# Large-scale NetFlow Information Management

Adrien Raulot, Shahrukh Zaidi

University of Amsterdam

*Supervisor: Wim Biemolt (SURFnet)*

February 5, 2018

# What is NetFlow?

- Traffic monitoring technology originaly developed by Cisco.

- *Flow*: "a set of IP packets passing an observation point in the network during a certain time interval. All packets belonging to a particular flow have a set of common properties."[4]

- Important differences with regular packet capture methods:
  - NetFlow considered to be less privacy sensitive
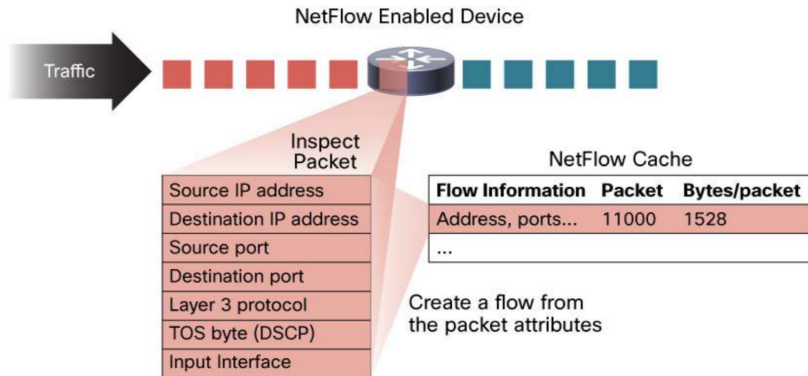  - NetFlow requires less computational resources for analysis

Figure 1: Schematic overview of the NetFlow export process.[2]

# NetFlow Analysis

Three main application areas[3]:

- Flow analysis and reporting

- Threat detection

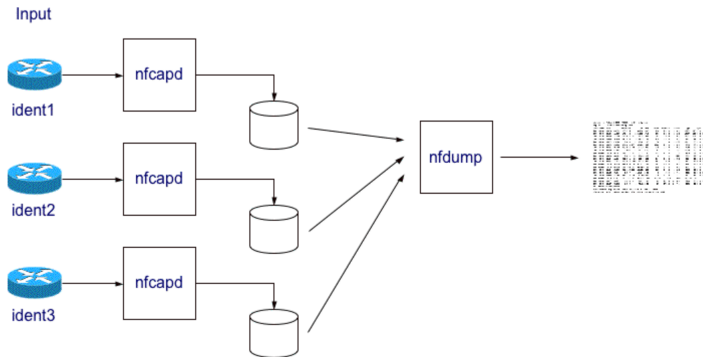- Performance monitoring

# NetFlow analysis techniques

- NfDump:



Figure 2: Schematic overview of the NfDump tool set.[1]

# Netflow Analysis techniques

- Limitations of this setup[5]:

    - **Inefficient file-based store:** NfDump typically stores NetFlow data in separate files for every 5 minutes time frame

    - **Very slow processing speed:** each file is read line by line from the beginning. Therefore, analysis of large amounts of NetFlow data takes a lot of time.

    - **Limited analysis methods:** as network situations are becoming more and more complex, new analysis approaches are required that allow for NetFlow data analysis.

*Which data analysis technique could be used in order to analyse the current SURFnet NetFlow data in a more time-efficient manner?*

# What is Apache Hadoop?

- Framework for large datasets processing
- Distributed, local computation & storage
- Hadoop Distributed File System (HDFS)
- YARN (Yet Another Resource Negotiator)
- Batch, interactive & real-time jobs
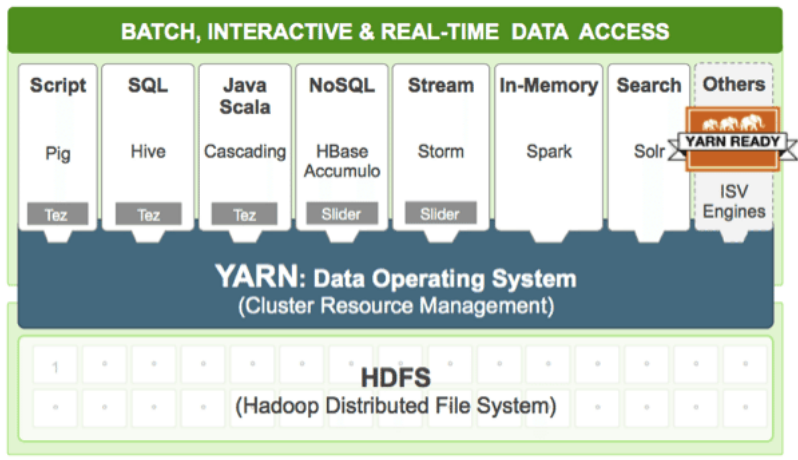- Designed to be scalable

# What is Apache Hadoop?



Figure 3: Schematic overview of Hadoop 2.0.

# What is Apache Spark?

- Hadoop-related project, but not only
- Powerful computing engine for Big Data processing
- In-memory
- Built-in modules for streaming, SQL, machine learning, etc.
- Binding for Java, Scala, Python and R
- Ease of use

# What is Apache Parquet?

- Data-store for Hadoop
- Column-oriented
- Fast access to data



Figure 4: Schematic overview of a row vs column-oriented database.

# Choice for analysis technique (summary)



Figure 5: Apache Parquet logo.



Figure 6: Apache Hadoop logo.



Figure 7: Apache Spark logo.

**To-Do list:**

1. Store NetFlow data into Parquet files on HDFS
2. Load Parquet files using PySpark (Python API)
3. Query the data using Spark SQL

# Experiments: test environment

**Hadoop cluster specifications:**

- $\sim$ 100 nodes
- $\sim$ 600 cores
- $\sim$ 4TB of memory
- $\sim$ 2PB of storage
- Apache Hadoop 2.7.2
- Apache Spark 2.1.1

**NfDump server specifications:**

- 1x Dell PowerEdge R230
- Intel Xeon CPU E3-1240L v5 @ 2.10GHz
- 4 cores
- 16GB of RAM
- $\sim$ 200GB of SSD storage
- NfDump v1.6.12

# Experiments: implementation

1. Convert NetFlow binary data to CSV

   ```
   nfdump -r nfcapd.201801011245 -o csv
   ```

2. Write two Spark jobs in Python:
   - Converter: Converts CSV data to Parquet format
   - Querier: Loads Parquet data & executes queries

3. Write SQL query

   ```
   query = 'SELECT ts, sa, da FROM nf_data'
   ```

4. Using the Querier, execute and cache the results

5. Proceed with next operations on the cached results

   ```
   print results.count()
   print results.show()
   ```

# Experiments: test queries

- Retrieve all flows containing a specific IP address

- Retrieve all flows with a byte count larger than 100MBs

- List the top 10 of Telnet connections with only the SYN flag set in the IP header ordered by the number of bits per second

- List the top 10 of IP addresses receiving the largest amount of traffic

- Retrieve all flows with only the SYN flag set in the IP header

Figure 8: Execution time of retrieving all flows containing a specific IP address.
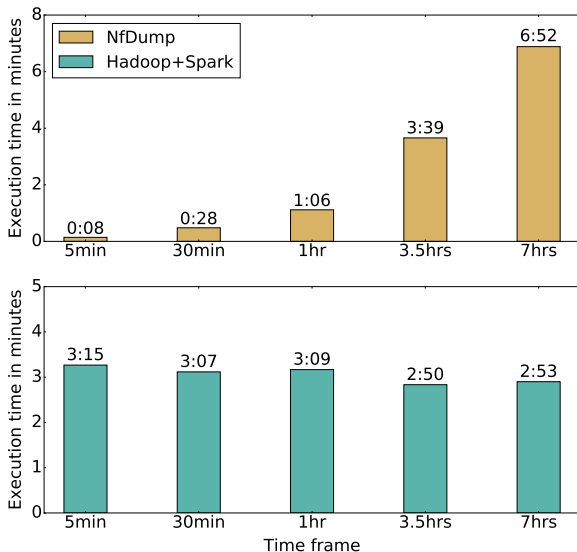
# Results: retrieve all flows with byte count >100MB



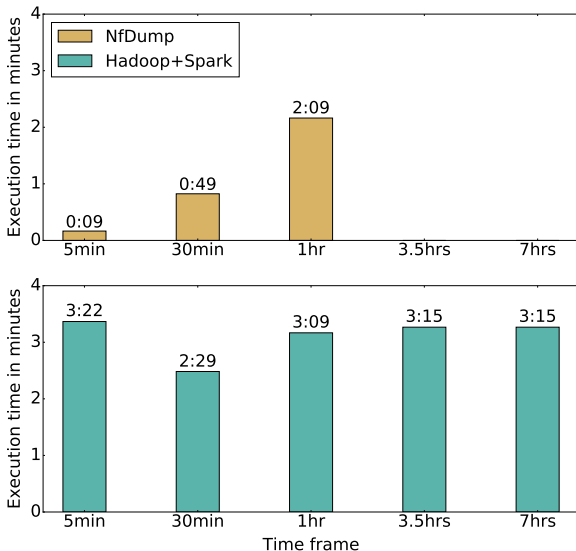Figure 9: Execution time of retrieving all flows with byte count larger than 100MB.

Figure 10: Execution time of retrieving the top 10 of Telnet connections with only the SYN flag set ordered by the number of bits per second.
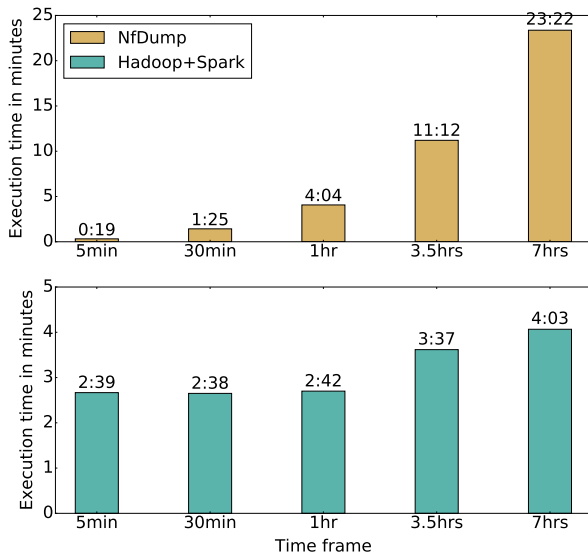
Figure 11: Execution time of Retrieving the top 10 IP addresses receiving the largest amount of traffic.

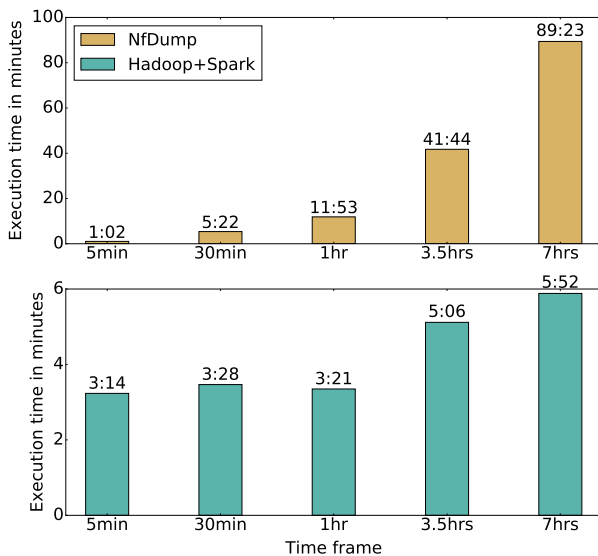# Results: Retrieve all flows with only SYN flag set



Figure 12: Execution time of retrieving all flows with only the SYN flag set in the IP header.

- Execution time of NfDump increases linearly with longer time frames.

- Hadoop scales very well:
  - Execution time of Spark with Hadoop does not increase significantly when dealing with larger amounts of data.

- NfDump struggles with executing more complex queries, whereas this is no problem for Spark and Hadoop.

# Conclusion and future work

- Combination of Hadoop and Apache Spark is a viable option for analyzing large-scale NetFlow data.

- Tuning and optimization to the Spark implementation and Hadoop cluster may lead to even better performance.

# Questions?

# References

📄 NfDump.
http://nfdump.sourceforge.net/.

📄 I. Cisco.
NetFlow. Introduction to Cisco IOS NetFlow C a technical overview, 2007.

📄 R. Hofstede, P. Čeleda, B. Trammell, I. Drago, R. Sadre, A. Sperotto, and A. Pras.
Flow monitoring explained: From packet capture to data analysis with netflow and ipfix.
*IEEE Communications Surveys & Tutorials*, 16(4):2037–2064, 2014.

📄 G. Sadasivan.
Architecture for ip flow information export.
*Architecture*, 2009.

📄 Z. Tian.
Management of large scale NetFlow data by distributed systems.