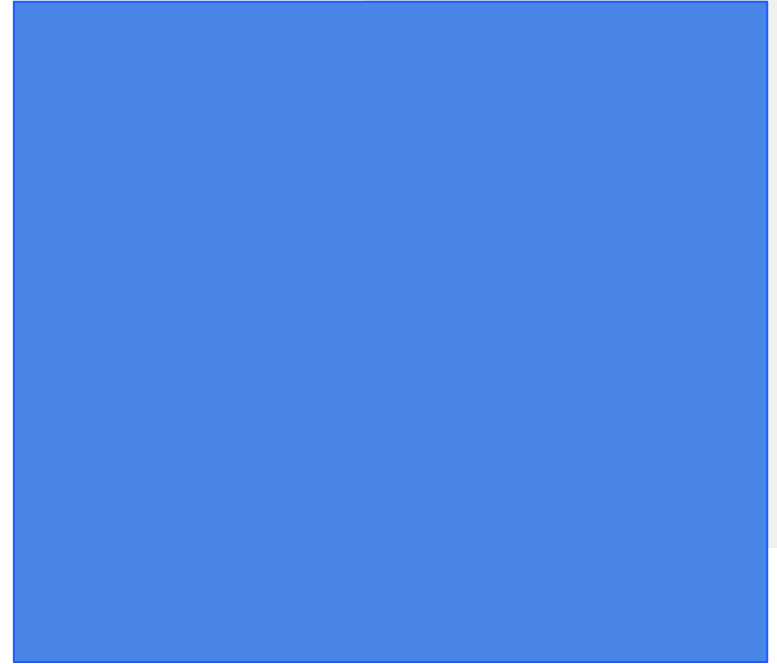# Big Data Analytics

Dr Sirintra Vaiwsri | Email: sirintra.v@itm.kmutnb.ac.th

# Streaming Processing

# Streaming Processing (Chambers and Zaharia, 2018)

- Streaming processing continuously receives input and computes a result.

- It responds quickly and is more efficiency than batch processing.

- However, it is challenging if we need to process input that is not in sequence. For example, to process 2, 10, and then 5 where receiving 2, 5, and 10.

Dr Sirintra Vaiwsri

# Streaming Processing (Antolínez García, 2023)

- Main characteristics of data streaming:
    - Continuous information
    - Unbounded information
    - High volume and velocity
    - Time sensitive
    - Heterogeneous data sources

# Streaming Processing (Antolínez García, 2023)

- Streaming processing uses timestamps to order events in sequence.

  - Event time is based on the event that is generated by a device.

  - Ingestion time is based on the time of the event's arrival.

  - Processing time is based on the beginning time of the event process.

Dr Sirintra Vaiwsri

# Spark Structured Streaming

# **Spark Structured Streaming** (Antolínez García, 2023)

- Spark Structured Streaming provides a high-level manipulation.

- It also provides scalable and near-real-time streaming processing.

- It is built on top of the Spark SQL library and based on the Dataframe and Dataset APIs.

Dr Sirintra Vaiwsri

# Spark Structured Streaming (Antolínez García, 2023)

- Concepts:

  - Input table: also called "unbounded input table", means arrived data is appended as a new row in the table.

  - Result table: also called "unbounded output table", means once new data input arrives, it is processed and the result table is updated.

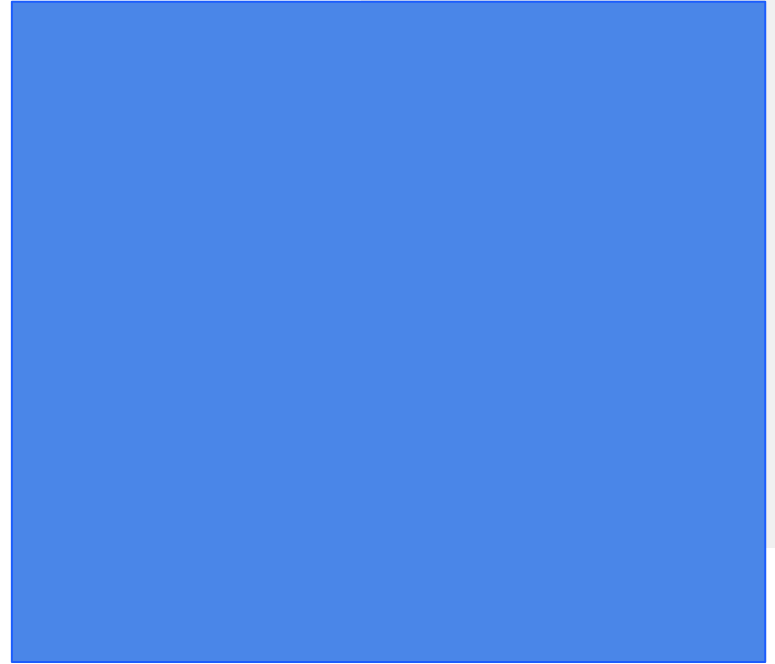# Spark Structured Streaming (Antolínez García, 2023)

○ Output modes:

- Append mode - only new rows in the result table are written to the output sink.

- Complete mode - write all rows in the result table every time that the data were processed.

- Update mode - only write the updated rows to the output sink.

Dr Sirintra Vaiwsri

# Programming Model for Structured Streaming
## (Spark, 2023)

- New row(s) is appended to the input table (unbounded table) in every trigger interval, e.g. every 1 second.

- Data is then queried resulting in an updated result table.

- The result table can be then written out to the external storage as an output.

Dr Sirintra Vaiwsri

# Dataframes Streaming API

# **Streaming Dataframes** (Antolínez García, 2023)

- Streaming Dataframes created using
  `SparkSession.readStream()`

- Input Sources:

  - Socket source

  - File source - a file such as CSV, JSON, or Parquet is read as a stream of data.

  - Kafka source - data is read from the Kafka source.

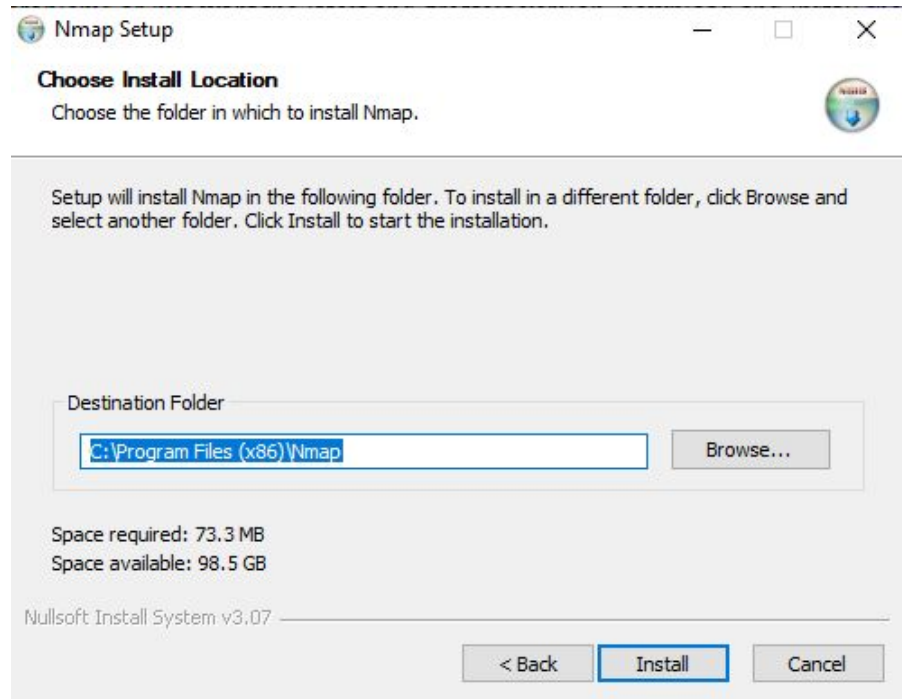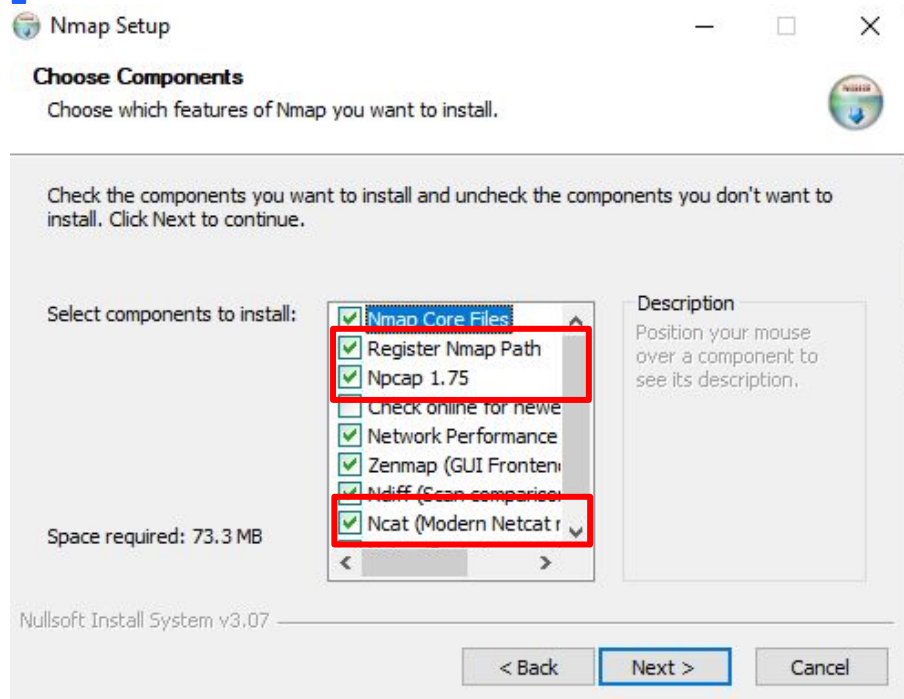Dr Sirintra Vaiwsri

# Socket Source

- Data can be ingested by listening to a socket connection.

- It is often used for testing.

- To execute, use NetCat which is in a Nmap package.

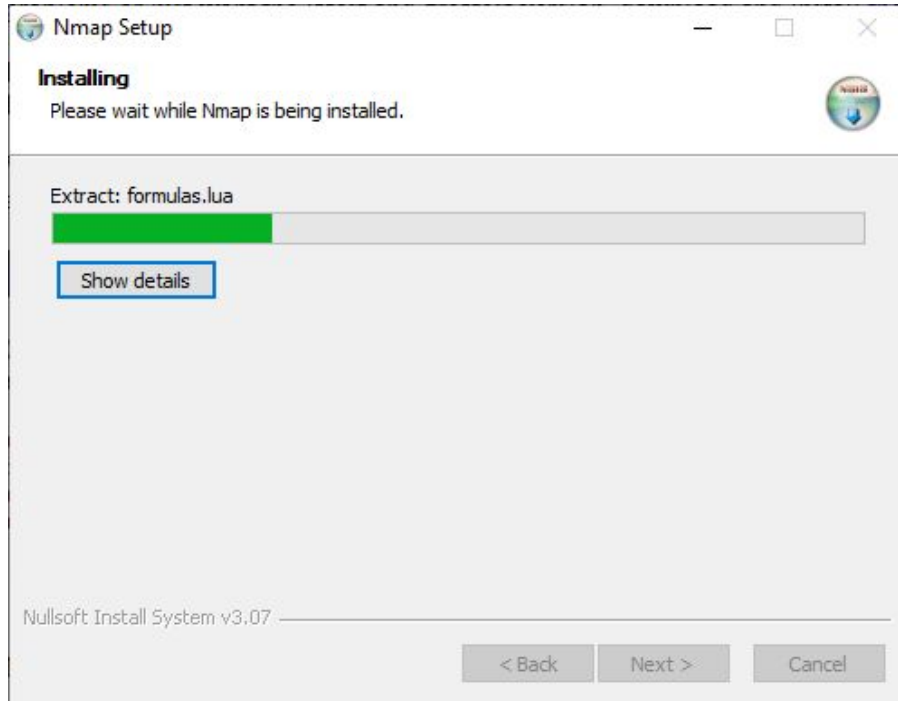- Download Nmap (.exe) for Windows here.

Dr Sirintra Vaiwsri

# Socket Source

- Install Nmap and Npcap

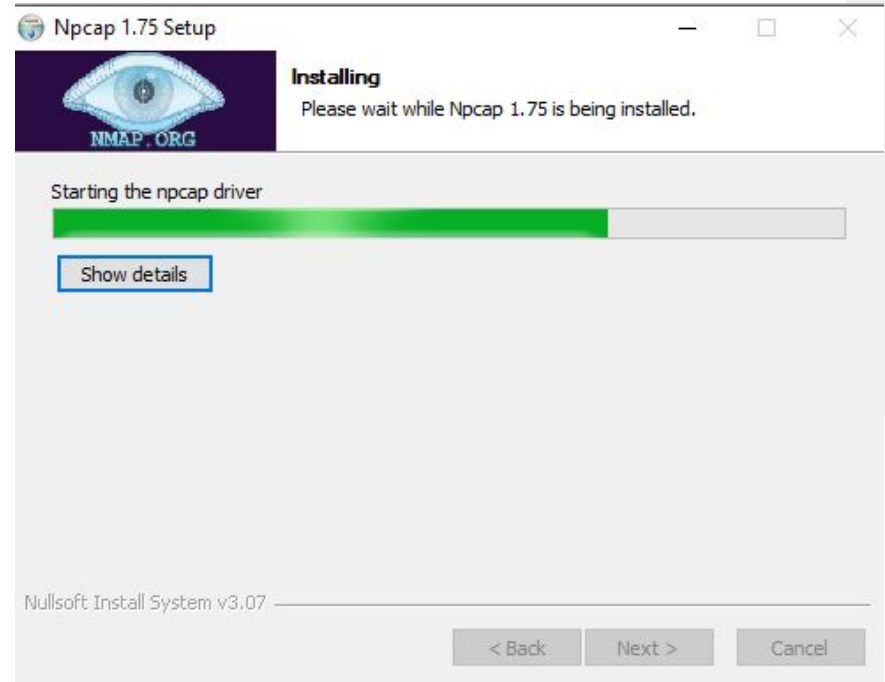- We also install Npcap because the Nmap uses the Npcap library.

Dr Sirintra Vaiwsri

# Socket Source

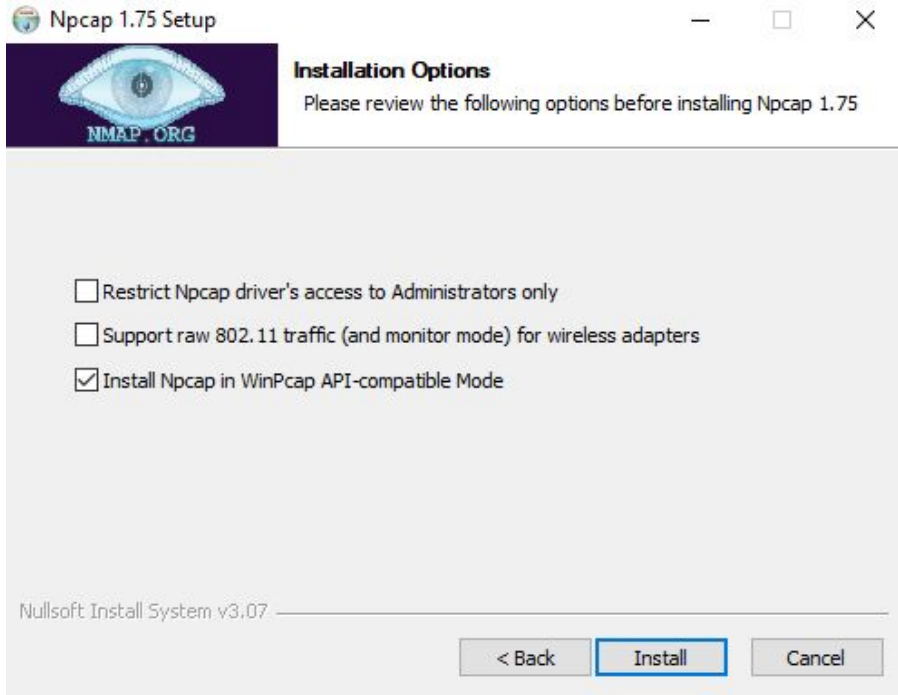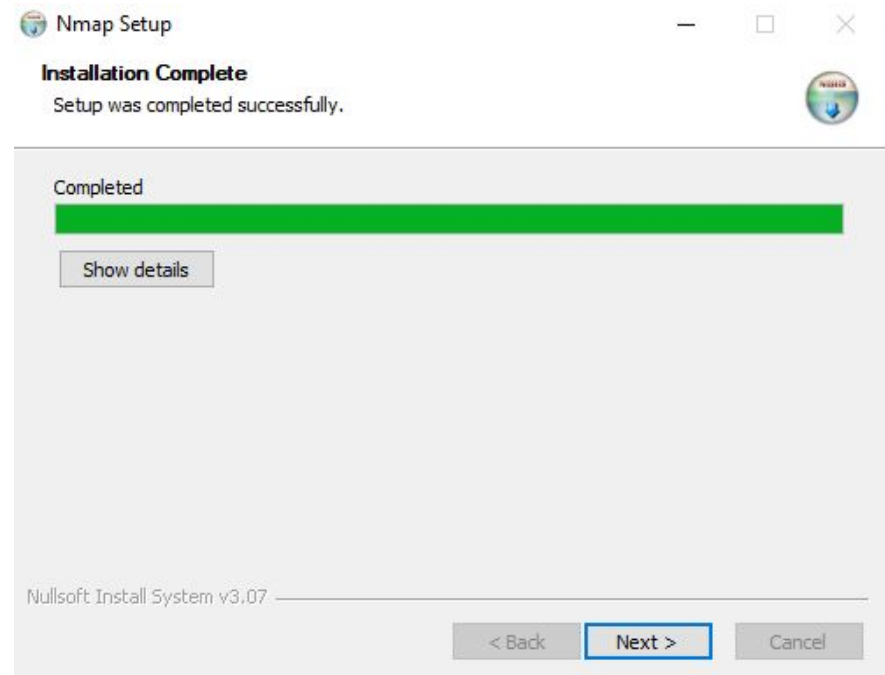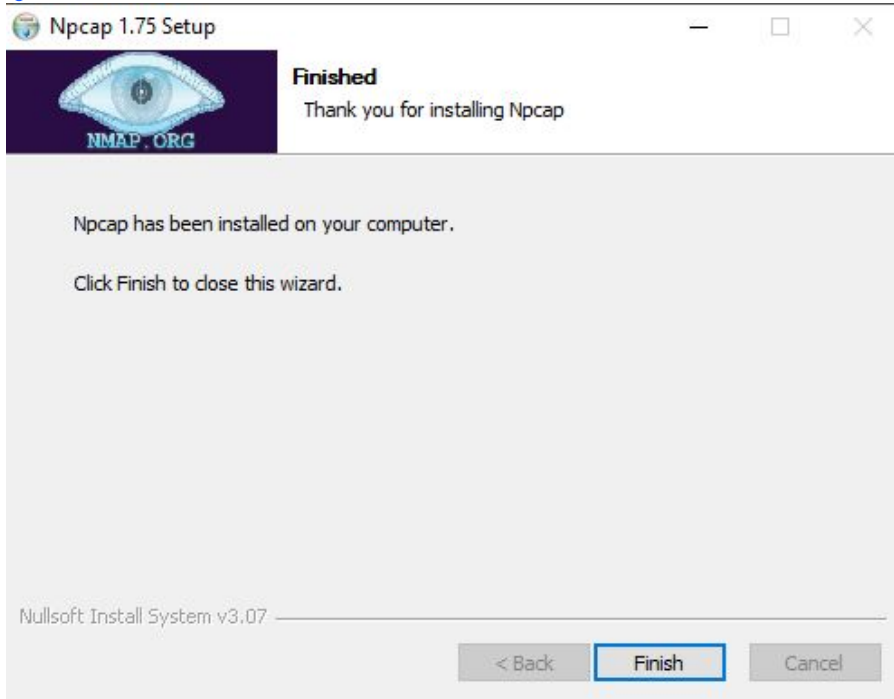# Socket Source
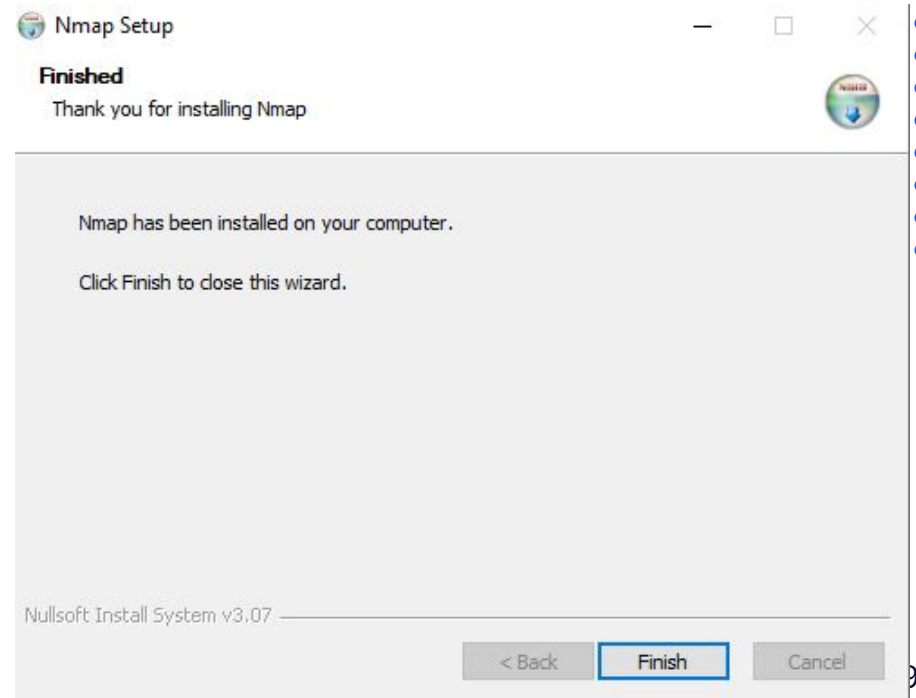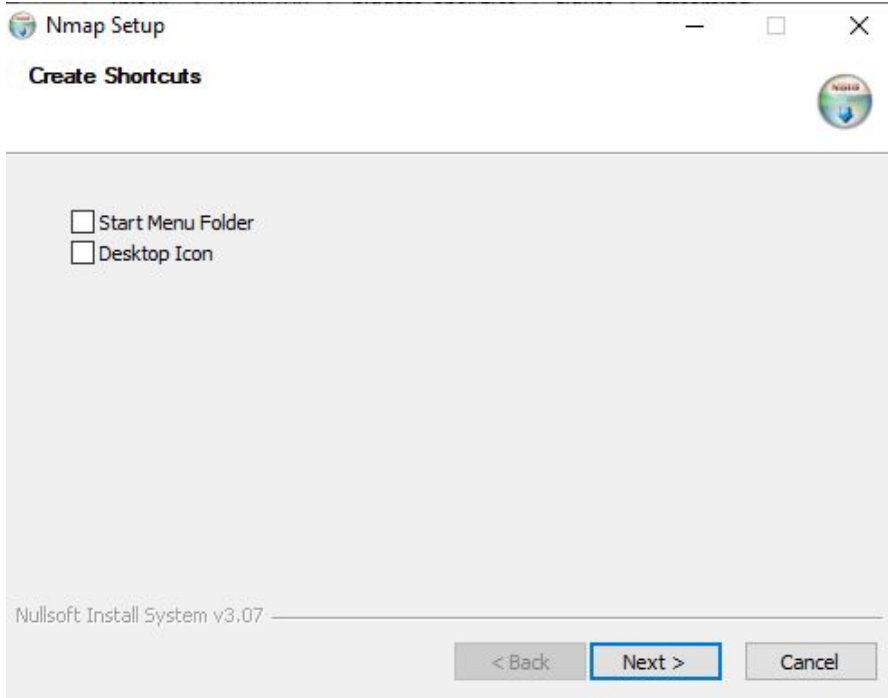
Dr Sirintra Vaiwsri

# Socket Source

Dr Sirintra Vaiwsri

# Socket Source

# Socket Source

# Socket Source

Socket Source Implementation:

```
# Import Spark Session
# Import explode and split from the pyspark.sql.functions
# Create Spark Session
lines = <spark session>.readStream.format("socket").\
        option("host", "localhost").\
        option("port", 9999).\
        load() #The function readStream() is used to return a
               # DataStreamReader for reading data streams as a
               # stream DataFrame.
```

20

Dr Sirintra Vaiwsri

# Socket Source

```
words = lines.select(
                    explode(
                        split(lines.value, " ")
                    ).alias("word")
                )
```

- The function explode() is used to return a new row for each value in a given array or map.

- The function split() is used to split value, where in this case we split value using a space (" ").

- The function alias() is used to name a new column.

# Socket Source

```
wordCounts = words.groupBy("word").count() # Group
                                    # and count each word
query = wordCounts.writeStream.\
        outputMode("<output mode>").\
        format("console").start()
```

- The function writeStream() is used to save the streaming Dataframe out to the external storage, where in this case we use "console".

```
query.awaitTermination()
```

- The function awaitTermination() is used to wait for the termination of the current query.

Dr Sirintra Vaiwsri

# Socket Source

- Open the first command prompt
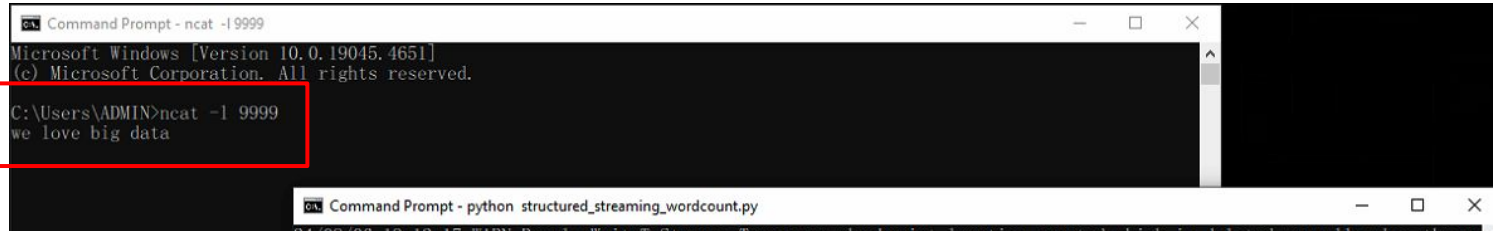
- Type the command ncat -l 9999 (-l is listen)

# Socket Source

- Open the second command prompt

- Run the created Structured Streaming Python file

- First, you will see only an empty Dataframe.

```
Batch: 0
-------------------------------------
+-----+-----+
|word |count|
+-----+-----+
+-----+-----+
```

Dr Sirintra Vaiwsri

# Socket Source



- Try typing something in the first command prompt

- You will see the words and their count in the second command prompt

Dr Sirintra Vaiwsri

# Socket Source



- Try typing something that contains the same words as your previous words

- You will see those word counts are incremented.

# **File Source** (Spark, 2023)

- A file is read as a stream of data.

- Once the file is modified, the file will be processed in the structured streaming

- Supported file formats are such as Text, CSV, JSON, etc.

- File source used in Structured Streaming requires a specified schema.

- This is to ensure a consistent schema being used for the streaming query.

# File Source

File Source Implementation:

```python
# Import Spark Session
# Import split from the pyspark.sql.functions
# Import StructType, StringType, and StructField from
    # pyspark.sql.types
# Create Spark Session
<your schema> = StructType([StructField(\
        "<column name>", StringType(), True), \
        .....................
        ]) # Create schema
```

Dr Sirintra Vaiwsri

# File Source

```
lines = <spark session>.readStream.format("csv").\
        option("maxFilesPerTrigger", 1).\
        option("header", True).\
        option("path", "<path to your folder>").\
        schema(<your schema>).load()
```

- **maxFilesPerTrigger** is the maximum number of new files to be considered in every trigger.

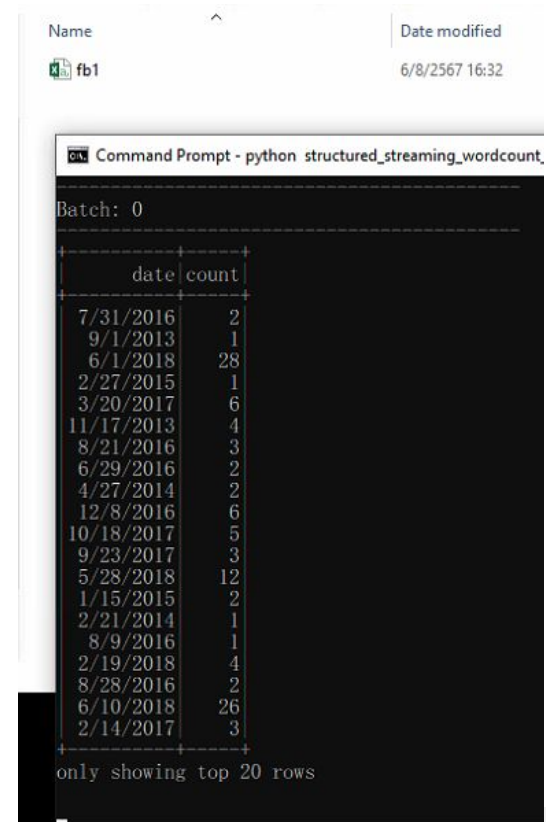- The option **path** is used to identify the folder containing data.

# File Source

```
# Print schema
words = lines.withColumn("date", \
          split(lines["status_published"], " ").\
          getItem(1)) # Get the date of status published
wordCounts = words.groupBy("date", \
          "status_type").count() # Count based on date
                                 # and status type

# Write Stream
# Wait for termination
```

Dr Sirintra Vaiwsri

# File Source

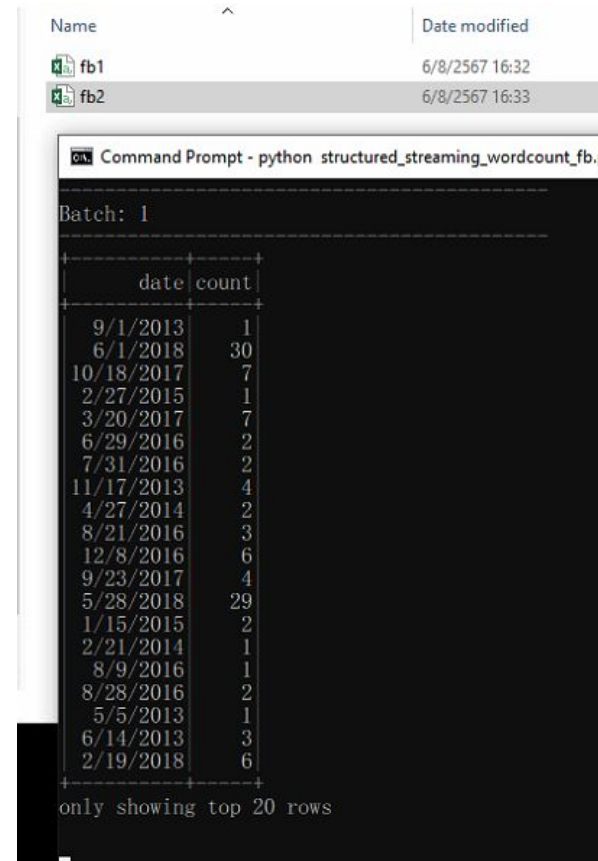- Create a folder and add fb1.csv into the created folder

- Run the Stream File Source Python code.

# File Source

- Add the file fb2.csv into the created folder

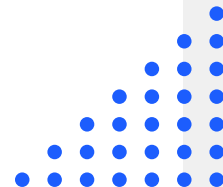- You will see the second Dataframe, where the numbers of counts are updated.

# Assignment (2 points)

- Please implement the socket and file codes.
- Please execute your code and show the results to get 2 points.

# References

- Chambers, B., & Zaharia, M. (2018). *Spark: The definitive guide: Big data processing made simple*. " O'Reilly Media, Inc.".
- Antolínez García, A. (2023). *Hands-on Guide to Apache Spark 3: Build Scalable Computing Engines for Batch and Stream Data Processing*. Berkeley, CA: Apress.
- Spark. https://spark.apache.org. Accessed: 2023-10-16.

Dr Sirintra Vaiwsri