# Practical Packet Deflection in Datacenters

Sepehr Abdous, Erfan Sharafzadeh, Soudeh Ghorbani

JOHNS HOPKINS
U N I V E R S I T Y

December 2023

# Datacenter Traffic is Busty!



**2** *Bursts* occur for different **workloads**

**1** *Bursts* are manifested at **small** timescales!

Legend: Data Mining, Web, Cache

CDF vs Burst Duration (μs)

[Zhang et al., IMC '17]

JOHNS HOPKINS UNIVERSITY

# Existing Techniques Fall Short Against Bursts



Congestion control

TCP Sawtooth

4 Packets

DCN

# Existing Techniques Fall Short Against Bursts

Swift: Delay is Simple and Effective for Congestion Control in the Datacenter

Gautam Kumar, Nandita Dukkipati, Keon Jang (MPI-SWS)', Hassan M. G. Wassel, Xian Wu, Behnam Montazeri,

TIMELY: RTT-based Congestion Control for the Datacenter

Congestion control

**Congestion control protocols are slow against bursts!**

*Alibaba Group♠, Harvard University♢, University of Cambridge◇, Massachusetts Institute of Technology♣*

Data Center TCP (DCTCP)

Mohammad Alizadeh†‡, Albert Greenberg†, David A. Maltz†, Jitendra Padhye†, Parveen Patel†, Balaji Prabhakar‡, Sudipta Sengupta†, Murari Sridharan†

†Microsoft Research   ‡Stanford University
{albert, dmaltz, padhye, parveenp, sudipta, muraris}@microsoft.com
{alizade, balaji}@stanford.edu

Deadline-Aware Datacenter TCP (D²TCP)

Balajee Vamanan          Jahangir Hasan          T. N. Vijaykumar
Purdue University          Google Inc.              Purdue University
bvamanan@ecn.purdue.edu   jahangir@google.com      vijay@ecn.purdue.edu

**~ few milliseconds    >> 50 μs**

**2 Packets**

TCP Sawtooth

[Zhang et al., IMC '17]

# Existing Techniques Fall Short Against Bursts

Load balancing paradigms are ineffective against drops at the **last hop**.

Congestion control

Load balancing

Last-hop Congestion

From Spine Blocks

Edge Aggregation Block

ToR fanin drops (35.9%)

ToR uplink oversub drops (1.3%)

ToRs

Hosts

Host fanin drops (62.8%)

[Singh et al., SIGCOMM '12]

Sender

Receiver

# Existing Techniques Fall Short Against Bursts



Priority Spectrum

Low priority:

High priority:

Packet schedulers are unable to address **bursts of high-priority packets**.

Drop

Congestion control

Load balancing

Packet scheduling

Programmable Packet Scheduling with a Single Queue

Zhuolong Yu
Johns Hopkins University

Xiao Sun
Stony Brook University

Chuheng Hu
Johns Hopkins University

Vladimir Braverman
Johns Hopkins University

Zhenhua Liu
Stony Brook University

Jingfeng Wu
Johns Hopkins University

Mosharaf Chowdhury
University of Michigan

Xin Jin
Peking University

Carousel: Scalable Traffic

Ahmed Saeed*          Nandita Du
...ia Institute of Technology    Google, ...

PI²: A Linearized AQM
for both Classic and Scalable TCP

Koen De Schepper*   Olga Bondarenko   Ing-Jyh Tsang‡   Bob Briscoe
*Nokia Bell Labs, Belgium                    ‡Simula Research Laboratory, Norway
‡{koen.de_schepper|ing-jyh.tsang}@nokia.com   {olgabo|bob}@simula.no
Google, Inc.

...ialendar Queues for High-speed Packet Scheduling

Zhenxingyu Zhao*       Ming Liu*       Pravein G Kannan†       Changhoon Kim‡

nan§

Programmable Packet Scheduling at Line Rate

ivinay Subramanian*, Mohammad Alizadeh*, Sharad Chole‡, Shang-Tse Chuang†, Anurag Agrawal†,
Hari Balakrishnan*, Tom Edsall‡, Sachin Katti*, Nick McKeown*
*MIT CSAIL, †Barefoot Networks, ‡Cisco Systems, ‡Stanford University

JOHNS HOPKINS
UNIVERSITY

# **Packet Deflection** Avoids Drops in the Hotspots!

- Deflection: **Re-routing** packets that arrive at a full buffer to a neighboring switch.

Deflection improves the query performance by up to **43x**!     [Vertigo, CoNEXT '21]

State-of-the-art **deflection** proposals are **not implementable** in existing programmable hardware!

# State-of-the-art Deflection Depends on **Two** Primitives



Clean-slate deflection requirements in **real-time**:

I. **Filtering congested ports**
II. Extracting packets from the queue

Filtering the potential destination ports requires many comparisons ⚠️

PISA pipeline

Parser | Ingress processing | Traffic Manager | Egress processing

Only deflect to free buffers!

JOHNS HOPKINS UNIVERSITY

# State-of-the-art Deflection Depends on **Two** Primitives



Clean-slate deflection requirements in **real-time**:

I.       Filtering congested ports

II.     **Extracting packets from the queue**

Pop a low-priority packet from the buffer for deflection! ⚠

JOHNS HOPKINS
UNIVERSITY

## Our Contribution:
## **Implementing Deflection** in Programmable Hardware

- Implementing two approaches to deflection:
  - **Simple Deflection**
  - Approximation of **Selective Deflection** called **Preemptive Deflection**
- Intuitions:
  - Using **packet recirculation** instead of expensive memory manipulation.
  - Using **admission control** instead of packet extraction from the queue.

- Using DCTCP Congestion control
- 100 Gbps 8-ary fat-tree cluster
- Incast size of 100 Requests per Query

Practical Packet Deflection in Datacenters, CoNEXT 2023

# Approaches to Deflection Suited for Different Needs

Just deflect to a non-congested port, **randomly**!

Simple

Selective

**Sort** packets in the queue and deflect lower-priority packets!

Priority Spectrum

Low priority          High priority

JOHNS HOPKINS
UNIVERSITY

# What Makes **Simple Deflection** Hard to Implement ❓



Queue Utilization Data

Selected Egress Deflection Port

Randomly Select

Packet Ingress

Congested Original Destination

✓ Simple Deflection is effective when congestion is infrequent
(71% lower QCT under moderate Incast)

CHALLENGE 2: **Filtering** non congested ports is non-trivial and computationally **expensive**!

CHALLENGE 1: **Queue utilization data** is not always available behind the traffic manager!

JOHNS HOPKINS
U N I V E R S I T Y

# Implementable Simple Deflection in PISA



Queue Utilization Data

Collect

Queue Utilization Data

Recirculate to Ingress

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 |

R=2

R=8

1. **Syncing** the egress and ingress pipelines with worker packets.
2. Selecting a non-congested port uniformly at random from **queue utilization bitmap**.

# What Makes **Selective Deflection** Hard to Implement ?



## Priority Spectrum

Low priority:
Longest remaining time

High priority:
Shortest remaining time

Packet Ingress

## Operation Steps

1. **High-priority** packet arrives
2. Extract and **deflect** a **low-priority** packet

JOHNS HOPKINS
UNIVERSITY

# What Makes **Selective Deflection** Hard to Implement ？

**Priority Spectrum**

Low priority:
Longest remaining time

High priority:
Shortest remaining time

Packet Ingress

Operation Steps

1. **High-priority** packet arrives
2. Extract and **deflect** a **low-priority** packet

# What Makes **Selective Deflection** Hard to Implement ❓



Operation Steps

1. **High-priority** packet arrives
2. Extract and **deflect** a **low-priority** packet
3. **Insert** the **new packet** at the head of its queue
4. The deflected packet experiences **extra hop latency** instead of retransmission!

# What Makes **Selective Deflection** Hard to Implement ?



Priority Spectrum

Low priority:
Longest remaining time

High priority:
Shortest remaining time

Packet Ingress

**CHALLENGE 1:** Switch traffic manager is limited to **FIFO** queues

**CHALLENGE 2: Packet extraction** is not feasible after inserting the packet into a queue

Operation Steps

1. **High-priority** packet arrives
2. Extract and **deflect** a **low-priority** packet
3. **Insert** the **new packet** at the head of its queue
4. The deflected packet experiences **extra hop latency** instead of retransmission!

Selective Deflection is effective, even under **extreme load**!

(up to **43x** lower query completion times)

[Vertigo, Conext '21]

Can we use **existing traffic management capabilities** to approximate Selective Deflection?

JOHNS HOPKINS
U N I V E R S I T Y

# Implementable Selective Deflection in PISA

Can we use existing traffic management capabilities to approximate **Selective Deflection**?

Yes, with **Preemptive Deflection!**

Admission policy: Should I **admit** the packet that just arrived?

Compare its priority to a **representation** of existing packets

Queue Capacity

Relative Packet Priority

Sensitivity Knob

Deflection Threshold =
($\tau$)
× [1 −                          × ]

JOHNS HOPKINS
UNIVERSITY

# Implementable Selective Deflection in PISA

Can we use existing traffic management capabilities to approximate **Selective Deflection**?

Yes, with **Preemptive Deflection!**

Admission policy: Should I **admit** the packet that just arrived?

Compare its priority to a **representation** of existing packets

Quantile-based

**1**

> > > > > ≤

$\tau = 8 \times [1 - 0.5 \times \frac{1}{6}] = $ **7.3** ✓ ADMIT

| 4 | 5 | 2 | 3 | 2 | 1 |

**4**

≤ > ≤ ≤ ≤ ≤

$\tau = 8 \times [1 - 0.5 \times \frac{5}{6}] = $ **4.6** ↩ DEFLECT

Queue Size     Sensitivity Parameter

# Implementable Selective Deflection in PISA

Can we use existing traffic management capabilities to approximate **Selective Deflection**?

Yes, with **Preemptive Deflection!**

Admission policy: Should I **admit** the packet that just arrived?

Compare its priority to a **representation** of existing packets

Quantile-based

Distribution-based

✓ Implementable using pre-filled tables
✓ Takes only 2 processing stages!
- Less accurate than quantile-based deflection

✓ Implementable using a sample rolling window for arriving packets
- Requires consecutive comparisons in limited PISA stages.

JOHNS HOPKINS
UNIVERSITY

# Putting it All Together

**Practical Deflection in Datacenters**

**Simple Deflection**

| Packet **recirculation** for syncing queue utilization data | **Bitmaps** for randomly selecting non-congested ports |
| --- | --- |

✓ Effective when congestion is infrequent

✓ Requires minimal resources, no external input

**Quantile-based Preemptive Deflection**

| **Admission vs deflection** policy on FIFO queues | **Priority comparators** parallelized in PISA pipelines |
| --- | --- |

✓ Can Handle Extreme degrees of congestion

✓ Accurately approximates selective deflection

**Distribution-based Preemptive Deflection**

Using tables prefilled with **statistical distribution** mean-values

✓ Can handle large degrees of congestion

✓ Requires few processing stages in Tofino

JOHNS HOPKINS
UNIVERSITY

# Implementable Deflection Improves the Performance

# Implementable Deflection Under Large-scale Incast

- Large-scale OMNET simulations
- 2-tier leaf-spine with 100 Gbps links
- 40 machines with all-to-all traffic
- Swift Congestion control

**Simple Deflection** is effective under moderate congestion!

### 50% Background + 5% Incast

### 50% Background + 35% Incast



Legend:
- ECMP (Baseline)
- Practical Simple Deflection
- Clean-slate Simple Deflection

JOHNS HOPKINS UNIVERSITY

# Implementable Deflection Under Large-scale Incast

- Large-scale OMNET simulations
- 2-tier leaf-spine with 100 Gbps links
- 40 machines with all-to-all traffic
- Swift Congestion control

**Simple Deflection** is effective under moderate congestion!

**Preemptive Deflection** offers **superior** performance under **large Incast**



50% Background + 5% Incast

50% Background + 35% Incast

21x lower QCT

6x lower QCT

Query Completion Time (s)

Query Completion Time (s)

CDF

ECMP (Baseline)

Practical Simple Deflection

Clean-slate Simple Deflection

Quantile-Preemptive Deflection

Selective Deflection

JOHNS HOPKINS
U N I V E R S I T Y

# Implementable Deflection Under Large-scale Incast

- Large-scale OMNET simulations
- 2-tier leaf-spine with 100 Gbps links
- 40 machines with all-to-all traffic
- Swift Congestion control

**Simple Deflection** is effective under moderate congestion!

**Preemptive Deflection** offers **superior** performance under **large Incast**

**Early drop** and packet prioritization **alone cannot recover** from short-term local bursts



50% Background + 5% Incast

50% Background + 35% Incast

3x higher QCT!

CDF

Query Completion Time (s)

Query Completion Time (s)

······· ECMP (Baseline)

—— Practical Simple Deflection

— — Clean-slate Simple Deflection

—■— Quantile-Preemptive Deflection

— — Selective Deflection

—— Admission Control [AIFO, SIGCOMM '21]

JOHNS HOPKINS
UNIVERSITY

# We Made Packet Deflection Practical

- We propose an accurate implementation of **Simple Deflection** on PISA architecture.

- We introduce **Preemptive Deflection**, an approximation of selective deflection on PISA.

- Choosing among deflection techniques depends on:

  - Network utilization & congestion intensity

  - Resource availability

  - Performance requirements

- Preemptive Deflection improves high-priority Flow Completion Times by **425x** in a physical testbed.

- Visit https://hopnets.github.io/practical_deflection for the codebase

- Contact Authors: sabdous1@jhu.edu, erfan@cs.jhu.edu

Illustrations generated using Bing AI

JOHNS HOPKINS
UNIVERSITY

# Backup slides