



From telemetry data to CSVs with Python, Spark and Azure Databricks

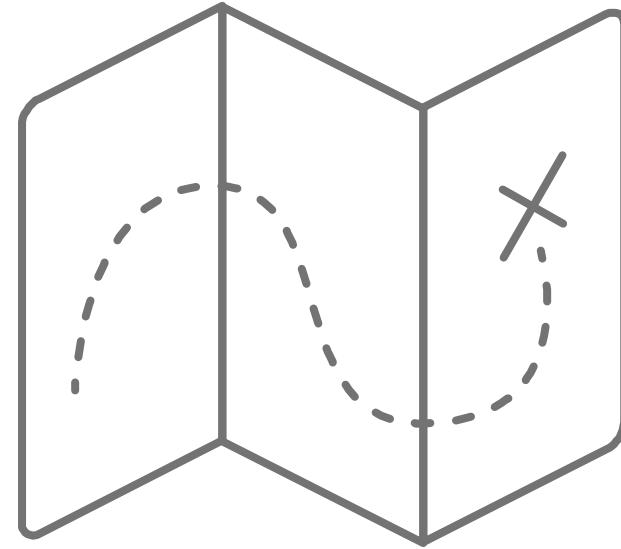
EUROPYTHON 2021

Nicolò Giso

Agenda



- Background
- Data Applications
- Architecture
- Code
 - Some lines of codes
 - Other lines
 - And another bunch of them
- Orchestration
- Conclusion



Background

Let me introduce myself



Nicolò Giso

About Me

Role: Data Scientist

Company: Tenova

Interests: books, tv series, data

Professional Skills

Programming

75%

PowerPoint

20%

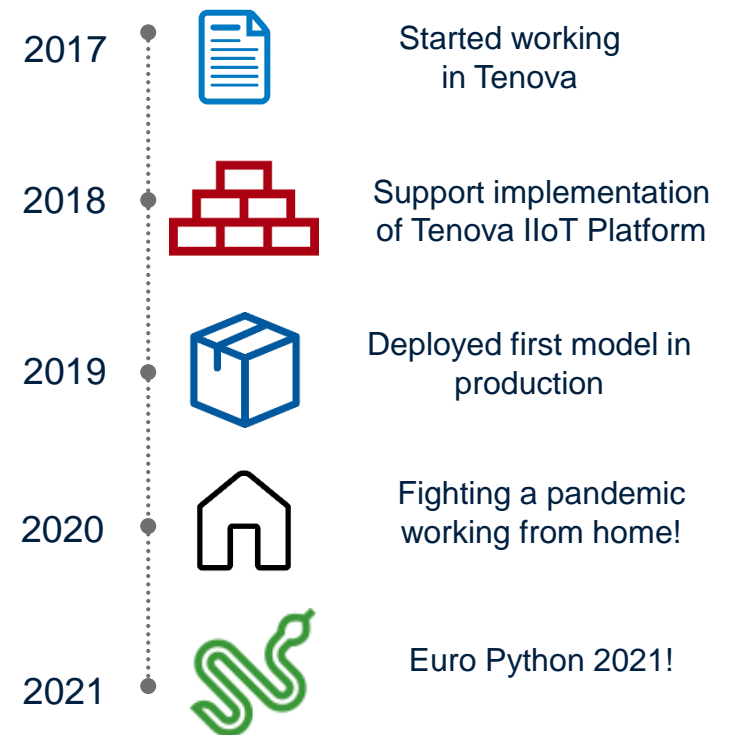
Bad Memes

91%

Ctrl+C – Ctrl+V

89%

My History



TENOVA, A TECHINT GROUP COMPANY, IS YOUR WORLDWIDE PARTNER
FOR **SUSTAINABLE, INNOVATIVE AND RELIABLE SOLUTIONS**
IN THE METALS AND –THROUGH THE WELL-KNOWN TAKRAF AND DELKOR BRANDS
– **IN THE MINING** INDUSTRIES.

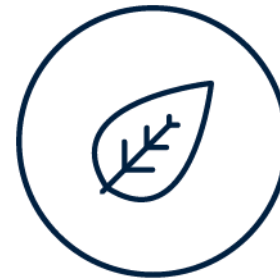
WE DESIGN AND DEVELOP SOLUTIONS THAT HELP COMPANIES TO:



REDUCE
COSTS



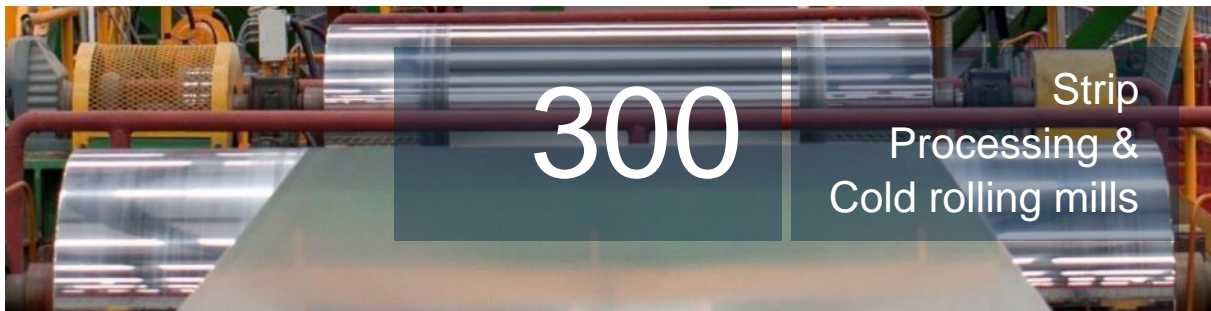
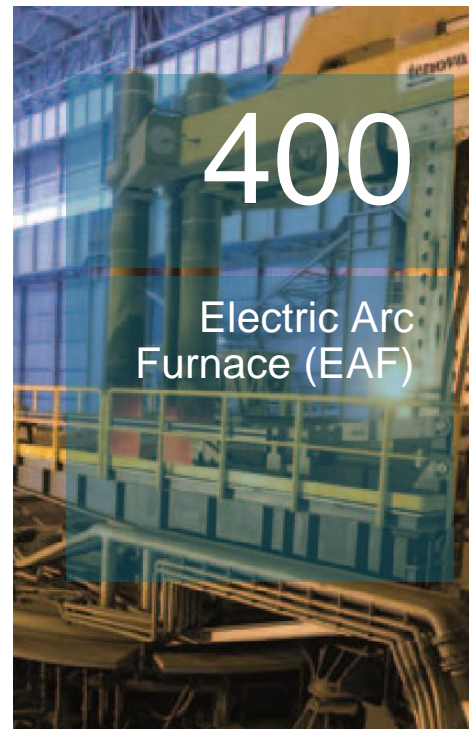
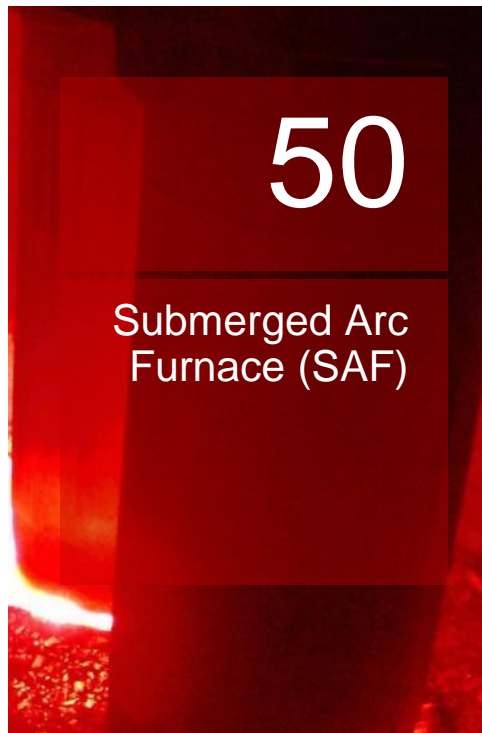
SAVE
ENERGY



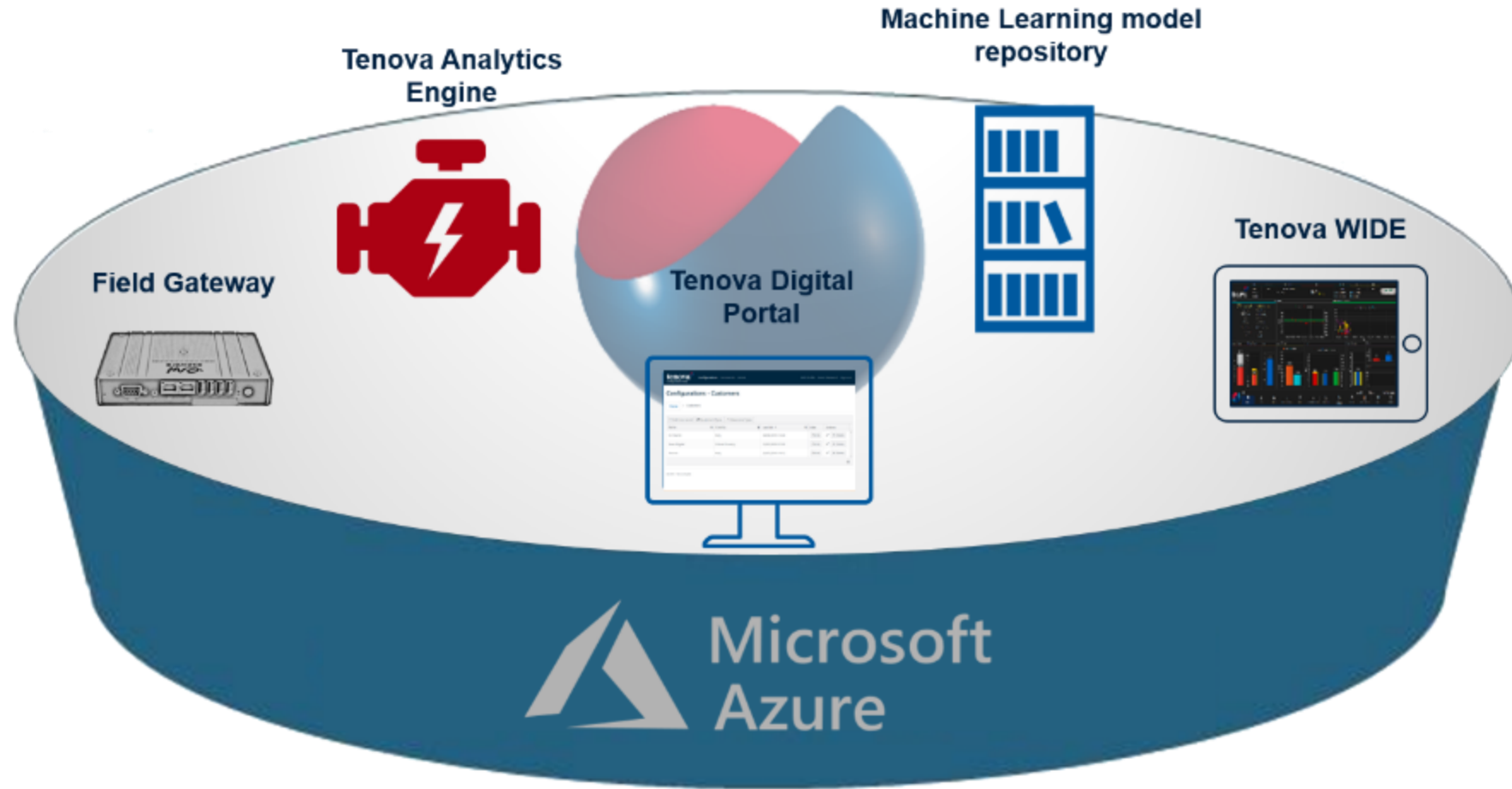
LIMIT
ENVIRONMENTAL
IMPACT



IMPROVE
WORKING
CONDITIONS

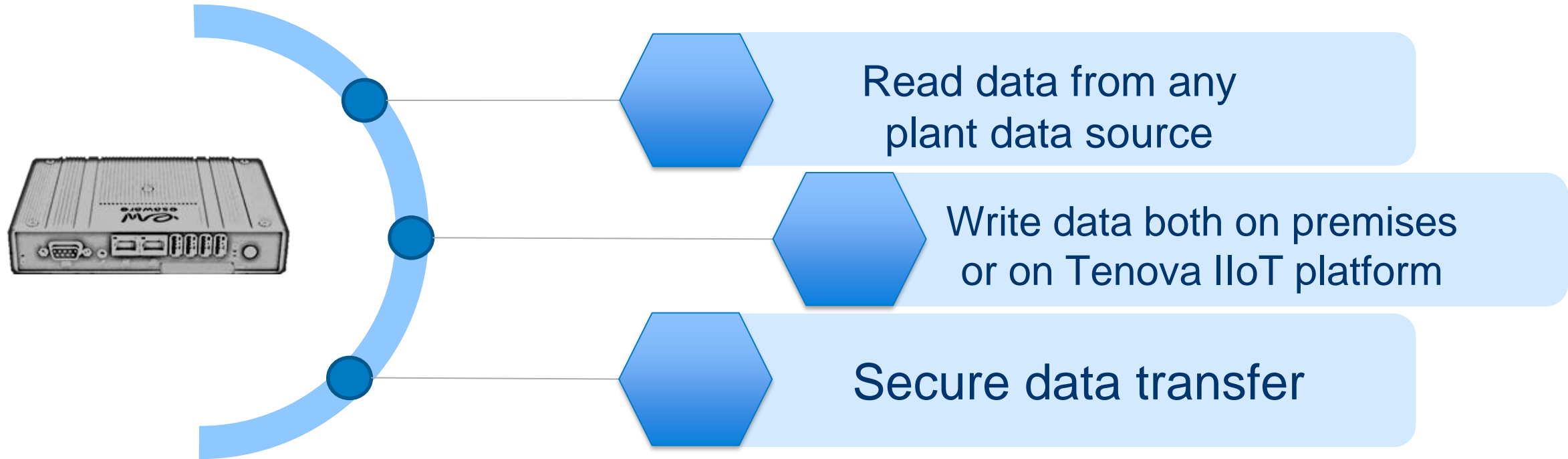


Tenova IloT platform



Tenova IloT platform

TENOVA EDGE



Data Applications

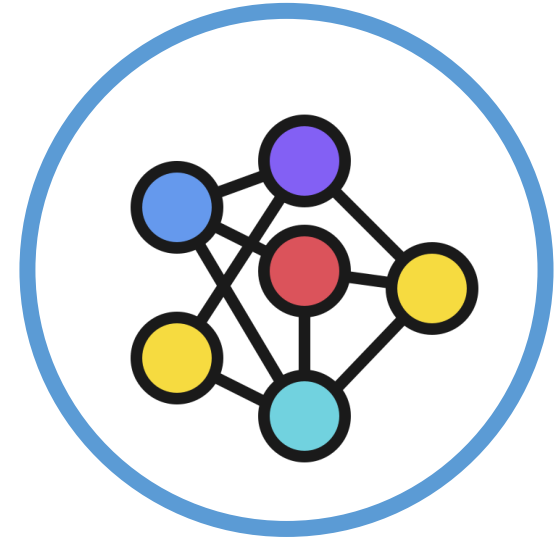
Data Applications



Monitoring



Data Analysis



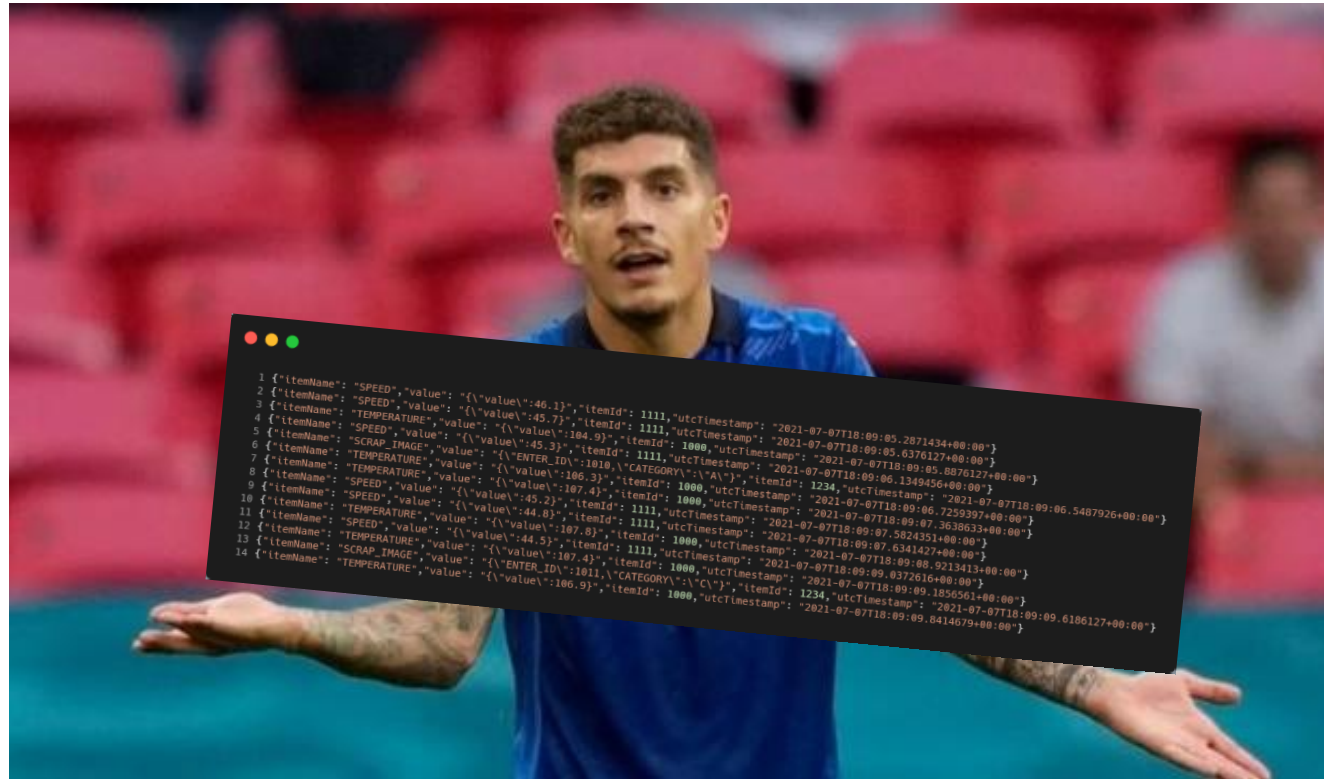
ML Models

Raw Data

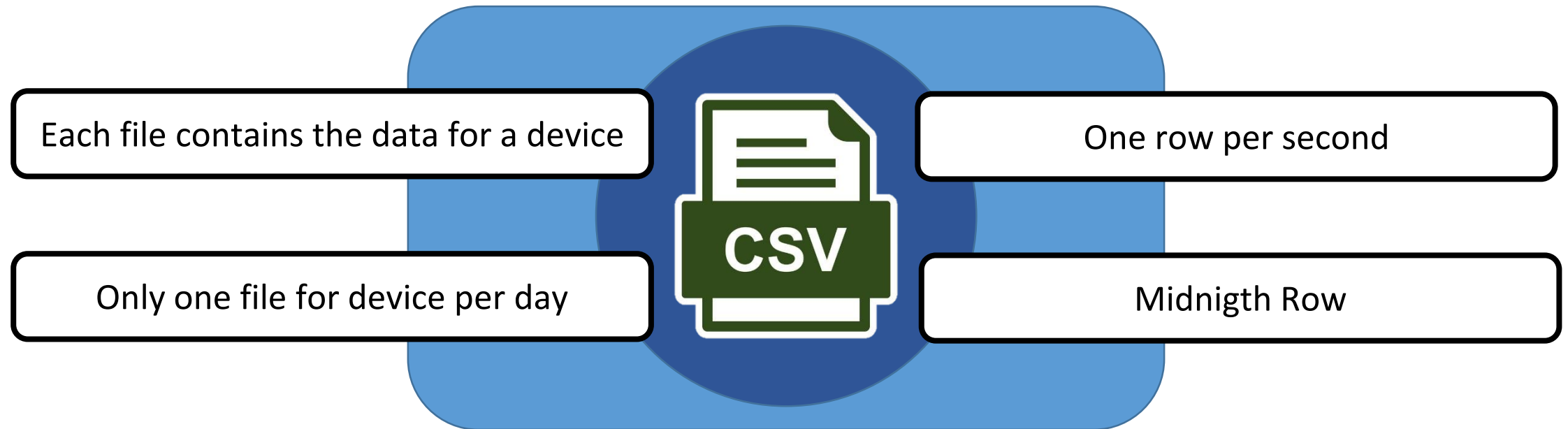


```
1 {"itemName": "SPEED", "value": "{\\"value\\":46.1}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:05.2871434+00:00"}
2 {"itemName": "SPEED", "value": "{\\"value\\":45.7}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:05.6376127+00:00"}
3 {"itemName": "TEMPERATURE", "value": "{\\"value\\":104.9}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:05.8876127+00:00"}
4 {"itemName": "SPEED", "value": "{\\"value\\":45.3}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:06.1349456+00:00"}
5 {"itemName": "SCRAP_IMAGE", "value": "{\\"ENTER_ID\\":1010,\\"CATEGORY\\":\\"A\\"}", "itemId": 1234, "utcTimestamp": "2021-07-07T18:09:06.5487926+00:00"}
6 {"itemName": "TEMPERATURE", "value": "{\\"value\\":106.3}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:06.7259397+00:00"}
7 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.4}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:07.3638633+00:00"}
8 {"itemName": "SPEED", "value": "{\\"value\\":45.2}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:07.5824351+00:00"}
9 {"itemName": "SPEED", "value": "{\\"value\\":44.8}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:07.6341427+00:00"}
10 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.8}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:08.9213413+00:00"}
11 {"itemName": "SPEED", "value": "{\\"value\\":44.5}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:09.0372616+00:00"}
12 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.4}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:09.1856561+00:00"}
13 {"itemName": "SCRAP_IMAGE", "value": "{\\"ENTER_ID\\":1011,\\"CATEGORY\\":\\"C\\"}", "itemId": 1234, "utcTimestamp": "2021-07-07T18:09:09.6186127+00:00"}
14 {"itemName": "TEMPERATURE", "value": "{\\"value\\":106.9}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:09.8414679+00:00"}
```

Our engineers looking at raw data

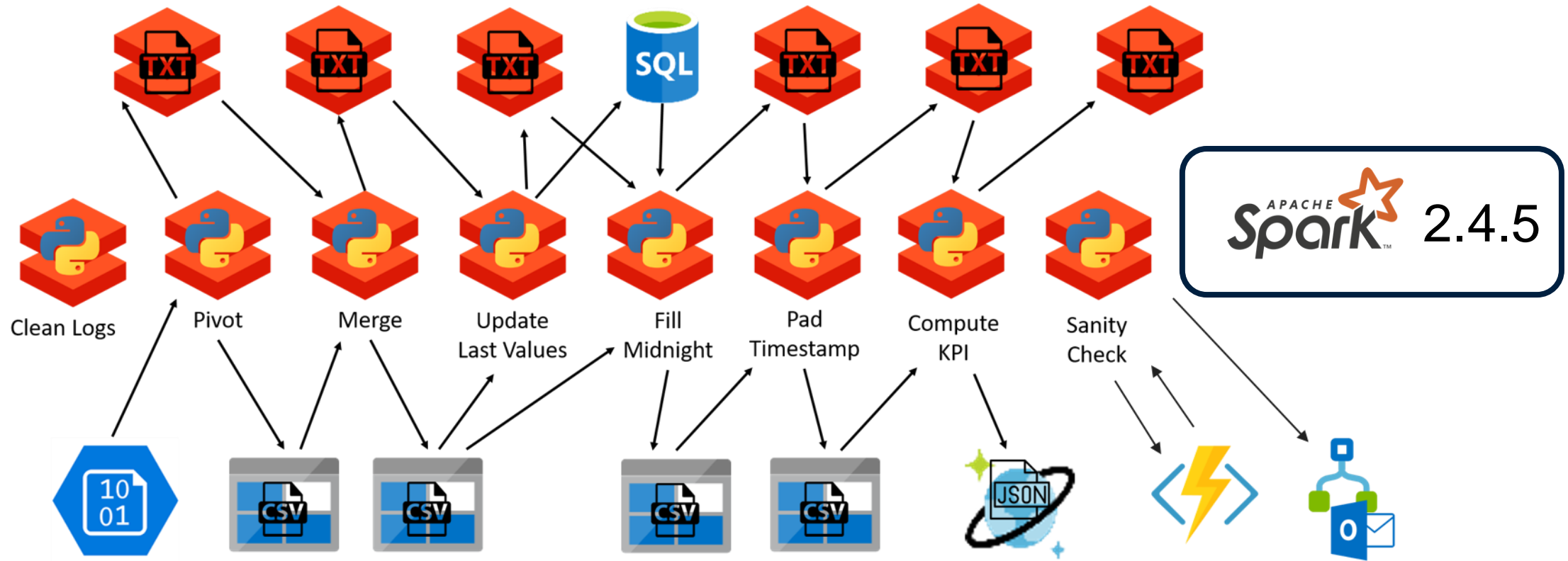


Requirements



Architecture

Architecture

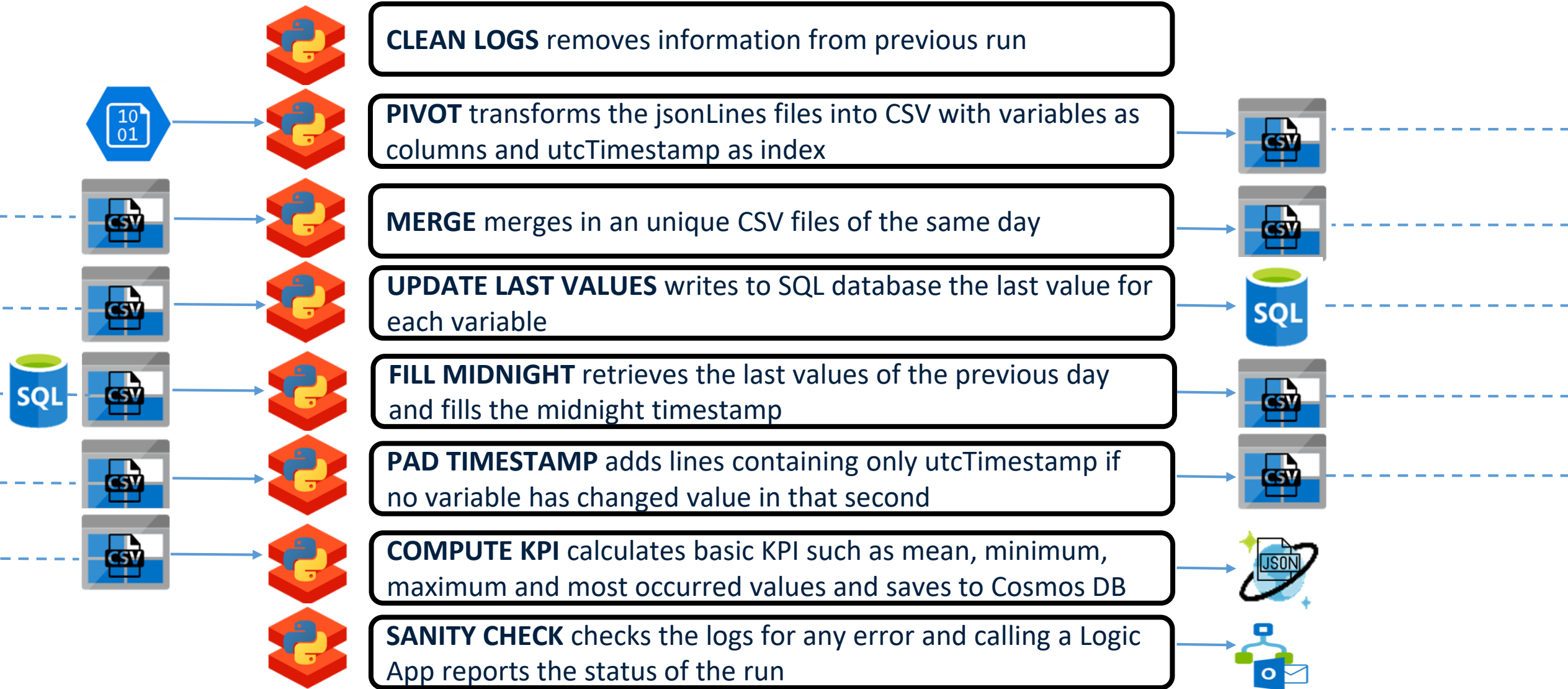


Our intern looking at the architecture



Architecture

I HOPE IT'S A LITTLE CLEARER



Code

Raw data



```
1 {"itemName": "SPEED", "value": "{\\"value\\":46.1}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:05.2871434+00:00"}
2 {"itemName": "SPEED", "value": "{\\"value\\":45.7}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:05.6376127+00:00"}
3 {"itemName": "TEMPERATURE", "value": "{\\"value\\":104.9}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:05.8876127+00:00"}
4 {"itemName": "SPEED", "value": "{\\"value\\":45.3}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:06.1349456+00:00"}
5 {"itemName": "SCRAP_IMAGE", "value": "{\\"ENTER_ID\\":1010,\\"CATEGORY\\":\\"A\\"}", "itemId": 1234, "utcTimestamp": "2021-07-07T18:09:06.5487926+00:00"}
6 {"itemName": "TEMPERATURE", "value": "{\\"value\\":106.3}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:06.7259397+00:00"}
7 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.4}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:07.3638633+00:00"}
8 {"itemName": "SPEED", "value": "{\\"value\\":45.2}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:07.5824351+00:00"}
9 {"itemName": "SPEED", "value": "{\\"value\\":44.8}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:07.6341427+00:00"}
10 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.8}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:08.9213413+00:00"}
11 {"itemName": "SPEED", "value": "{\\"value\\":44.5}", "itemId": 1111, "utcTimestamp": "2021-07-07T18:09:09.0372616+00:00"}
12 {"itemName": "TEMPERATURE", "value": "{\\"value\\":107.4}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:09.1856561+00:00"}
13 {"itemName": "SCRAP_IMAGE", "value": "{\\"ENTER_ID\\":1011,\\"CATEGORY\\":\\"C\\"}", "itemId": 1234, "utcTimestamp": "2021-07-07T18:09:09.6186127+00:00"}
14 {"itemName": "TEMPERATURE", "value": "{\\"value\\":106.9}", "itemId": 1000, "utcTimestamp": "2021-07-07T18:09:09.8414679+00:00"}
```

CODE



PIVOT transforms the jsonLines files into CSV with variables as columns and utcTimestamp as index



```
1 from pyspark.sql.functions import to_timestamp, explode, col, desc, last, date_trunc, row_number
2 from pyspark.sql.window import Window
3
4
5 def transform_dataframe(raw_df):
6     df = raw_df.select(
7         to_timestamp(raw_df.utcTimestamp).alias("utcTimestamp"),
8         raw_df.itemName,
9         explode(parse(raw_df.value)).alias("key", "value"),
10        raw_df.itemId,
11    )
12    df = df.withColumn("variable", merge_columns("key", "itemName", "itemId"))
13    df = df.select(df.utcTimestamp, df.variable, df.value)
14
15    # truncatedUtc is introduced to select just the last value (desc) for each pair of [utcTimestamp, variable]
16    df = df.withColumn("truncatedUtc", date_trunc("second", "utcTimestamp"))
17    w = Window.partitionBy("variable", "truncatedUtc").orderBy(desc("utcTimestamp"))
18    # for each variable at a given truncatedUtc the first row is selected.
19    # Since the window is in descending order the last value is selected.
20    df = df.withColumn("row_number", row_number().over(w)).where(col("row_number") == 1).drop("rownumber")
21
22    # there is just one value per timestamp, per variable. last is used just to take the unique value.
23    df = df.groupby("truncatedUtc").pivot("variable").agg(last("value"))
24    df = df.drop("utcTimestamp").withColumnRenamed("truncatedUtc", "utcTimestamp").orderBy("utcTimestamp")
25
26    return df
```

utcTimestamp	1000 TEMPERATURE	1111 SPEED	1234 SCRAP_IMAGE_ENTER_ID	1234 SCRAP_IMAGE_CATEGORY
2021-07-07T18:09:05	104.9	45.7		
2021-07-07T18:09:06	106.3	45.3	1010	A
2021-07-07T18:09:07	107.4	44.8		
2021-07-07T18:09:08	107.8			
2021-07-07T18:09:09	106.9	44.5	1011	C

Merge

CODE



MERGE merges in an unique CSV files of the same day



```
1 def merge(day_files, conf):
2     dfs = [spark.read.csv(_, header=True, multiLine=True) for _ in day_files]
3     df_raw = reduce(
4         lambda df1, df2: df1.join(
5             df2, list(set(df1.columns).intersection(df2.columns)), "outer"
6         ),
7         dfs,
8     )
9     df_raw = sort_columns(df_raw)
10    schema = df_raw.schema
11    df_no_duplicates = df_raw.dropDuplicates(subset=["utcTimestamp"])
12    if df_no_duplicates.count() == df_raw.count():
13        return df_raw.sort(asc("utcTimestamp"))
14    else:
15        df_raw = df_raw.toPandas()
16        df_duplicated = df_raw[df_raw.duplicated(subset="utcTimestamp", keep=False)]
17        df_no_duplicates = df_raw[~df_raw.duplicated(subset="utcTimestamp", keep=False)]
18
19        resolved_duplicates = df_duplicated.groupby("utcTimestamp").agg(unique_udf)
20        resolved_duplicates["utcTimestamp"] = resolved_duplicates.index
21        resolved_duplicates.index = range(len(resolved_duplicates))
22
23        result = pd.concat([df_no_duplicates, resolved_duplicates]).sort_values(
24            "utcTimestamp"
25        )
26        result = result[schema.fieldNames()]
27        result = spark.createDataFrame(result, schema)
28        return result
```


Update Last Values



CODE



UPDATE LAST VALUES writes to SQL database the last value for each variable



```
1 def find_last_value_column(obj):
2     last_non_null = obj.dropna().tail(1)
3     try:
4         last_non_null_value = last_non_null[0]
5         last_non_null_timestamp = last_non_null.index[0]
6         return {"value": last_non_null_value, "timestamp": last_non_null_timestamp}
7     except:
8         return None
9
10
11 def generate_last_values_dict(df):
12     last_values_dict = {
13         column: find_last_value_column(df[column])
14         for column in df.columns
15         if column != "utcTimestamp"
16     }
17     last_values_dict = {
18         key: value for key, value in last_values_dict.items() if value is not None
19     }
20     return last_values_dict
```

Update Last Values



QUERY

Timestamp	ItemId	Value	DeviceId	ItemName
2021-07-26 23:59:20	1000	103.5	1	TEMPERATURE
2021-07-26 23:59:56	1111	46.8	1	SPEED
2021-07-26 23:55:15	1234	1982	1	SCRAP_IMAGE_ENTER_ID
2021-07-26 23:55:15	1234	A	1	SCRAP_IMAGE_CATEGORY

```
1 /*get timestamp*/
2 SELECT Timestamp FROM lastValues WHERE ItemId={itemId} AND CONVERT(DATE, Timestamp)=CONVERT(DATE, '{last_timestamp}') AND ItemName='{itemName}'
3
4 /*update last value*/
5 UPDATE lastValues SET ItemId = {itemId}, Timestamp = '{last_timestamp}', Value = '{value}', ItemName='{itemName}', DeviceId= {device} WHERE ItemId={itemId}
  AND ItemName='{itemName}' AND CONVERT(DATE, Timestamp)=CONVERT(DATE, '{last_timestamp}')
6
7 /*insert last value*/
8 INSERT INTO lastValues(ItemId, Timestamp, Value, DeviceId, ItemName) VALUES ('{itemId}', '{last_timestamp}', '{value}', {device}, '{itemName}')
```

Fill Midnight

RETRIEVE LAST VALUES



FILL MIDNIGHT retrieves the last values of the previous day and fills the midnight timestamp



```
1 from pandas import read_sql
2
3
4 def retrieve_last_values(device, day):
5     data = read_sql(
6         f"""SELECT itemid, itemname, value FROM LastValues
7             WHERE DATEADD(DAY, +1, CONVERT(DATE, Timestamp))='{day}'
8             and deviceid={device}""",
9         cnxn,
10    )
11    yesterday_column_names = [
12        f'{value}|{data["itemname"][key]}' for key, value in enumerate(data["itemid"])
13    ]
14    yesterday_last_values = dict(zip(yesterday_column_names, data["value"]))
15    return yesterday_last_values
```

Fill Midnight

UPDATE DATAFRAME



```
1 def midnight_exists(df):
2     first_timestamp_string = df.select("utcTimestamp").first()["utcTimestamp"]
3     # very bad trick
4     hour_minute_seconds_string = first_timestamp_string.split("T")[1][:8]
5     if hour_minute_seconds_string == "00:00:00":
6         return True
7     else:
8         return False
9
10
11 def update_dataframe(df, info):
12     variables_df = [_ for _ in df.columns if _ != "utcTimestamp"]
13     if info == {}:
14         return df
15     for _ in variables_df:
16         if _ not in info.keys():
17             info[_] = ""
18     if midnight_exists(df):
19         df = fill_midnight_timestamp(df, info)
20         print("fill midnight")
21     else:
22         print("insert midnight")
23         df = insert_midnight_timestamp(df, info)
24     return df
```

Pad Timestamp



PAD TIMESTAMP adds lines containing only utcTimestamp if no variable has changed value in that second



```
1 from pyspark.sql.functions import col, min as min_, max as max_
2
3
4 def pad_timestamp(df, step=1):
5     minp, maxp = df.select(
6         min_("utcTimestamp").cast("timestamp").cast("int"),
7         max_("utcTimestamp").cast("timestamp").cast("int"),
8     ).first()
9     reference = spark.range(
10        (minp / step) * step, ((maxp / step) + 1) * step, step
11    ).select(col("id").cast("timestamp").alias("utcTimestamp"))
12    result = reference.join(df, ["utcTimestamp"], "leftouter")
13    return result
```

Compute KPI



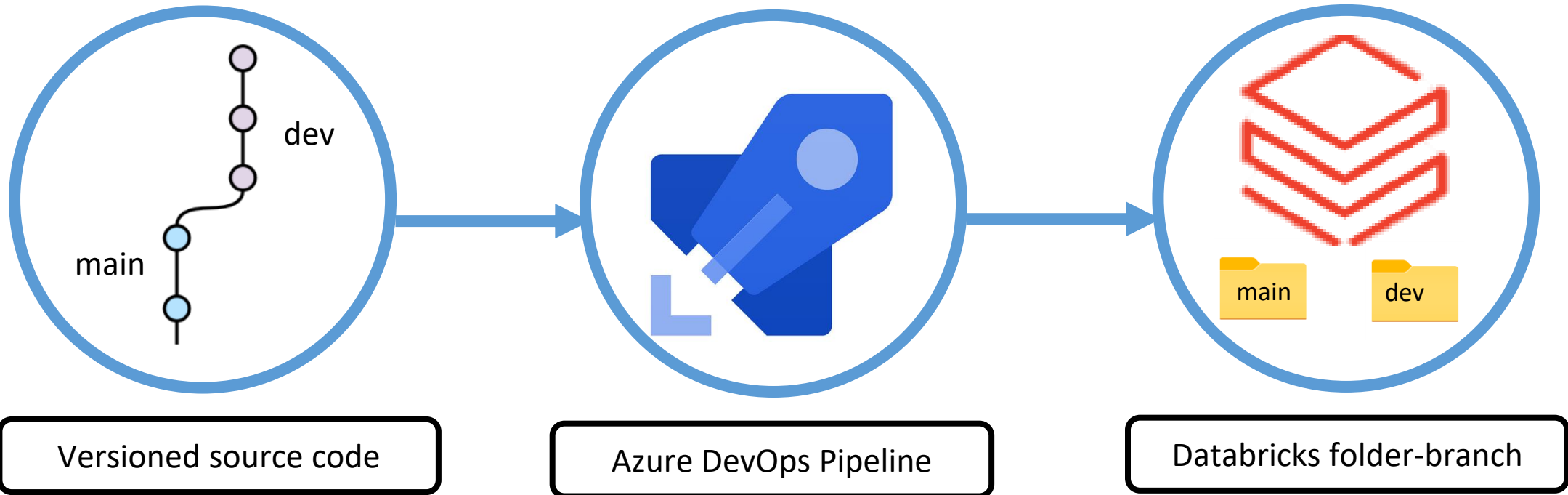
COMPUTE KPI calculates basic KPI such as mean, minimum, maximum and most occurred values and saves to Cosmos DB



```
1 import pandas as pd
2
3
4 def compute_basic_KPIs(df):
5     df = df.apply(pd.to_numeric, errors="ignore")
6     summary = df.describe(include="all")
7     summary_pd = summary.drop(["utcTimestamp"], axis=1)
8     df_filled = df.ffill()
9     summary_filled = df_filled.describe(include="all")
10    summary_filled_pd = summary_filled.drop(["utcTimestamp"], axis=1)
11
12    result_pd = []
13    for column in summary_pd.columns:
14        item = {}
15        tag_id, tag_name = column.split("|")
16        item["TagId"] = tag_id
17        item["KPI"] = build_KPI_dict_for_column(
18            summary_not_filled=summary_pd,
19            summary_filled=summary_filled_pd,
20            column=column,
21        )
22        time_occ = df_filled[column].value_counts()
23        item["TimeOccurencies"] = dict(time_occ[time_occ > 1])
24        item["TagName"] = tag_name
25        serialized_item = pd.Series(item).to_json()
26        result_pd.append(item)
27    return result_pd
```

Orchestration

Continuous Deployment



The logo for Azure Databricks, consisting of three stacked, slightly offset red squares.

Azure Databricks

Same environment as the notebooks

Little flexibility for future developments



Azure Data Factory

Easy to customize

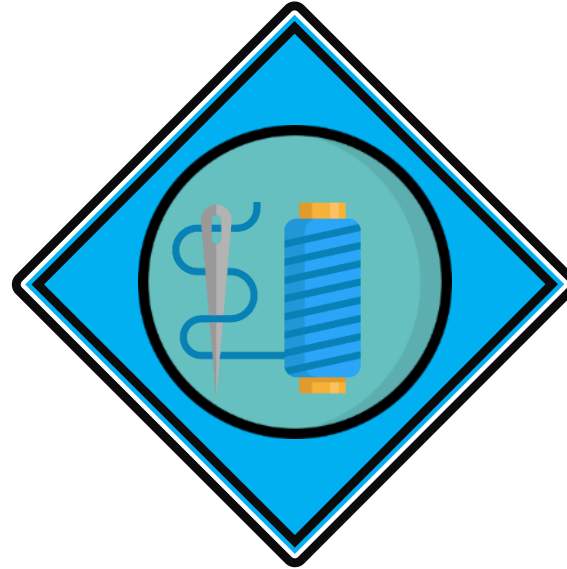
Flexibility in replacing some notebooks with other services

Conclusion

Next Steps



UPGRADE TO SPARK 3



COMPUTE AD HOC KPI



**IMPROVE
PERFORMANCE
COMPUTE KPI**

The End!

