

,

()

()

()

Big Data

()

:

II

,

-20-1

(,)

123 «

'

»

()

-

(- -)

,

()

:

(, ,)

()

(,)

,

()

123 « ' »

()

-

(- -)

,

()

:

“ ” 20 .

(, ,)

1.
Big Data

“ 05 ” 2021 . 1656

2. 13 2021 .

3. 1) : ; 2) :

Apache Hadoop; 3) : , .

4. , _____ , ;

1) _____ ;

2) _____ ;

3) _____ ;

4) _____ ;

6) _____ .

5. _____ , _____ , _____ , _____ , _____
 () _____
 - -12

6. _____ , _____ .1) (_____)

	(_____ , _____ , _____ , _____)		

1		09.11.21-12.11.21	
2	NTMA	13.11.21-18.11.21	
3		19.11.21-29.11.21	
4		30.11.21-03.12.21	
5		04.12.21-07.12.21	
6		08.12.21-09.12.21	
7		10.12.21-11.12.21	

08 2021 .

_____ () _____
 | _____ () _____ (; ,) _____

: 72 ., 22 ., 10 ., 1

., 24 .

, , ,
,

, .

.

Big Data.

.

,

TE,

.

ABSTRACT

Master's thesis: 72 pages, 22 figures, 10 tables, 1 appendix, 24 sources.

BIG DATA, BIG DATA PLATFORMS, MACHINE LEARNING,
NETWORK MEASUREMENTS, TRAFFIC ANALYSIS.

The purpose of the qualification work is to develop a method of monitoring computer network traffic based on big data technologies.

During the qualification work, a detailed method of monitoring for the Big Data collection and generation of traffic statistics using measured values is proposed. To confirm the effectiveness of this method, the implementation was performed using special tools Big Data. The Big Data analysis methods for real-time analysis of a large amount of stored historical data is also investigated. The results of Big Data analysis can help solve TE problems such as switching load balancing and data plane fault tolerance.

		,	,	,	
					8
					9
1					
					11
1.1	NTMA	Big Data			11
1.2					12
1.2.1					12
1.2.2					13
1.2.3					15
1.3	NTMA				18
1.3.1					18
1.3.2	NTMA				19
1.3.3			NTMA		22
2					
	SDN				24
2.1		SDN			24
2.2					26
2.3					27
3				SDN	
					29
3.1					29
3.2					30
3.2.1					33
3.2.2					35
3.2.3					35
3.2.1					36

3.2.2	38
3.2.3	40
3.2.4	41
3.2.5	41
3.3	43
3.3.1	43
3.3.2	44
3.4	45
4	47
4.1	50
4.2	51
4.3	52
4.4	54
4.5	55
4.6	57
	60
	61
	64

ER – « – » (., Entity-Relation)

HDFS – Hadoop (., Hadoop Distributed File System)

KDD – (., Knowledge Discovery in Data)

NTMA – (., Network Traffic Monitoring and Analysis)

OD – « – » (., Origin–Destination)

OF – OpenFlow

QoE – (., Quality of Experience)

QoS – (., Quality of Service)

SDN – - (., Software-defined networking)

TE – (., Traffic Engineering)

TM – (., Traffic Matrix)

(., NTMA – Network Traffic Monitoring and Analysis)

, . - , NTMA.

. , , : , NTMA . NTMA.

NTMA, (,) NTMA. (., SDN – Software-defined networking).

. , , . OpenFlow (OF),

(., TE – Traffic Engineering). TE ,

TE

«

»

TE.

1

1.1 NTMA Big Data

(NTMA).

NTMA

NTMA.

Big Data),

(„

NTMA [1, 2, 3].

HDFS – Hadoop Distributed File System;

Spark;

MLlib Apache Giraph.

NTMA [4, 5, 6],

NTMA (NTMA)

[7, 8], NTMA [9, 10, 11],

NTMA

1.2

1.2.1

[12]

«

».

[13, 14]

«V's», – (Volume, Velocity, Variety, Veracity, Value).

[6, 12, 15]

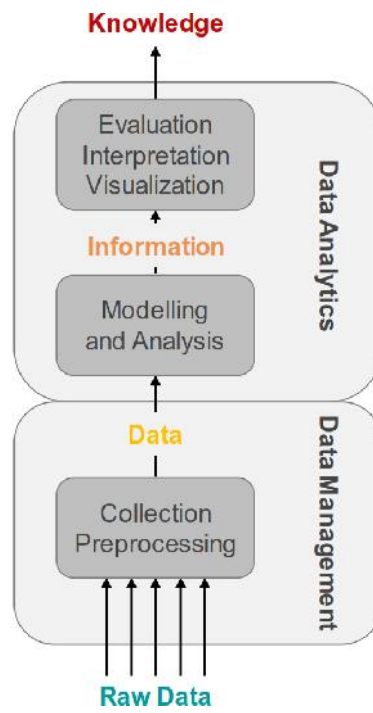
[4, 5]

«5-V»,

NTMA.

1.2.2

(., KDD – Knowledge Discovery in Data)



1.1 –

KDD

(, «5-V»)

KDD

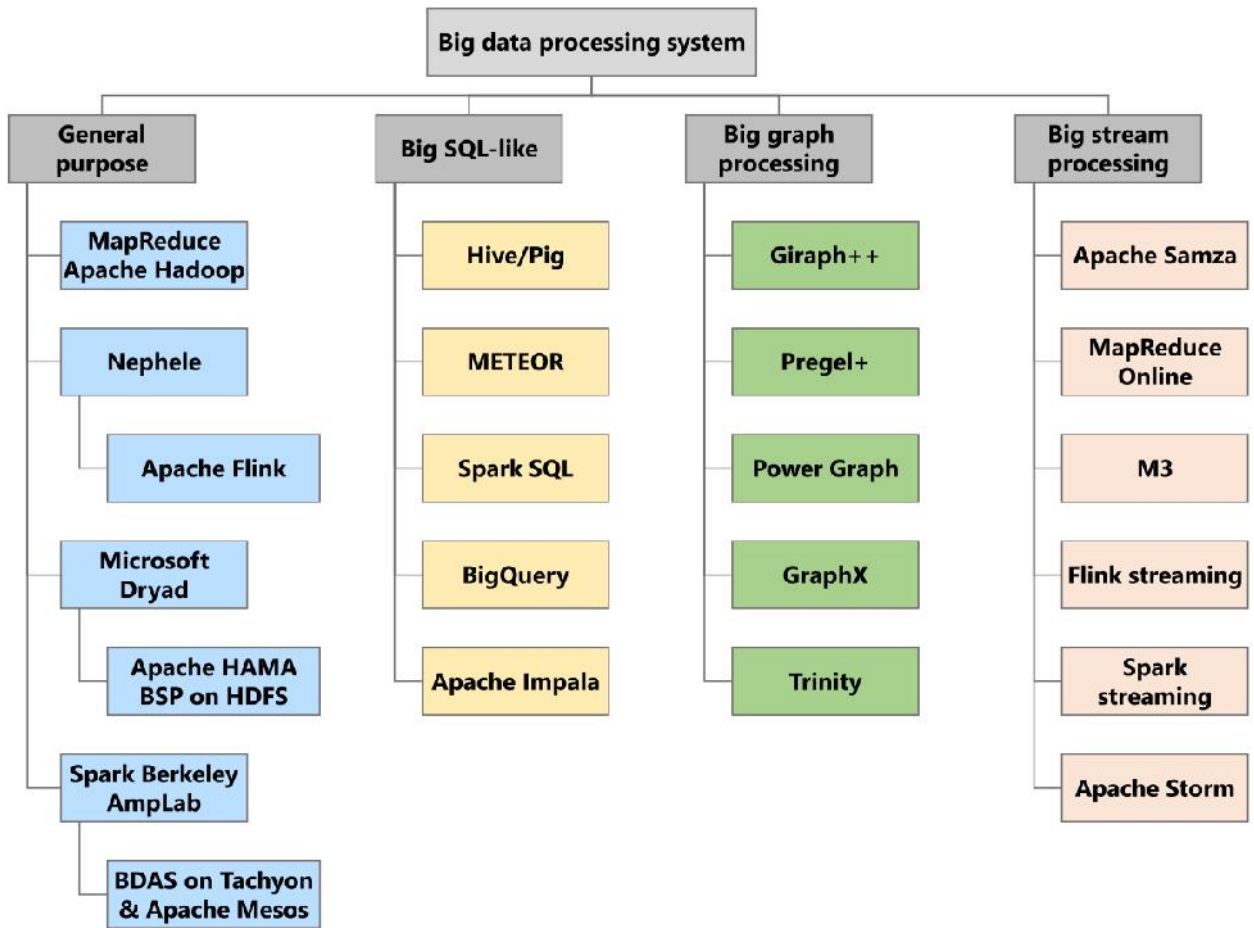
KDD.

1.2.3

[6, 15],

- SQL-

1.2
 MapReduce, Dryad, Flink Spark . Hive, HAWQ,
 Apache Drill Tajo SQL. Pregel, GraphLab
 , Storm S4

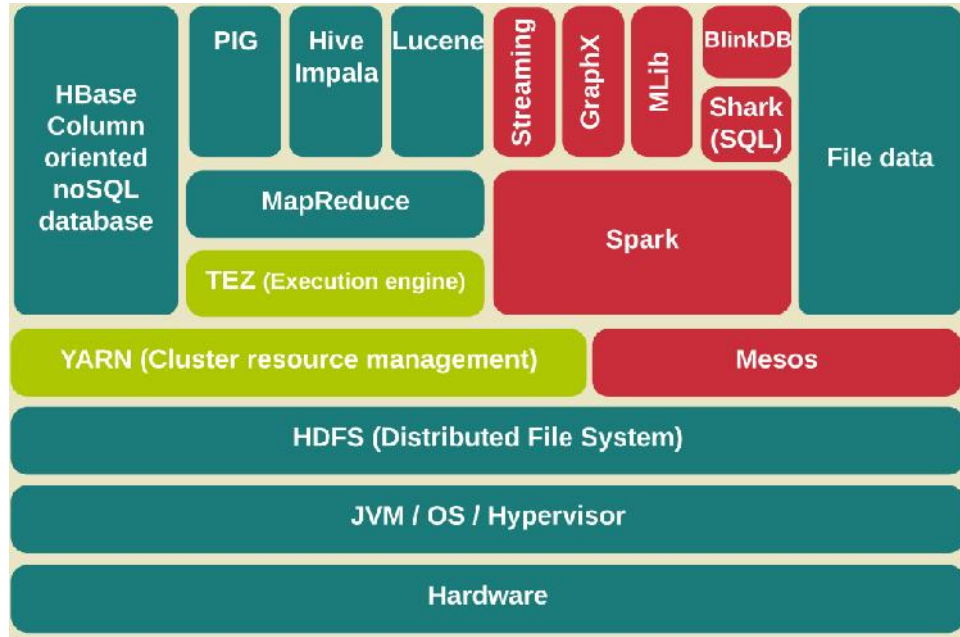


1.2 –

Hadoop.

1.3

MapReduce Hadoop Hadoop, Hadoop (HDFS). YARN – Mesos) (, TEZ) Spark.



1.3 – Hadoop 2.0 Spark ([14])

Spark Hadoop, 2.0, MapReduce. Spark («Map» «Reduce» –), . Spark Hadoop

Facebook, Sawzall Google, Pig Latin Yahoo, Hive
 SCOPE Microsoft, Mahout MapReduce MLlib
 Spark, Hadoop. GraphX Spark
 Streaming, Apache Hadoop,
 Cloudera CDH, Hortonworks HDP MapR Converged Data
 Platform.

1.3 NTMA

1.3.1

NTMA,

[16].

QoS/QoE.

[16]

NTMA

QoS/QoE

1.3.2 NTMA

1.1

QoS.

[16].

(, ARIMA SARIMA),
(,).

(,
)

TCP/UDP.

, - . -
 , ,
 . - ,
 , -
 (, HTTPS TCP). ,
 - .

1.1 – NTMA

	.
	, QoE, .
	.
	, ().

. ,
 :
 - ,
 « »
 ;
 - , ,
 ;
 - , ,
 . ,
 ;

[17].

() -

[16].

NTMA.

,
/ .

(., Variety):

,
, , .
— ,

DNS

,
(, , ,
).

(., Veracity),

.
. ,
. ,

(., Value)

. , , ,
,

2

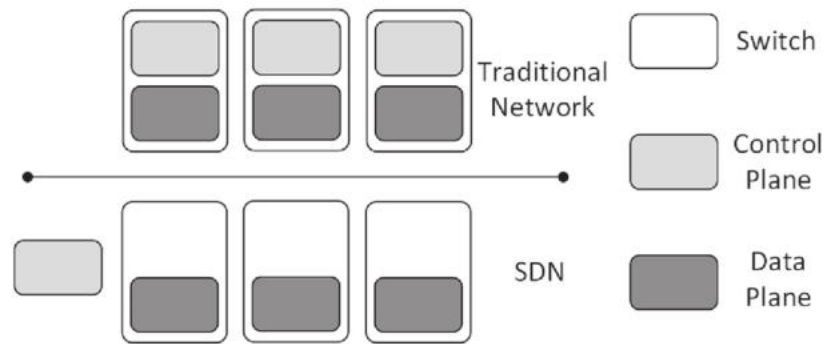
SDN

2.1 SDN

(SDN) ,
 ,
 ,
 ,
 ,

2.1

SDN.



2.1 –

SDN

:
 - - ;
 , ;

SDN

OpenFlow (OF) [18]. OF API,

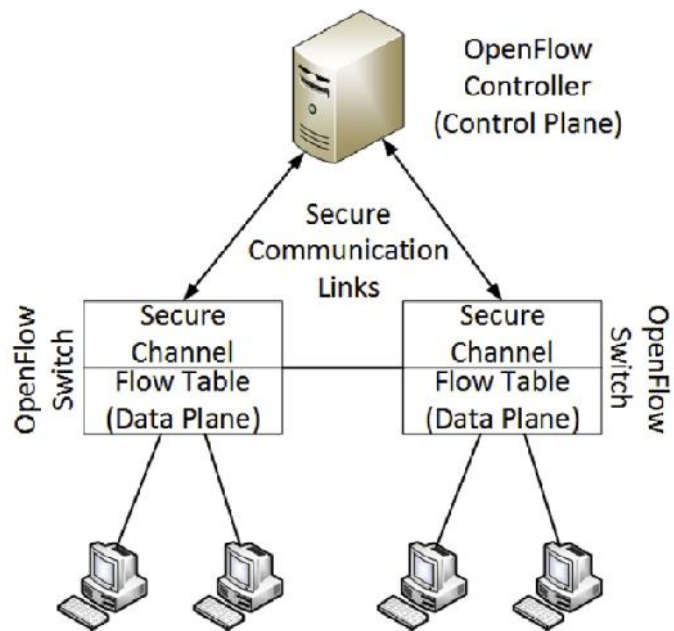
OF OF.

OF (2.2). OF

Open Networking Foundation (ONF)

« OpenFlow» [18].

[19].



2.2 –

OpenFlow

2.2

Traffic Engineering (TE)

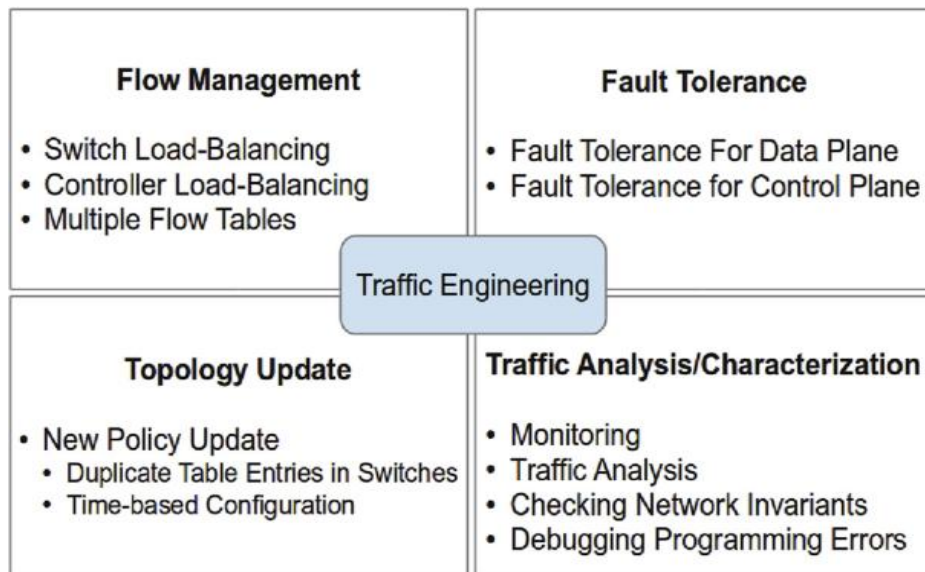
[20].

SDN

SDN

2.3:

[21].



,
.
.
.
.

API

Big Data,

3.1

OF

SDN

[18],

:

-
-
-

OF,

3.1.

1 -

SDN.

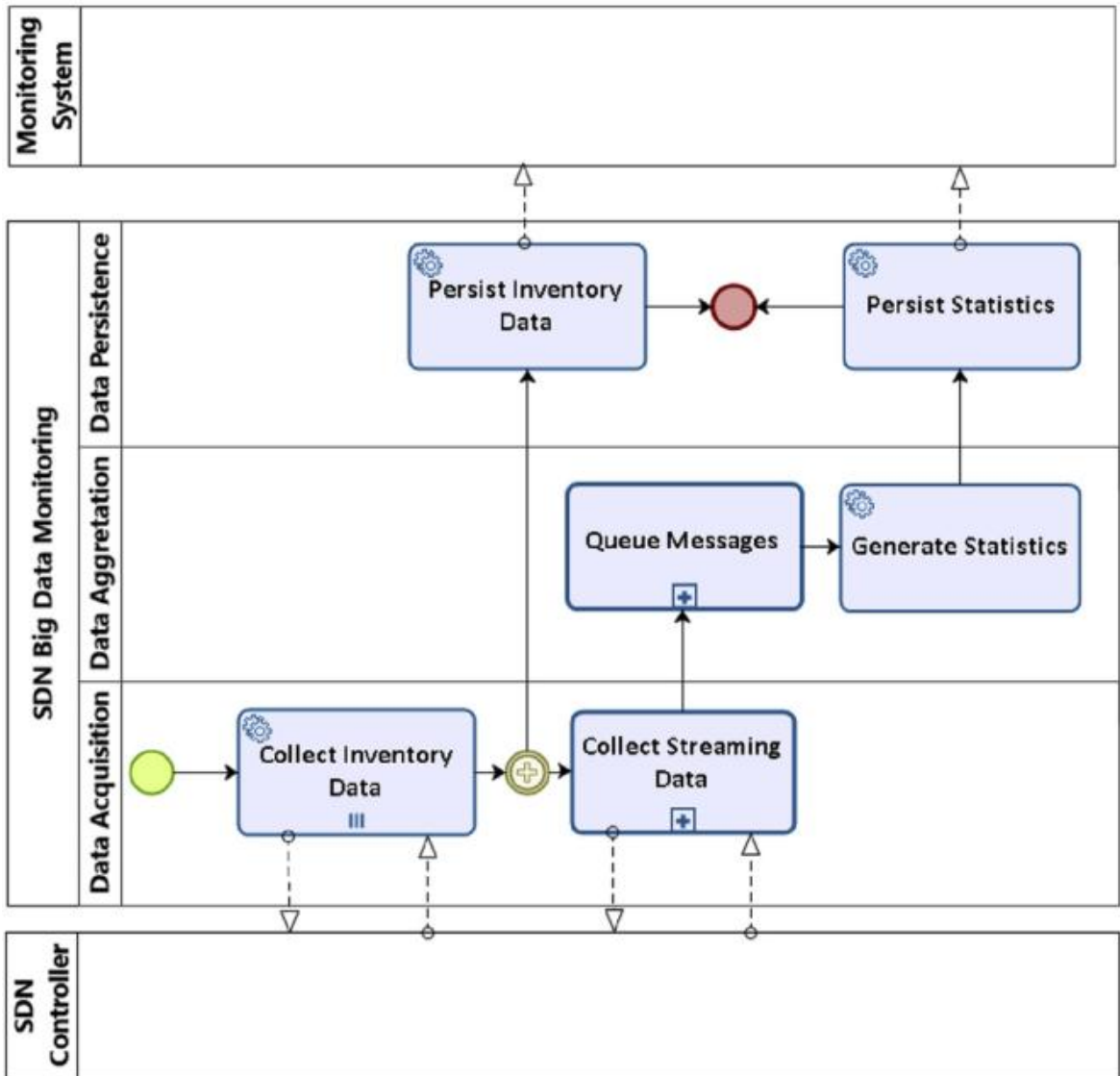
2 -

[22]

(

).

3 -



3.1 –

3.2

()

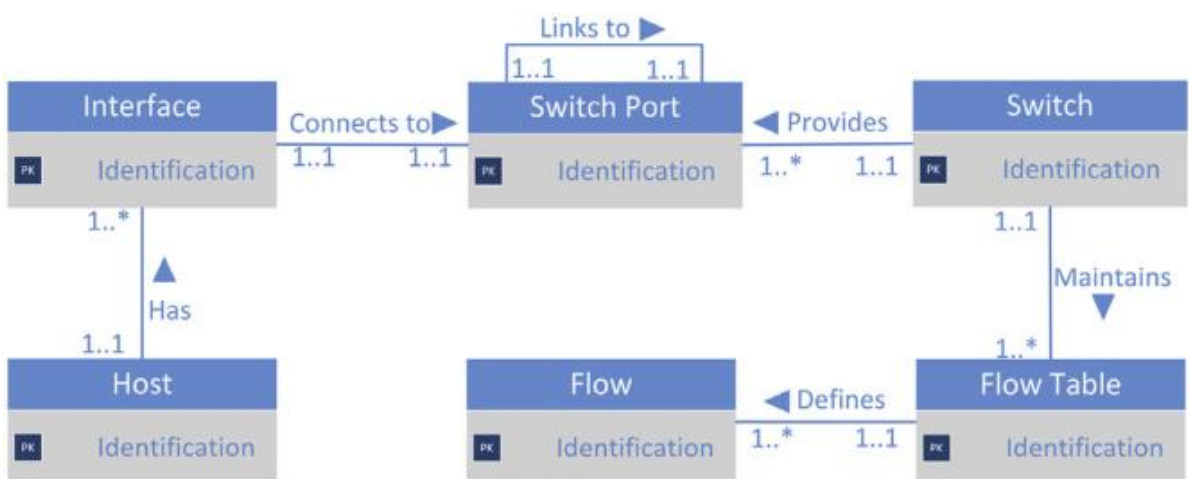
3.2 « - » (ER)

Host, Interface, Switch Port, Switch, FlowTable, Flow.

- 3.2:
- Host (« ») – ;
 - Interface (« ») – ;
 - Switch Port (« ») – ;
 - Switch (« ») – ;
 - FlowTable (« ») – ;
 - Flow (« ») – .

3.2,

- Has (« »). Host Interface
- Connects to (« »). Interface Interface;
- Links to (« ' »). Port ;
- Provides (« »). Switch Switch Port , Switch Port Switch;
- Maintains (« »). Switch Flow Table, Flow Table Switch;



3.2 – ER-

- Defines (« »). Flow Table
 Flow, Flow Flow Table.
 , Host Interface,
 OF. SDN ,
 ,
 , identification,
 .
 Host Interface. 3.1 ,
 Switch, OF.

3.1 – Switch

Identification	
MAC address	MAC-
Buffers	,
Tables	,
Capabilities [1..9]	(, ,)

, Switch Port
 , 3.2.
 Flow Table Flow .
 - , SDN.
 3.2.1
 ,
 . OF .

Connects to Links

to (3.2).

3.2 – Switch Port

Identification	
Hardware Address	MAC-
Name	' ,
Configuration	
State	
Current Feature	
Supported Features	
Peer Features	' ,
Advertised Features	' ,
Current Speed	/
Maximum Speed	/

, 3.2, Provides
Switch Switch Port.

Switch Port Switch. Maintains
Switch Flow Table. Defines Flow Table Flow
Flow Flow Table.

3.1.

SDN

(3-5 3.1).

n (7 3.1).

3.1 –

```

1 Procedure CollectData()
2   while TRUE do
3     CollectHostInventory()
4     CollectSwitchInventory()
5     CollectLinkInventory()
6     update Host, Switch, Switch Port, Flow
       Table, Flow repositories
7     sleep(n units of time)
8   end

```

3.2.2

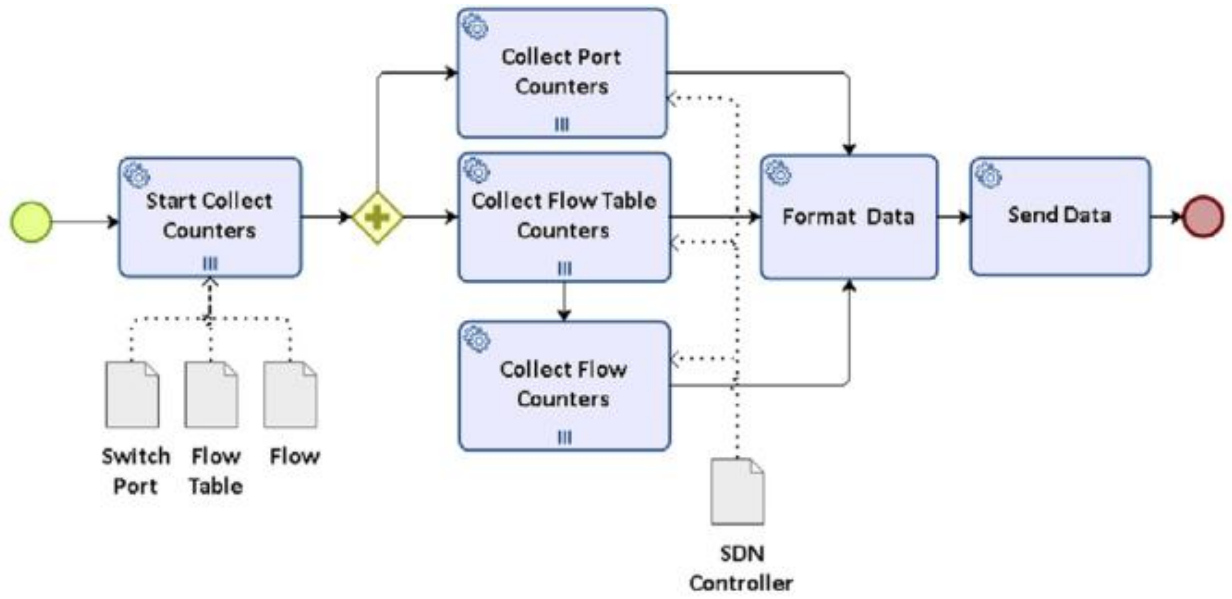
3.3,

API

3.2.3

$$C = \sum_{S=1}^n (p + ft), \tag{4.1}$$

C - ; n - ; p -
, ft - .



3.3 -

3.2.

5 11

eth0

if,

3.2.1

3.3.

3.3.

3.2 –

```

1 Procedure StartCollectCounters()
2   while TRUE do
3     read Switch Port repository
4     for each switch port do
5       if collect counters not started then
6         start collect port counters task
7       end
8     end
9     read Flow Table repository
10    for each flow table do
11      if collect counters not started then
12        start collect flow table counters
13        task
14      end
15    end
16    sleep (n units of time)
17  end

```

3.3 –

```

1 Procedure CollectPortCounters()
2   while TRUE do
3     CollectPortCounters(switch,port)
4     send counters to Format Data
5     sleep(n units of time)
6   end

```

CollectPortCounters(switch, port) (3)

SDN n (5) ,
 (switch) (port)

SDN.

3.3 –

Timestamp	,
Packets Received	
Packets Transmitted	
Bytes Received	
Bytes Transmitted	
Collision Count	
Over RunError Received	RX
Drops Transmitted	, TX
Drops Received	, RX
Frame Error Received	
CRC Error Received	
Seconds	,
Nanoseconds	

3.2.2

3.4.

3.4 –

Timestamp	,
Active entries	
Packet Lookups	
Packet Matches	,

(3.2),

$$Cf = \sum_{S=1}^n \sum_{ft=1}^m f, \tag{4.2}$$

Cf – ; n –
; m – ; f –

3.4

CollectFlowTableCounters(switch, flow table) (3)

SDN n (10) , (switch)
(table flow)

SDN.

3.4 –

```

1 Procedure CollectFlowTableCounters()
2   while TRUE do
3     CollectFlowTableCounters(switch,flow
      table)
4     send counters to Format Data
5     for each flow in the flow repository do
6       if collect counters not started then
7         start collect flow counters task
8       end
9     end
10    sleep(n units of time)
11  end

```

3.2.3

3.5.

3.5.

3.5 –

Timestamp	,
Received Packets	
Received Bytes	
Seconds	,
Nanoseconds	

3.5 –

```

1 Procedure CollectFlowCounters()
2   while TRUE do
3     CollectFlowCounters(switch,flow
4       table,flow)
5     send counters to Format Data
6     sleep(n units of time)
7   end

```

CollectFlowCounters (switch, flow table, flow) (3)
 SDN n (5)
 . (switch),
 (flow table) (flow)
 SDN.

3.2.4

, , , ,
 , ,
 , .

3.6.

3.6 –

```

1 Procedure FormatData()
2   while TRUE do
3     receive counters data
4     format data
5     forward data to the send data task
6     sleep(n units of time)
7   end
    
```

3.2.5

, (,)

3.4,

(Collect Port Counters, Collect Flow

Table Counters, Collect Flow Counters)

SDN

3

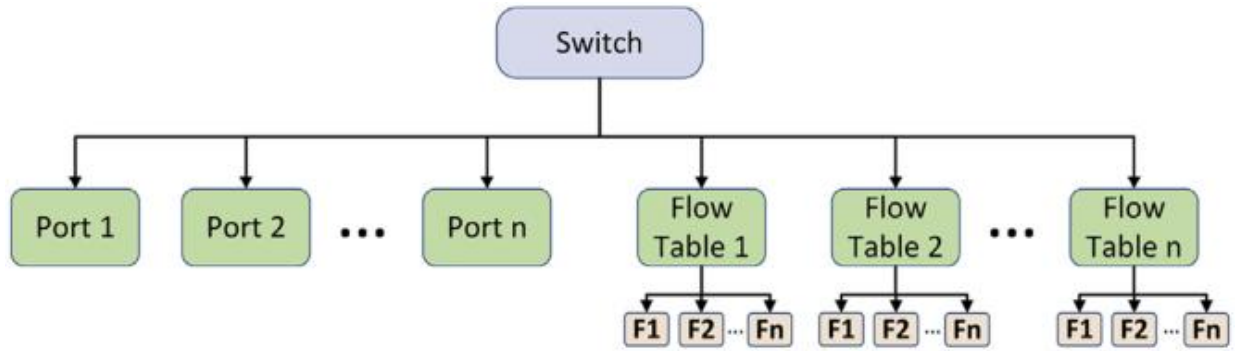
3.3, 3.4

3.5.

3, 4 5

3.1.

SDN



3.4 –

sleep,

n 3.1

SDN

while

3.1

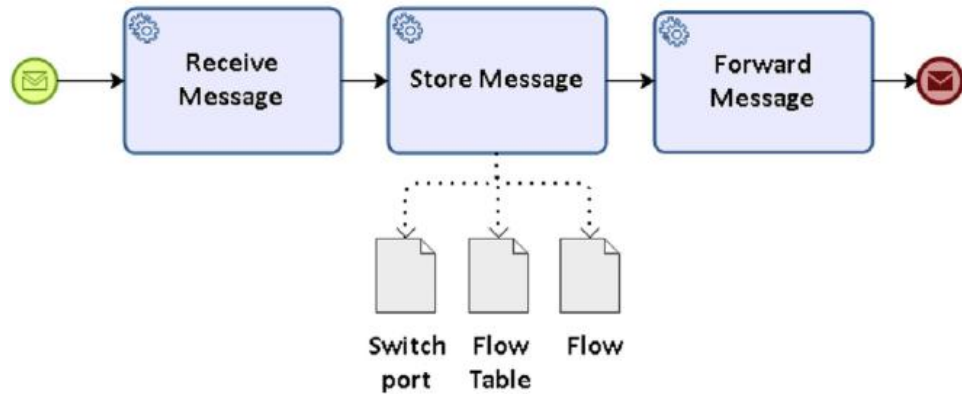
SDN

(Host, Switch Link),

API,

3.3

3.3.1



3.5

3.5 –

Receive message

Store message

(,) , , , .

Forward message

, .

3.3.2

, (,), , .

3.7.

3.7 –

```

1 Procedure GenerateStatistics()
2   while TRUE do
3     request data
4     ComputeStatistics()
5     send statistics
6     sleep(n units of time)
7   end

```

3.6

GenerateStatistics.

1

(Batch 1)

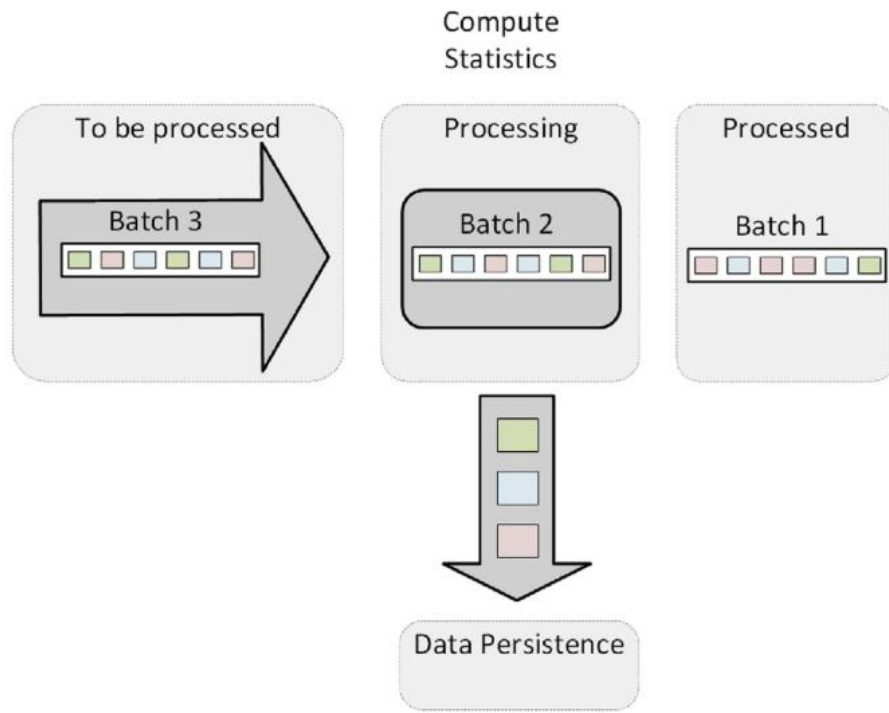
,

2 (Batch 2),

1

2

3 (Batch 3), ComputeStatistics



3.6 –

3.4

Persist inventory data

Persist statistics.

’
.
,
.
.
-
,
,
.

4

OpenDaylight

SDN

4.1.

4.1 –

id			
1	2	Intel(R) Xeon(R) E5-2630 v2 @ 2,60	96
2	2	Intel(R) Xeon(R) E5-2630 v2 @ 2,60	96
3	2	Intel(R) Xeon(R) E5-2620 v3 @ 2,40	96
4	1	Intel(R) Core(TM) i7-4710MQ @ 2,50	16

1, 2 3

HP MSA 2040

SAN 12

2,5

Linux Ubuntu Server 16.04,

Multipath,

4.1

S1-S6

h1-h8 –

eth1,eth2,

Mininet.

4.2,

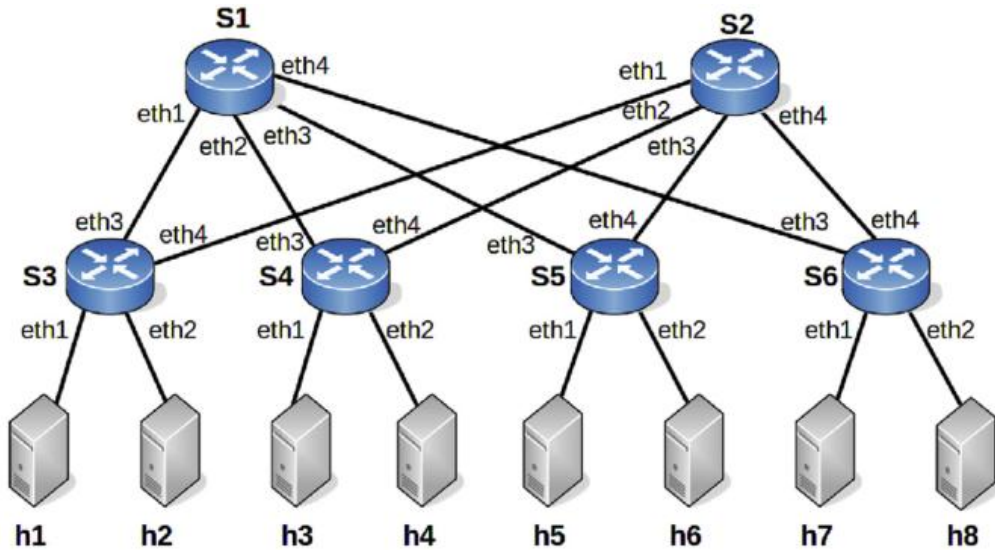
TCP

Iperf3.

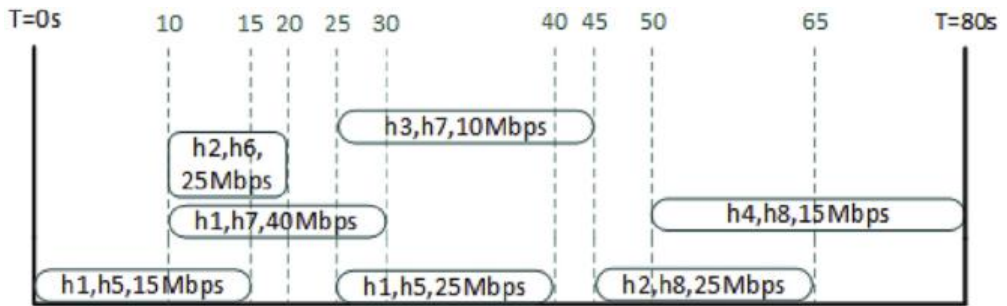
OpenDaylight,

3,

4.2.



4.1 –



4.2 –

SDN

Spark,

Elasticsearch.

4.2

4.2 – 3-

	IP-	IP-	
S1	10.0.0.1	10.0.0.5	: eth3
	10.0.0.2	10.0.0.8	: eth4
	10.0.0.3	10.0.0.7	: eth4
S2	10.0.0.1	10.0.0.7	: eth4
	10.0.0.2	10.0.0.6	: eth3
	10.0.0.4	10.0.0.8	: eth4
S3	10.0.0.1	10.0.0.5	: eth3
	10.0.0.1	10.0.0.7	: eth4
	10.0.0.2	10.0.0.6	: eth4
	10.0.0.2	10.0.0.8	: eth3
S4	10.0.0.3	10.0.0.7	: eth3
	10.0.0.4	10.0.0.8	: eth4
S5	*	10.0.0.5	: eth1
	*	10.0.0.6	: eth2
S6	*	10.0.0.7	: eth1
	*	10.0.0.8	: eth2

tcpdump

tshark

python,

;

,

tshark. tshark

OF,

tshark,

4.1

4.3

S6

(eth1-eth4)

SDN

:

$$S = \sum_{n=1}^N T_n, \tag{4.1}$$

S –

; T_n –

n; N –

4.2, 20-

h1 h7 10-

eth1,

10 30.

30,

25

h3 h7.

eth2

20-

h2 h8,

30-

h4 h8.

eth3

20-

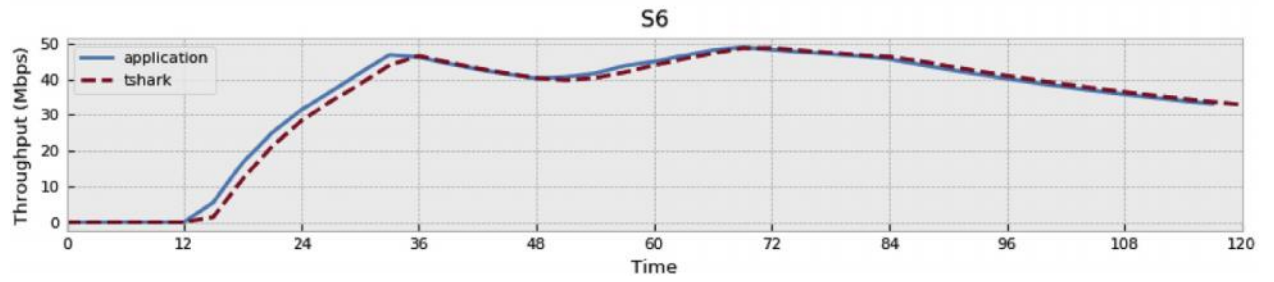
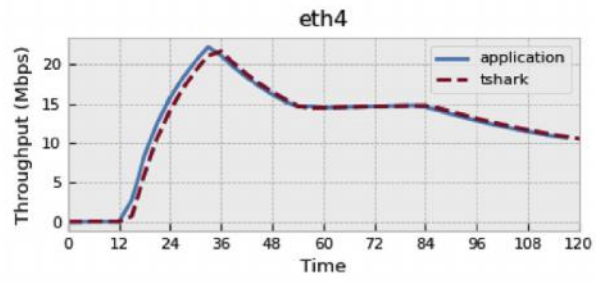
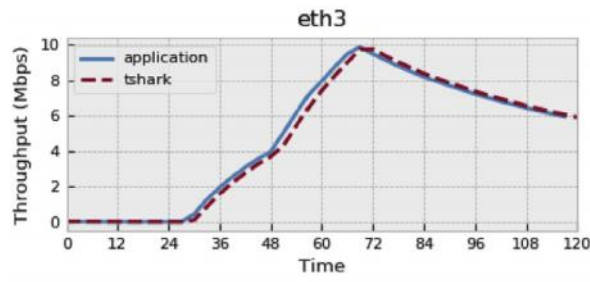
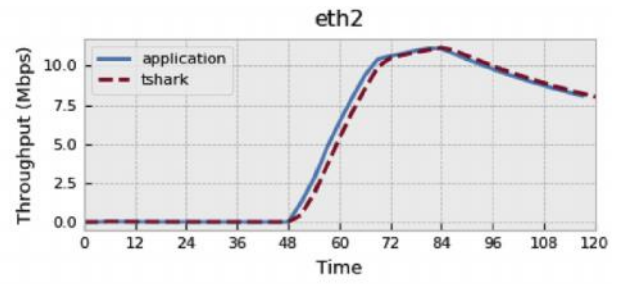
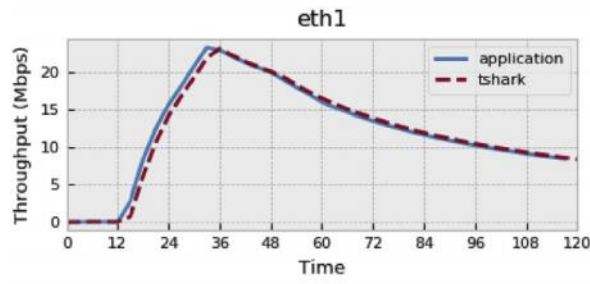
h3 h7 20-

h2

h8. eth4
h4 h8.

20-

h1 h7 30-



4.3 –

4.2

()

4.4.

:

$$P = \frac{T_p}{T_s} \times 100, \tag{4.2}$$

T_p - ; T_s -

4.3 , udp 10 S6.

4.4, 0 10- ,

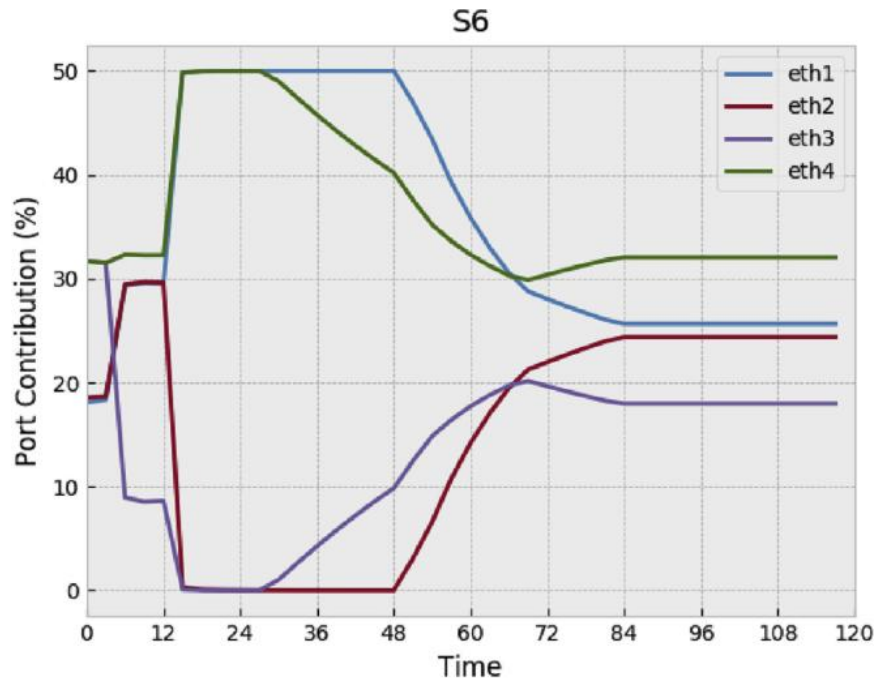
MDNS, LLDP, ARP ICMP,

ping

4.4 , 10-25 eth1 eth4

4.3.

eth2 50



4.4 -

4.3

(4.1).

(S1-S6). 4.5

tshark S1, S3 S4. 15 S1 S3

. S1

h1 h5 15 . 15 S1

10 , 25

h1 h5 45 ()

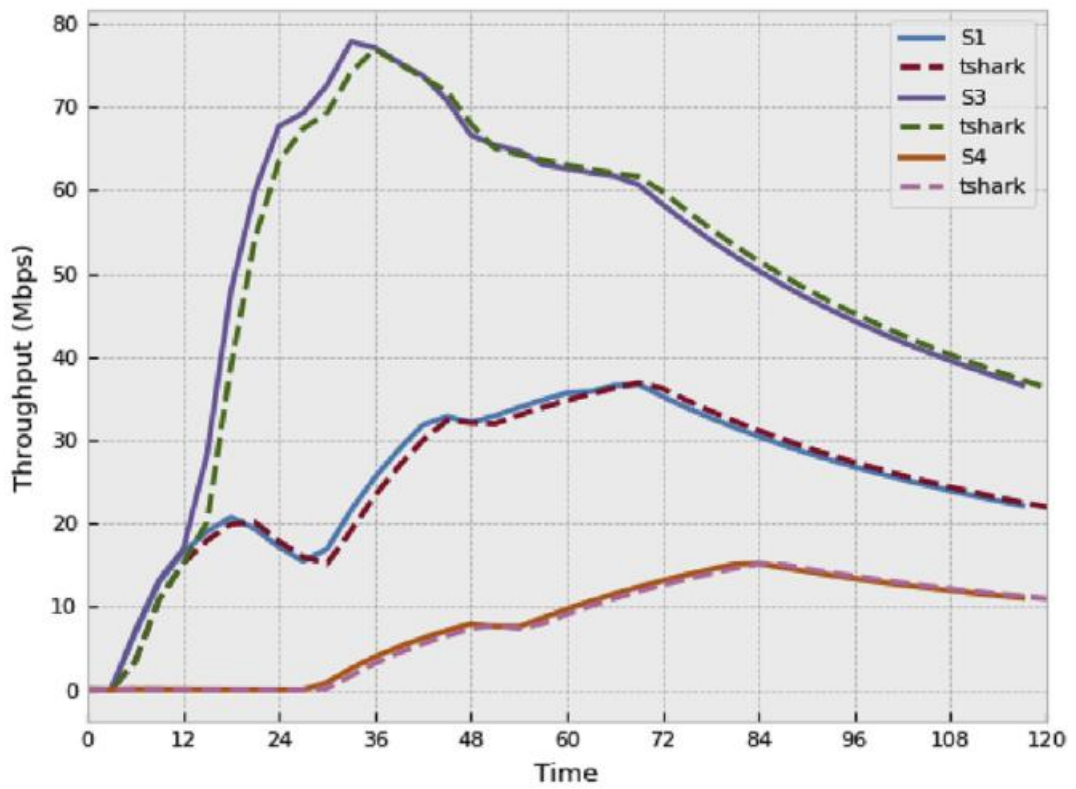
h2 h8 20 . S3 ,

h1 h5 h1 h7 10

h1 h5 25 . S4

25- h3 h7 50-

h4 h8.



4.5 –

4.4

4.6.

$$P = \frac{T_p}{P_c} \times 100, \tag{4.3}$$

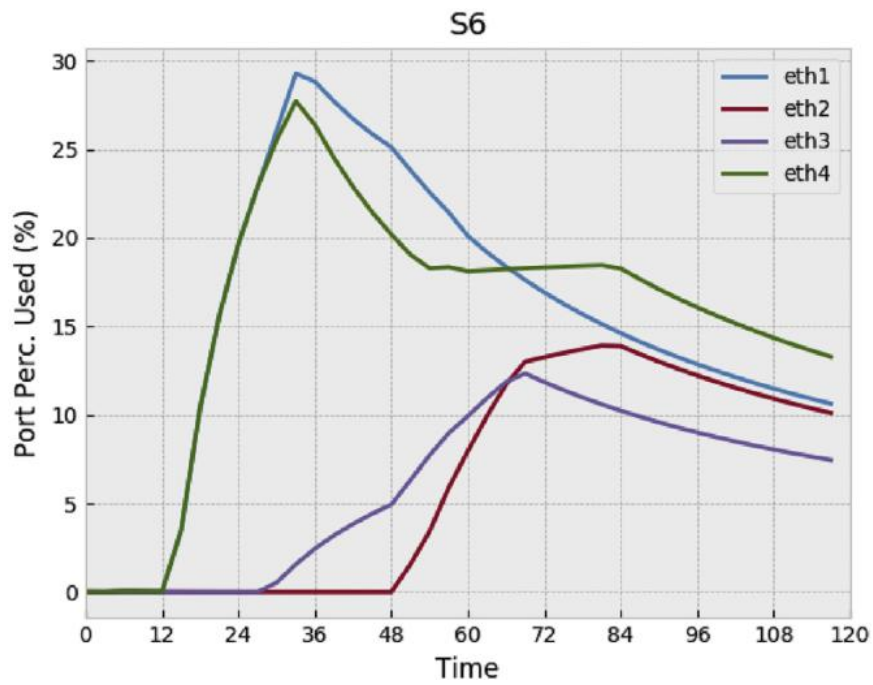
T_p – ; P_c –

P_c Switch Port 3.2

Current Feature. , OpenDaylight ,

1 / .

100 / .



4.6 –

4.6 P, (eth1-eth4)

S6. eth1 25 30

h1 h7 (10) h3 h7 (25).

eth2 , , ,

45 (h2-h8 h4-h8).

eth3 h3 h7 h2 h8. , eth4

h7 h4 h8. , eth3, h1

4.5

4.7.

SDN

3.3,

OF

Bytes Received Bytes

Transmitted

(Second Nanosecond)

(4.3),

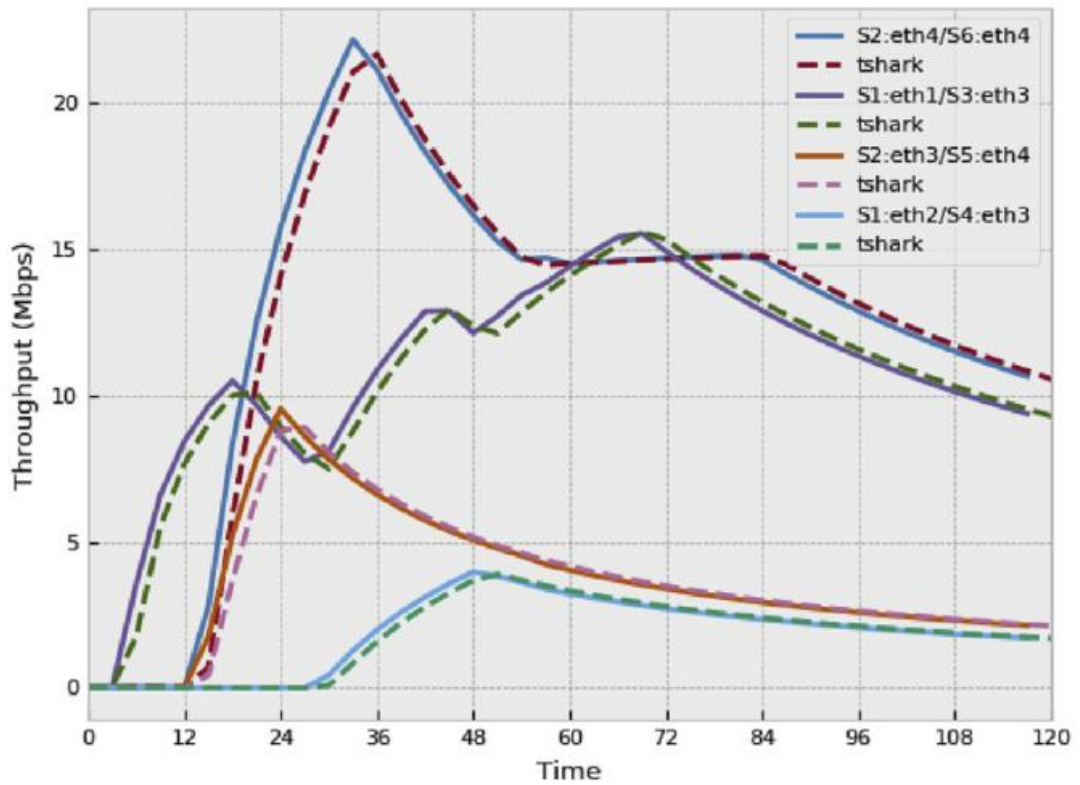
Bytes Received

, Bytes Transmitted

Second Nanosecond

S_{a_i} S_{b_j} .

1. S_{a_i} i a.
2. S_{b_j} j b.
3. BR_{a_i} BT_{a_i} Bytes
Received Bytes Transmitted, S_{a_i} .
4. BR_{b_j} BT_{b_j} Bytes
Received Bytes Transmitted, S_{b_j} .
5. $(BR_{a_i} + BT_{a_i}) > (BR_{b_j} + BT_{b_j}),$
 $S_{a_i}; - S_{b_j}.$



4.7 –

4.7
S2:eth4 S6:eth4, S1:eth1 S3:eth3, S2:eth3 S5:eth4 S1:eth2

S4:eth3.

S2:eth4/S6:eth4

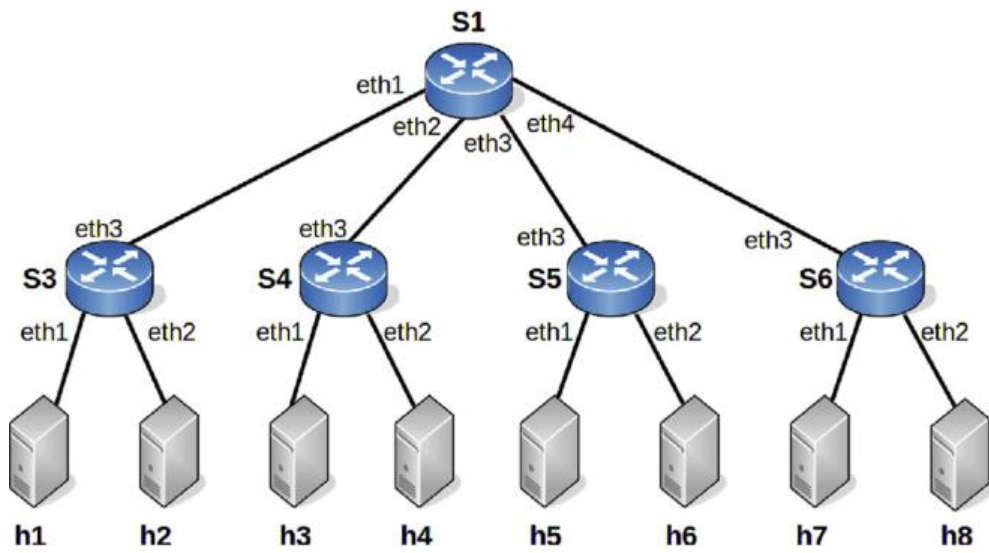
h1 h7, ,
 h4 h8. ,
 S1:eth1/S3:eth3, h1 h5
 h2 h8. h2 h6 ,
 S2:eth3/S5:eth4. S1:eth2/S4:eth3
 h3 h7.

4.6

(TM),
 (OD) [23].
 (4.1), , - ,
 . ,
 : S1 ,
 S2 (4.8, 4.9).
 OD, S1 S2, 4.2. 4.10

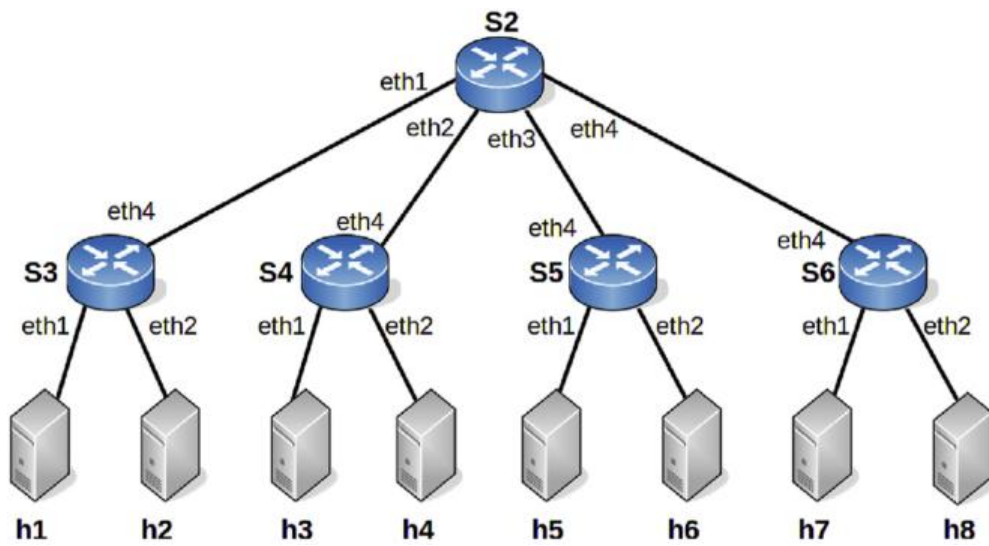
4.2 –

S1	S2
h1 – h2	h1 – h2
h1 – h3	h1 – h3
h1 – h4	h1 – h4
h7 – h8	h7 – h8



4.8 –

S1



4.9 –

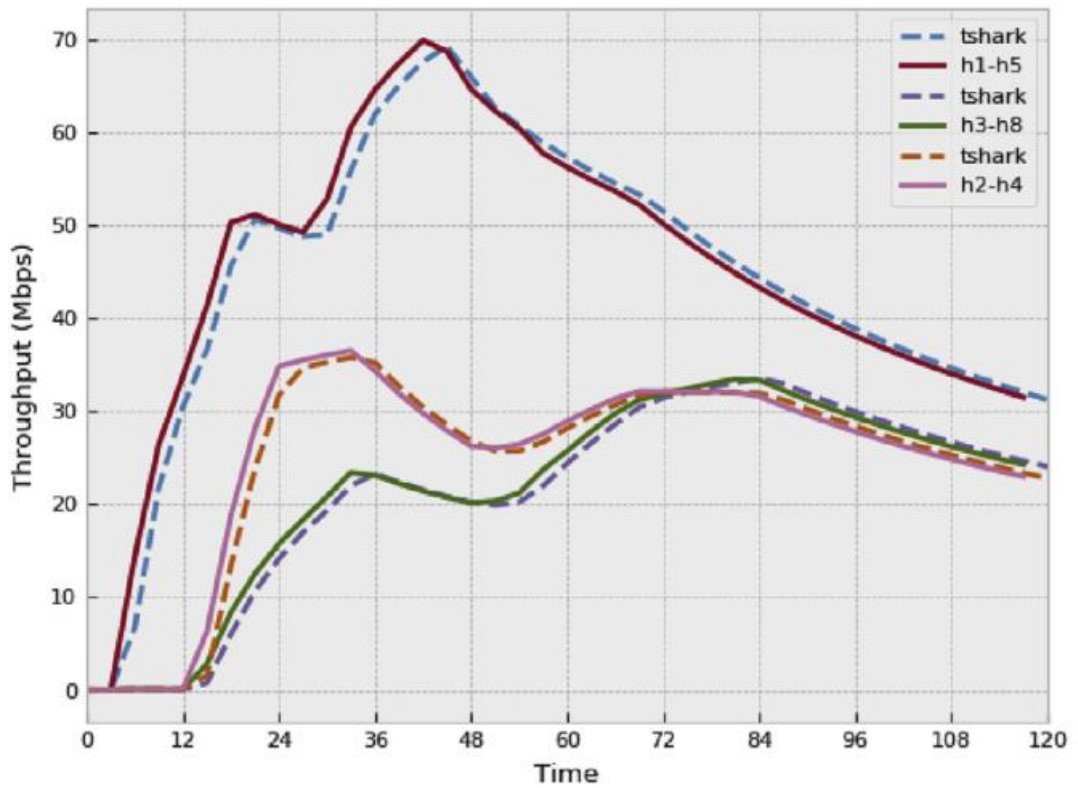
S2

OD,

:

$$T_{OD} = \sum_{n=1}^N T_{Ln}, \tag{4.4}$$

T_{OD} – OD; T_{Ln} –
 n.
 4.10
 : h1-h5 (root S1), h3-h8 (root S2) h2-h4 (root S2). 4.1
 OD.
 , h3 h8, S2
 : S4:eth1, S4:eth4/S2:eth2, s2:eth4/S6:eth4 S6:eth2. OD h1-h5
 0 25 .
 h3-h8,
 : h1-
 h7, h3-h7 h4-h8. OD h2-h4.
 h1-h7, h2-h6 h4-h8
 OD.



SDN.

Big Data.

TE,

OD,

NoSQL,

TE,

3 1 ,

SDN.

» [24].

1. Y. Lee and Y. Lee. Toward Scalable Internet Traffic Measurement and Analysis with Hadoop. *SIGCOMM Comput. Commun. Rev.*, 43(1):5-13, 2013.
2. C. Orsini, A. King, D. Giordano, V. Giotsas, and A. Dainotti. BG-PStream: A Software Framework for Live and Historical BGP Data Analysis. *Proc. of the IMC*, pages 429-444, 2016.
3. S. Wassermann, P. Casas, T. Cuvelier, and B. Donnet. NETPerfTrace: Predicting Internet Path Dynamics and Performance with Machine Learning. *Proc. of the Big-DAMA*, pages 31-36, 2017.
4. A. Fahad, N. Alshatri, Z. Tari, A. Alamri, I. Khalil, A. Y. Zomaya, S. Fofou, and A. Bouras. A Survey of Clustering Algorithms for Big Data: Taxonomy and Empirical Analysis. *IEEE Trans. Emerg. Topics Comput.*, 2(3):267-279, 2014.
5. C.-W. Tsai, C.-F. Lai, H.-C. Chao, and A. V. Vasilakos. Big Data Analytics: A Survey. *J. Big Data*, 2(1):21, 2015.
6. Y. Zhang, T. Cao, S. Li, X. Tian, L. Yuan, H. Jia, and A. V. Vasilakos. Parallel Processing Systems for Big Data: A Survey. *Proc. IEEE*, 104(11):2114-2136, 2016.
7. A. Callado, C. Kamienski, G. Szabo, B. P. Gero, J. Kelner, S. Fernandes, and D. Sadok. A Survey on Internet Traffic Identification. *Commun. Surveys Tuts.*, 11(3):37-52, 2009.
8. T. T. Nguyen and G. Armitage. A Survey of Techniques for Internet Traffic Classification Using Machine Learning. *Commun. Surveys Tuts.*, 10(4):56-76, 2008.
9. M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita. Network Anomaly Detection: Methods, Systems and Tools. *Commun. Surveys Tuts.*, 16(1):303-336, 2014.
10. R. Hofstede, P. Celeda, B. Trammell, I. Drago, R. Sadre, A. Sperotto,

and A. Pras. Flow Monitoring Explained: From Packet Capture to Data Analysis with NetFlow and IPFIX. *Commun. Surveys Tuts.*, 16(4):2037-2064, 2014.

11. S. Valenti, D. Rossi, A. Dainotti, A. Pescapé, A. Finamore, and M. Mellia. Reviewing Traffic Classification. In *Data Traffic Monitoring and Analysis – From Measurement, Classification, and Anomaly Detection to Quality of Experience*. Springer, Heidelberg, 1 edition, 2013.

12. H. Hu, Y. Wen, T.-S. Chua, and X. Li. Toward Scalable Systems for Big Data Analytics: A Technology Tutorial. *IEEE Access*, 2:652-687, 2014.

13. D. Laney. 3d Data Management: Controlling Data Volume, Velocity, and Variety. Technical report, META Group, 2001.

14. J. Manyika and others. *Big Data: The Next Frontier for Innovation, Competition, and Productivity*, 2011.

15. F. Bajaber, R. Elshawi, O. Batarfi, A. Altalhi, A. Barnawi, and S. Sakr. Big Data 2.0 Processing Systems: Taxonomy and Open Challenges. *J. Grid Comput.*, 14(3):379-405, 2016.

16. R. Boutaba, M. A. Salahuddin, N. Limam, S. Ayoubi, N. Shahriar, F. Estrada-Solano, and O. M. Caicedo. A Comprehensive Survey on Machine Learning for Networking: Evolution, Applications and Research Opportunities. *J. Internet Serv. Appl.*, 9(16), 2018.

17. T. Karagiannis, K. Papagiannaki, and M. Faloutsos. BLINC: Multilevel Traffic Classification in the Dark. *Proc. of the SIGCOMM*, pages 229-240, 2005.

18. OpenFlow Switch Specification 1.5.0, 2014 [Online] Available: <https://www.opennetworking.org/images/stories/downloads/sdn-resources/onf-specifications/openflow/openflow-switch-v1.5.0.noipr.pdf>.

19. F. Volpato, M.B. Castro, M.A.R. Dantas, 2019. OFQuality: quality of service management module for software-defined networking. *Int. J. Grid Comput. Util. – IJGUC* 10 (2), 187-198.

20. I.F. Akyildiz, A. Lee, P. Wang, M. Luo, W. Chou. Research challenges for traffic engineering in software defined networks. *IEEE Network* 30 (3), 52-58, Jun.2016.

21. I.F. Akyildiz, A. Lee, P. Wang, M. Luo, W. Chou. A roadmap for traffic engineering in SDN-OpenFlow networks. *Comput. Network.* 71, 1-30, 2014.

22. Z. Shu, J. Wan, J. Lin, S. Wang, D. Li, S. Rho, C. Yang. Traffic engineering in software-defined networking: measurement and management. *IEEE Access* 4, 3246-3256, 2016.

23. N.L.M. van Adrichem, C. Doerr, F.A. Kuipers. OpenNetMon: network monitoring in OpenFlow software-defined networks. In: 2014 IEEE Network Operations and Management Symposium (NOMS), pp. 1-8, 2014.

24. . . , . . .
 Big Data //
 . , -
 , 18-19 2021 ., : , 2021. – . 28.