



**HAL**  
open science

# Network Configuration and Flow Scheduling for Big Data Applications

Lautaro Dolberg, Jérôme François, Shihabur Rahman Chowdhury, Reaz Ahmed, Raouf Boutaba, Thomas Engel

► **To cite this version:**

Lautaro Dolberg, Jérôme François, Shihabur Rahman Chowdhury, Reaz Ahmed, Raouf Boutaba, et al.. Network Configuration and Flow Scheduling for Big Data Applications. Networking for Big Data, Chapman and Hall/CRC , 2015, Big Data Series, 9781482263497. hal-01244585

**HAL Id: hal-01244585**

**<https://inria.hal.science/hal-01244585>**

Submitted on 16 Dec 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Network Configuration and Flow Scheduling for Big Data Applications

Lautaro Dolberg      Thomas Engel      Jérôme François  
*SnT. University of Luxembourg*      *INRIA Nancy Grand Est*

Shihabur Rahman Chowdhury      Reaz Ahmed      Raouf Boutaba  
*William G. Davis Computer Research Centre. Waterloo University*

## 1 Introduction

Big Data applications play a crucial role in our evolving society. They represent a large proportion of the usage of the cloud [5, 22, 3] because the latter offers distributed and on-line storage and elastic computing services. Indeed, Big Data applications require to scale computing and storage requirements on the fly. With the recent improvements of virtual computing, data-centers can thus offer a virtualized infrastructure in order to fit custom requirements. This flexibility has been a decisive enabler for the Big Data application success of the recent years. As an example, many Big Data applications rely, directly or indirectly, on Apache Hadoop<sup>1</sup> which is the most popular implementation of the Map-Reduce programming model [24]. From a general perspective, it consists in distributing computing tasks between *mappers* and *reducers*. Mappers produce intermediate results which are aggregated in a second stage by the reducers. This process is illustrated in Figure 1(a), where the mappers send partial results (values) to specific reducers based on some keys. The reducers are then in charge of applying a function (like sum, average or other aggregation function) to the whole set of values corresponding to a single key. This architectural pattern is fault tolerant and scalable. Another interesting feature of this paradigm is the execution environment of the code. In Hadoop, the code is directly executed near the data it operates on, in order to limit the data transfer within the cluster. However, large chunks of data are still transferred between the mappers and reducers (shuffle phase) which thus necessitates an efficient underlying network infrastructure. It is important to note that the shuffle phase does not wait for the completion of the mappers to start as the latter already emits (*key,value*) pairs based on partial data it has read from the source (for example, for each line). Since some failures or bottlenecks can occur, Hadoop tasks are constantly monitored. If one of the components (*i.e.*, mappers or reducers) is not functioning well (*i.e.* it does not progress as fast as others for example), it can be

---

<sup>1</sup>hadoop.apache.org

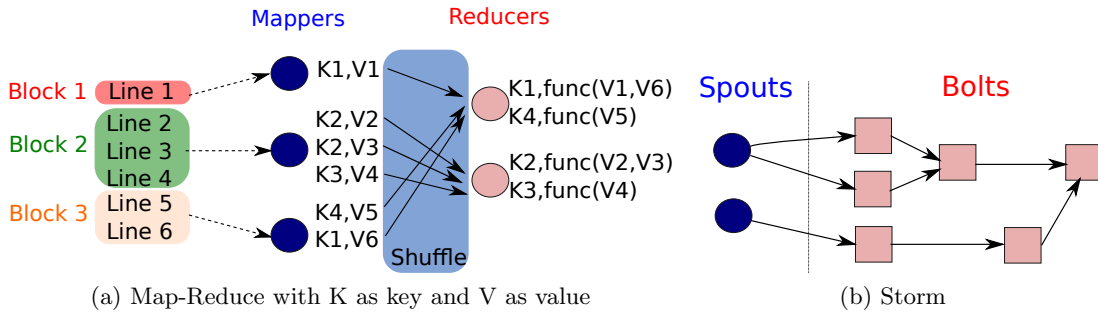


Figure 1: Big data computational model and the underlying network traffic as plain arrows

duplicated into another node for balancing load. In such a case, this leads also to additional data transfers.

Storm [34] is another approach that aims at streaming data analytics, while Hadoop was originally designed for batch processing. Storm consists of *spouts* and *bolts*. Spouts read a data source to generate tuples and emits them towards bolts. Bolts are responsible of processing the tuples and eventually emit new tuples towards other *bolts*. Therefore, a Storm application is generally represented by a graph as shown in Figure 1(b). The main difference between Storm and Map-Reduce is that data transfers occur all the time (streaming) and so are not limited to a specific phase (shuffle phase in Hadoop). As a result, among the diversity in big data applications, there are common problems, whose the probably more aimed at is optimizing the data transfer rate between hosts.

Therefore, while big data technological improvements were mainly highlighted by new computing design and approaches, like Hadoop, network optimizations are primordial to guarantee high performances. This chapter reviews existing approaches to configure network and schedule flows in such a context. In the following sections, we will cover the diverse optimization methods grouped according to their intrinsic features and their contributions. In particular, recent network technologies such as Software Defined Networking (SDN) empowered the programmability of switching devices. Consequently, more complex network scheduling algorithms can be afforded to leverage the performance of Map-Reduce jobs. That is why this chapter focuses on SDN-based solutions but also introduces common networking approaches which could be applied as well as virtualization techniques. The latter are strongly coupled with the network design. For example, end-hosts in a data-center are virtual machines which can be assigned to different tasks and so would lead to various traffic types, which can be better handled if the network is adaptive and so reconfigurable easily.

This chapter is structured as follows:

1. Optimization of the VM placement: even not dealing with network configuration, it has a significant impact on the same;

2. Topology design: it is an important topic as the way the machines are wired have an impact on performance;
3. Conventional networking, in particular routing and QoS scheduling: these might be customized to support Big Data as well.
4. Software-Defined Networking (SDN): this highlights recent approaches that leverages a global view of the network to implement efficient traffic management policies.

## 2 VM Placement for reducing Elephant Flow impact

Very large flows, normally associated to long Map-Reduce jobs are often called Elephant Flows[30]. Since any VM can be potentially hosted in any physical server, grouping VM that are involved in large data transfers can reduce the impact on the overall bandwidth usage of the network. This approach is based on the internal routing of Hypervisor systems used in virtualized data centers such as XEN[6], KVM[23] or VMWare[32] solutions. From a more general point of view, virtual machines can be colocated in a certain region of a network, even if on different physical machines. This is illustrated in Figure 2 where in the case of the original allocations (Figure 2(a), the tasks of the same job are scattered in the network and so the traffic between them has to go through many hops eventually resulting in network congestion. In Figure 2(b), by moving only two tasks (one from J1 and one from J3), each job is isolated in a single rack (under a single switch) and so no congestion occurs at higher level switches while improving the data transfer efficiency between tasks of the same job since these are connected through a single switch.

VM placement is basically related to virtual machine allocation problems which are optimization problems under certain criteria. One of the criterion should be the usage of network resources. Because this is not the focus of this chapter, we recommend the reader to read [37] for more details about network-aware VM placement.

The down side of existing network-aware VM placement approaches is the lack the reactivity. Normally, given the nature of Map-Reduce phases, it is not possible to match exactly in advance Map-Reduce jobs and needed network resources (for example, how large the data transfer will be during the shuffle phase is depending on the underlying data and applications). To cope with this practical issue, virtualized data-centres may estimate the VM-to-VM traffic matrix but such a method works well only with known batch job only. Another solution is to migrate VMs during their execution but this might be also resource consuming and negatively impact the finishing time of the Big Data jobs if this occurs too frequently.

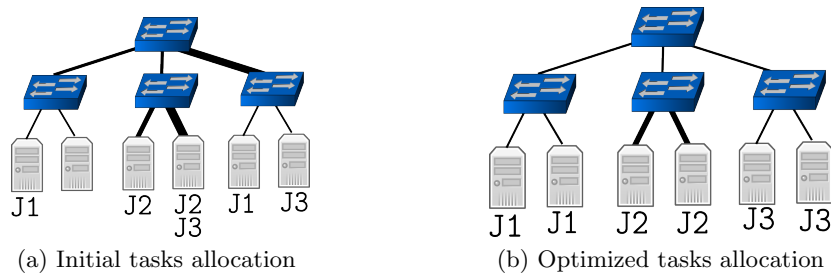


Figure 2: VM placement with 3 jobs (each job  $J_i$  has two tasks). The width of a link represents its load.

### 3 Topology design

Data-centers networks are usually organized in a tree topology [1, 29] with three defined layers:

- Core layer: This layer is the backbone of the network where high-end switches and fibers are deployed. In this layer only L2 forwarding takes place without any packet manipulation. The equipment for this layer is the more expensive among the hierarchical network model.
- Aggregation or distribution layer: In this layer takes place most of the L3 routing.
- Access layer: This layer provides connectivity to the end nodes and so are located at the top of the racks. They performs the last step of L3 packet routing and packet manipulation. Normally, those are the cheapest devices in the hierarchical network model.

Thanks to this hierarchical model, a low latency is achieved for traffic between two nodes in the same rack. This explains why approaches like Hadoop leverage rack awareness to ensure fast replication of data by selecting nodes in the same rack for copying data (but also others out of the rack in order to guarantee data availability under a rack failure). In addition, this type of configuration supports a large number of ports at the access layer.

A specific instance of the hierarchical model is the fat tree proposed in [3] and illustrated in Figure 3 which enables fault-tolerance by ensuring redundant paths in a deterministic manner.

The fat tree or Clos topology was introduced more than 25 years ago [25] to reduce the cost of telephony switched networks. The topology layout is organized as  $k$ -ary trees, where in every branch of the tree there are  $k$  switches, grouped in pods. Actually, a pod consists in  $(k/2)^2$  end-hosts and  $k/2$  switches. At the edge level, switches must have at least  $k$  ports connected as follows: half of the ports are assigned to end nodes and the

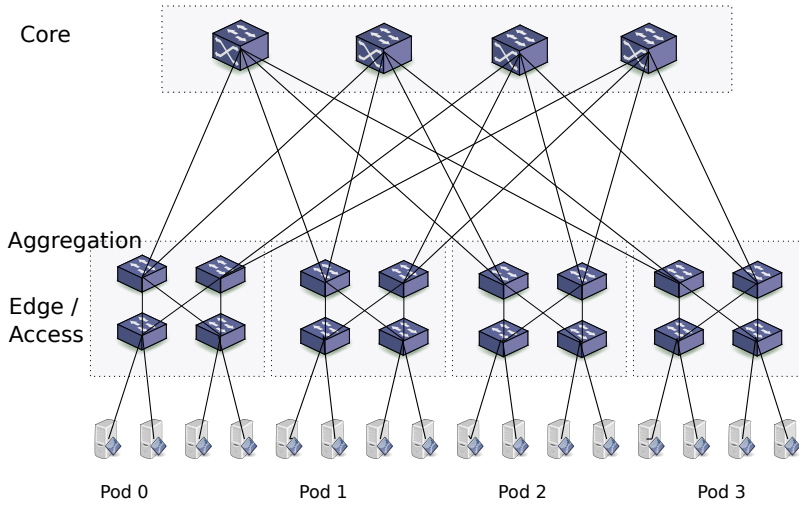


Figure 3: Example of a Hierarchical Network Model: Multi-rooted Network Topology

other half is connected to the upper aggregation layer of switches. In total, the topology supports  $(k^2/2)$   $k$ -port switches for connecting host nodes.

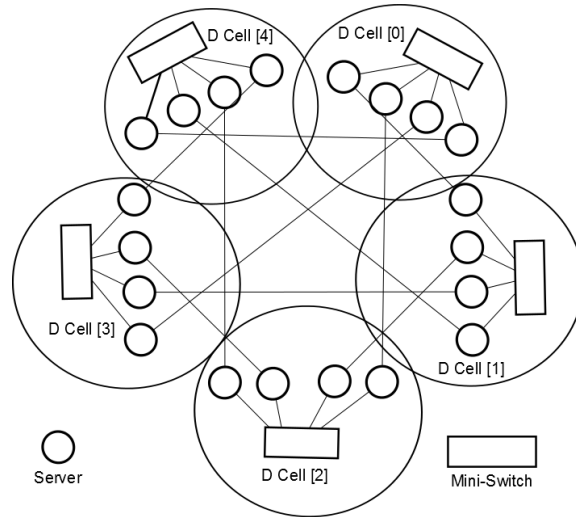


Figure 4: A DCell topology for 5 Cells of level 0, each containing 4 servers (src: [20])

DCell [20] is a recursively interconnected architecture proposed by Microsoft. Compared to a Fat Tree Topology, DCell is a fully interconnected graph in order to be largely fault tolerant even under several link failures. In fact, high level DCell nodes are recursively connected to low level ones, implemented with mini switches to scale out as showed in

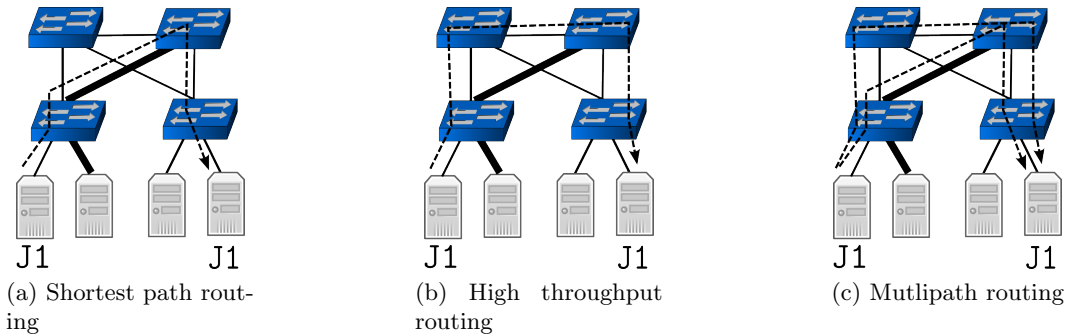


Figure 5: Routing decisions for one job with two tasks. The width of a link represents its load.

Figure 4). Experimental results have shown that with a 20 nodes network can outperform by two times a large data-center used for Map-Reduce. As a downside, DCell requires a full degree of connectivity, making it in practice costly to maintain and deploy. To enhance network connectivity between servers, CamCube [2] is a torus topology where each server is interconnected to other 6 servers and all communications are going through them, without any switch for internal communication. Finally, recent propositions like [33] promote a high flexibility by alleviating the need for a well-defined fixed graph structure, as the fat-trees are, and so by introducing some randomness in the topology bounded by some criteria.

## 4 Conventional Networking

### 4.1 Routing

Data-center network topologies like fat trees imply a large number of links leading to redundant paths. Therefore, routing algorithms can take benefit of that to achieve a higher bandwidth. As an illustrative example in Figure 5(a), the shortest path is used to route the traffic between the two tasks of the job  $J1$ . Unfortunately, it goes through a congested link. Hence, a redundant path can be used (Figure 5(b)) and even multiple of them conjointly (Figure 5(c)). Although these approaches have been proposed for routing in general, they are also used in data-centers to improve the performance of the Big Data applications. This is the reason why this section covers some propositions about how to use these principles in case of Big Data. However, the general issues are (1) to predict the traffic patterns and (2) to be able to rapidly change the configuration of the routing when the traffic suddenly changes, which is the case in a cloud infrastructure.

Nowadays, a major representative of such an approach is the Equal Cost Multi Path (ECMP) algorithm [21]. ECMP leverages the opportunity to route flows among multiple paths. Unlike traditional routing algorithms like OSPF which consider a single best path,

ECMP consider all the best multi-paths according to any metric (as for example the number of hop) among which a single one is selected for a given flow through a load balancer. The number of multiple paths is dependent on the router implementation but usually bounded to 16. Hence, this may yield to a lower performance than expected for large data-centers. In fact, the amount of entries in the routing tables grows at exponential rate, increasing the latency of the routing algorithm. Commercial solutions promoting multi-path routing include FabricPath by Cisco Systems, BCube, VL2 and Oracle Sun data-centre InfiniBand.

In addition to promoting the fat tree topology usage for data-centers, the authors of [3] proposed a dedicated routing algorithm based on an approach called *Two-Level Routing Tables*, where the routing tables are split into two hierarchical tables linked on the prefix length of the network address. A two layer table approach aims at leveraging the routing algorithm speed for establishing a route. This is possible because the authors introduced a private addressing system respecting a pre-established pattern like *8.pod.switch.host* assuming a class A network. The first table index entries use a left handed prefix length (eg, 8.1.2.0/24, 8.1.1.0/24, etc). The entries of the first table are linked to a smaller secondary table indexed by a right handed suffix (eg, 0.0.0.1/4, 0.0.0.4/4). For example, to find the route to the address 8.8.8.8, the algorithm will lookup the first table, find the corresponding entry for the first part of the network address 8.8.8.0/24, then jumps to the secondary table and find the remaining of the route. Since each switch of the aggregation layer in a fat tree topology has always a  $k/2$  degree of connectivity to the access layer, two-Level routing tables are bounded in the worst case to  $k/2$  entries for suffixes and prefixes. Moreover, flows can be actually classified by duration and size. Then, the proposed algorithm in [3] minimizes the overlap between the paths of voluminous flows. To achieve that, a central scheduler is in charge of keeping track of used links in the network in order to assign a new flow to a non used path. From this perspective, it falls into the category of centralized networking (see section 5.1) where a server acts as the controller by informing other ones about the link to use to forward specific packets of a flow.

The flow establishment is also leveraged by the previously described route lookup. In this approach, instead of routing traffic at a packet level, streams of data are grouped into flows and routed as a whole entity. One of the benefits of this approach is a faster route computation as it is reduced in a similar fashion as in circuit switching legacy technology. For example, if a host node requires to transfer a large data file as a part of a Big Data job, the whole stream will follow a pre-established route, reducing the latency of establishing a different route for each packet of the stream.

In order to enhance routing and network speed, hardware plays a core role. Therefore, there have been propositions to replace standard hardware. In particular, the authors in [18] argue for an hybrid optical-electric switch as optical links achieve higher throughput but are not well adapted to bursty traffic. Combining both technologies thus helps in obtaining good trade-off between accuracy and cost. Moreover, the technological availability of programmable circuits also lead to the possibility of implementing switching devices, specially in the aggregation and core layer using ASIC and FPGA devices. Authors of [27]



propose an approach for implementing switching cards with a PCI-E interface.

A recent proposal [9] addresses dynamic routing by replacing the traditional DHCP address configuration by an another automated address configuration system. In this approach, the network is automatically blue printed as a graph. Then, by interpreting a set of labels assigned to each computing node, the system tries to find an isomorphism that minimize the traffic at the aggregation layer. From the preliminary results, this approach has yielded promising results. However, it actually runs only over BCube or DCell because they have a fully connected topology.

## 4.2 Flow Scheduling

Network operators perform various traffic engineering operations in order to provide different network services on a shared network. This consists in classifying the traffic according to the intrinsic characteristics of each service or application using the network. For example, it is possible to define policies to specially treat Big Data applications. Similarly, the IPv6 Traffic Class includes the possibility of injecting information specific to applications in the packet stream. Another types of support for enabling network infrastructure to perform management of the traffic is proposed in RFCs [7] and [8]. The first (DiffServ) proposes a protocol for differentiating services and its network behavior. The latter, RSVP (Resource Reservation protocol), specifies also a protocol, that enables application to reserve network resources in advance of initiating a data transfer.

As highlighted in the introduction, Big Data applications includes both batch processing and streaming analytics, which are different by nature. In particular, batch processing jobs are more prone to use the network heavily during certain phases while streaming uses the network constantly with various rates. Therefore, the apparition of a batch job (Hadoop) may suddenly impact the network and so the other underlying applications. The authors in [17] have proposed to schedule flows from BigData applications in a data center using a variation of FIFO scheduling that allows some level of multiplexing between the flows. The authors propose to schedule flows in the order of arrival with a certain degree of freedom and allow multiplexing over a limited number of flows which in turn allows small flows to be processed alongside large flows. This approach allows the co-execution of batch and streaming Big Data applications.

### Limitations

It is worth mentioning that, in traditional data center networks, only aggregation and core layer switches have the capability of scheduling flows. This is a limitation imposed by the hardware. To be able to exploit the full potential of flow scheduling, an additional network function is required. This is often implemented in a central controller, this way allowing core and aggregation switches to be replaced by simple switches. One of the main advantages of using this approach is the reduced cost of switching and forwarding (L2) devices.

Another disadvantage of traditional networking is that the network configuration remains static and so impacts on the maintenance cost of the infrastructure because any modification of the topology must be wired manually by the network administrators. Virtualized networks come into play for coping with the lack of flexibility in traditional networks, and became popular over the last years, thanks to the emerging virtualization technologies and computing power to support them. As a result, data-center owners offering their clients not only virtual machines (know as Virtual Private Servers (VPS)) but also virtual network infrastructure. This allows VPS users to create customized topologies. Virtual LANs (VLAN) have been popular in the past decades for splitting large organizational networks into smaller ones. However, this approach fails to segregate application traffic because of the coarse routing granularity inside a VLAN. A possible solution to this issue is to use a dynamic topology, that adapts to the specific needs of each application. In such a scope, the following Section 5.1 covers emerging technologies facilitating dynamic network configuration using a centralized control plane implemented in software.

## 5 Software-Defined Networking

This section covers both theoretical approaches as well as practical implementations. Solutions highlighted in the following paragraphs combine three aspects: computational patterns present in most of Big Data services, data-centers network architectural improvements such as hierarchical topologies (*e.g.* Fat-Trees) and dynamic routing algorithms leveraged by the adoption of technologies such as SDN. These three aspects combined together allow to adapt the network configuration from the core to the aggregation infrastructure layer to suit better Big Data application needs.

Routing and scheduling decisions rely on the traffic matrix. Such a matrix can be observed in real-time at the network level but can also be predicted in order to plan next course of actions. The traffic matrix usually reflects the flow's size, duration and frequency for each pair of nodes and eventually application instances or even between multiple tasks of a single job. Alternatively, Big Data applications can interact with a central controller to expose their current usage and needs. These two types of approaches are differentiated in Figures 6(a) and 6(b). In every cases, there is a Big Data application controller or manager (*e.g.* the *jobtracker* or the *resource manager* in Hadoop) which is in charge of triggering and monitoring the tasks. In Figure 6(a), a monitoring service is gathering traffic from forwarding devices and sends the information to the network controller itself which is in charge of taking routing decisions. The monitoring can even be done by OpenFlow [12] as an OpenFlow controller can request such statistics from OpenFlow switches. In this a case, both the monitor and controller are merged in a single entity. In a second scenario (Figure 6(b)), the Big Data controller sends itself information about the running jobs to the network controller which can thus take proper configuration actions. Finally, it is also possible to imagine an hybrid approach (Figure 6(c)) where both types of information are

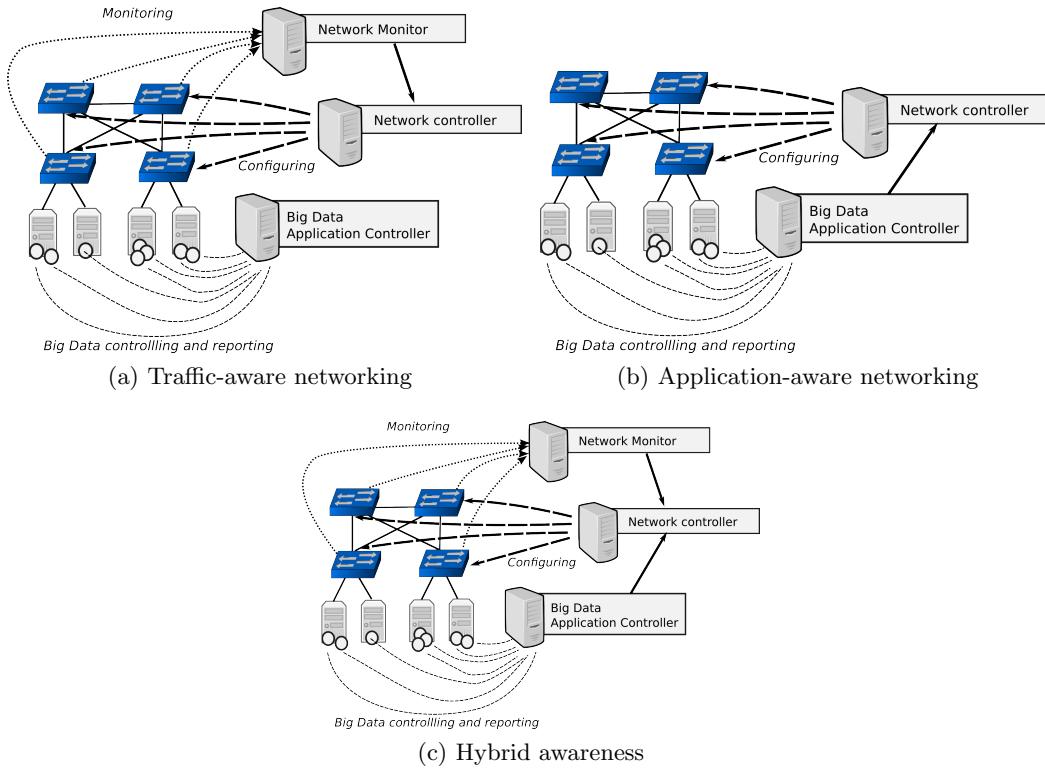


Figure 6: The different type of \*-aware networking (Small circles represent a task of a Big Data process)

made available to the controller. It might be useful if the level of details from the Big Data controller is coarse-grained.

To summarize, the different methods covered in the following subsections are, actually, similar to conventional networking (select better paths, minimizing congestion, etc.) but they rely on a higher and more dynamic coupling between the network configuration and applications (or the corresponding traffic).

## 5.1 Software Defined Networks

In recent years, Software-Defined Networking (SDN) emerged introducing a new layer of abstraction for more flexible network management. Under this approach, switches are just forwarding devices while most of the control (*e.g.* routing decisions) is performed in a central controller. As a result, network can be built with merchant silicone and can be programmatically controlled by the central control plane. This eventually results in reduction of both CAPEX and OPEX.

SDN decouples the data and the control plane as shown in Figure 7, where:

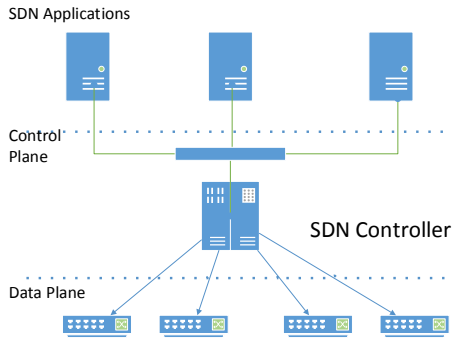


Figure 7: Software Defined Network Architecture Example

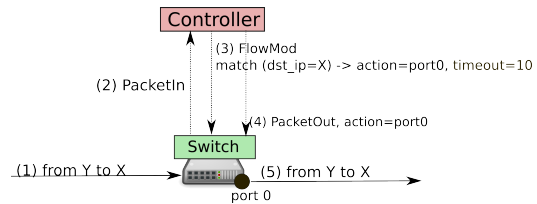


Figure 8: Software Defined Network with Open Flow rules

- **Control Plane:** The concept of the control plane is to have a dedicated communication channel for exchanging signalization messages among forwarding and management devices. Most of the available products for SDN expose a North Bound API for applications to subscribe to real time statistics and service usage.
- **Data Plane:** This layer, also refereed as the Forwarding Plane, performs the actual switching/forwarding of the network traffic. The traffic in this plane is accounted and measured but not interpreted by any decisional algorithms.

Additionally, the application layer is composed of custom made applications. The latter subscribe to the North Bound API of the SDN controller to enable extra functionality not provided by out of the box controller. For example, these applications might be security oriented[31] or for routing purposes[15].

OpenFlow[28] is adopted as de facto standard control protocol. OpenFlow acts as the communication protocol between switches and controllers (*e.g.* NOX, Floodlight, POX, etc). An OpenFlow rule consists of two parts: a match field, that filters packet headers, and instructions, indicating what actions to take with the matched packets.

Upon arrival of a packet at a switch, the controller decides on the route of the packet and sends the corresponding rule to the switch. This event is known as *FlowMod*. Finally, the packet is sent (*PacketOut*). Figure 8, illustrate an example where a routing action is taken upon arrival of a packet with destination *X* and source *Y*. Additionally, a controller can provision switches with flow tables entries in advance. Hence, a *PacketIn* message is not required to emit an event *FlowMod*. The rules also have soft (last seen packet) and

hard (maximum absolute value) timeouts, after expiration of these timeouts the rule is removed.

While originally proposed for campus networks, the modification proposed by authors of [13] consists of reducing the overhead induced by OpenFlow to enable a more efficient flow management for Big Data analytics applications networking through the extensive use of wildcard rules within the switches to avoid invoking the OpenFlow controller for each new flow. However, the extensive use of wildcards on OpenFlow might cause loss of granularity in the statistics derived from the counters on the controller and evidently on routing and scheduling decisions. As mentioned in [13], DevoFlow aims to devolve control by cloning rules whenever a flow is created using wildcards. The cloned rule, will replace the wildcard fields using the clone’s specific information. Additionally, DevoFlow enriches OpenFlow rules by including local routing actions (without relying on the OpenFlow controller), such as, multi path routing. This last feature allows to rapidly reconfigure the route for a given flow leveraging the flow scheduling.

## 5.2 Traffic-aware networking

The Topology Switching approach [36] proposes to expose several adaptive logical topologies on top of a single physical one. It is similar to the allocations problem in VM placement introduced in Section 2 by trying to assign every individual flow to a specific path to optimize an objective. The optimization objectives can be multiple in case of Big Data applications, the most important one is the total capacity, *i.e.* trying to use the available bandwidth as much as possible in order to reduce the job completion time. For example, considering a fat-tree topology as showed in Figure 3, every Map-Reduce typical bisection traffic is considered as a separate *routing task*. Thus, each task runs an instance of a particular routing system. For every routing system, a pre-allocated bandwidth is established in the physical topology to maximize the bandwidth. Topology Switching is implemented in a central topology server, responsible for allocating resources but also for substracting unused resources and collecting metrics. The two metrics used in this approach are the bisection bandwidth and the all-to-all transfer. Bisection bandwidth is used to measure the topology ability to handle concurrent transfers at the physical layer. The all-to-all metric is used to evaluate how the logical topologies react under a worst case scenario. Based on both metrics, the Topology Switching approach runs an adaptive algorithm for readjusting the logical configurations for the virtual networks. Topology Switching offers an alternative for "one-size fit all" data-center design, providing a good trade off between performance and isolation.

Hedera[4] scheduler assigns the flows to non-conflicting paths similarly to [3], especially by aiming at not allocating more than one flow on routes that cannot satisfy its network requirements in terms of aggregate bandwidth of all flows. Hedera works by collecting flow information from the aggregation layer switches, then computing non-conflicting paths, and re-programming the aggregation layer to accommodate the network topology in order

to fulfill the Map-Reduce jobs requirements. More especially, bottlenecks can be predicted based on a global overview of path states and traffic bisection requirements in order to change the network configuration.

### 5.3 Application-aware networking

The methods described in this section improves the network performance by scheduling flows according to an application-level inputs and requirements. At the transport layer, flows are not distinguishable from each other but groups of computing nodes in Big Data Application usually expose an application semantic. For example, an application can be composed of several shuffle phases and each of them corresponds to a specific set of flows. Furthermore, a Big Data application can evaluate its current stage. For instance, in a Map Reduce task, the mapper status (completion time) is computed from the proportion of the data, from the source, which has been read and such a completion time can approximate the remaining data to transfer. Therefore, a mapper having read 50% of its data source and having already send 1GB of data should approximatively send another 1GB. This is an approximation and it cannot be guaranteed that the mapper will send as much information for the remaining data it has to read. For example, a usual example where a mapper sends a  $\langle key, value \rangle$  pair for each read line can also apply some filtering and so may emit nothing based on the input data.

Therefore, some methods build a semantic model reflecting the Big Data application needs. The semantic model used for these approaches associates the network traffic to be managed with the characteristics and the current state of the application it originates from. This model might differ among the different proposed works but generally aims at assessing the state of the Big Data applications and their related flows.

In this context, the authors in [19] propose to optimize network performance by arranging QoS policies according to applications requests. Host nodes running Big Data applications can exchange messages within their proposed framework called PANE to submit QoS policies similarly to what can be done with conventional networks (Section 4.2). Naturally, this approach will lead to traffic over subscription under high traffic demand circumstances. To solve this issue, users have also to provide conflict resolution rules for each QoS rule they submit into the system. Also, this approach can be employed for implementing security policies such as denial of service prevention by setting a top hierarchy policy triggered at the SDN controller.

OFScheduler[26] is a scheduler which assesses the network traffic while executing Map-Reduce jobs and then load-balance the traffic among the links in order to decrease the finishing time of jobs based on the estimated demand matrix of Map-Reduce jobs. OF-Scheduler assumes that Map-Reduce flows can be marked (for example by Hadoop itself) to distinguish those related to the shuffle from those related to the load balancing (when a task is duplicated). The scheduling first searches for heavily loaded links and then selects flows to be offloaded by giving the preference to (1) load-balancing flows and (2) larger

flows in order to limit the impact on performance (cost of the offloading due to OpenFlow rule installation). The reason for (1) is that it corresponds to a duplicated task the original of which may finish somewhere else in the data-center unlike the others. The rationale behind (2) is to minimize the global cost of offloading and so by moving big flows, there are more chances to remedy the problem of the link load without re-scheduling additional ones.

Assuming optical links, authors of [35] describe an application-aware SDN controller that configures optical switches in real-time based on the traffic demand of Big Data applications. By enabling the Hadoop Job Scheduler to interact with the SDN controller, they propose an aggregation methodology to optimize the use of optical links by leveraging intermediate nodes in the aggregation. In the simplest case, when a single aggregate has to gather data through  $N$  switches whereas the number of optical links is lower, it has to go through multiple rounds (optical switching) in order to complete the job. The other switches only using a single connection to the aggregating switch can also be connected together to act as intermediate nodes to form a spanning tree rooted in the aggregator and so to avoid the multiple rounds. Such a principle (Many to One) is extended towards general case with Many to Many jobs or when multiple single aggregation overlaps (*e.g.*, different sources overlap their aggregators). This requires more complex topologies such as torus. Other data center network topologies discussed in this chapter such as DCell or CamCube also make use of high redundancy to build similar shaped topologies. Building a torus topology is more complicated than a tree because the search space for suitable neighbors is larger, a greedy heuristic is used to support as much as possible the traffic demand. The routing algorithm within the torus topology is meant to exploit all possible optical paths. Authors also propose to assign weights to the optical links for load-balancing purposes on the torus topology.

FlowComb[14] is a combination of proactive and reactive methods for flow scheduling. It allows the Hadoop controller to specify requirements but also promotes the use of a statistic-based method that predicts based on the network load of previous runs. Hence, this approach lies between application-aware and traffic-aware. Based on that, any routing or scheduling approach described in section 5.2 could be applied, especially Hedera[4] which has been chosen by the authors. The central decision engine gathers all the job pertinent data and creates a set of Open Flow rules to be installed temporarily and erased after job completion. However, the main drawback of the proactive method using estimation is that about 30% of jobs are detected after they start, and 56% before they finish.

CoFlow[10] proposes a full reactive method, that only after receiving the Hadoop Job Scheduler network requirements is able to yield results. Its implementation exposes an API for declaring flows at application level. This API can be used for example from the Hadoop Job Scheduler as it is mentioned by the authors to express on demand bandwidth requirements at the different phases of a Map-Reduce job. Actually, CoFlow introduced an abstraction layer to model all dependencies between flows in order to schedule an entire application, i.e. a set of flows, and not only a single flow.

In contrast with the methods described previously, the authors of [16] propose an approach for routing on a packet basis by splitting the flows in chunks similarly to TCP. These chunks are distributed to the available ports of a switch using different strategies: random, round robin and counter based. However, the main limitation of this approach is the necessity to reorder the chunks.

## 6 Conclusions

Big data applications are a major representative in today's cloud services which have also guided the network design and configuration for performance purposes. For example, the fat-tree network topology is a popular choice among data-centers hosting Big Data applications. Also, the usage of ECMP as a routing algorithm, leverages the notion of flow routing for a better efficiency in redundant linked networks. Complementary to the Fat-Tree approach, the DCell and BCube design patterns propose a high degree or almost full connectivity between the nodes of the data-centre. The usage of these kind of topologies is tightly related to the type of applications running over the network. Therefore, one size (network architecture/topology) does not fit all applications and some will experience degraded performance. To cope with this situation, alternatives in the field of dynamic routing and flow scheduling have been proposed.

The network topology can be adapted dynamically to meet the application bandwidth needs in terms of data transfer but also to reduce the latency and improve the Big Data job's finishing time. Many of the solutions proposed in this field consist in regrouping application nodes (VMs) that concentrate a high volume of data to be transferred.

Programmable networks are more flexible in having a central controller that can take a lead role in flow scheduling. Many Big Data applications have an observable traffic pattern which is exploited by several works to propose specific scheduling to make a more efficient network usage (*e.g.*, load balancing, traffic management and resources allocation). In this direction, several authors have highlighted the notion of "network awareness". In general, two kinds of application state-full controllers and network architectures have been proposed: Passive application controllers (traffic-awareness) are those that take as input the traffic matrix; On the active controllers, there is an interface that allows the application, for instance the Hadoop Job Scheduler, to interact with the network controller about the job status.

Furthermore, applications can also leverage network awareness such that they adapt themselves to the network conditions like for instance bandwidth usage and topology. This has been demonstrated in [11] for different types of applications including Big Data ones.

In summary, network awareness seems to be a very promising direction for Big Data applications and its early adoption has already shown improvements. Programmable networks are a fundamental enabler for leveraging the statefulness of the controllers, and accordingly provide customized support for Big Data applications.



## References

- [1] *Cisco Data Center Infrastructure 2.5 Design Guide*. 2008.
- [2] ABU-LIBDEH, H., COSTA, P., ROWSTRON, A., O’SHEA, G., AND DONNELLY, A. Symbiotic routing in future data centers. *Computer Communication Review* 40, 4 (Aug. 2010), 51–62.
- [3] AL-FARES, M., LOUKISSAS, A., AND VAHDAT, A. A scalable, commodity data center network architecture. In *SIGCOMM Computer Communication Review* (2008), vol. 38, ACM, pp. 63–74.
- [4] AL-FARES, M., RADHAKRISHNAN, S., RAGHAVAN, B., HUANG, N., AND VAHDAT, A. Hedera: Dynamic flow scheduling for data center networks. In *Symposium on Networked Systems Design and Implementation*, NSDI.
- [5] ARMBRUST, M., FOX, A., GRIFFITH, ET AL. A view of cloud computing. *Communications of the ACM* 53, 4 (2010), 50–58.
- [6] BARHAM, P., DRAGOVIC, B., FRASER, K., HAND, S., HARRIS, T., HO, A., NEUGEBAUER, R., PRATT, I., AND WARFIELD, A. Xen and the art of virtualization. *ACM SIGOPS Operating Systems Review* 37, 5 (2003), 164–177.
- [7] BLAKE, S., BLACK, D., CARLSON, M., DAVIES, E., WANG, Z., AND WEISS, W. RFC 2475: An Architecture for Differentiated Service, 1998.
- [8] BRADEN, R., ZHANG, L., BERSON, S., HERZOG, S., AND JAMIN, S. RFC 2205: Resource ReSerVation Protocol (RSVP) – Version 1 Functional Specification, September 1997.
- [9] CHEN, K., GUO, C., WU, H., YUAN, J., FENG, Z., CHEN, Y., LU, S., AND WU, W. Dac: Generic and automatic address configuration for data center networks. *Transactions on Networking* 20, 1 (Feb 2012), 84–99.
- [10] CHOWDHURY, M., AND STOICA, I. Coflow: A networking abstraction for cluster applications. In *Workshop on Hot Topics in Networks* (Redmond, VA, USA, 2012), HotNets, ACM, pp. 31–36.
- [11] CHOWDHURY, M., ZAHARIA, M., MA, J., AND JORDAN. Managing data transfers in computer clusters with orchestra. In *SIGCOMM Computer Communication Review* (2011), vol. 41, ACM, pp. 98–109.
- [12] CHOWDHURY, S. R., BARI, M. F., AHMED, R., AND BOUTABA, R. PayLess: A Low Cost Network Monitoring Framework for Software Defined Networks. In *Network Operations and Management Symposium* (2014), NOMS, IEEE/IFIP.

- [13] CURTIS, A. R., MOGUL, J. C., TOURRILHES, J., YALAGANDULA, P., SHARMA, P., AND BANERJEE, S. Devoflow: scaling flow management for high-performance networks. In *Computer Communication Review* (2011), vol. 41, ACM SIGCOMM, pp. 254–265.
- [14] DAS, A., LUMEZANU, C., ZHANG, Y., SINGH, V., JIANG, G., AND YU, C. Transparent and flexible network management for big data processing in the cloud. In *Workshop on Hot Topics in Cloud Computing* (Berkeley, CA, 2013), USENIX.
- [15] DINH, H. T., LEE, C., NIYATO, D., AND WANG, P. A survey of mobile cloud computing: architecture, applications, and approaches. *Wireless communications and mobile computing* 13, 18 (2013), 1587–1611.
- [16] DIXIT, A., PRAKASH, P., AND KOMPELLA, R. R. On the efficacy of fine-grained traffic splitting protocols in data center networks. In *SIGCOMM* (2011), ACM.
- [17] DOGAR, F. R., KARAGIANNIS, T., BALLANI, H., AND ROWSTRON, A. Decentralized task-aware scheduling for data center networks. In *SIGCOMM* (New York, NY, USA, 2014), ACM.
- [18] FARRINGTON, N., PORTER, G., RADHAKRISHNAN, S., BAZZAZ, H. H., SUBRAMANYA, V., FAINMAN, Y., PAPEN, G., AND VAHDAT, A. Helios: A hybrid electrical/optical switch architecture for modular data centers. In *SIGCOMM* (2010), ACM, pp. 339–350.
- [19] FERGUSON, A. D., GUHA, A., LIANG, C., FONSECA, R., AND KRISHNAMURTHI, S. Participatory networking: An api for application control of sdns. In *SIGCOMM* (2013), ACM, pp. 327–338.
- [20] GUO, C., WU, H., TAN, K., SHI, L., ZHANG, Y., AND LU, S. Dcell: A scalable and fault-tolerant network structure for data centers. In *Conference on Data Communication* (2008), SIGCOMM, ACM.
- [21] ISELT, A., KIRSTADTER, A., PARDIGON, A., AND SCHWABE, T. Resilient routing using mpls and ecmp. In *Workshop on High Performance Switching and Routing* (2004), HPSR, IEEE, pp. 345–349.
- [22] KAVULYA, S., TAN, J., GANDHI, R., AND NARASIMHAN, P. In *International Conference on Cluster, Cloud and Grid Computing*, CCGrid, IEEE/ACM.
- [23] KIVITY, A., KAMAY, Y., LAOR, D., LUBLIN, U., AND LIGUORI, A. kvm: the linux virtual machine monitor. In *Proceedings of the Linux Symposium* (2007), vol. 1, pp. 225–230.

- [24] LEE, K.-H., LEE, Y.-J., CHOI, H., CHUNG, Y. D., AND MOON, B. Parallel data processing with mapreduce: A survey. *SIGMOD Rec.* 40, 4 (Jan. 2012), 11–20.
- [25] LEISERSON, C. E. Fat-trees: Universal networks for hardware-efficient supercomputing. vol. 1985, IEEE, pp. 892–901.
- [26] LI, Z., SHEN, Y., YAO, B., AND GUO, M. Ofscheduler: a dynamic network optimizer for mapreduce in heterogeneous cluster. Springer, pp. 1–17.
- [27] LU, G., GUO, C., LI, Y., ZHOU, Z., YUAN, T., WU, H., XIONG, Y., GAO, R., AND ZHANG, Y. Serverswitch: A programmable and high performance platform for data center networks. In *Conference on Networked Systems Design and Implementation* (2011), vol. 11 of *NSDI*, Usenix, pp. 2–2.
- [28] MCKEOWN, N., ANDERSON, T., BALAKRISHNAN, H., PARULKAR, G., PETERSON, L., REXFORD, J., SHENKER, S., AND TURNER, J. Openflow: enabling innovation in campus networks. *SIGCOMM Computer Communication Review* 38, 2 (2008), 69–74.
- [29] NIRANJAN MYSORE, R., PAMBORIS, A., FARRINGTON, N., HUANG, N., MIRI, P., RADHAKRISHNAN, S., SUBRAMANYA, V., AND VAHDAT, A. Portland: A scalable fault-tolerant layer 2 data center network fabric. In *Conference on Data Communication* (2009), SIGCOMM, ACM, pp. 39–50.
- [30] PANDEY, S., WU, L., GURU, S. M., AND BUYYA, R. A particle swarm optimization-based heuristic for scheduling workflow applications in cloud computing environments. In *International Conference on Advanced Information Networking and Applications* (2010), AINA, IEEE, pp. 400–407.
- [31] ROSCHKE, S., CHENG, F., AND MEINEL, C. Intrusion detection in the cloud. In *Dependable, Autonomic and Secure Computing, 2009. DASC'09. Eighth IEEE International Conference on* (2009), IEEE, pp. 729–734.
- [32] ROSENBLUM, M. Vmwares virtual platform. In *Proceedings of hot chips* (1999), vol. 1999, pp. 185–196.
- [33] SINGLA, A., HONG, C.-Y., POPA, L., AND GODFREY, P. B. Jellyfish: Networking data centers randomly. In *Conference on Networked Systems Design and Implementation* (2012), NSDI, USENIX Association.
- [34] TOSHNIWAL, A., TANEJA, S., SHUKLA, A., RAMASAMY, K., PATEL, J. M., KULKARNI, S., JACKSON, J., GADE, K., FU, M., DONHAM, J., BHAGAT, N., MITTAL, S., AND RYABOY, D. Storm@twitter. In *SIGMOD International Conference on Management of Data* (2014), ACM.

- [35] WANG, G., NG, T. E., AND SHAIKH, A. Programming your network at run-time for big data applications. In *First Workshop on Hot Topics in Software Defined Networks* (2012), HotSDN, ACM, pp. 103–108.
- [36] WEBB, K. C., SNOEREN, A. C., AND YOCUM, K. Topology switching for data center networks. In *Hot Topics in Management of Internet, Cloud, and Enterprise Networks and Services* (2011), Hot-ICE, USENIX, pp. 14–14.
- [37] YAO, Y., CAO, J., AND LI, M. A network-aware virtual machine allocation in cloud datacenter. In *Network and Parallel Computing*, vol. 8147 of *Lecture Notes in Computer Science*. Springer, 2013.