



Thus, DCTCP avoids long-term congestion, while DIBS ensures that there is no packet loss during large and transient bursts created by severe incast-like traffic. DIBS can also be paired with other congestion control schemes such as QCN [4] as well as the classic Random Early marking (RED/REM).

**Spare capacity:** If all links in the network are fully utilized, DIBS will not work, since there will be no capacity to handle the extra load caused by detoured packets. However, in DCNs, congestion is usually transient and local [5]. Thus, when we detour packets away from one congestion hotspot, the network will typically have capacity elsewhere to handle the traffic. We are working to derive a limit on network utilization, beyond which detouring can be detrimental.

We note that DIBS is particularly suited for deployment in DCNs. Many popular DCN topologies offer multiple paths [2], which detouring can effectively leverage. The link bandwidths in DCNs are very high and the link delays are quite small. Thus, the additional delay of a detour is quite low. Current DCN switches do not implement the random detouring, but the change required to do so is minimal. Indeed, our NetFPGA implementation is less than 25 lines of code. DIBS does not come into play until there is extreme congestion – it has no impact whatsoever when things are “normal”.

## 2. DISCUSSION

**Packet reordering:** Detouring will cause packet reordering, which can adversely impact TCP’s performance. A simple fix is to disable fast retransmission on end hosts.

**CIOQ switches:** Many modern switches use a combined input-output queued architecture (CIOQ), with shared buffer to store packets. DIBS can be implemented on these switches quite easily by defining a threshold. The switch fabric controller keeps track of packets on a per-port basis. When the number of packets buffered for an output port exceeds the threshold, the switch detours subsequent to other ports.

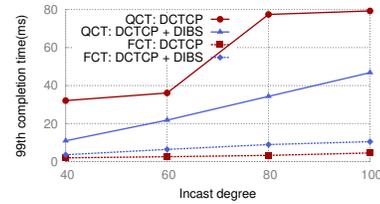
**L2 routing:** DIBS cannot work with self-learning L2 routing schemes such as the spanning tree protocol. However, most modern DCN topologies such as FatTree rely on pre-computed routing tables.

**Rogue flows and congestion collapse:** Rogue flows that do not implement congestion control can degrade performance of any network. DIBS may exacerbate this problem by constantly detouring packets of such flow. Today, DCNs use a variety mechanisms to detect and deal with such rogue flows. We expect that these mechanisms will be sufficient for our purposes as well. We are studying this issue in more detail. In the absence of such rogue flows, and as long as DIBS is paired with a congestion control protocol such as DCTCP, DIBS will not lead to congestion collapse.

**Collateral damage:** When congestion occurs at a particular switch, DIBS detours excess packets to other switches. These detoured packets interfere with other flows that would have otherwise not have been impacted by the traffic at the congested switch. We term this collateral damage. Our simulations show that collateral damage is low as long as a DIBS is used with a scheme like DCTCP and there is sufficient spare capacity. We are currently exploring the extreme conditions where collateral damage can be significant.

**Routing loops:** DIBS can appear to cause routing loops, as seen in Figure 1. However, these are merely temporary artifacts and do not affect actual routing.

**Early results:** We have implemented a preliminary version of DIBS in a NetFPGA router, in a Click modular router, and in NS-



**Figure 2: Mixed traffic: Variable incast degree. Compared to DCTCP alone, DIBS+DCTCP improves query completion time (QCT), with little collateral damage to background flows (FCT). (Background inter-arrival time: 120ms; query arrival rate: 300 qps; response size: 20KB)**

3. Our initial experiments are quite encouraging. Here we present only a sample result.

Figure 2 shows performance of DIBS with a complex mix of traffic, derived from traffic traces from a large search engine [3]. The traffic consists of both query traffic, which causes incast flows, and background flows. As the degree of incast increases, the 99th percentile of query completion time (QCT) is much worse with DCTCP alone. DIBS improves performance by avoiding packet losses. At the same time, the collateral damage to background flows (FCT) is low, since the network has spare capacity.

**Ongoing work:** We are currently evaluating the performance of DIBS with a variety of traffic patterns using detailed simulations and experiments. We are also considering modifications to the basic DIBS concept such as limiting the number of times a packet can be detoured, probabilistically detouring, and priority-based detouring.

**Related work:** DIBS can be implemented alongside traffic spreading schemes such as ECMP (flow or packet level) or MPTCP [6]. These schemes reduce the possibility of extreme congestion, but do not eliminate it. Technologies such as Ethernet flow control and its modern variant, Priority Flow Control (PFC) [1], as well as the Infiniband link layer, guarantee loss-free L2 networks. Unlike these technologies, DIBS does not guarantee a loss-free L2 network. However, compared to these technologies DIBS is significantly easier to configure [1]. Also, unlike Ethernet flow control, DIBS is deadlock free [1]. We are exploring this relationship further. [7] proposed reducing TCP RTomin value to mitigate the impact of packet losses. This proposal is orthogonal to DIBS, and the two can be implemented together. Centralizing scheduling and rate allocation [8] can *in theory* avoid all packet losses. However, in practice, such systems face scalability issues.

## 3. REFERENCES

- [1] Priority flow control. [http://www.cisco.com/en/US/prod/collateral/switches/ps9441/ps9670/white\\_paper\\_c11-542809.pdf](http://www.cisco.com/en/US/prod/collateral/switches/ps9441/ps9670/white_paper_c11-542809.pdf).
- [2] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. In *SIGCOMM*, 2008.
- [3] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan. Data center TCP (DCTCP). In *SIGCOMM*, 2010.
- [4] M. Alizadeh, A. Kabbani, B. Atikoglu, and B. Prabhakar. Stability analysis of QCN: The averaging principle. In *SIGMETRICS*, 2011.
- [5] S. Kandula, J. Padhye, and P. Bahl. Flyways to de-congest data center networks. In *HotNets*, 2009.
- [6] C. Raiciu, S. Barré, C. Pluntke, A. Greenhalgh, D. Wischik, and M. Handley. Improving datacenter performance and robustness with multipath TCP. In *SIGCOMM*, 2011.
- [7] V. Vasudevan, A. Phanishayee, H. Shah, E. Krevat, D. Andersen, G. Ganger, G. Gibson, and B. Mueller. Safe and effective fine-grained TCP retransmissions for datacenter communication. In *SIGCOMM*, 2009.
- [8] D. Zats, T. Das, P. Mohan, D. Borthakur, and R. H. Katz. DeTail: reducing the flow completion time tail in datacenter networks. In *SIGCOMM*, 2012.