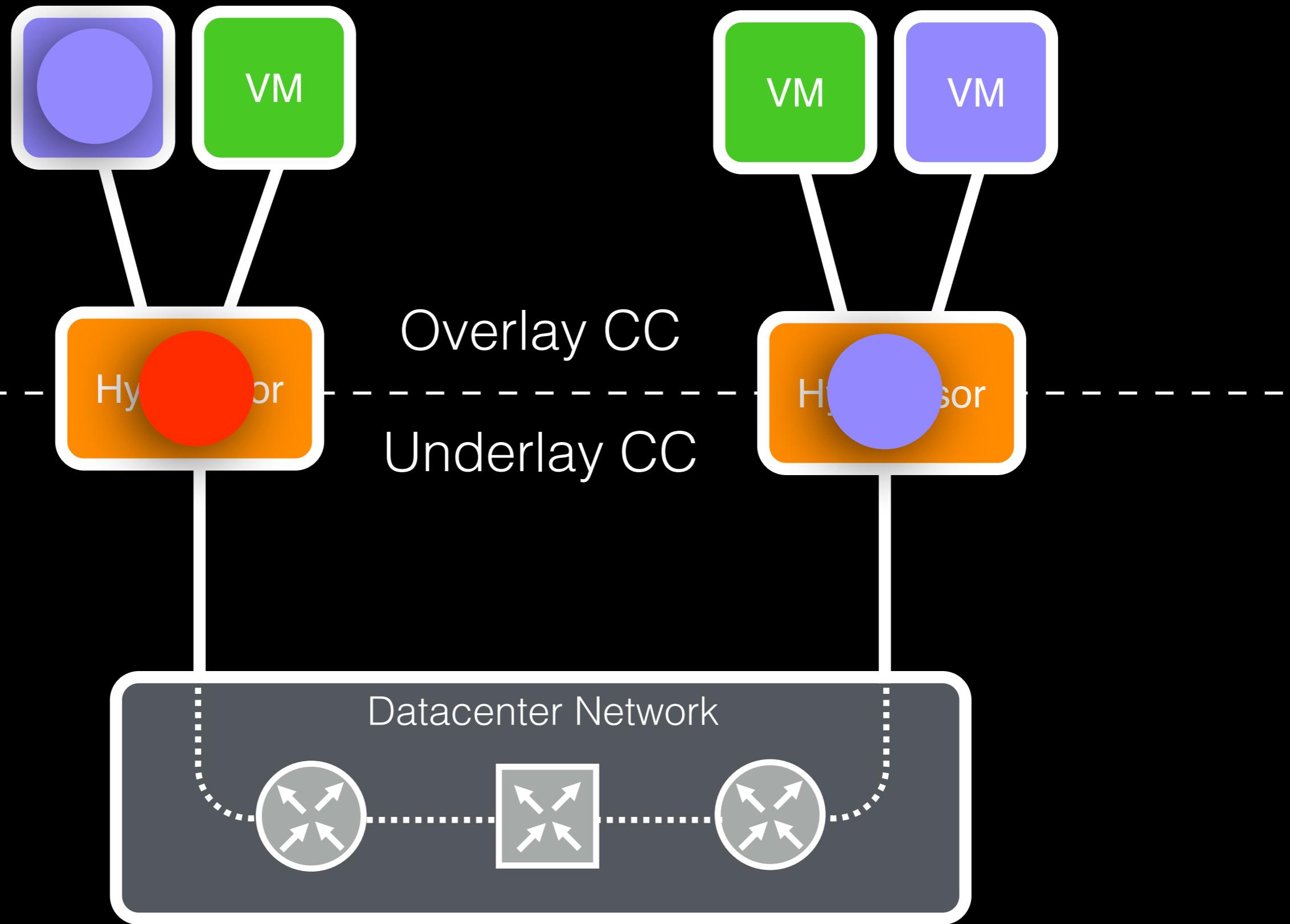
Multitenant &
Enterprise Datacenters

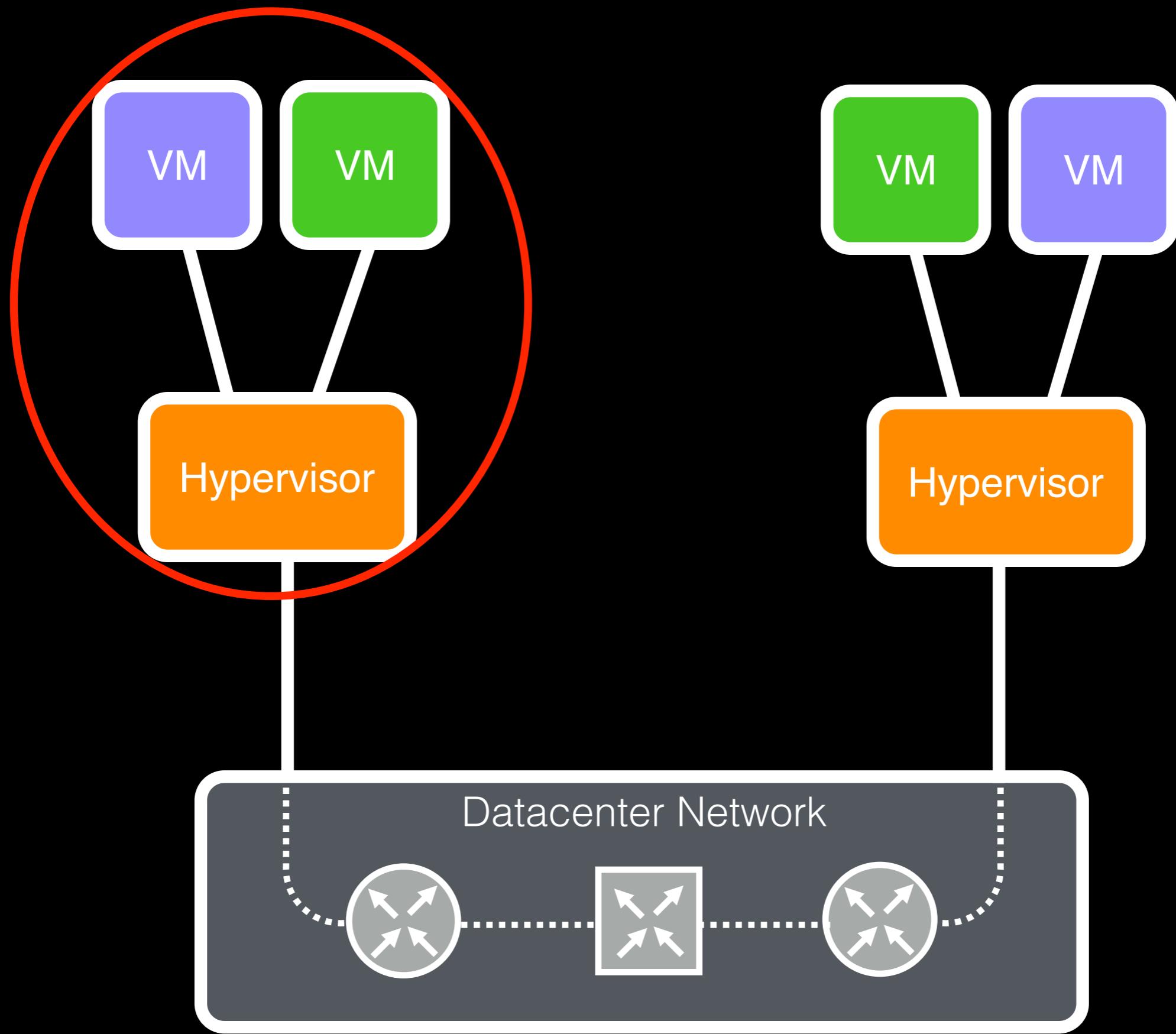




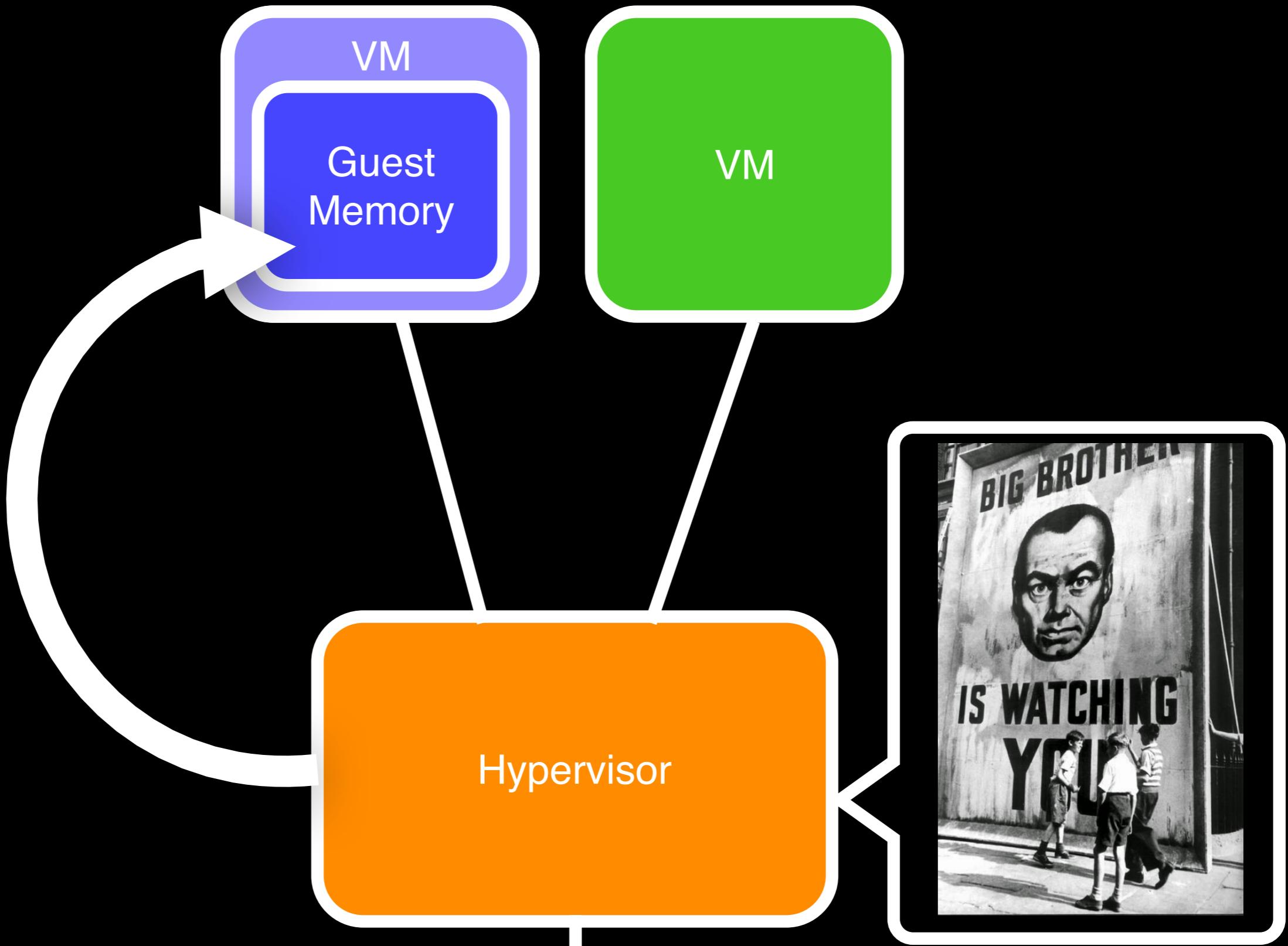




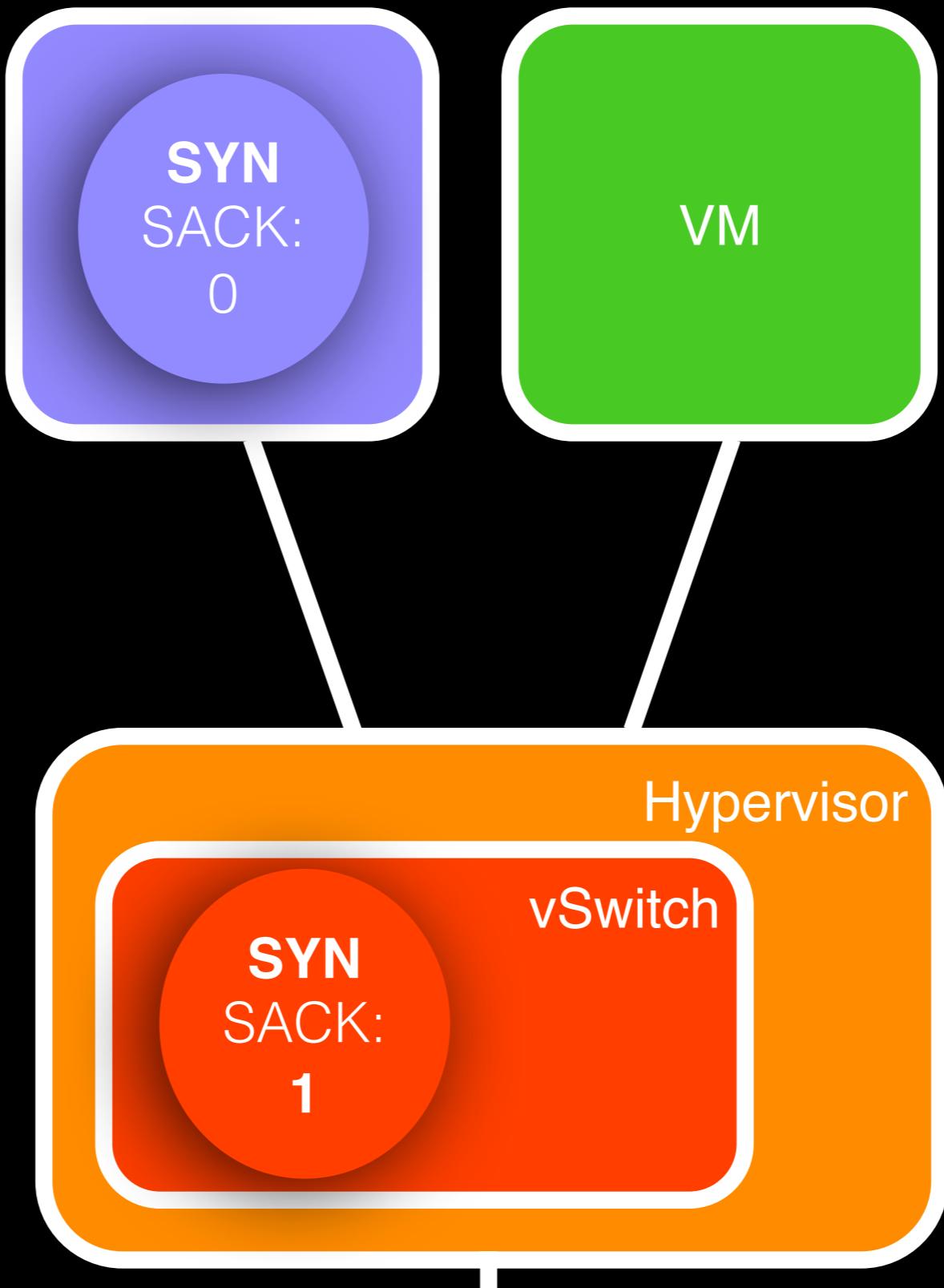




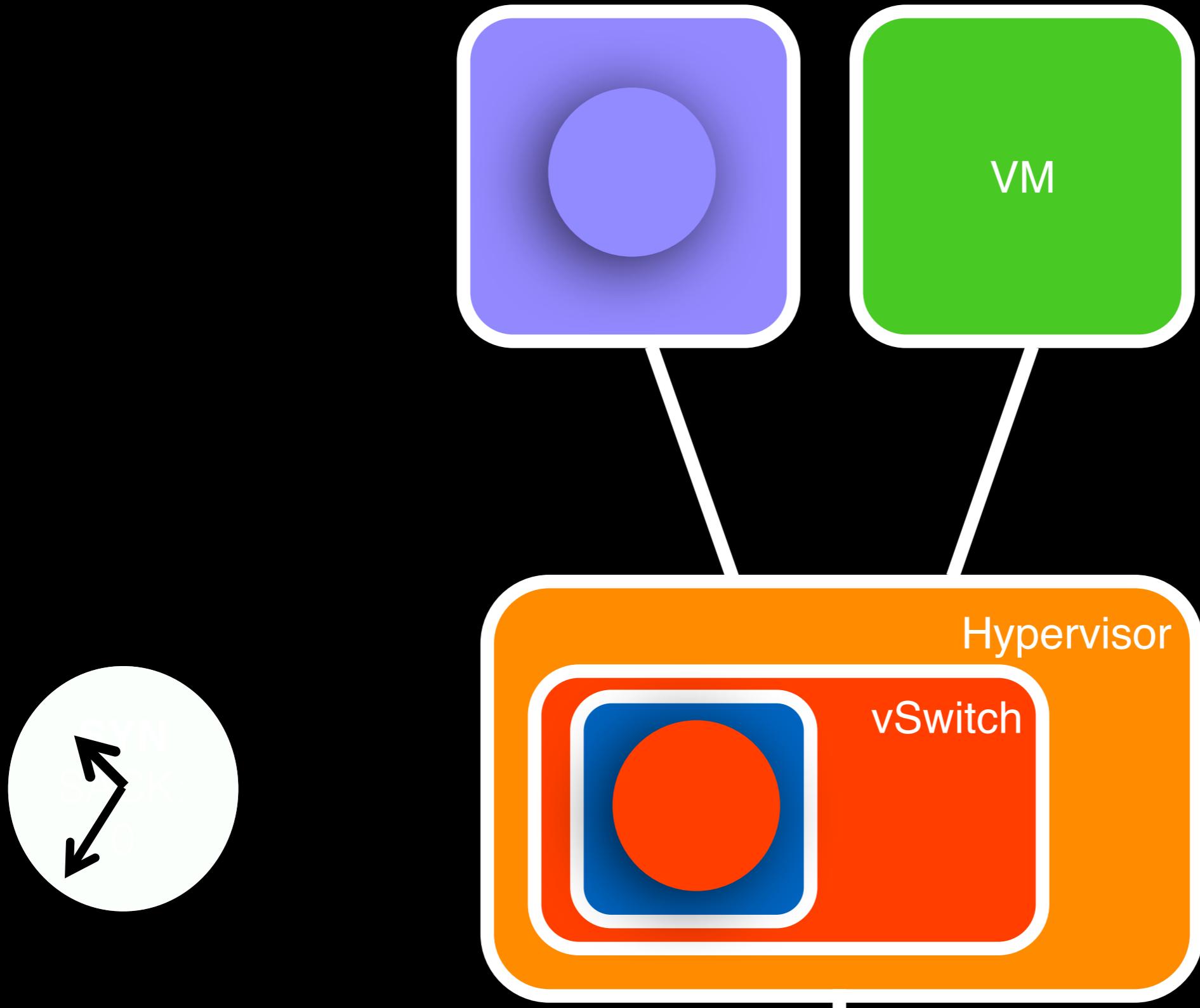
Guest Introspection



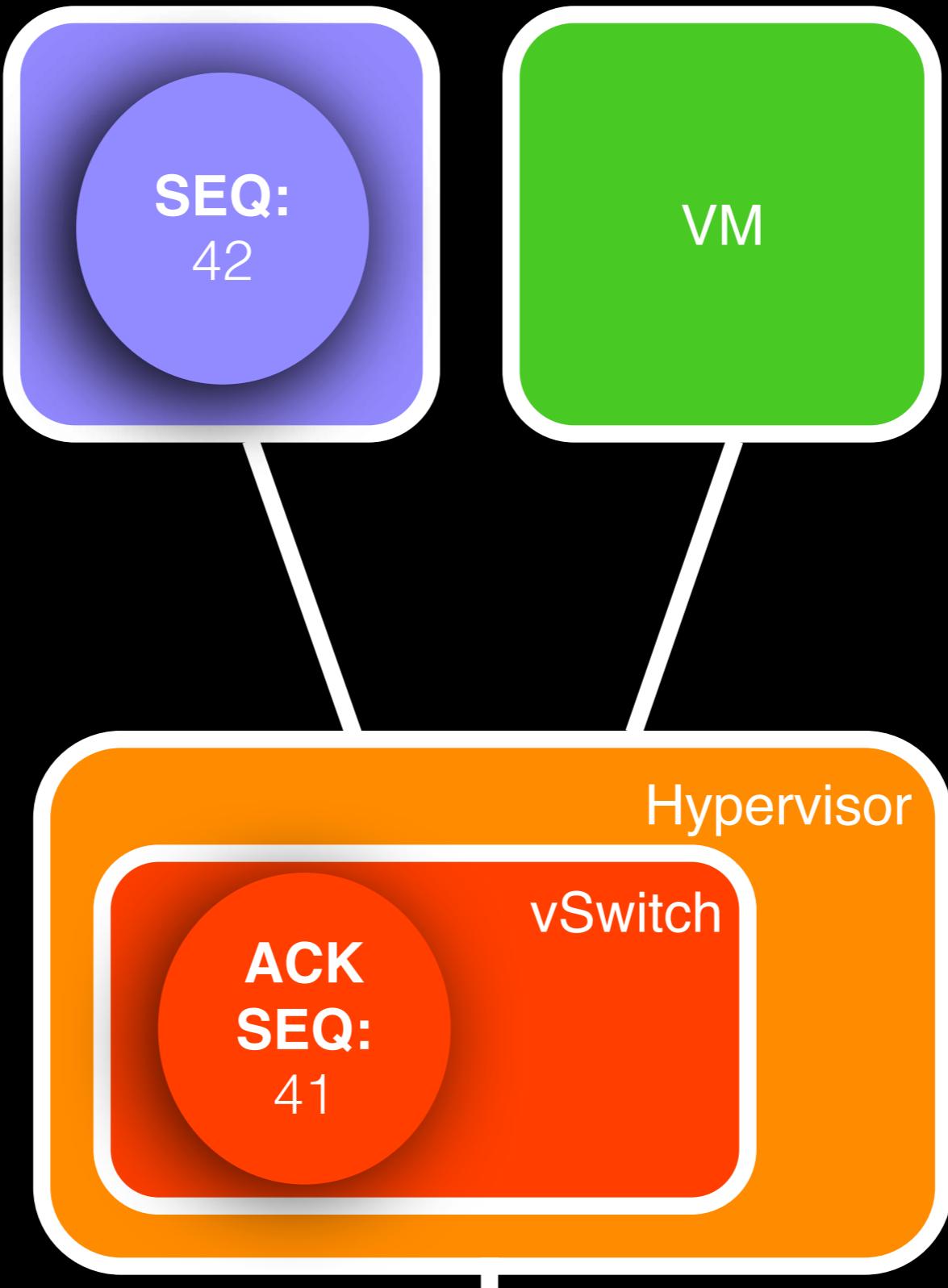
TCP Header Modification



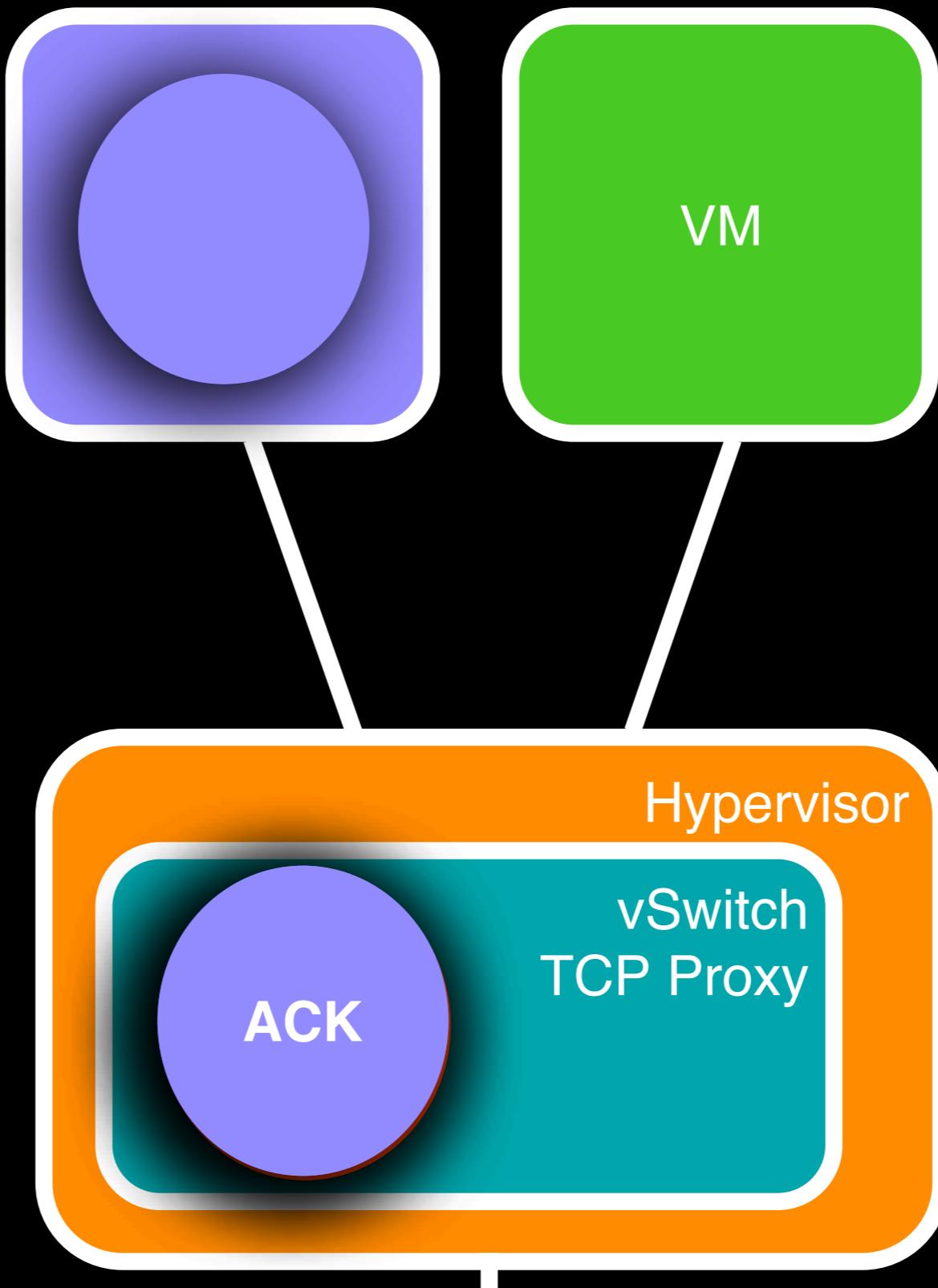
Buffering



Fake ACKs



TCP Proxy



vCC Flow Modification Techniques

- Guest Introspection
- Buffering
- TCP Header Modification
- Fake ACK generation
- TCP Proxy

A Binary Feedback Scheme for Congestion Avoidance in Computer Networks

K. K. RAMAKRISHNAN and RAJ JAIN
Digital Equipment Corporation

We propose a scheme for *congestion avoidance* in the network layer. The scheme uses a minimal amount of feedback from the network to the source to adjust the amount of traffic allowed into the network. It uses a binary feedback scheme to indicate when the network is congested and set a *congestion-indication* bit on packets flowing through the network. This information is communicated back to the users through the network. The scheme is distributed, adapts to the dynamic state of the network, and is robust. It is quite simple to implement, and has low overhead. It can be applied to multiple sources. This paper presents the basic idea and choice of the various decision mechanisms. We also discuss transient changes in the network and pathological cases.

Categories and Subject Descriptors: C.2.1 [Computer Systems Organization]: Architecture and Design—*network communications*; C.2.3 [Computer Communications]: *forward networks*; C.2.3 [Computer Communications]: *monitoring*; C.4 [Computer System Organization]: *operating systems*

General Terms: Algorithms, Design, Performance

Additional Key Words and Phrases: Computer communication, congestion avoidance, congestion control, congestion indication, performance, network power

1. INTRODUCTION

Congestion in computer networks is a significant problem in today's networks and increased link speeds. Flow control problems that have been addressed by several mechanisms in the increasing range of speeds of links, distributed computing, effective control of traffic, and important. The lack of control may result in lost packets, retransmissions, may ultimately lead to *congestion collapse*.

The control mechanisms adopted to combat congestion may be categorized into two distinct types:

An earlier version of this paper was presented at the "Workshop on Distributed Systems Architectures and Protocols," Stanford, CA, October 1989. Authors' address: Distributed Systems Architecture Group, Digital Equipment Corporation, Littleton, MA 01460-1289.

Permission to copy without fee all or part of this material is granted by the copyright holders for research and educational purposes only. The material must not be sold in whole or in part without the written permission of the copyright holders. The copyright holders make no representations about the accuracy of the information contained in this document. The copyright holders will not be liable in any way for any damages resulting from its use.

© 1990 ACM 0734-2071/90/0500-0158 \$01.50

ACM Transactions on Computer Systems, Vol. 8, No. 2, May 1990, pp. 158–179.

TCP and Explicit Congestion Notification

Sally Floyd*

Lawrence Berkeley Laboratory
One Cyclotron Road, Berkeley, CA 94704
floyd@ee.lbl.gov

Abstract

This paper discusses the use of Explicit Congestion Notification (ECN) mechanisms in the TCP/IP protocol. The first part proposes new guidelines for TCP's response to ECN mechanisms (e.g., Source Quench packets, ECN fields in packet headers). Next, using simulations, we explore the benefits and drawbacks of ECN in TCP/IP networks. Our simulations use RED gateways modified to set an ECN bit in the IP packet header as an indication of congestion, with Reno-style TCP modified to respond to ECN as well as to packet drops as indications of congestion. The simulations show that one advantage of ECN mechanisms is in avoiding unnecessary packet drops, and therefore avoiding unnecessary delay for packets from low-bandwidth delay-sensitive TCP connections. A second advantage of ECN mechanisms is in networks (generally LANs) where the effectiveness of TCP retransmit timers is limited by the coarse granularity of the TCP clock. The paper also discusses some implementation issues concerning specific ECN mechanisms in TCP/IP networks.

1 Introduction

This paper proposes guidelines for TCP's response to ECN (Explicit Congestion Notification) mechanisms, and explores the effect upon performance of ECN mechanisms in TCP/IP networks. The paper discusses some implementation issues concerning ECN mechanisms, but does not make specific recommendations concerning the use of ECN mechanisms in TCP/IP networks.

In current TCP/IP networks, TCP relies on packet drops as the indication of congestion. The TCP source detects dropped packets either from the receipt of three duplicate acknowledgements (ACKs) or after the time-out of a retransmit timer, and responds to a dropped packet by reducing the congestion window [J88]. TCP implementations also respond to ICMP Source Quench

*This work was supported by the Director, Office of Energy Research, Scientific Computing Staff, of the U.S. Department of Energy under Contract No. DE-AC03-76SF00098.

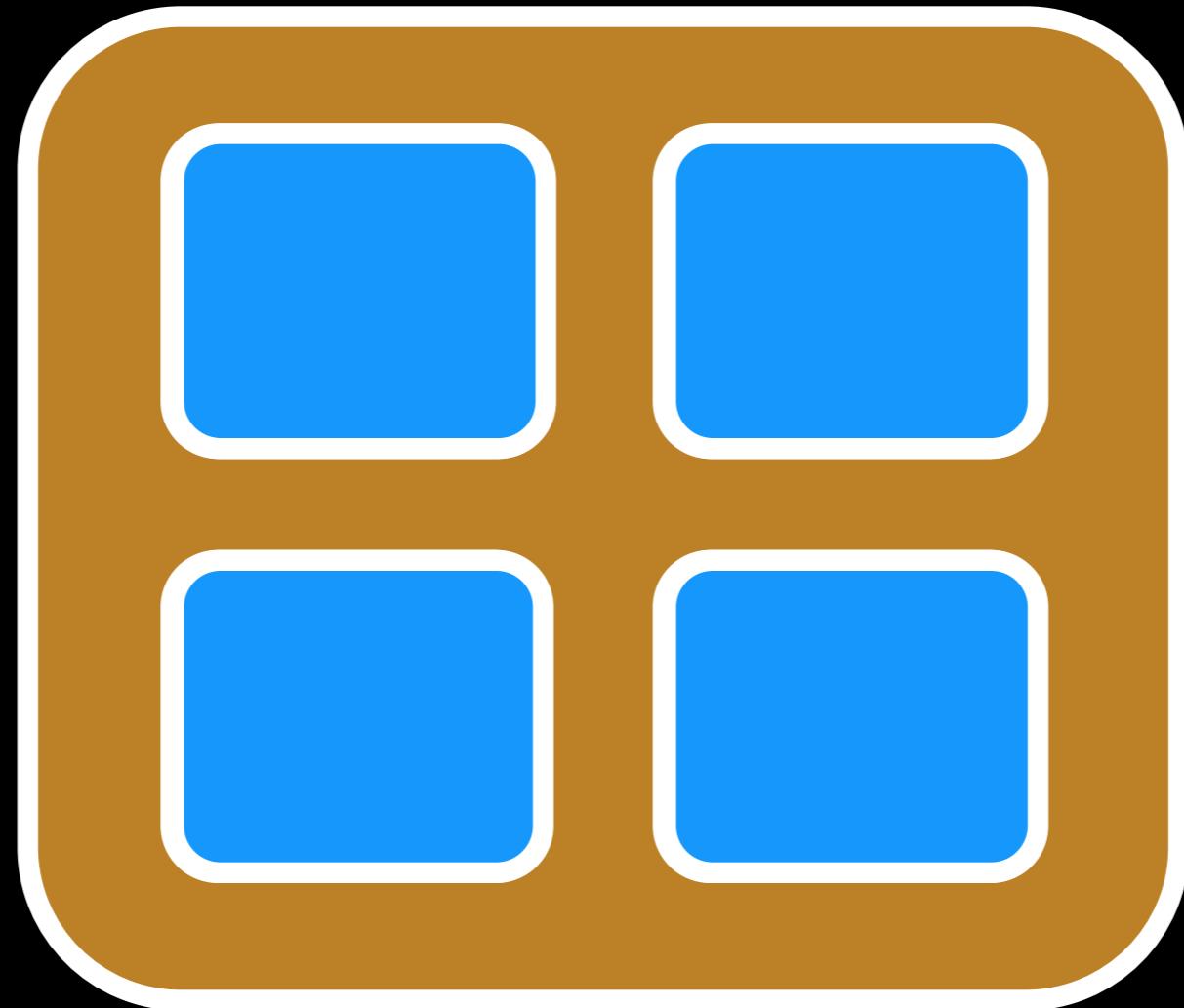
messages, but Source Quench messages are rarely used, in part because they can consume network bandwidth in times of congestion.

The reliance on packet drops as the indication of congestion is perfectly appropriate for a network with routers whose main function is to route packets to the appropriate output port. Most current routers in TCP/IP networks have no provision for the detection of incipient congestion. When a queue overflows, packets are dropped. When the TCP source detects this packet drop, the TCP source infers the presence of congestion in the network.

Future routers are likely to have more developed mechanisms for the detection of incipient congestion. With the DECbit scheme, for example, routers detect incipient congestion by computing the average queue size, and set the ECN bit in packet headers when the average queue size exceeds a certain threshold [RJ90]. Recently-proposed Random Early Detection (RED) gateways have a similar ability to detect incipient congestion [FJ93]. Gateways with mechanisms for detecting incipient congestion before the queue overflows are not limited to packet drops as the method of informing sources of congestion.

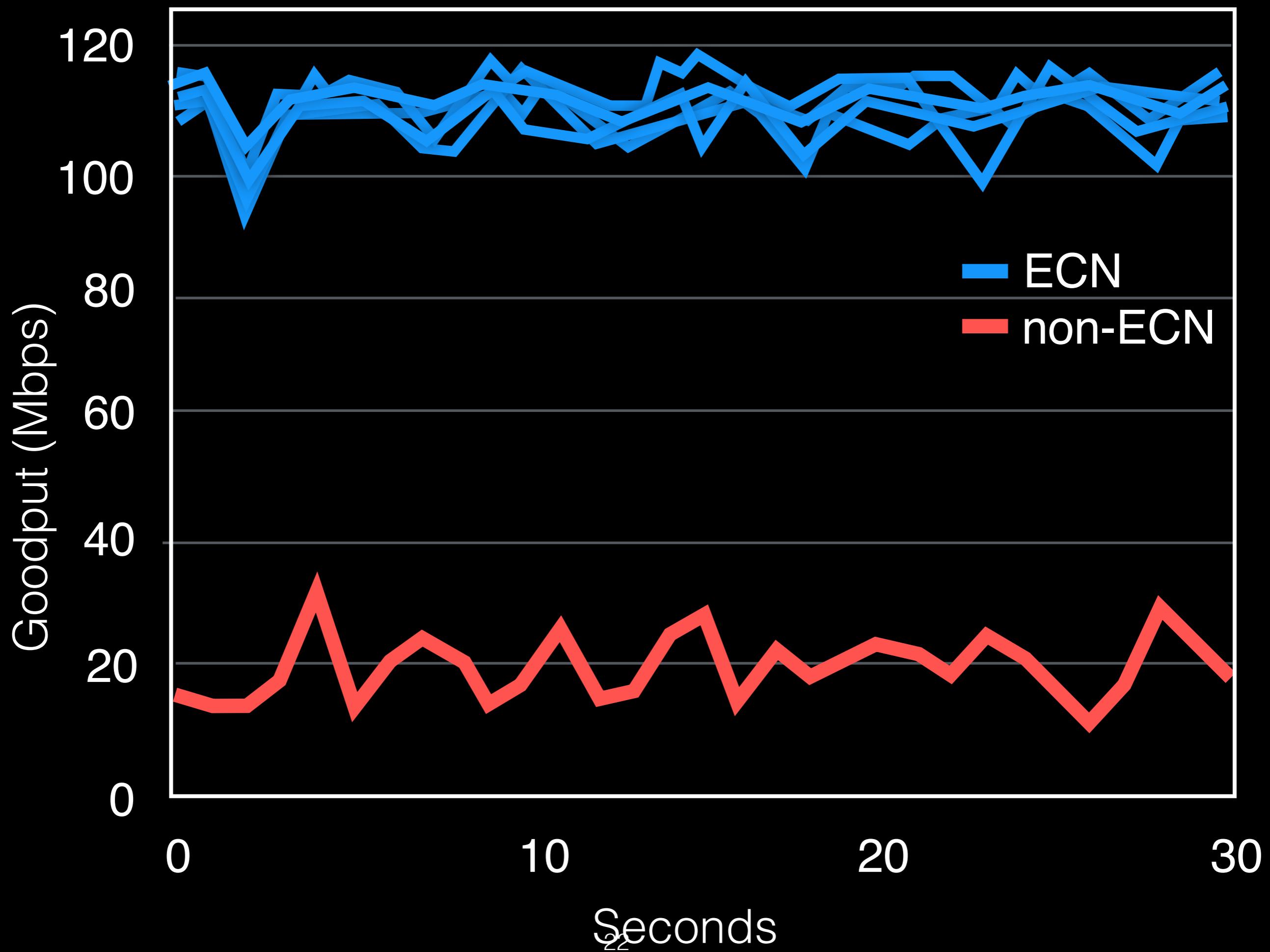
For networks with mechanisms for the detection of incipient congestion, the use of ECN mechanisms for the notification of congestion to the end nodes prevents unnecessary packet drops. For bulk-data connections, the user is concerned only with the arrival time of the last packet of data, and delays of individual packets are of no concern. For some interactive traffic, however, such as telnet traffic, the user is sensitive to the delay of individual packets. For such low-bandwidth delay-sensitive TCP traffic, unnecessary packet drops and packet retransmissions can result in noticeable and unnecessary delays for the user. For some connections, these delays can be exacerbated by a coarse-granularity TCP timer that delays the source's retransmission of the packet.

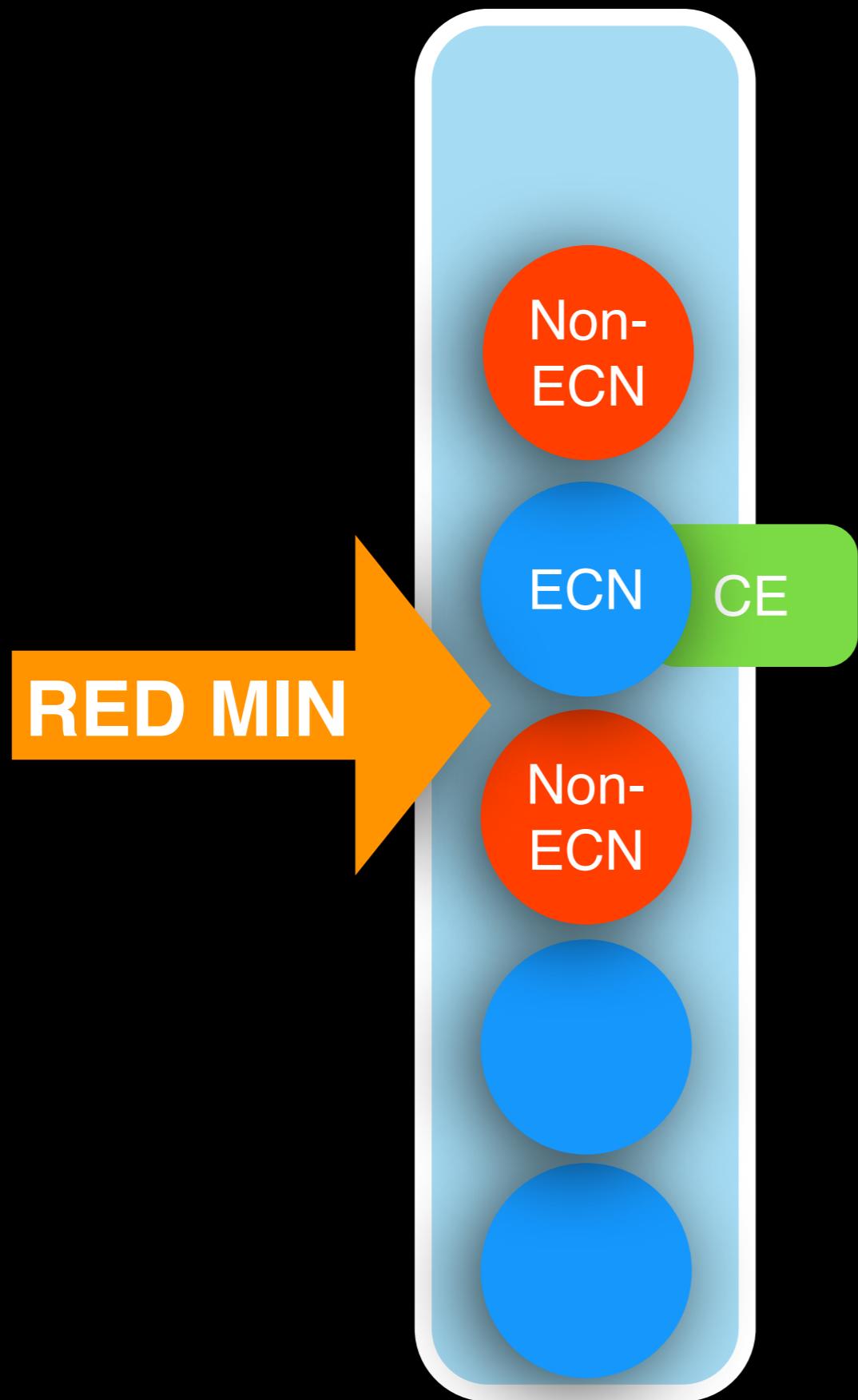
A second benefit of ECN mechanisms is that with ECN, sources can be informed of congestion quickly and unambiguously, without the source having to wait for either a retransmit timer or three duplicate ACKs to infer a dropped packet. For bulk-data TCP connections,



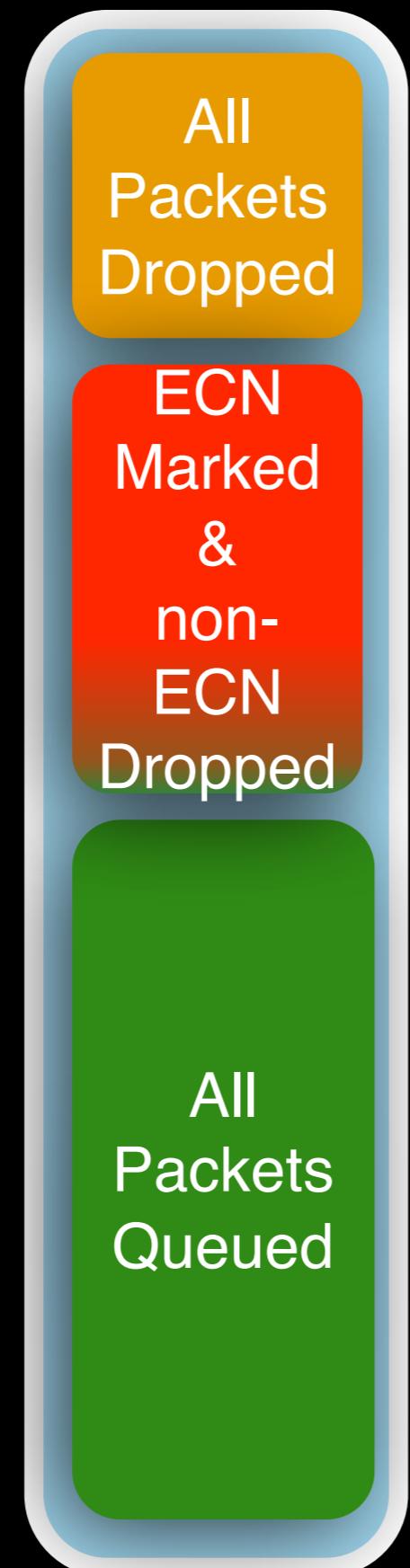








3 Dropped
2 Marked
1 Queued



Sender

Bottleneck Link

ECN

ECN

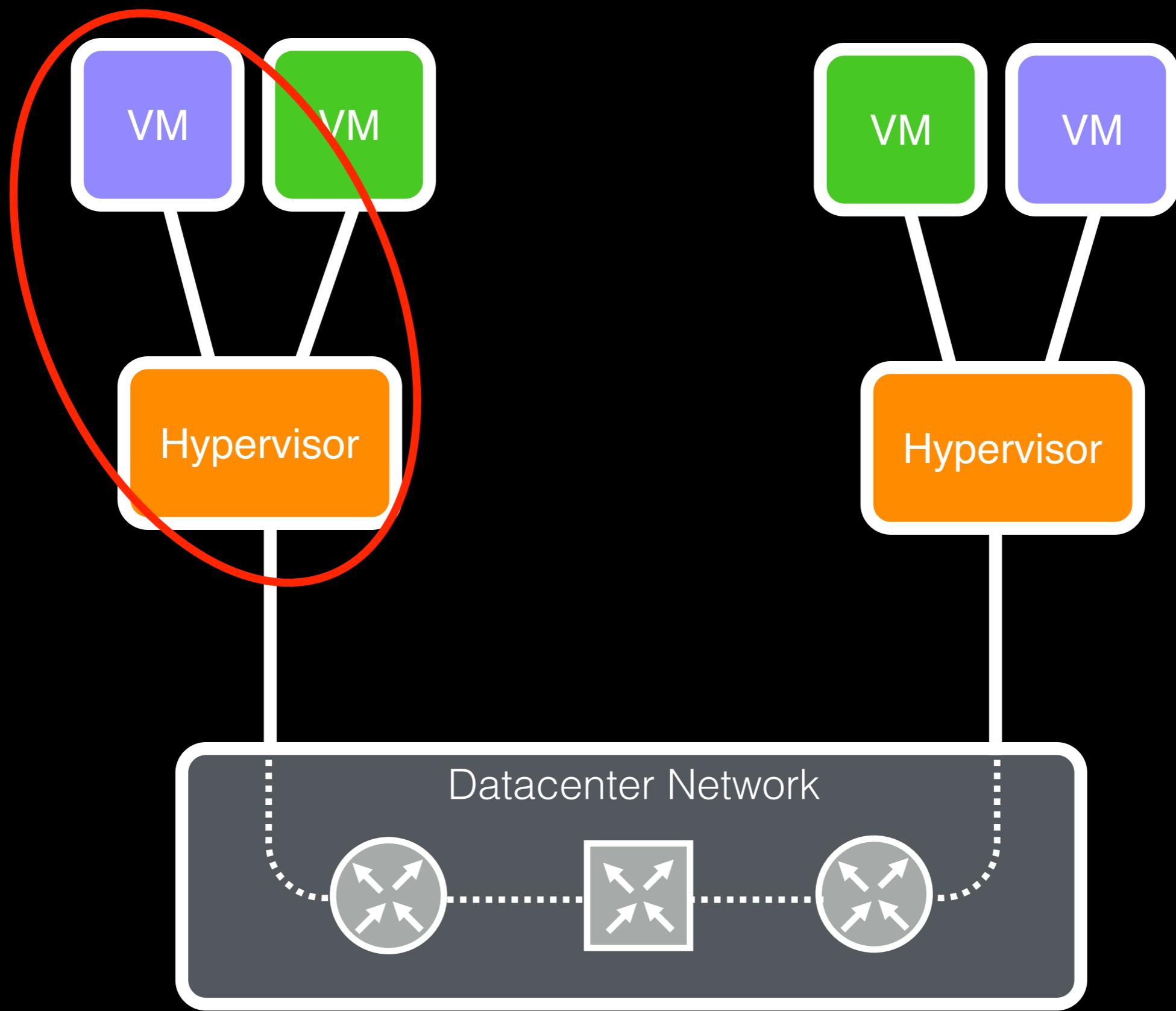
ECN

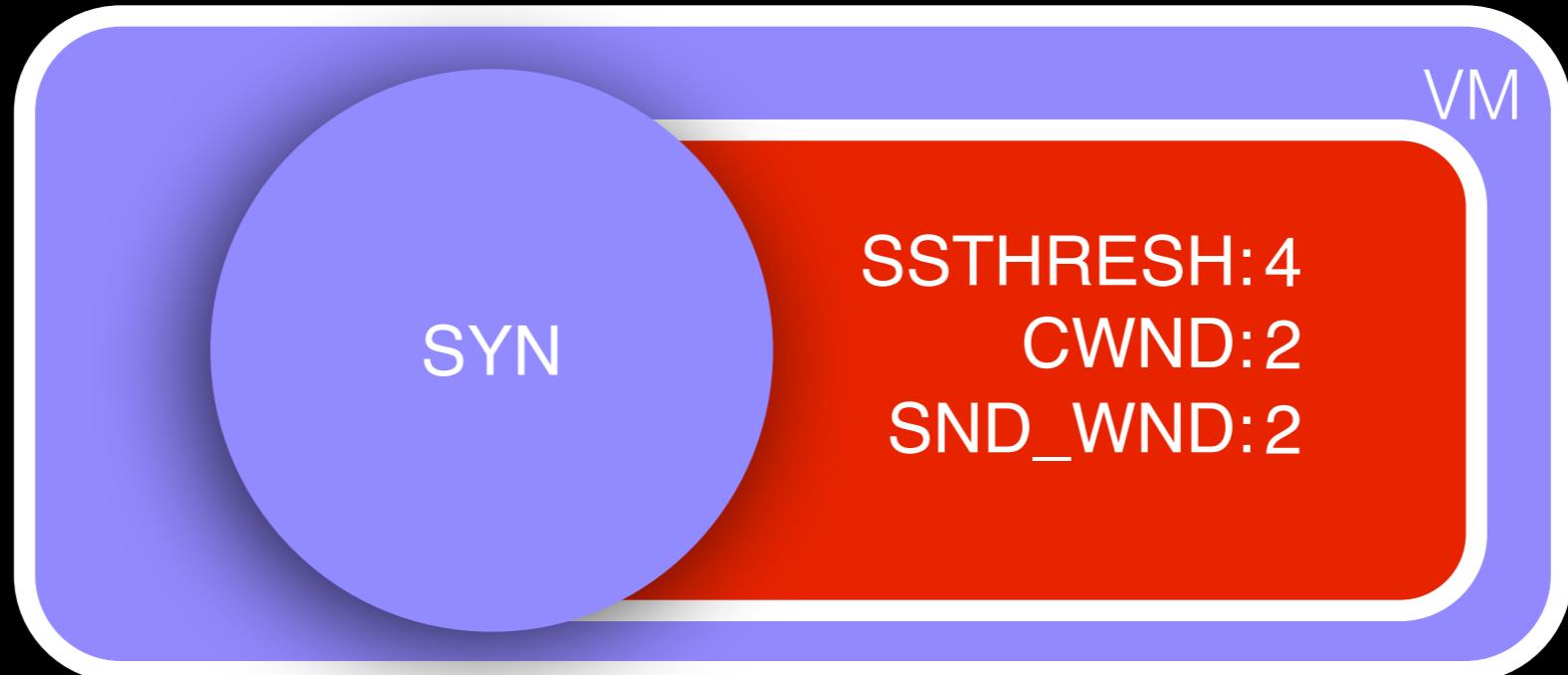
translated to ECN

Receiver

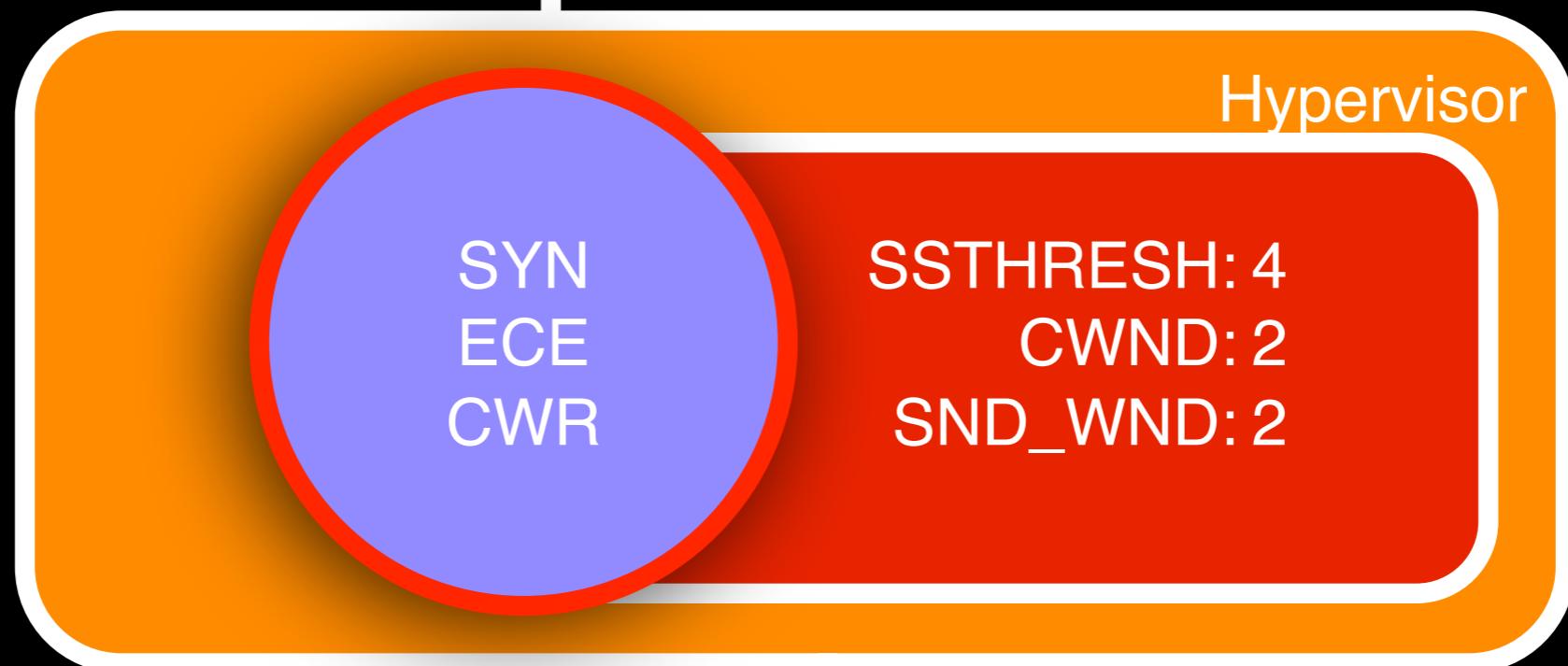
vCC

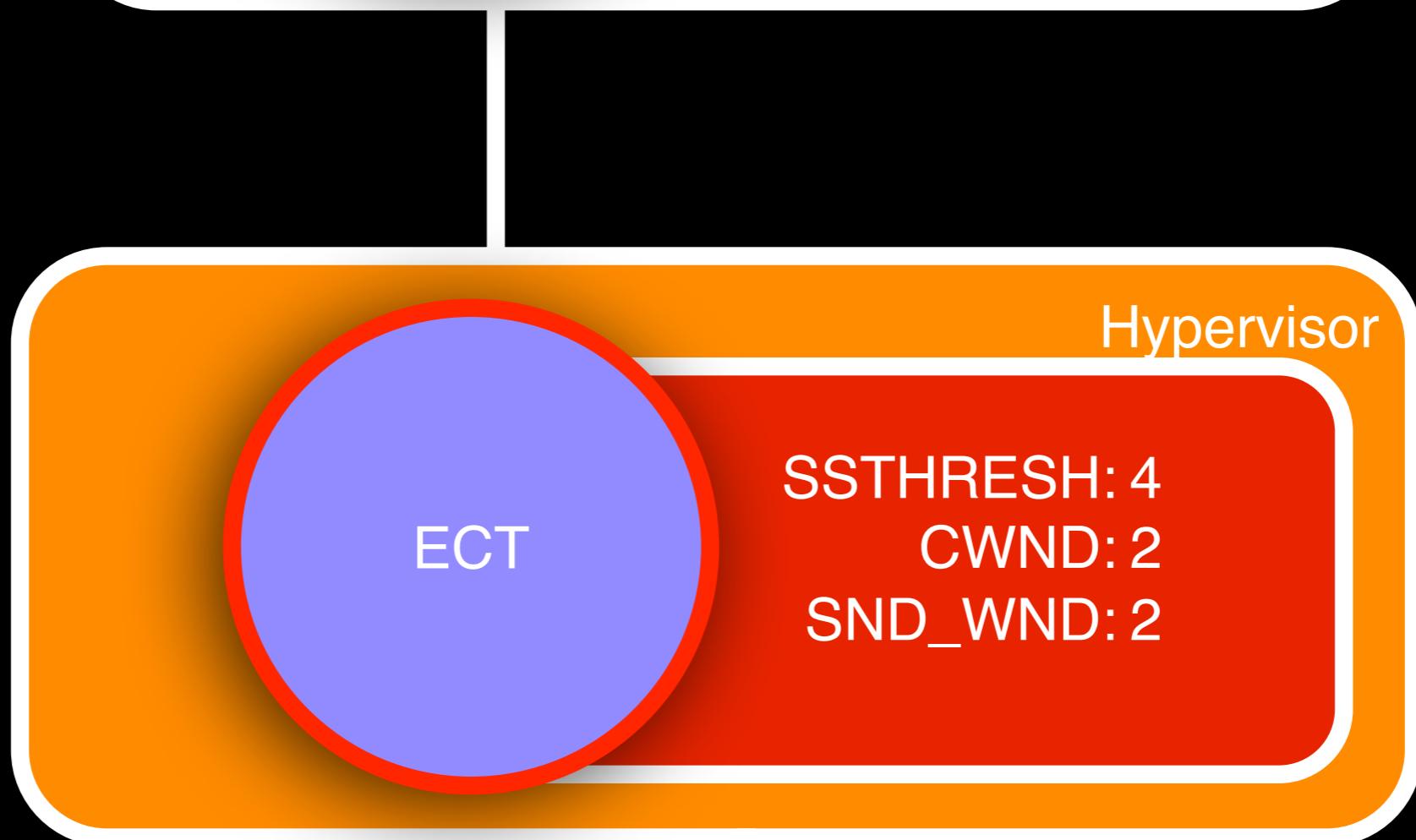
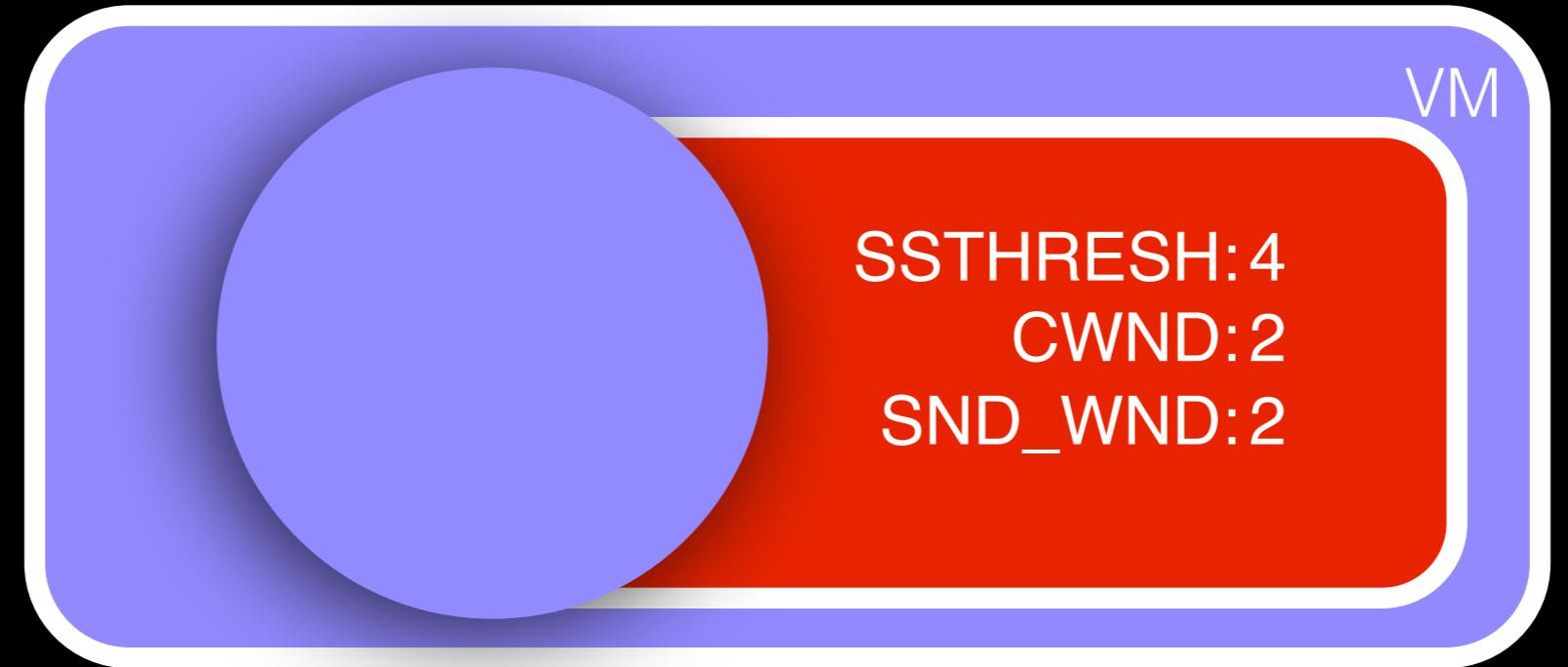
Virtualized Congestion Control

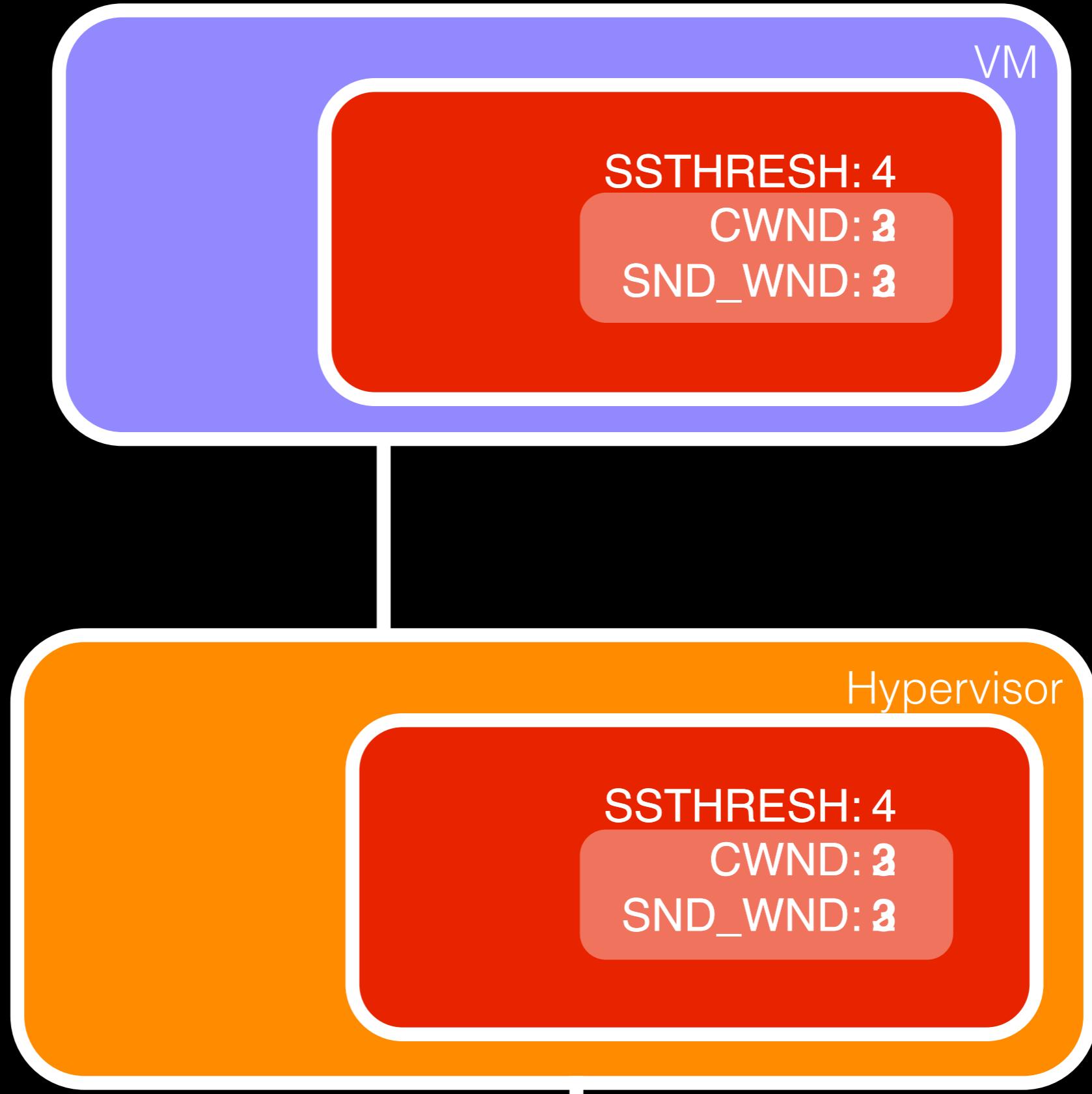


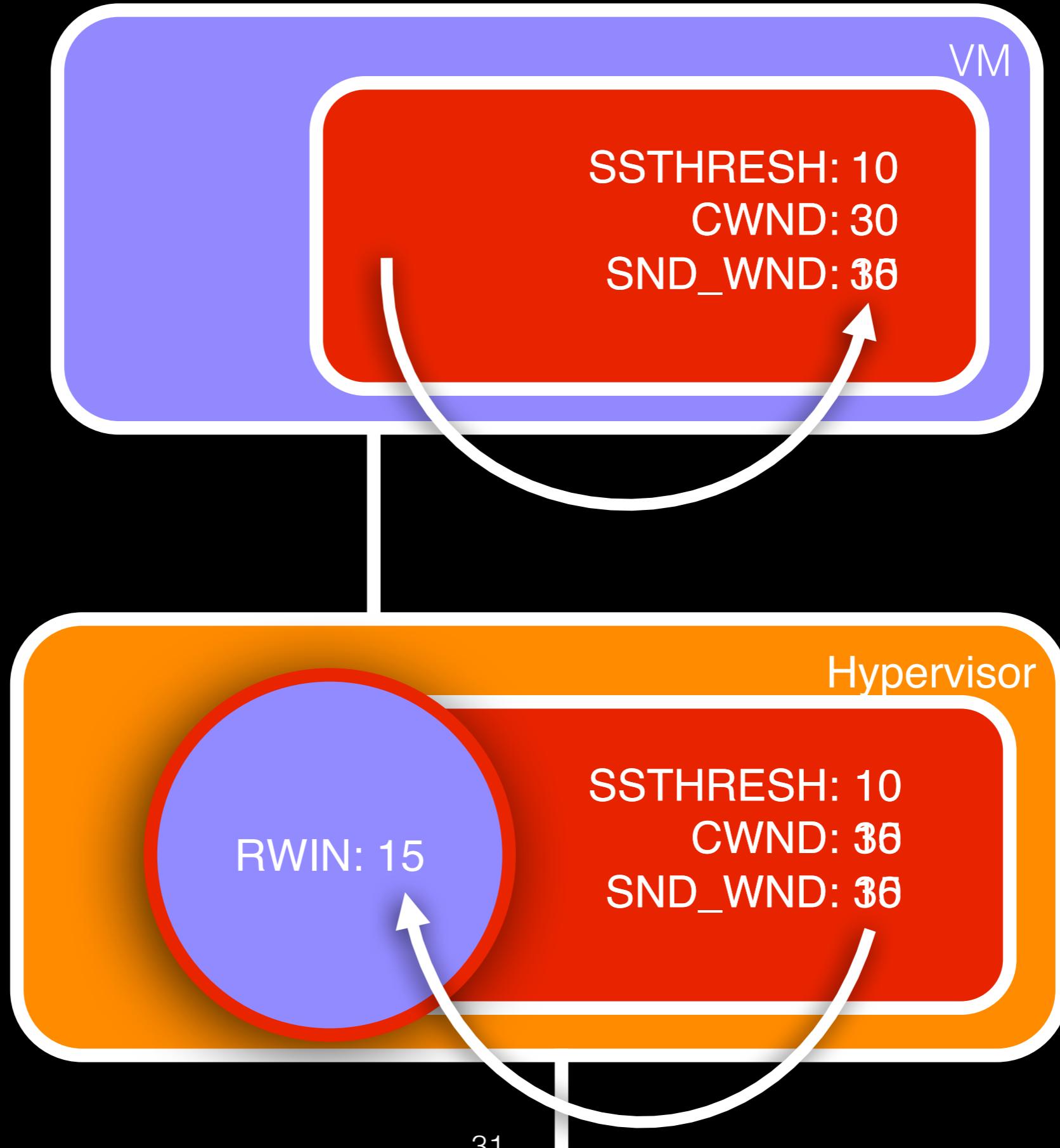


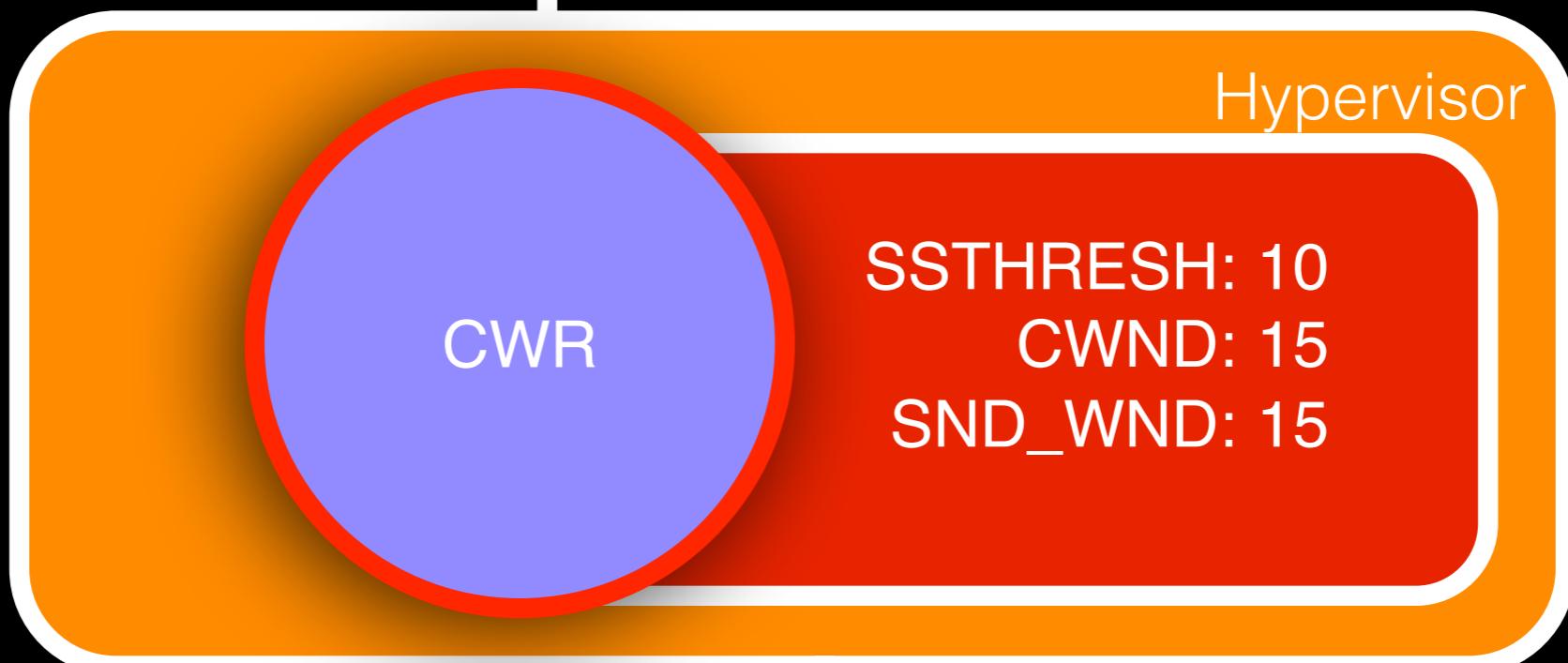
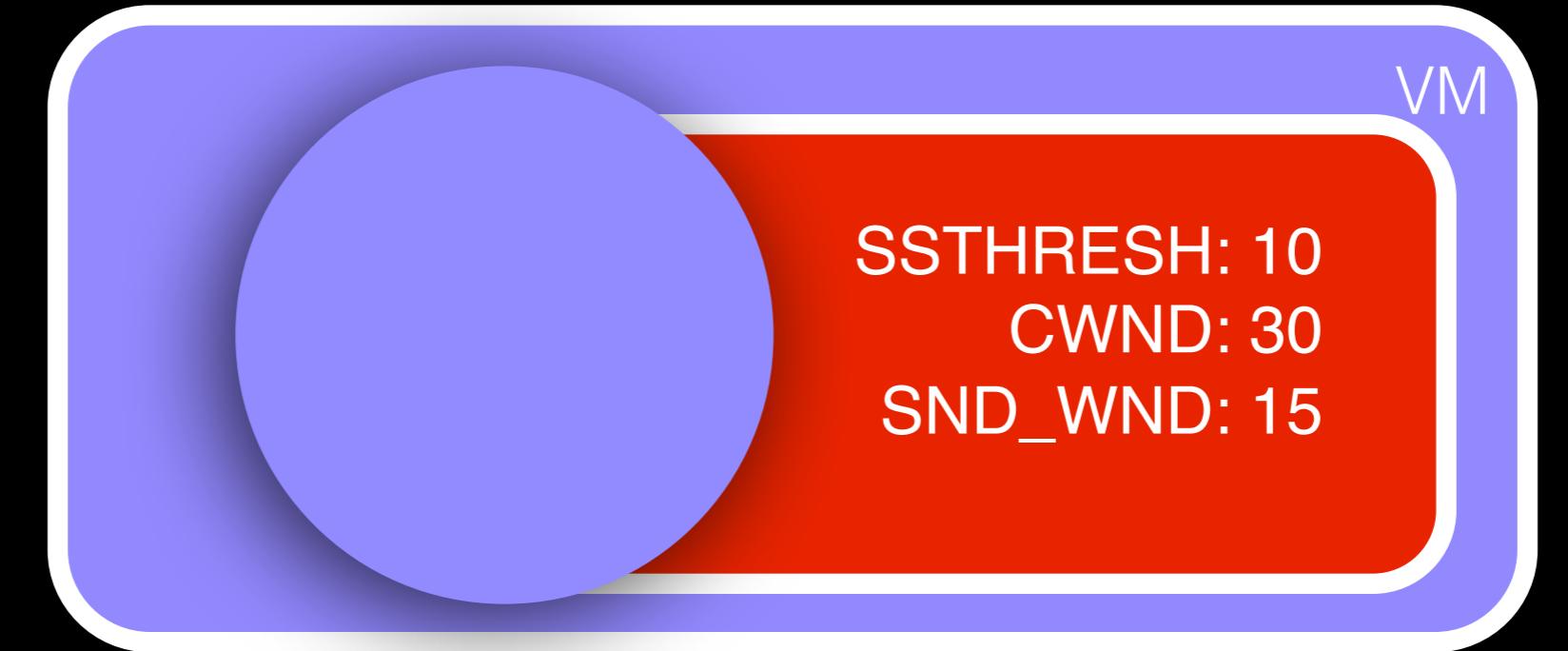
SSTHRESH: Slow Start Threshold
CWND: Congestion Window
SND_WND: Send Window

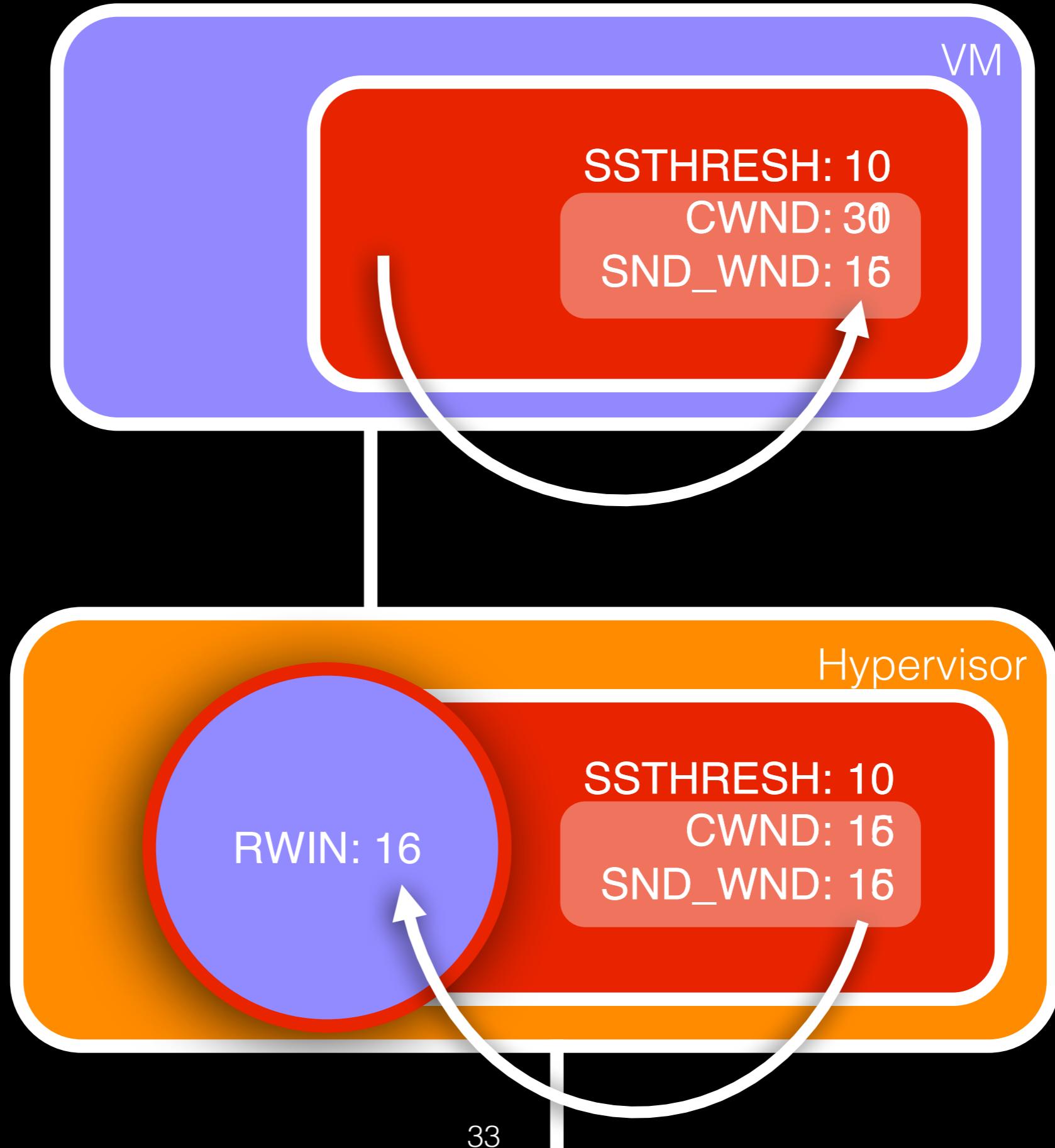


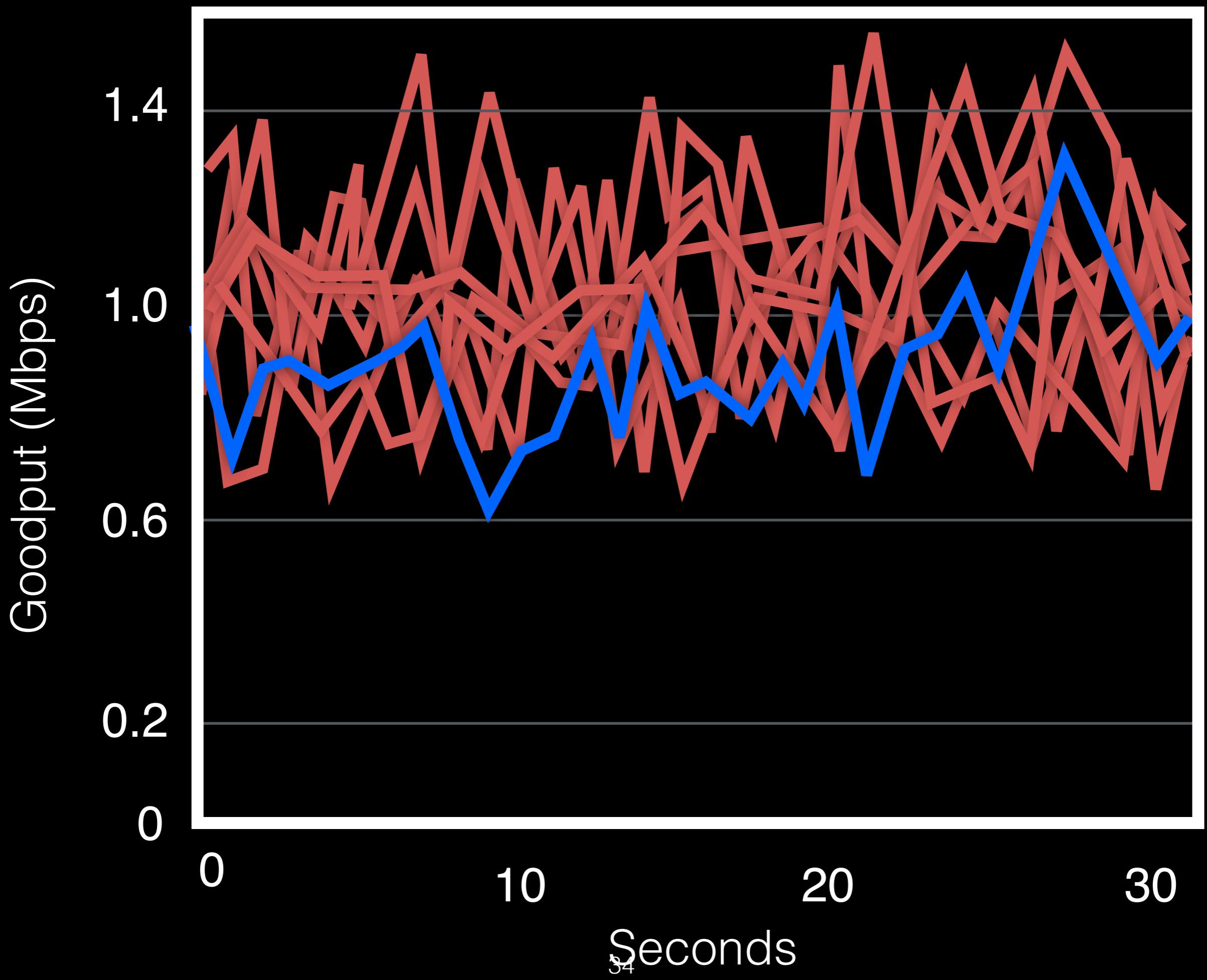


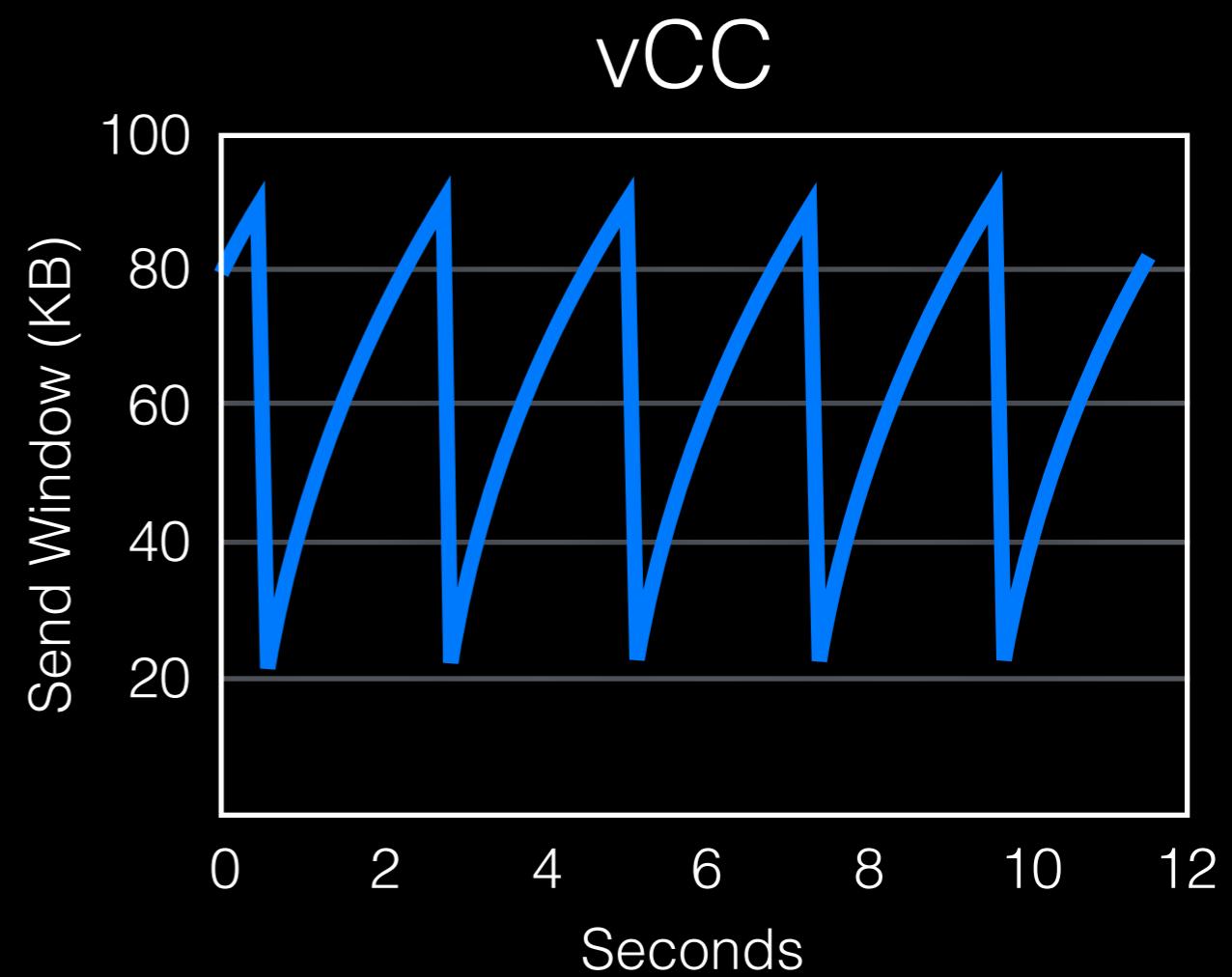
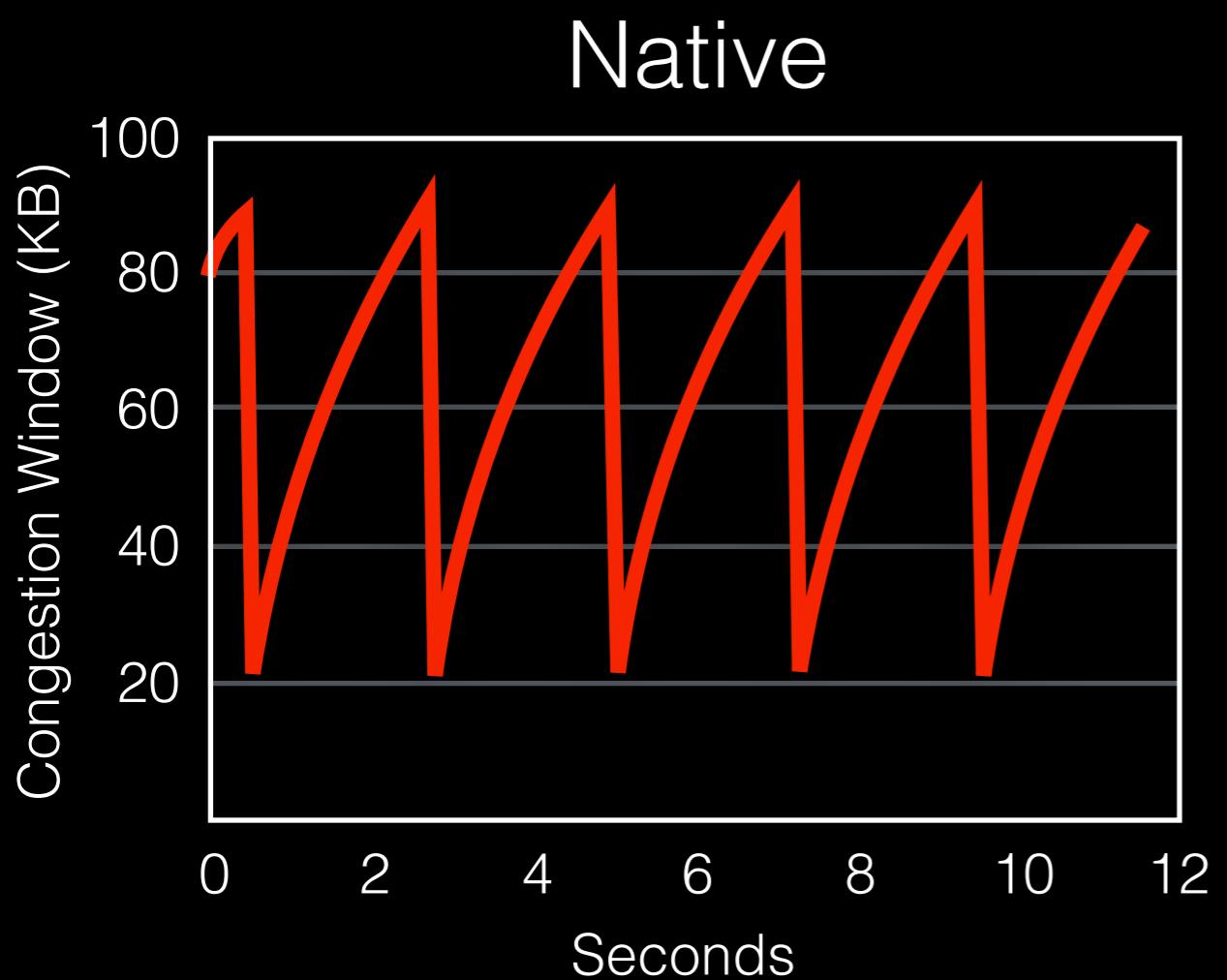


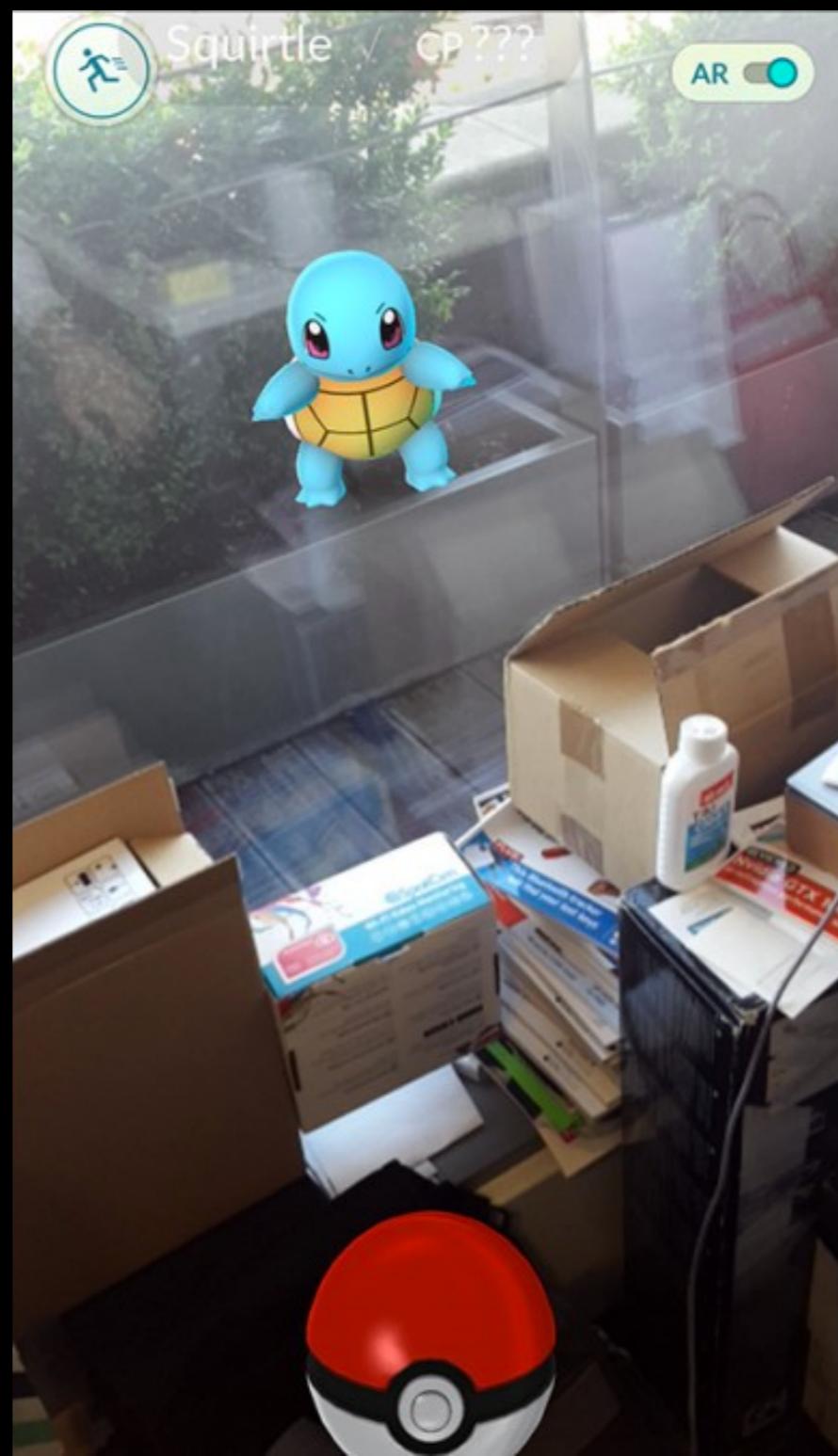


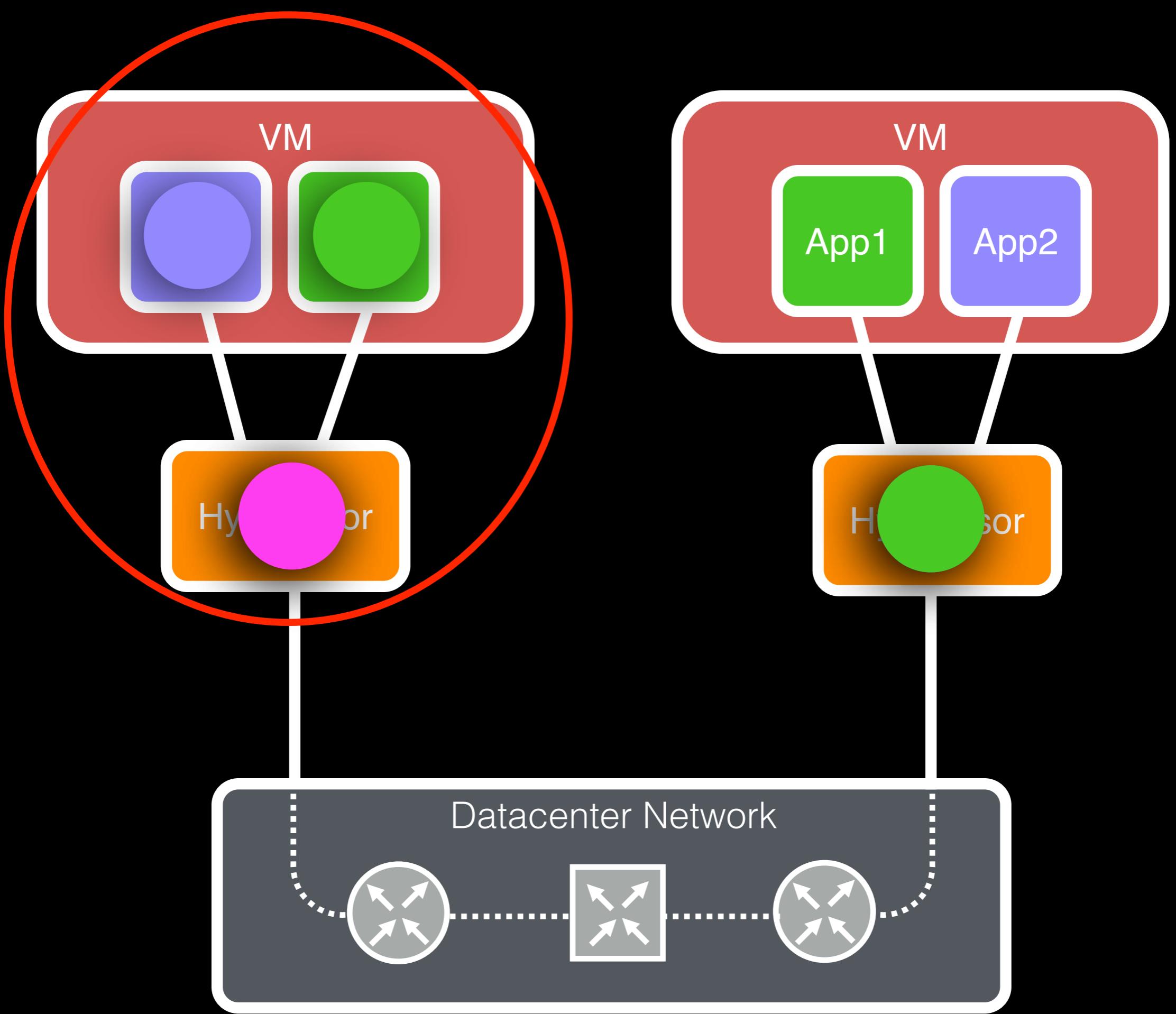


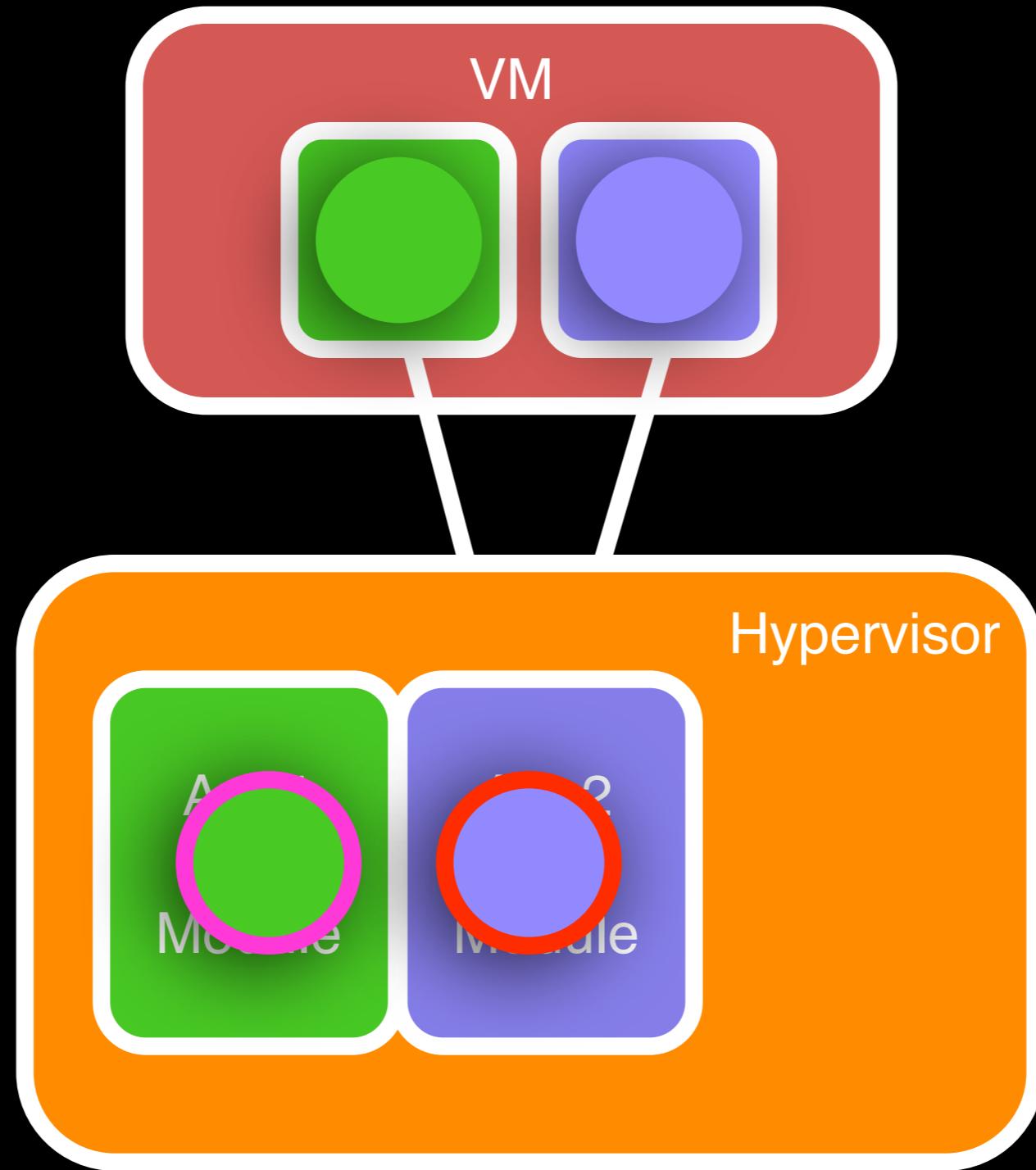












Throughput (% link capacity)

100

80

60

40

20

0

Unpreferred

Unpreferred

Unpreferred

Preferred

Preferred

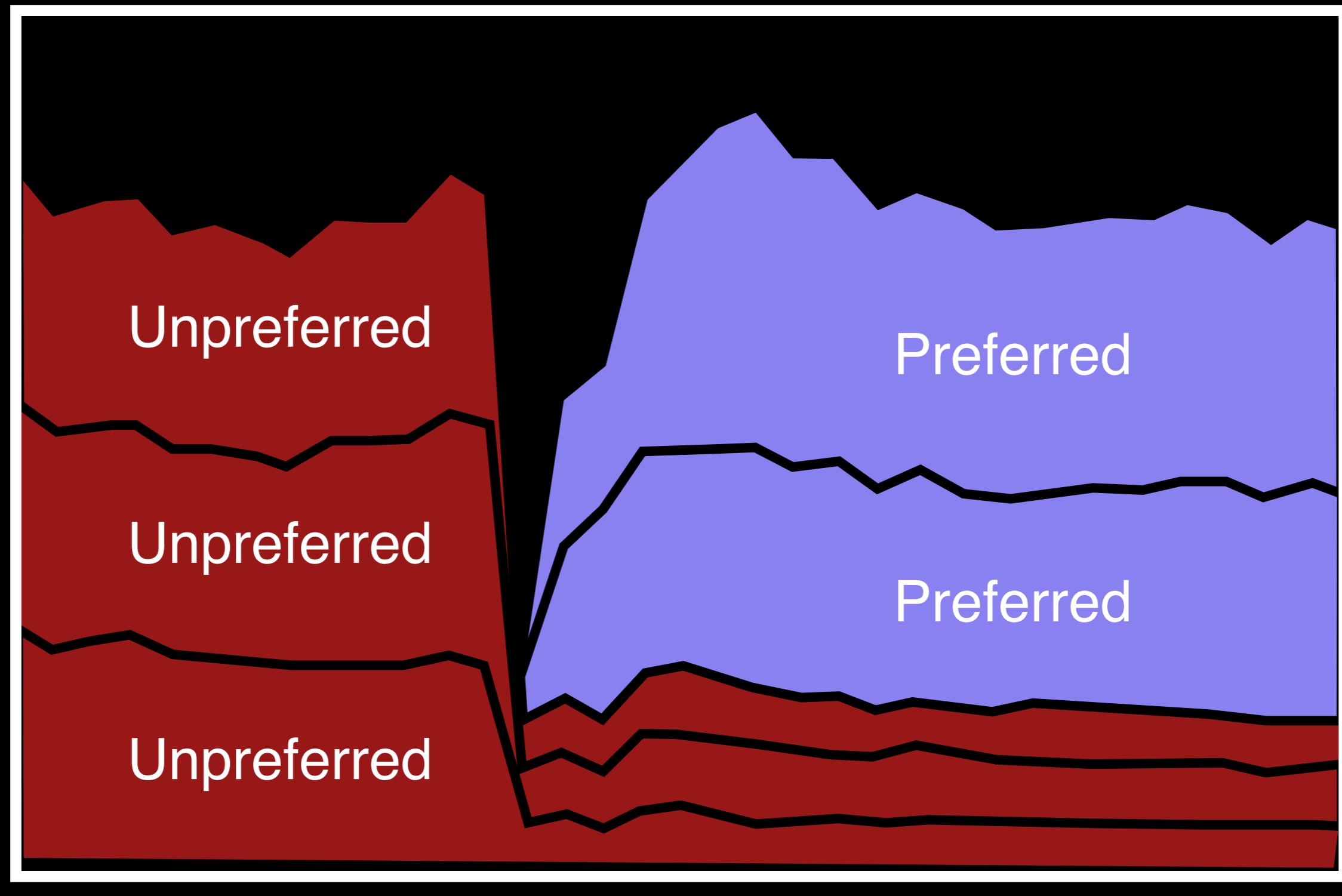
10

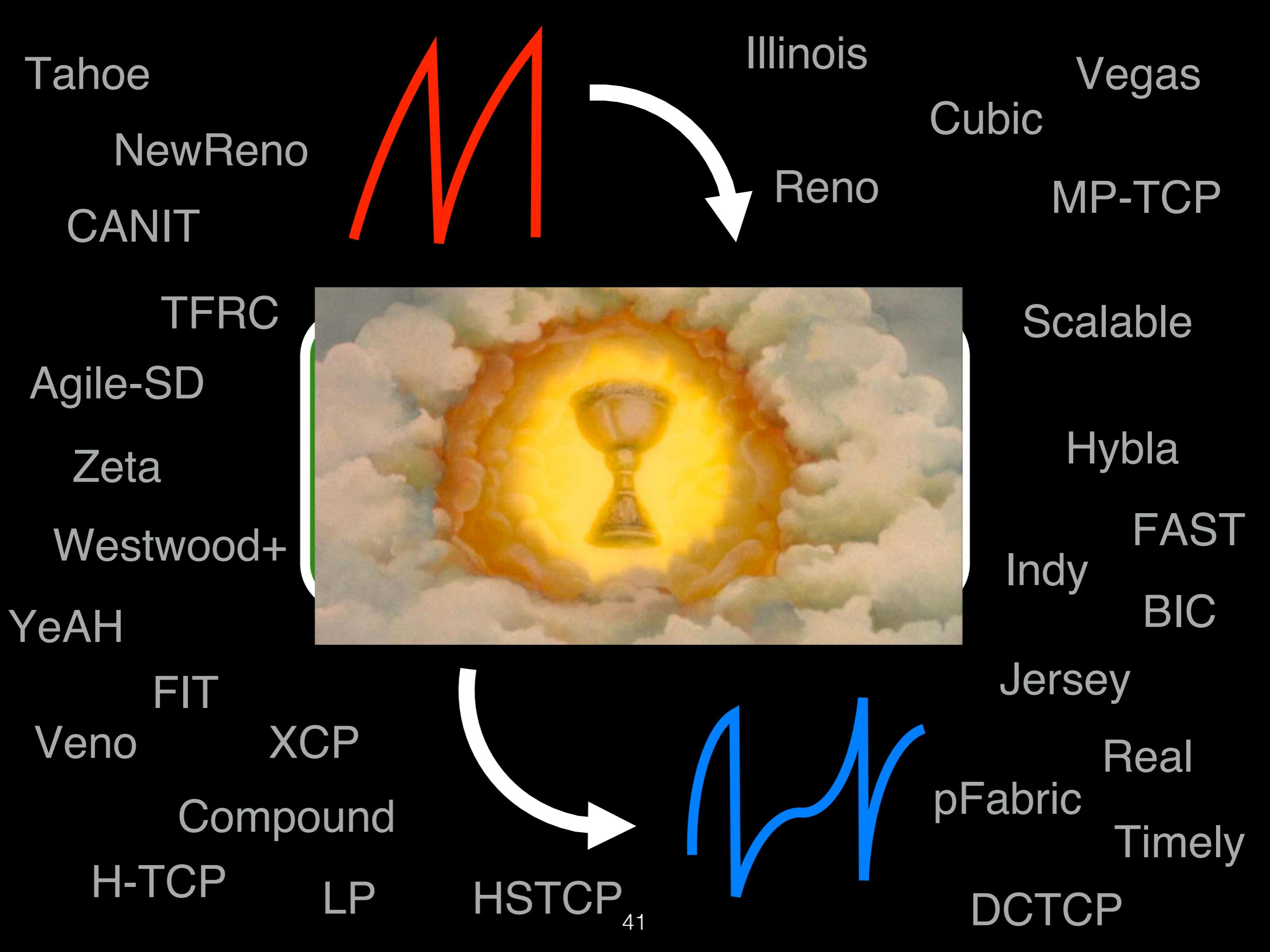
20

30

Time (seconds)

40





ECN Experiments in this Presentation:

Linux 3.19 mininet testbed
vCC implemented as a kernel patch

Public vCC Code:

Linux patch and mininet test suite available on Github
Public Code and Extended Paper:

<http://webee.technion.ac.il/~isaac/vcc/>

Our Paper:

ECN unfairness, virtual-ECN vCC, ESX app-level
throttling vCC, technique survey, future
implementation analysis

Thank You!



Stanford
PLATFORMLAB

 **TECHNION**
Israel Institute of Technology

vmware®
RESEARCH

Virtualized Congestion Control

Translate between Congestion Control Algorithms
in the Hypervisor

in order to

- Achieve uniform congestion control in multitenant and enterprise datacenters
- Programmatically assign congestion control to fine-grained flow signatures
- Simplify rollout of new congestion control algorithms