

15-712 Project Proposal

Introduce Indirection Scheduling to Hybrid Switch

Conglong Li, Anbang Zhao
{conglonl, anbangz}@andrew.cmu.edu

1 Introduction

Traditionally, circuit switching and packet switching have been considered as alternative design choices, each with their own advantages and disadvantages. Packet switches are efficient at multiplexing traffic across a large number of ports, while circuit switches can often service higher line rates in a more cost-effective manner, especially when combined with optical links. However, as optical circuit switching technology advances, the reconfiguration time (i.e., the amount of time it takes to alter the input-to-output mapping) is swiftly decreasing, thus increasingly blurring the division between the packet and circuit regimes.

Indeed, a range of new datacenter switch designs take advantage of this trend, and propose to schedule appropriately large traffic demands via a high-bandwidth circuit switch and handle any remaining traffic with a slower packet switch. However, all recent proposals for such hybrid designs presume the existence of an omniscient scheduling oracle that can compute switch configurations and map traffic to them in an optimal fashion. Recently, Liu *et al.* [6] recently proposed a hybrid switch scheduling algorithm, Solstice, that it is both highly effective at scheduling datacenter-like traffic workloads and has practical computational overheads when doing so.

In this proposal, we propose to improve the Solstice scheduling algorithm by introducing an indirection heuristic that reduce the number of configuration needed for the scheduling. This proposed heuristic aims to redirect the small traffic demands to the host with large traffic demand. In this way, the small demands are indirectly transfered by a third intermediate host and won't require additional configuration time. We plan to evaluate this heuristic by both simulations and comparisons with integer linear programming (ILP) formulations.

The rest of this proposal is organized as follows. Section 2 discusses Solstice and other related work. Section 3 describes the proposed indirection heuristic. Finally, Section 4 describes the proposed evaluation plan and goals.

2 Related Work

Recently, researchers proposes hybrid datacenter network architectures that offer higher throughputs at lower cost by combining switching technologies. In particular, recent proposals suggest employing highspeed optical [2, 3, 8] or wireless [4, 5, 9] networks configured to service the heavy flows, while passing the remainder of the traffic through a traditional, relatively underprovisioned packet-switched network.

As technology trends usher in dramatically faster reconfiguration times, the distinction between packet and circuit is blurred, and ever smaller flows can take advantage of a hybrid fabric. This trend will soon allow servicing the bulk of the traffic through a rapidly reconfigurable optical switch [7], leaving a relatively minor portion to be serviced by the packet network [6]. While the potential cost savings that hybrid technologies could realize is large, the design space for scheduling resources in the hybrid regime is not yet well understood. What range of traffic demands can be scheduled for a given switch design and how can this schedule be computed efficiently? These questions are not addressed by the existing switch scheduling literature.

Many of the classical approaches to scheduling for switches with non-trivial reconfiguration delays divide the offered demand into two parts: an initial, heavy-weight component that is served by $O(N)$ highly utilized configurations with significant durations, and a second, residual component that is serviced by a similar number of short, under-utilized schedules. The recent Solstice scheduling algorithm [6] exploits the skewed nature of datacenter traffic patterns to create a small number of configurations with long durations that minimize the penalty for reconfiguration and leaves only a small amount of residual demand to be serviced by a low-speed (and lowcost) unconstrained packet switch.

3 The Indirection Heuristic

In this proposal, we propose to improve the Solstice scheduling algorithm by introducing an indirection heuristic that reduce the number of configuration needed for the scheduling. Originally, Solstice will schedule each traffic demand by a direct path, which means that a new configuration is required to fulfill a demand that has different sources or destinations. Our heuristic proposes to schedule the small traffic demands indirectly. For example, port a has 5 demand to port b and 20 demand to port c; port c has 20 demand to port b. Solstice will require 3 configurations for these demands. However, we could let port a send all 25 demand to port c and let port c send 25 demand to port b. In this way we reduce one configuration time, which is nontrivial for optical switches.

4 Evaluation Plan and Goals

We plan to evaluate our heuristic first by simulations. We received the simulator of the Solstice algorithm from the authors. We plan to introduce our heuristic into the simulator and compare it with the original algorithm. The simulation results will be comparable in terms of number of configurations, demand served on circuit switch, and demand served on hybrid switch. This will be our 75% goal.

In addition, we plan to evaluate our heuristic by comparing with integer linear programming (ILP) formulations. We plan to write a ILP problem generator which will generate scheduling problems with any size of network and any arbitrary demand between each ports. We will then use the SCIP Optimization Suite [1] to solve these ILP problems and compare the results with our heuristic's simulation results. We believe that this comparison will show the gap between our heuristic and an optimal answer. This will be our 100% goal. Furthermore, we may deploy our heuristic into a real optical switch and evaluate it with real world workload. This will be our 125% goal.

References

- [1] SCIP Optimization Suite. URL <http://scip.zib.de/>.

- [2] K. Chen, A. Singla, A. Singh, K. Ramachandran, L. Xu, Y. Zhang, X. Wen, and Y. Chen. OSA: An optical switching architecture for data center networks with unprecedented flexibility. In *USENIX NSDI*, 2012.
- [3] N. Farrington, G. Porter, S. Radhakrishnan, H. H. Bazzaz, V. Subramanya, Y. Fainman, G. Papen, and A. Vahdat. Helios: a hybrid electrical/optical switch architecture for modular data centers. In *SIGCOMM 2010*.
- [4] D. Halperin, S. Kandula, J. Padhye, P. Bahl, and D. Wetherall. Augmenting data center networks with multi-gigabit wireless links. In *SIGCOMM 2011*.
- [5] S. Kandula, J. Padhye, and P. Bahl. Flyways to de-congest data center networks. In *HotNets 2009*.
- [6] H. Liu, F. Lu, A. Forencich, R. Kapoor, M. Tewari, G. M. Voelker, G. Papen, A. C. Snoeren, and G. Porter. Circuit Switching Under the Radar with REACToR. In *USENIX NSDI*, 2014.
- [7] G. Porter, R. Strong, N. Farrington, A. Forencich, P. Chen-Sun, T. Rosing, Y. Fainman, G. Papen, and A. Vahdat. Integrating microsecond circuit switching into the data center. In *SIGCOMM 2013*.
- [8] G. Wang, D. G. Andersen, M. Kaminsky, K. Papagiannaki, T. Ng, M. Kozuch, and M. Ryan. c-Through: Part-time optics in data centers. In *SIGCOMM 2010*.
- [9] X. Zhou, Z. Zhang, Y. Zhu, Y. Li, S. Kumar, A. Vahdat, B. Y. Zhao, and H. Zheng. Mirror mirror on the ceiling: flexible wireless links for data centers. In *SIGCOMM 2012*.