

Keddah: Capture Hadoop Networking Behaviour

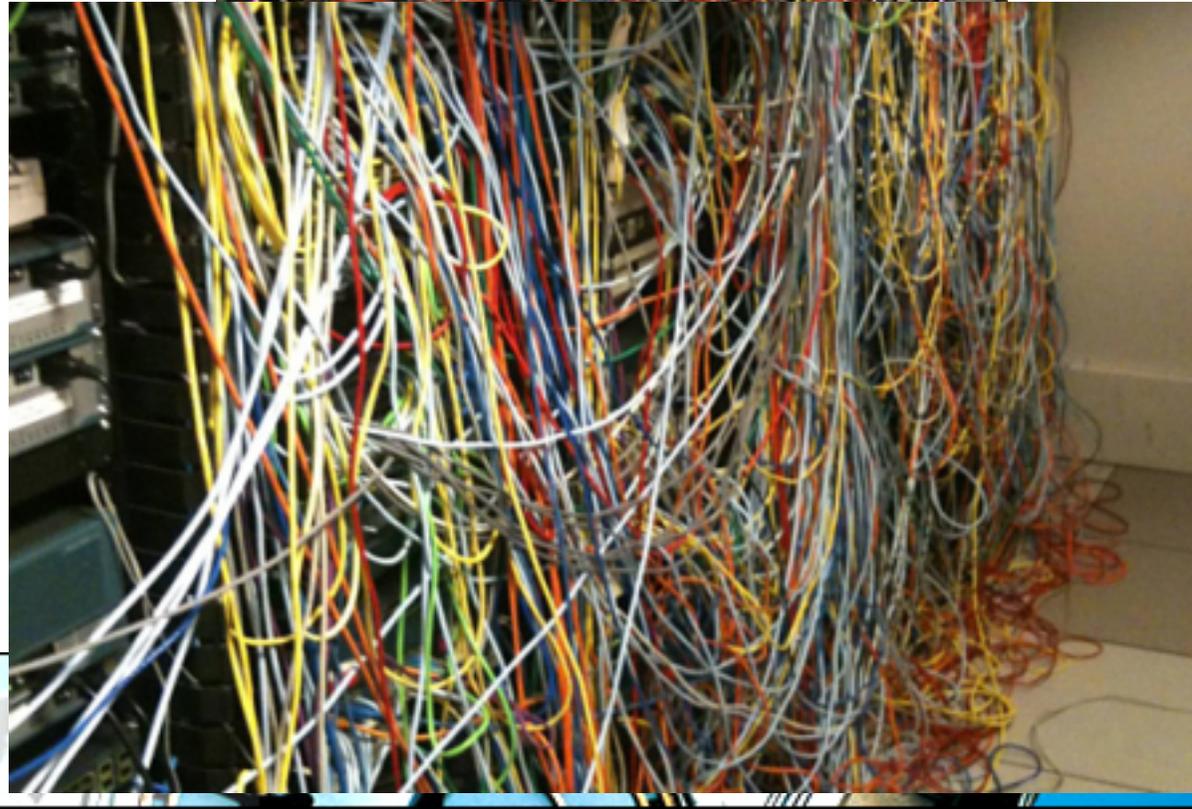
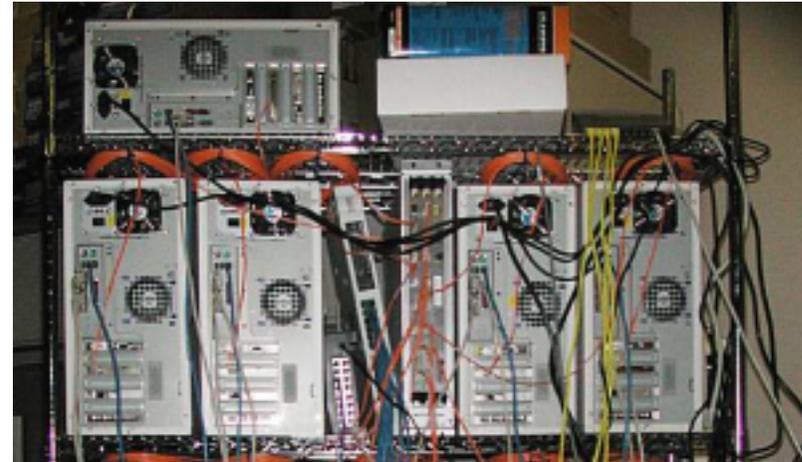


Presenter: Jie Deng

Supervisor: Steve Uhlig, Felix Cuadrado, Gareth Tyson

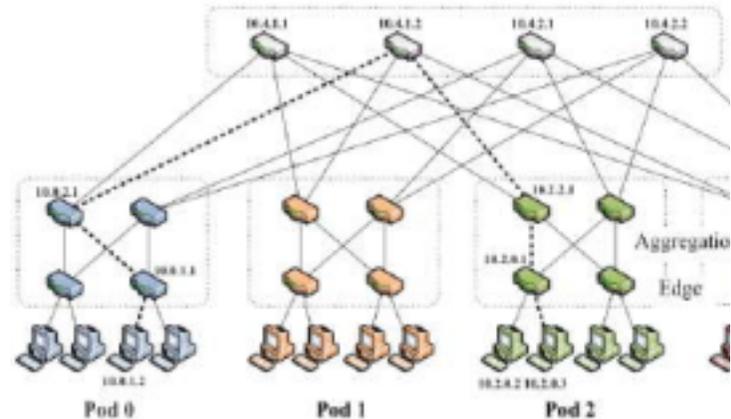
Distributed systems

- Network is essential
- But not well studied yet

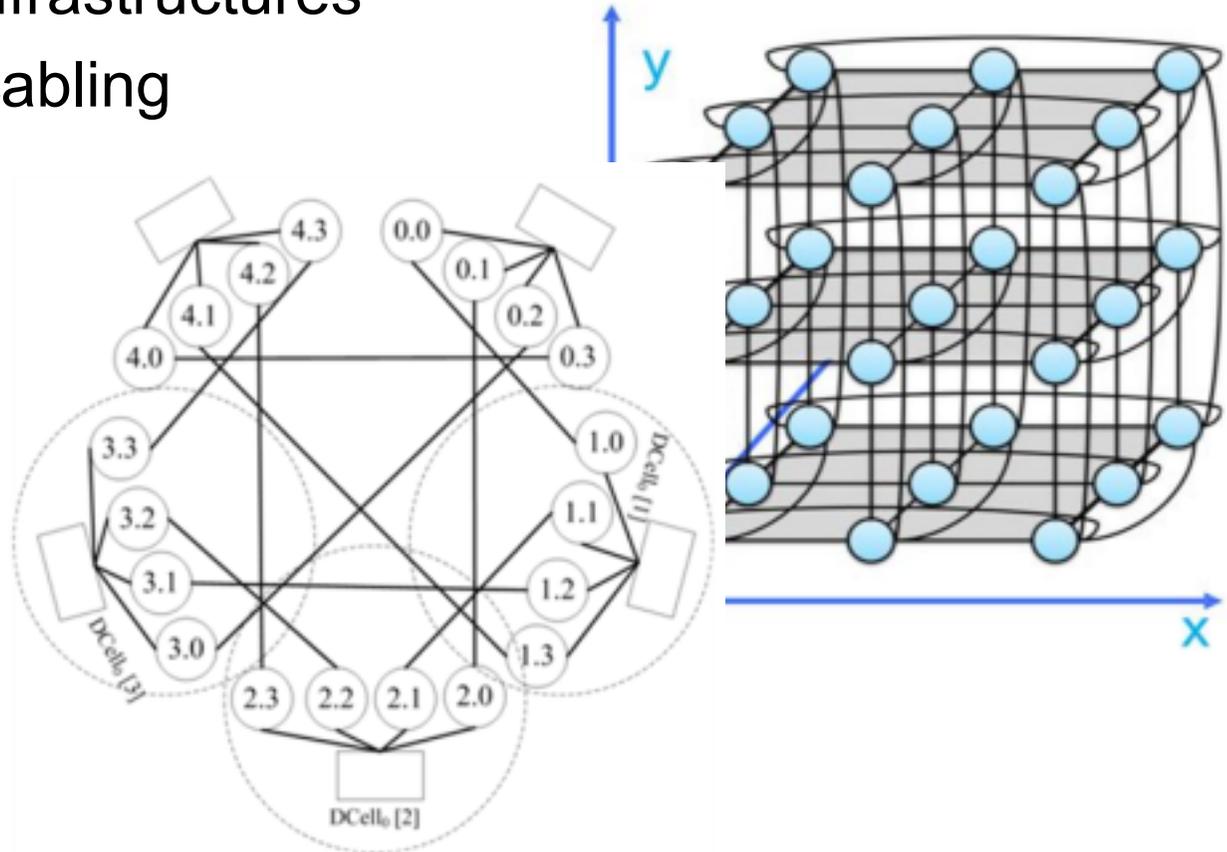


Networks are complicated and costly...

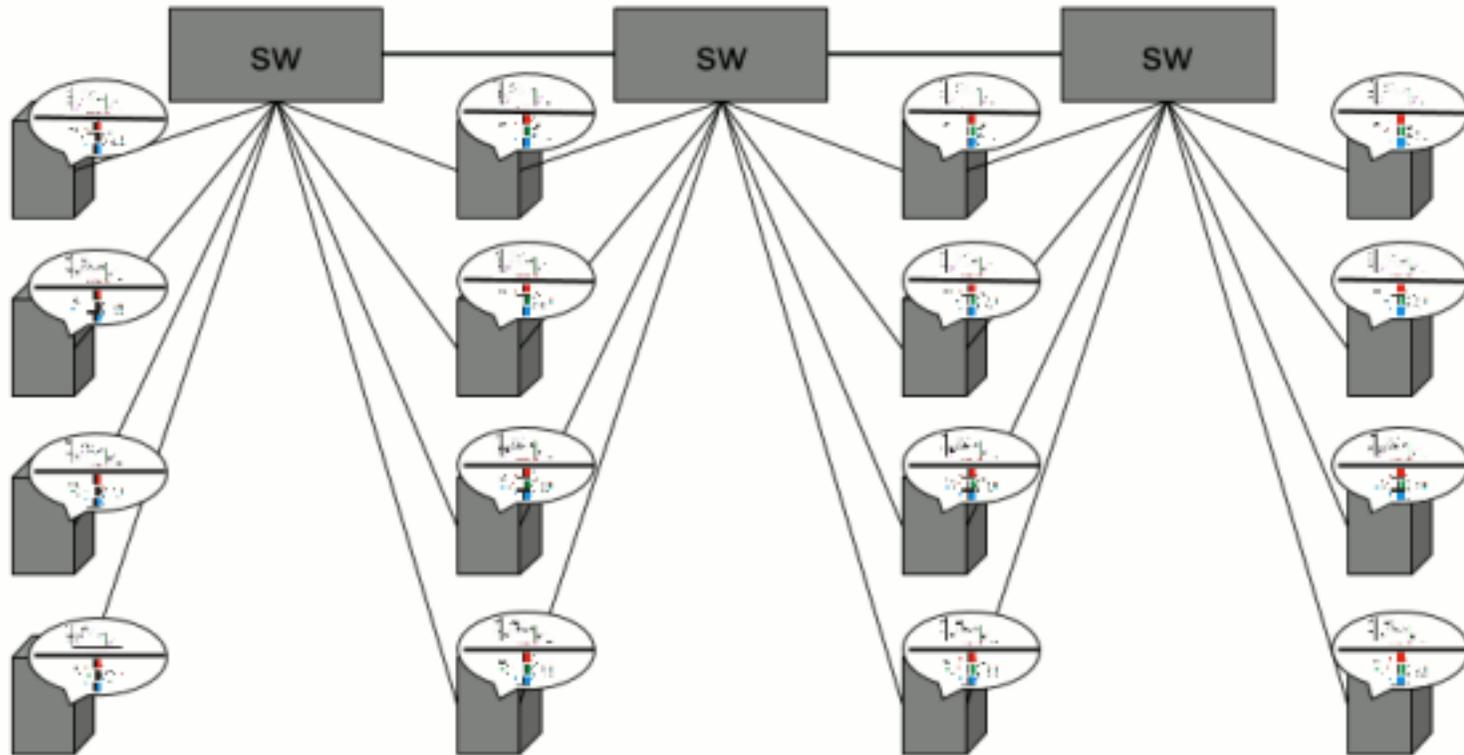
- Cost consuming infrastructures
- Time consuming cabling



Simple fat-tree topology. Using the two-level routing tables packets from source destination 10.2.0.3 would take the dashed path.



What we want... Traffic generator!



Methodology

Testbed:

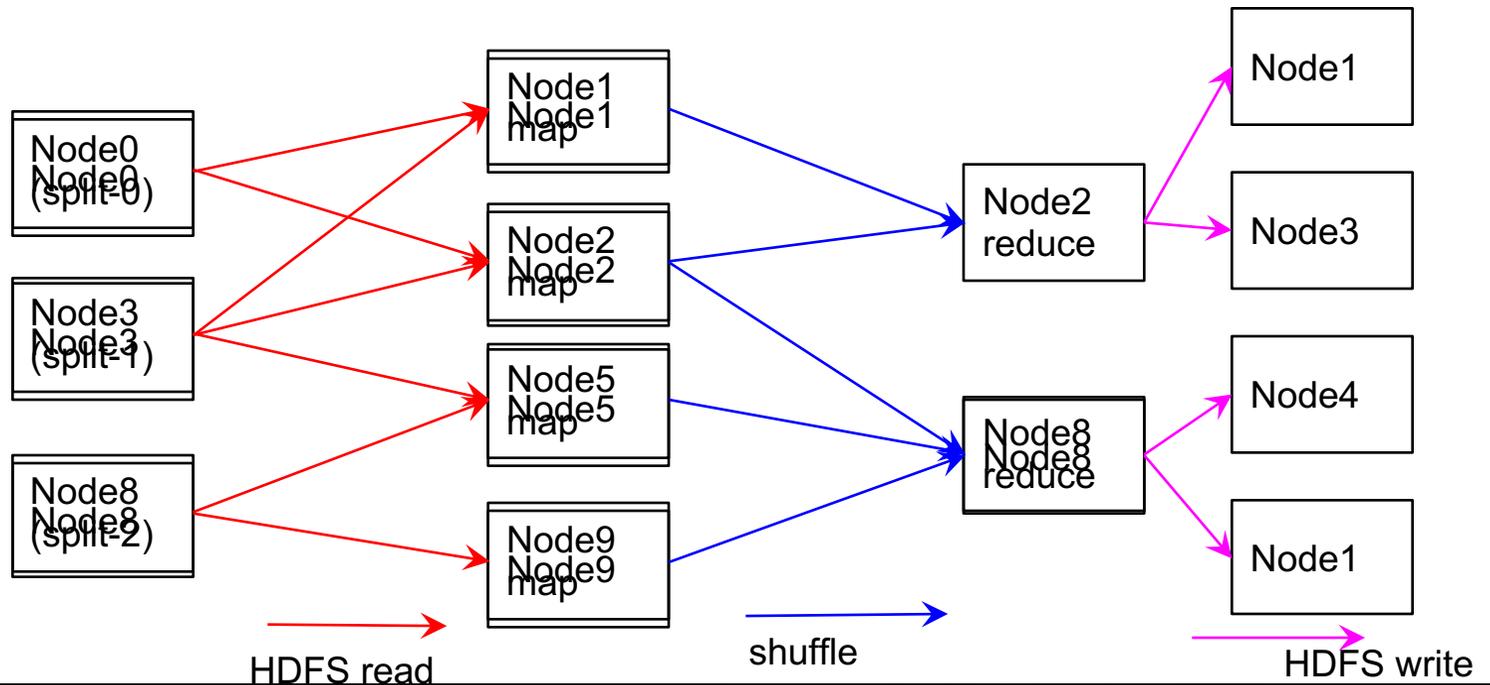
- 16 nodes w/ Hadoop & sflow installed

Reproduce Hadoop traffic:

- General over jobs
- Extendable over parameters
- Replicable over clusters

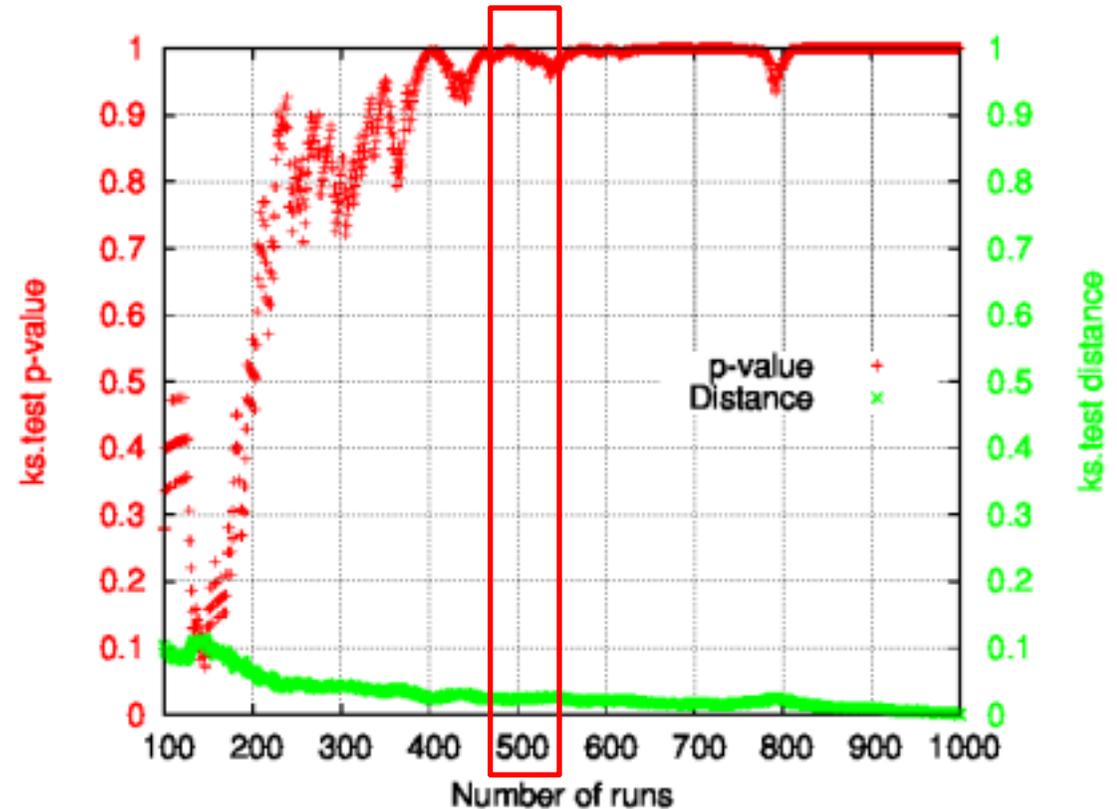
Hadoop

- Storage: HDFS
- Processing: MapReduce jobs



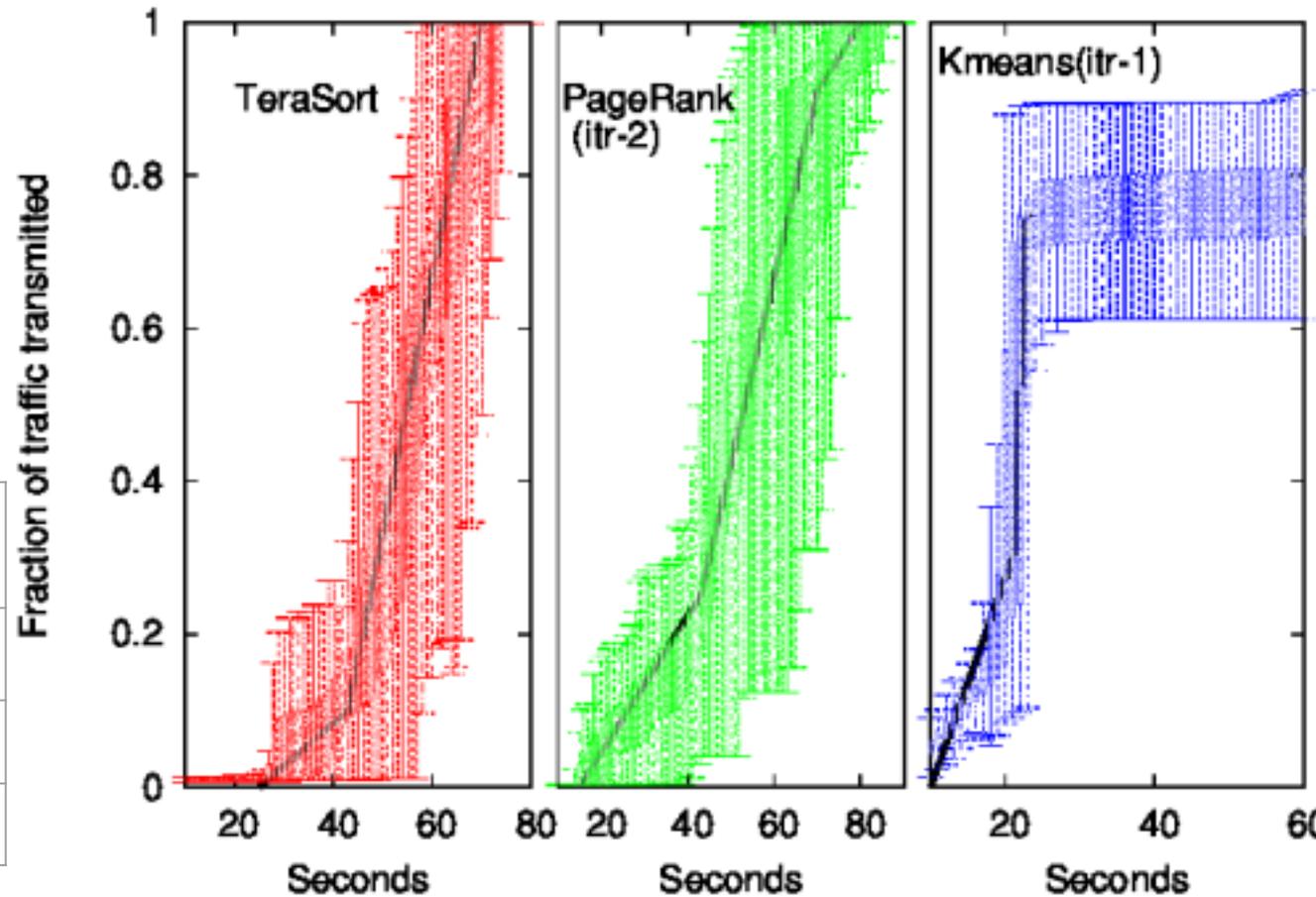
Capture the randomness...

- Nodes are randomly selected in Hadoop, thus the distribution.



Traffic Model

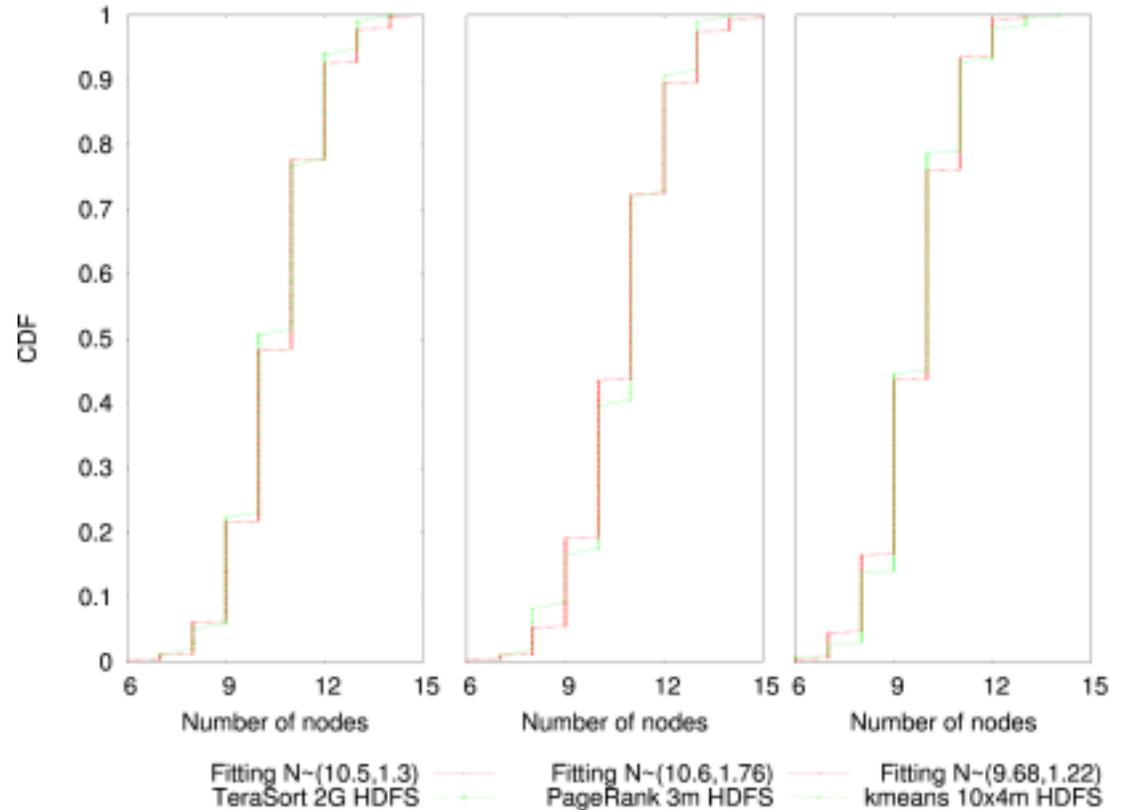
HDFS import
 -> Shuffle
 -> HDFS export



Job	HDFS import	Shuffle	HDFS export
TeraSort	0.1	0.9	0
PageRank	0.25	0.66	0.09
kmeans	0.31	0	0.67

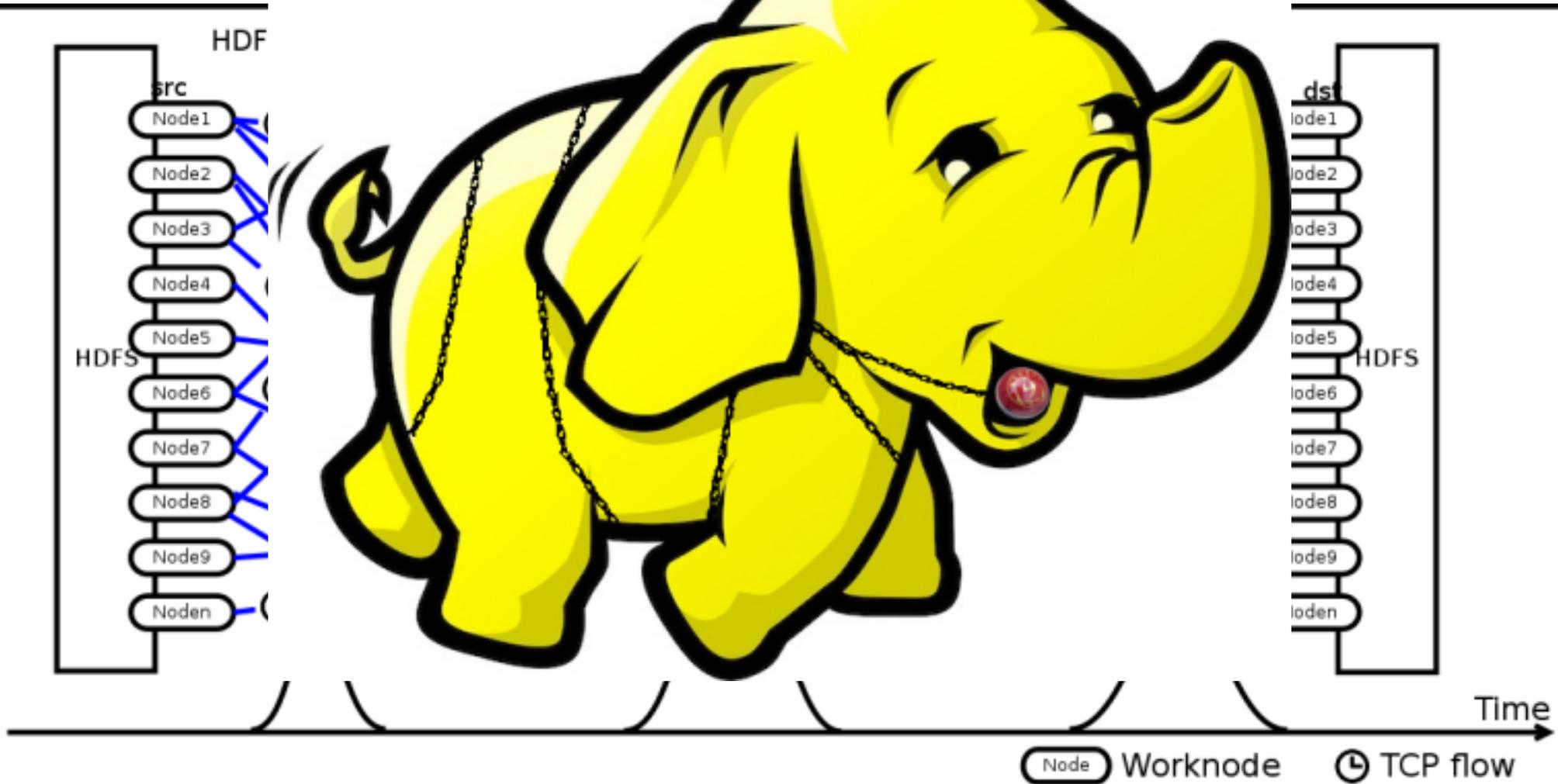
Modelling

- Number of nodes
- Number of flows
- Flow size
- Flow start time



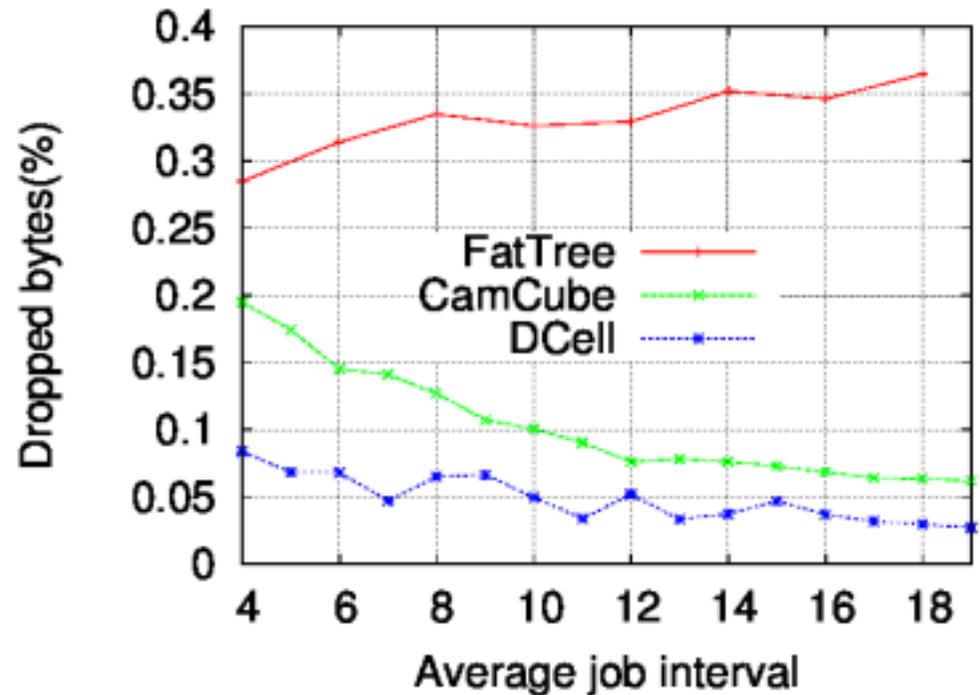
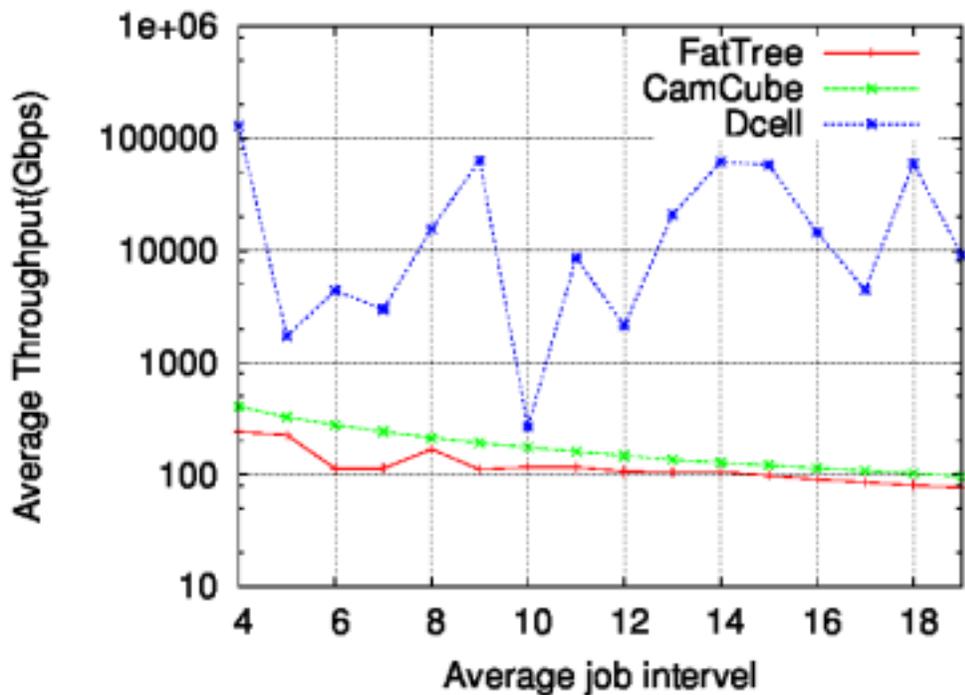
Keddah in ns3

- HDFS imp



Usefulness validation

- DCN topologies: FatTree(8x8) CamCube(4) Dcell(2,3)



- GitHub: <https://git.io/vKelv>
- Current work
 - Produce fine-grained workload traffic
 - Network(bandwidth) as a resource
- Thanks!
- Q&A

Resistance from systems

- Time consuming software engineering

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/types.h>
#include <arpa/inet.h>

void serveur1(portServ ports)
{
    int sockServ1, sockServ2, sockClient;
    struct sockaddr_in monAddr, addrClient;
    socklen_t lenAddrClient;

    if ((sockServ1 = socket(AF_INET, SOCK_STREAM, 0)) < 0)
        perror("Erreur socket");
    exit(1);
    if ((sockServ2 = socket(AF_INET, SOCK_STREAM, 0)) < 0)
        perror("Erreur socket");
    exit(1);
}

bzero(&monAddr, sizeof(monAddr));
monAddr.sin_family = AF_INET;
monAddr.sin_port = htons(ports.port);
monAddr.sin_addr.s_addr = INADDR_ANY;
bzero(&addrServ2, sizeof(addrServ2));
```

```
Kernel-o-Matic - vagrant@vagrant-ubuntu-precise-32: ~ - ssh - 80x24
CC [M] drivers/staging/rtl8712/rtl8712_cmd.o
CC [M] fs/nls/nls_utf8.o
CC [M] net/mac80211/mesh_sync.o
LD fs/nls/built-in.o
LD fs/ntfs/built-in.o
CC [M] fs/ntfs/aops.o
CC [M] drivers/net/wireless/ath/ath
CC [M] drivers/staging/rtl8712/rtl8
CC [M] drivers/media/usb/dvb-usb/te
CC [M] net/mac80211/mesh_ps.o
CC [M] drivers/net/wireless/b43/tab
CC [M] drivers/net/wireless/b43/sys
CC [M] drivers/media/usb/dvb-usb-v2
CC [M] drivers/net/wireless/ath/ath
CC [M] drivers/media/usb/dvb-usb/tt
CC [M] net/mac80211/pm.o
CC [M] drivers/media/usb/dvb-usb/um
CC [M] drivers/net/wireless/b43/xml
CC [M] fs/ntfs/attrib.o
CC [M] drivers/staging/rtl8712/rtl8
CC [M] drivers/media/usb/dvb-usb/vp
CC [M] drivers/media/usb/dvb-usb-v2
CC [M] drivers/net/wireless/ath/ath
```



MapReduce Job

job_1400191378494_0002

Logged in as: drhiv

Job Overview

Job Name:	scopedTarget wordcountJob
User Name:	hivert
Queue:	default
State:	SUCCEEDED
Urgency:	false
Started:	Thu May 15 22:09:22 EDT 2014
Finished:	Thu May 15 22:09:38 EDT 2014
Elapsed:	15sec
Diagnostics:	
Average Map Time:	0sec
Average Reduce Time:	0sec
Average Shuffle Time:	4sec
Average Merge Time:	0sec

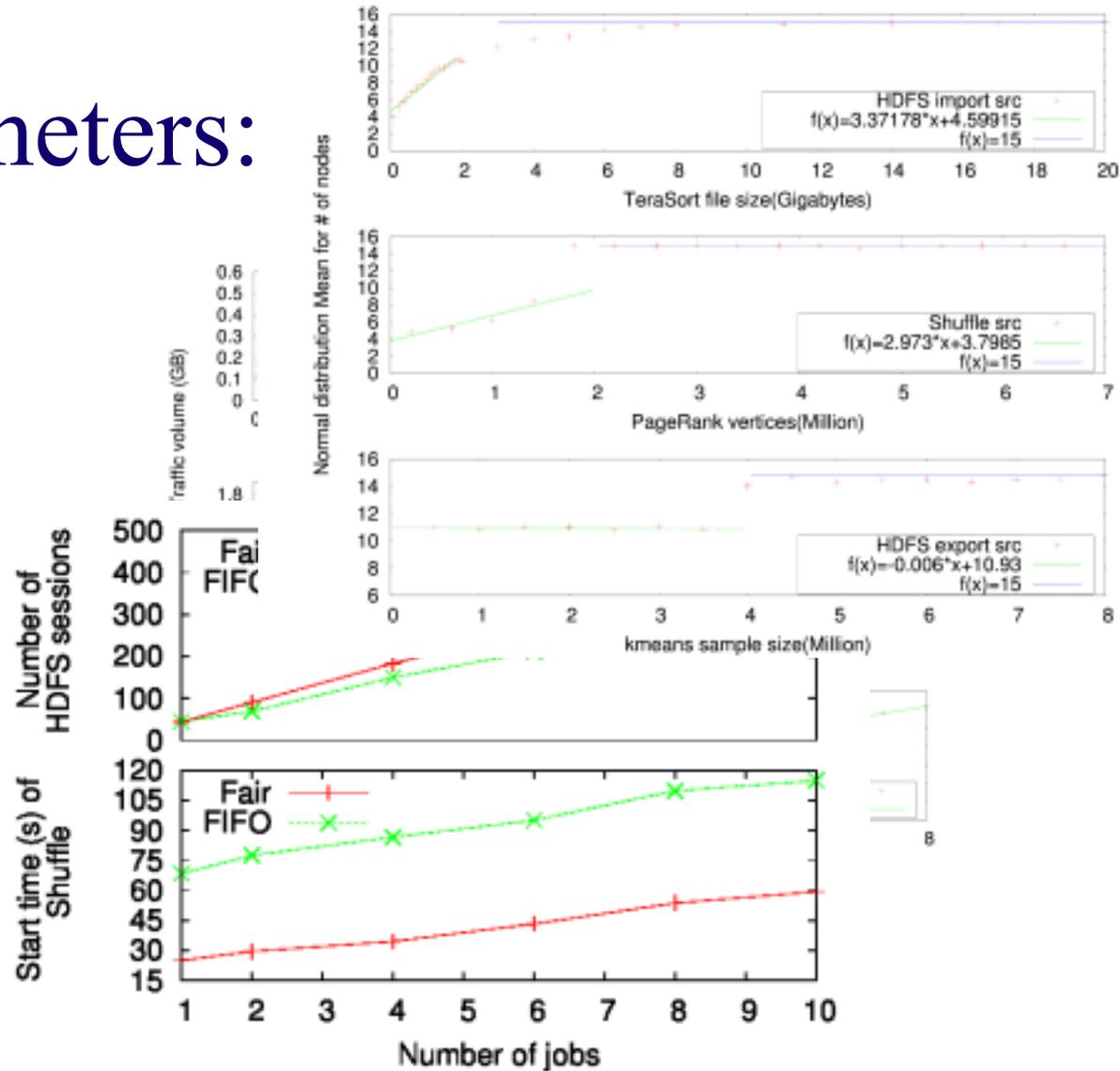
ApplicationMaster			
Attempt Number	Start Time	Node	Logs
1	Thu May 15 22:09:18 EDT 2014	10.0.1.4:8042	logs

Task Type	Total	Complete
Map	1	1
Reduce	1	1

Attempt Type	Failed	Killed	Successful
Maps	0	0	1
Reduces	0	0	1

Extend the parameters:

- Job parameters
- Scheduler



Use cases2

- TCP vs DCTCP

