

Data center resource management for in-network processing

Marco Micera

Politecnico di Torino, Technische Universität Darmstadt

1. Introduction
2. Analysis
3. FIXME Requirements? It is necessary?
4. Design
5. Conclusions

Introduction

Introduction

- (NetCache introduction: "modern Internet services, such as search, social networking and e-commerce, critically depend on high-performance key-value stores. Rendering even a single web page often requires hundreds or even thousands of storage accesses." [3])
- need to scale up (SHArP introduction, *justifying INP*: "As the number of compute elements grows, and the need to expose and utilize higher levels of parallelism grows, it is essential to reconsider system architectures, and focus on developing architectures that lend themselves better to providing extreme-scale simulation capabilities.") [1])
- However, modern-day data centers only exploit servers to perform computation

In-Network Processing (INP)

- Offloading computation to network devices (e.g., programmable switches, network accelerators, middleboxes, etc.), hence reducing load on servers
- (Daiet introduction: "The functionality of networks can now be enriched without hardware changes while retaining the capability of processing packets at very high rates, even above Terabits per second") [5]
- Few solutions out there already: Daiet [5], SHArP [1], NetChain [3], IncBricks [4]

Problem statement

- there is no Resource Manager (RM) that considering server and switch resources **conjunctly**...

(abstract)

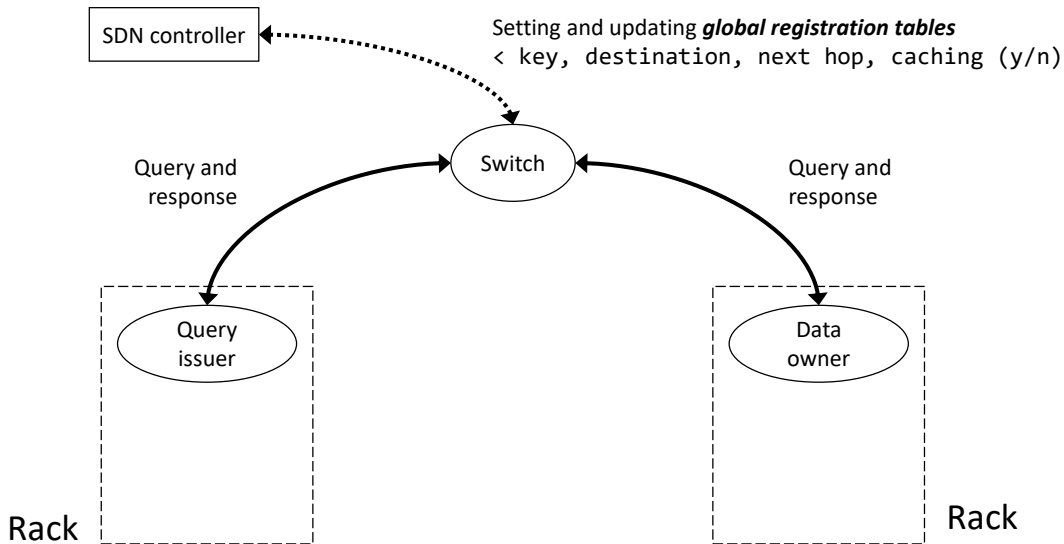
1. Model and evaluate an API through which applications can ask for INP resources
2. Discuss the importance of a scheduler which can reject INP requests and propose their server-only equivalent when needed (e.g., high switch utilization)

Analysis

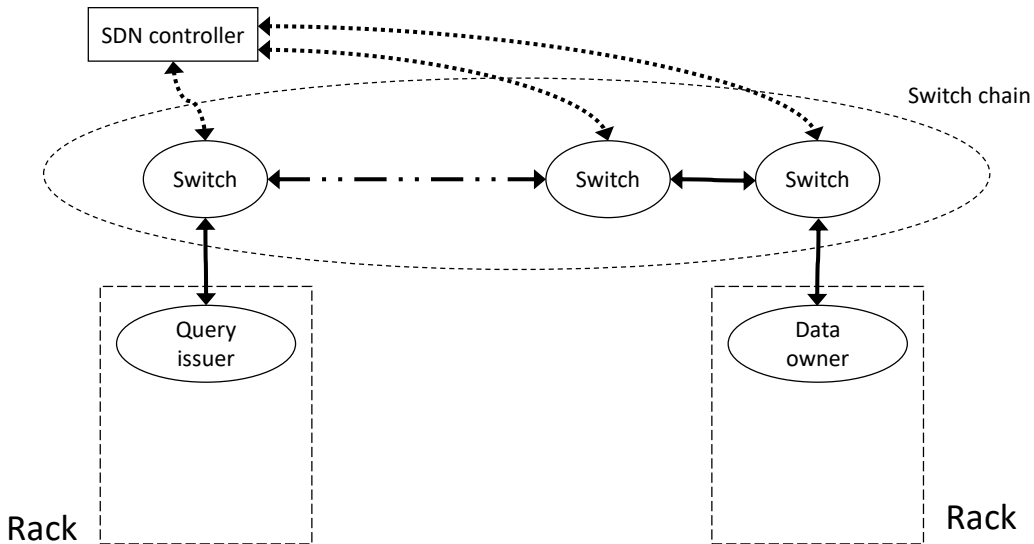
Currently existing In-Network Processing (INP) solutions

- In-network **storage**
 - Switches must
 - dedicate part of their local memory to store a distributed map
 - form a chain
 - IncBricks [4], NetChain [3]
- In-network **data aggregation**
 - Switches must
 - form a tree whose root is connected to data consumers and whose leaves are connected to data producers
 - dedicate part of their local memory to store a key-value map
 - be able to perform basic operations on data, such as writing and hashing
 - wait for all its children to send aggregated data
 - Dalet [5], SHArP [1]

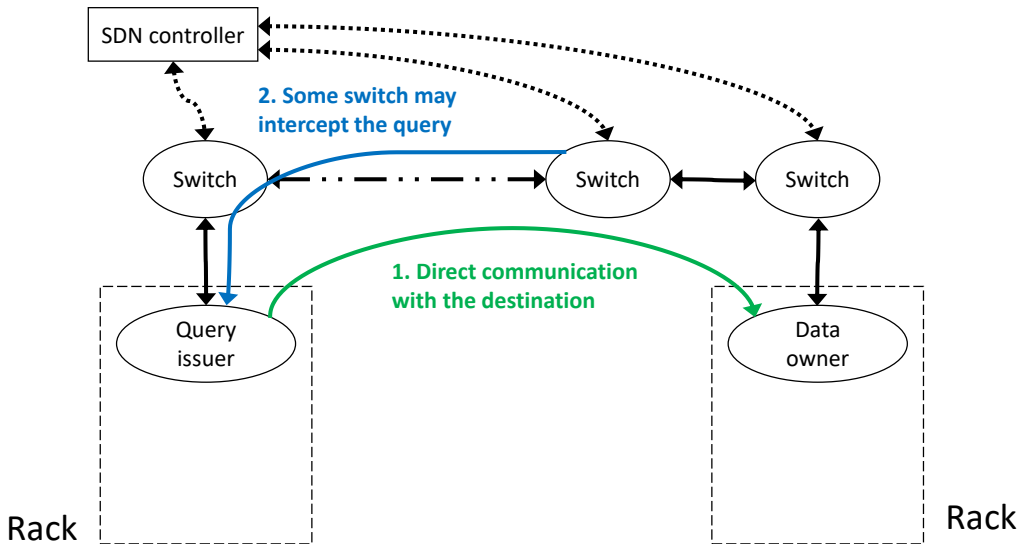
In-network caching system: IncBricks



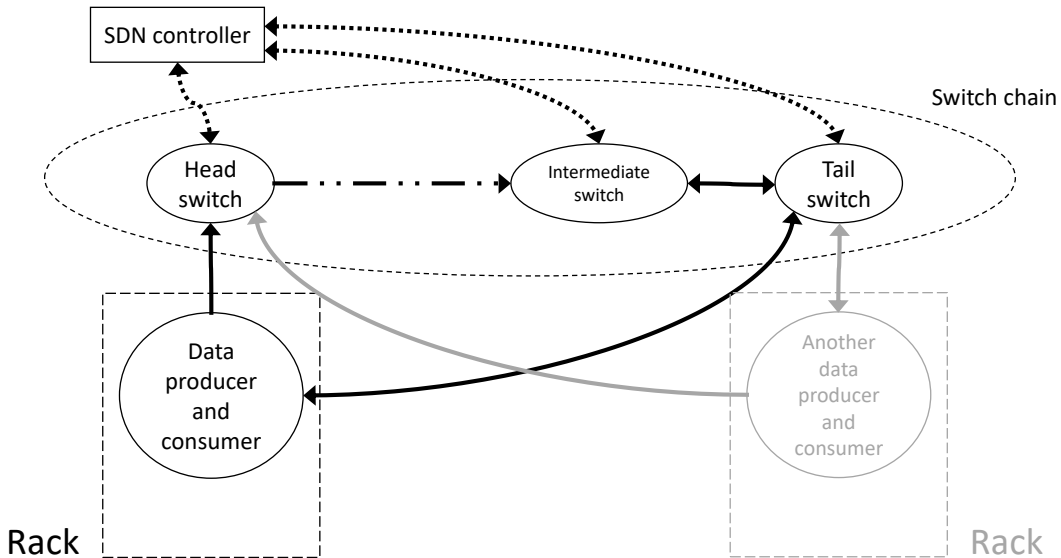
In-network caching system: IncBricks



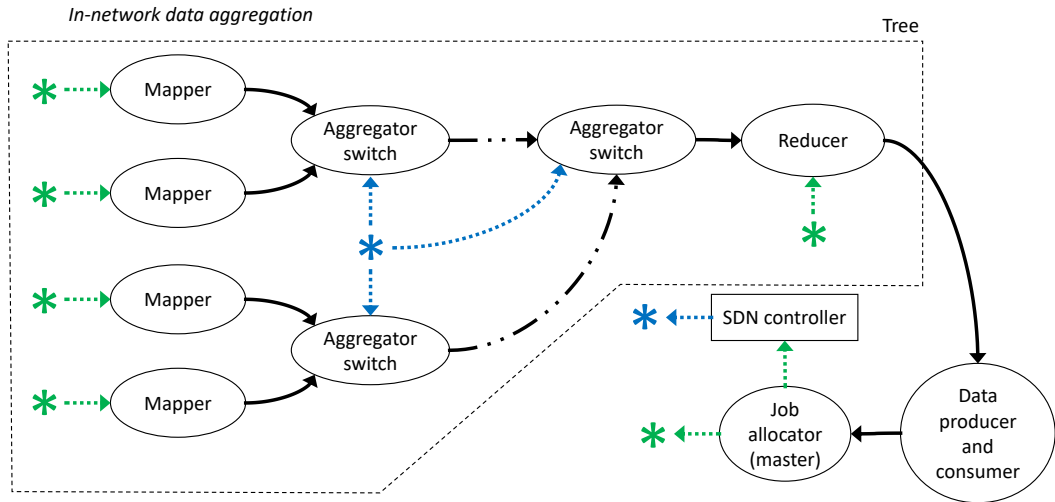
In-network caching system: IncBricks



Coordination services: NetChain

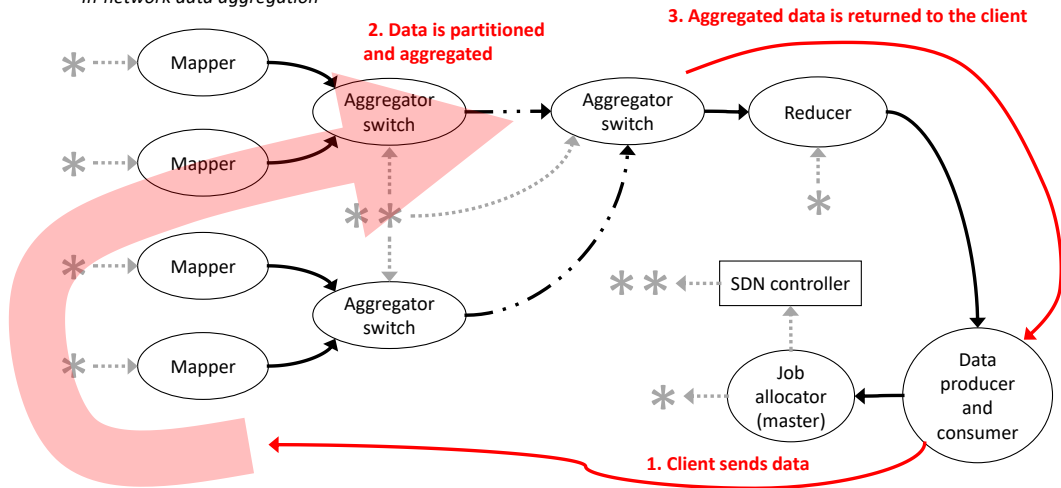


In-network aggregation: Daiet

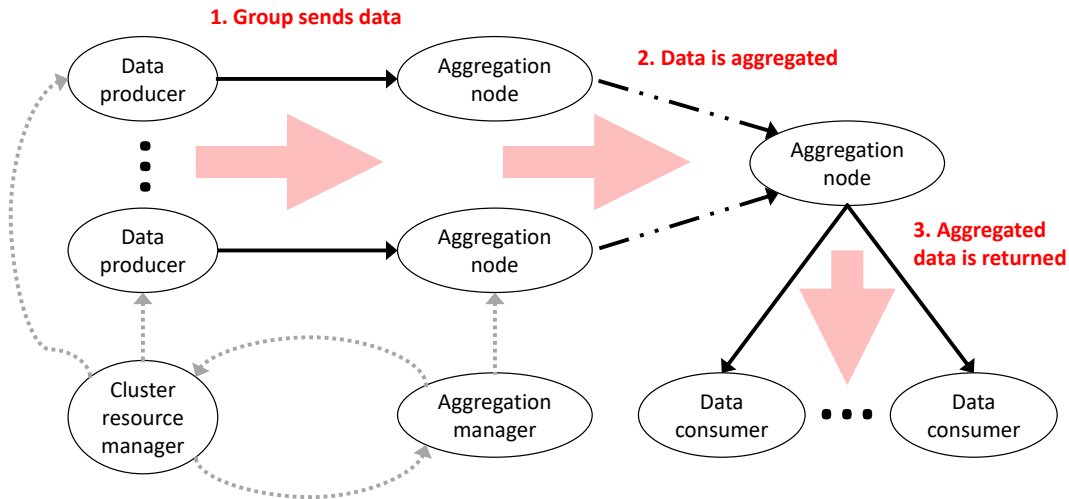


In-network aggregation: Daiet

In-network data aggregation

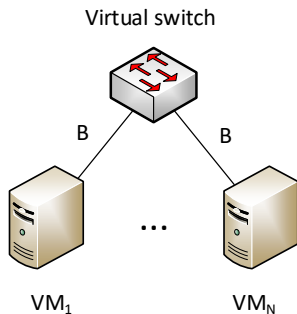


Aggregation protocol: SHArP



Resource models (3.4)

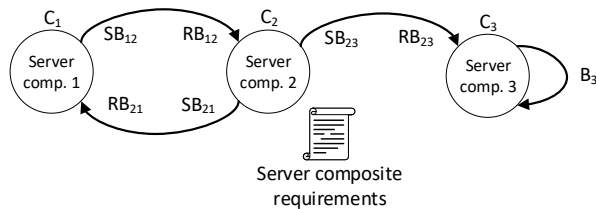
1. Virtual Cluster (VC)



2. Virtual Oversubscribed Cluster (VOC)

- N VCs connected to a single root virtual switch

3. Tenant Application Graph (TAG)



4. Fine-grained resource requests

- List of server-only resource demands

5. High-level goals

- E.g., job completion time (Bazaar [2])

(3.3)

FIXME: the only “network-aware” RM + its problems

(3.3.2)

**FIXME Requirements? It is
necessary?**

Design

(Thesis repository `dock/thesis/figures/design/model/presentation.pdf`)

The extended-Tenant Application Graph (eTag)

- (5.2.1 why existing resource models do not satisfy all requirements)
- (5.2.2 eTag)

The template database



- (5.3 generic groups)
- (5.1.2 template database role)(5.1.2 template database role)



FIXME I should introduce composites earlier

Conclusions

Questions?

Thank you

-  R. L. Graham, D. Bureddy, P. Lui, H. Rosenstock, G. Shainer, G. Bloch, D. Goldenberg, M. Dubman, S. Kotchubievsky, V. Koushnir, et al.
Scalable hierarchical aggregation protocol (sharp): a hardware architecture for efficient data reduction.
In Proceedings of the First Workshop on Optimization of Communication in HPC, pages 1–10. IEEE Press, 2016.
-  V. Jalaparti, H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron.
Bridging the tenant-provider gap in cloud services.
In Proceedings of the Third ACM Symposium on Cloud Computing, SoCC '12, pages 10:1–10:14, New York, NY, USA, 2012. ACM.

-  X. Jin, X. Li, H. Zhang, N. Foster, J. Lee, R. Soulé, C. Kim, and I. Stoica.
Netchain: Scale-free sub-rtt coordination.
In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 35–49, Renton, WA, 2018. USENIX Association.
-  M. Liu, L. Luo, J. Nelson, L. Ceze, A. Krishnamurthy, and K. Atreya.
Incbricks: Toward in-network computation with an in-network cache.
In *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '17*, pages 795–809, New York, NY, USA, 2017. ACM.



A. Sapio, I. Abdelaziz, A. Aldilaijan, M. Canini, and P. Kalnis.

In-network computation is a dumb idea whose time has come.

In *Proceedings of the 16th ACM Workshop on Hot Topics in Networks*, HotNets-XVI, pages 150–156, New York, NY, USA, 2017. ACM.