



Building an E2E Analytics Pipeline with PyFlink

Marta Paes (@morsapaes)
Developer Advocate

About Ververica



Original Creators of
Apache Flink®



Enterprise Stream Processing
With Ververica Platform



Part of
Alibaba Group



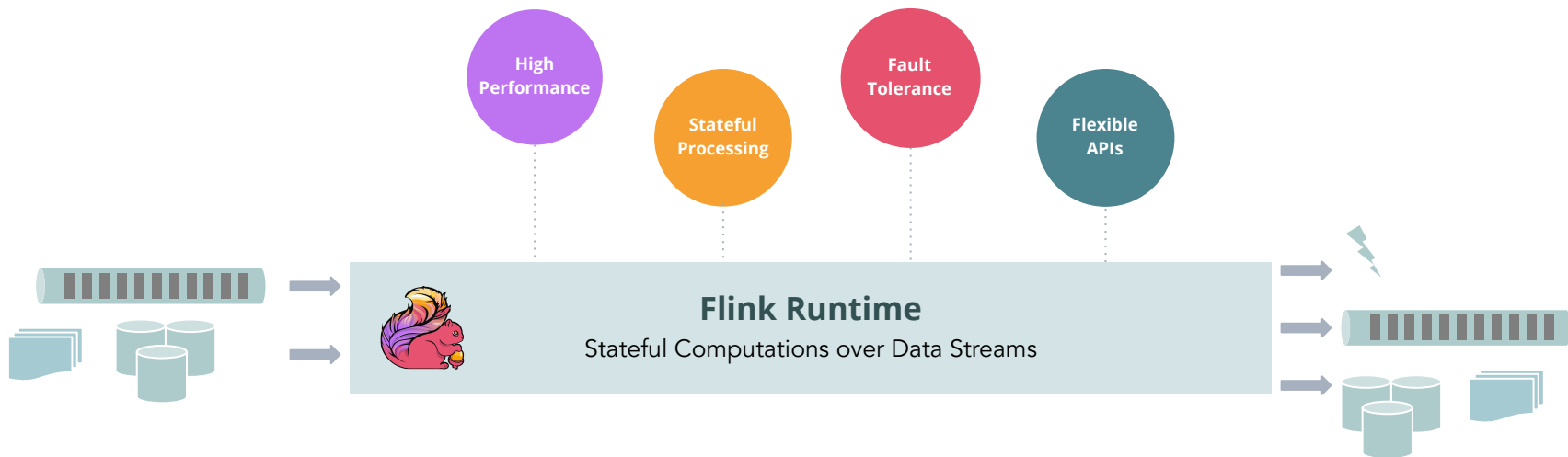
Apache Flink

Flink is an **open source** framework and **distributed** engine for **unified batch and stream processing**.



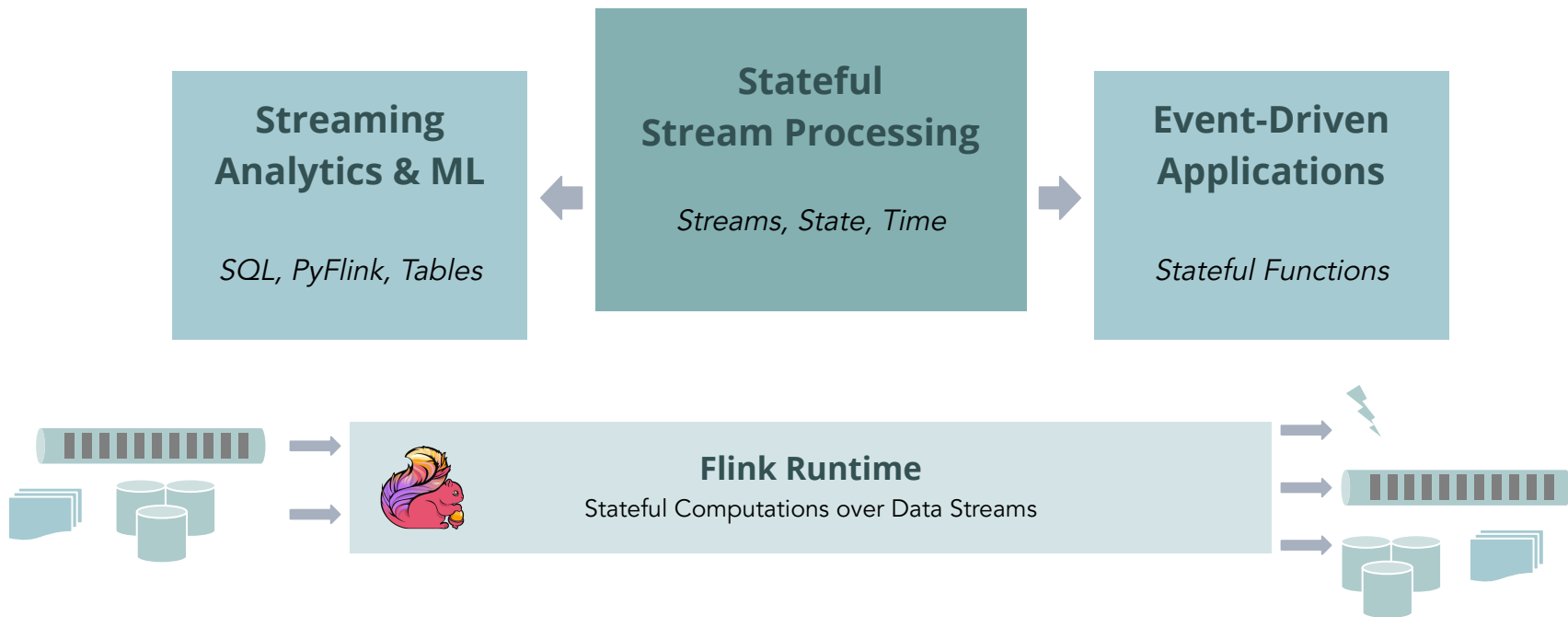
Apache Flink

Flink is an **open source** framework and **distributed** engine for **unified batch and stream processing**.



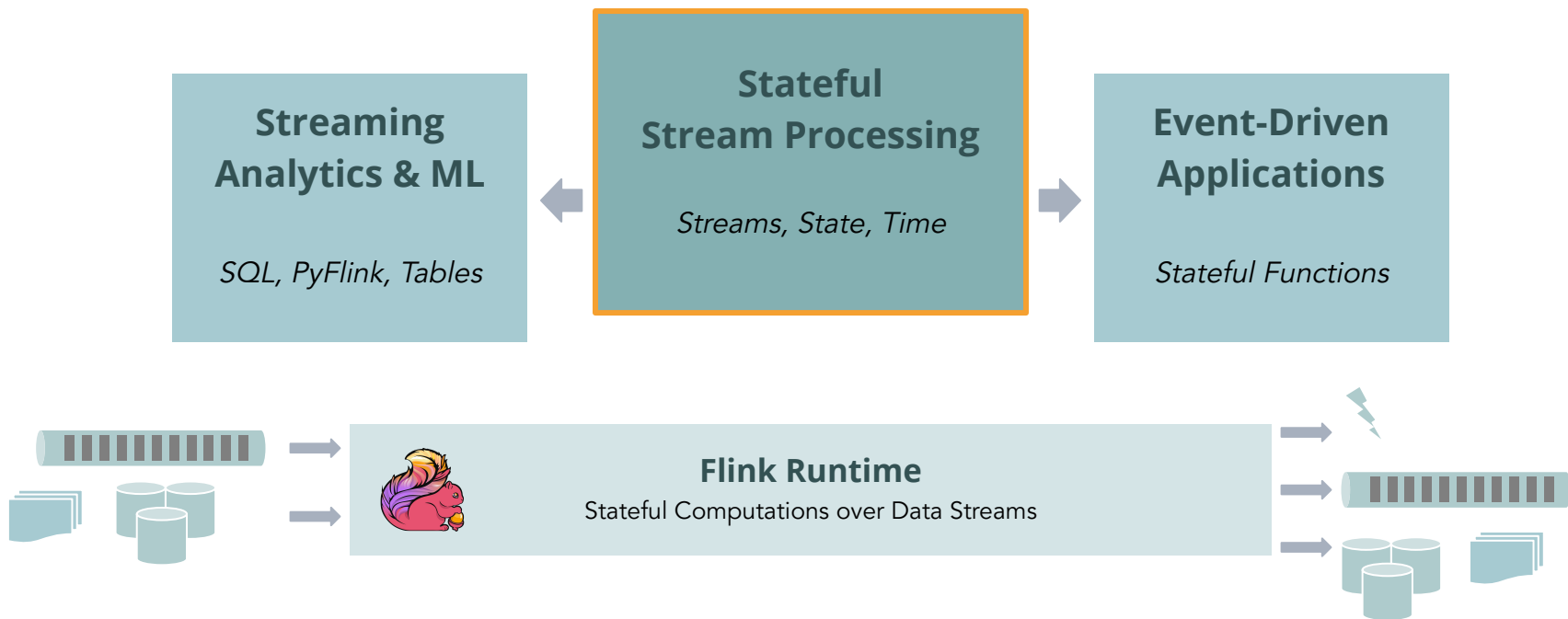
Use Cases

This gives you a robust foundation for a wide range of use cases:



Use Cases

Classical, core stream processing use cases that build on the primitives of **streams**, **state** and **time**.



Stateful Stream Processing

Classical, core stream processing use cases that build on the primitives of **streams**, **state** and **time**.

- Explicit control over these primitives
- Complex computations and customization
- Maximize **performance** and **reliability**

Example Use Cases

NETFLIX

[Large-scale Data Pipelines](#)

ING 

[ML-Based Fraud Detection](#)

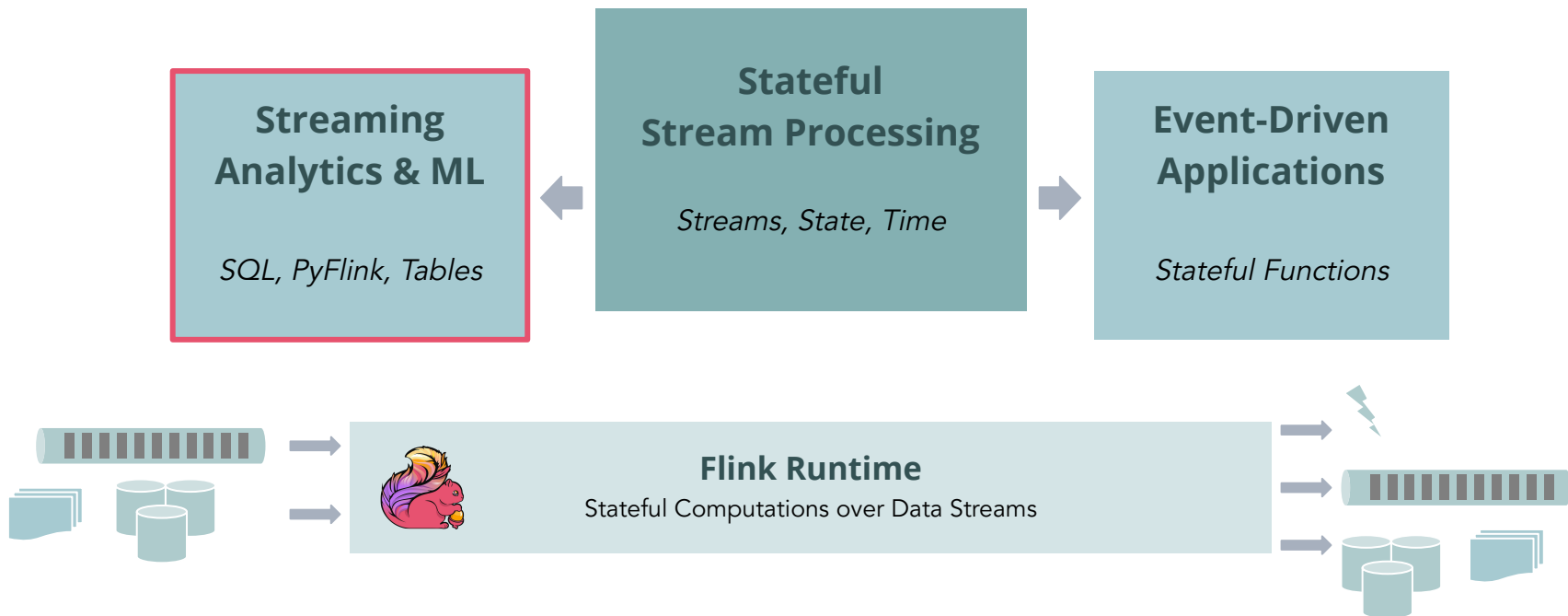
aws 

[Service Monitoring & Anomaly Detection](#)



Use Cases

More high-level or domain-specific use cases that can be modeled with SQL or Python and dynamic tables.



Streaming Analytics & ML

More high-level or domain-specific use cases that can be modeled with SQL or Python and dynamic tables.

- Focus on logic, not implementation
- Mixed workloads (batch and streaming)
- Maximize **developer speed** and **autonomy**

Example Use Cases



Uber

criteo.

[Unified Online/Offline Model Training](#)

[E2E Streaming Analytics Pipelines](#)

[ML Feature Generation](#)



More Flink Users



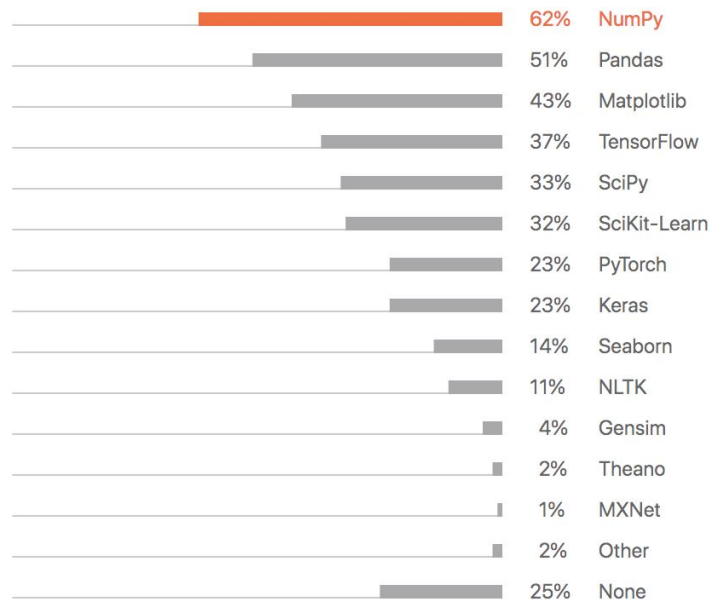


Why PyFlink?



Python is...pretty stacked?

What data science frameworks do you use in addition to Python?

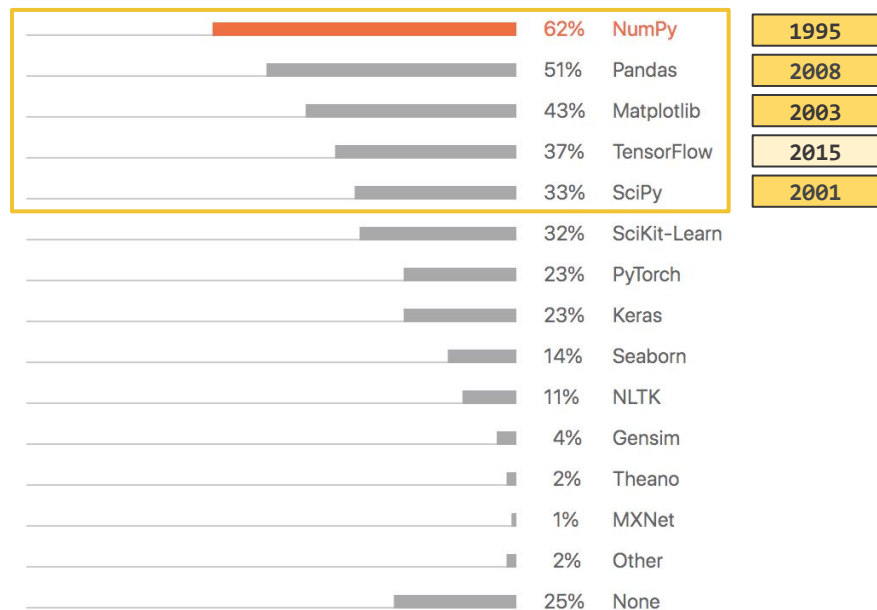


Mature analytics stack, with libraries that are **fast** and **intuitive**.



...and also timeless!

What data science frameworks do you use in addition to Python?

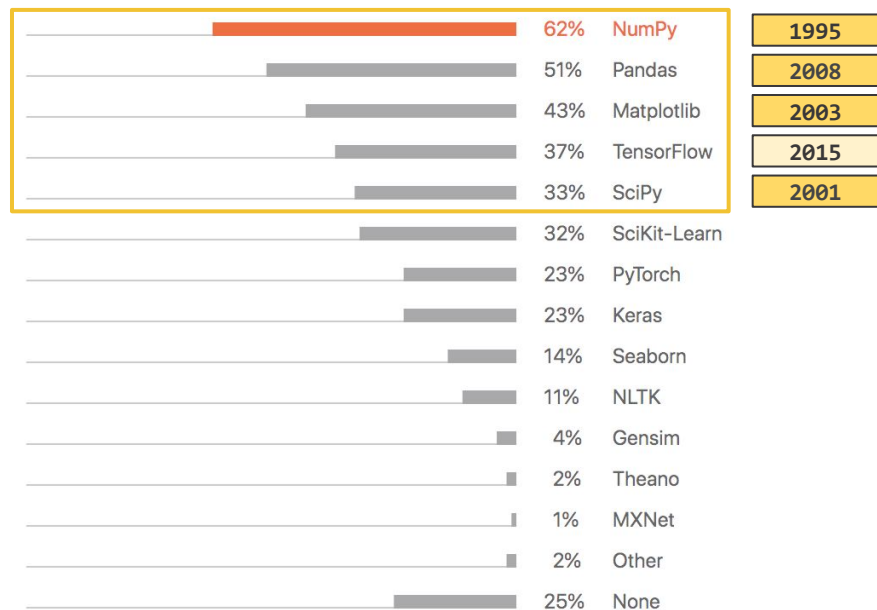


Mature analytics stack, with libraries that are **fast** and **intuitive**.



...and also timeless!

What data science frameworks do you use in addition to Python?



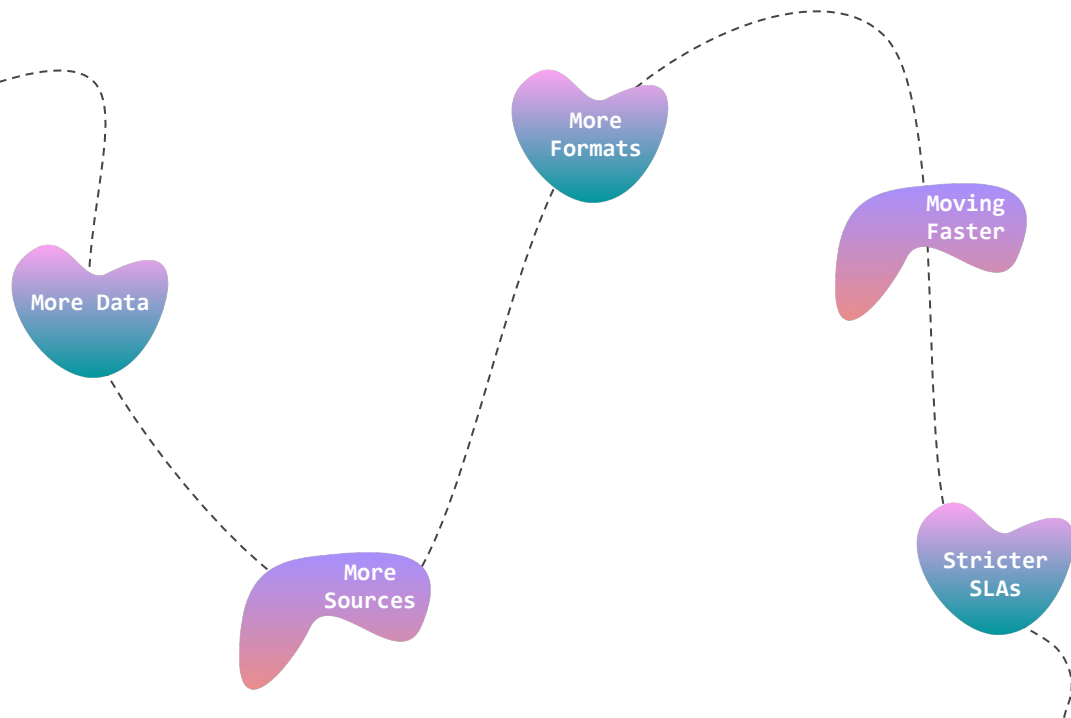
Mature analytics stack, with libraries that are **fast** and **intuitive**.



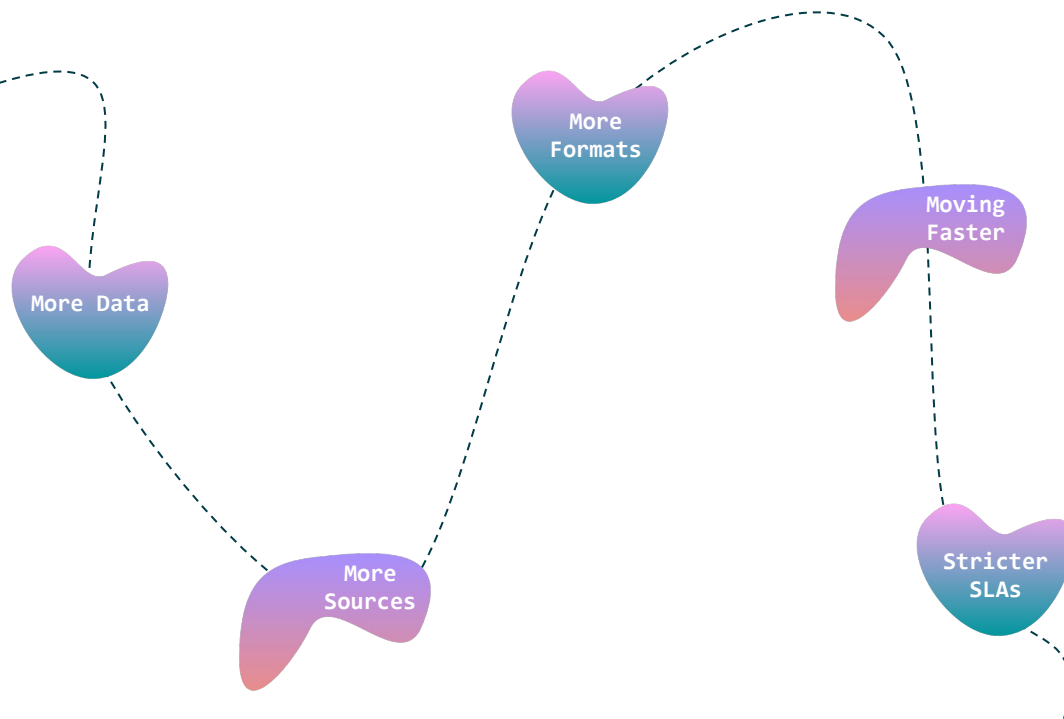
Older libraries are mostly **restricted** to a data size that **fits in memory** (RAM), and designed to run on a **single core** (CPU).



This is a problem, because



This is a problem, because



But you still want to use these powerful libraries, right?



Why PyFlink?

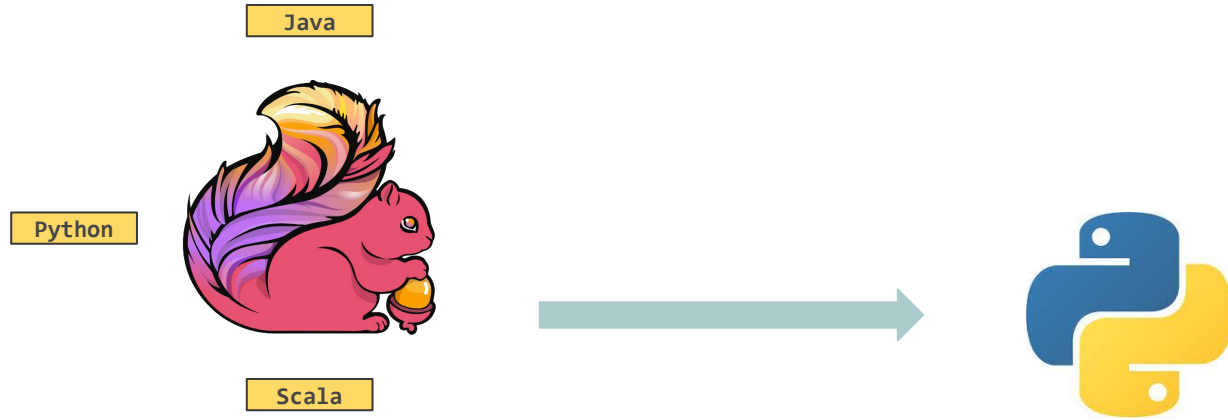
Java



Scala



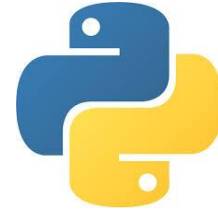
Why PyFlink?



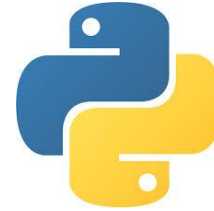
Expose the functionality of Flink beyond the JVM



Why PyFlink?



Why PyFlink?



Distribute and **scale** the functionality of Python through Flink



Flink at Alibaba scale

Double 11 / Singles Day 

Real-time Data Applications incl. sub-second updates to the GMV dashboard

Search

Recomm.

Ads

BI

Security



Data Size



1.7EB

Throughput (Peak)



4B
events/sec

State Size (Biggest)



100TB

Latency



Sub-sec





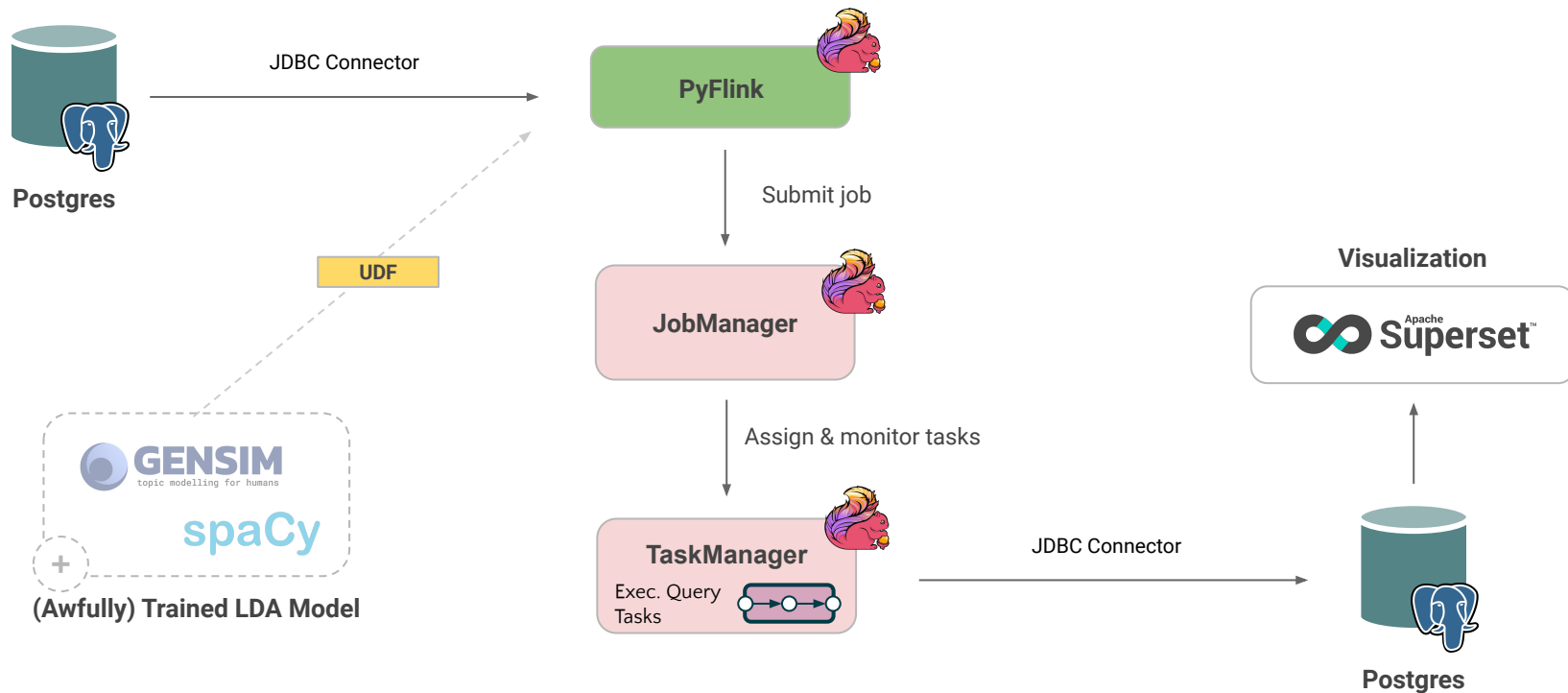
Demo



Can we use PyFlink to identify the most frequent topics in the Flink User Mailing List?



The Demo Environment



Step 1. Create the source and sink tables.

```
ddl_jdbc_source = """CREATE TABLE user_ml_messages (  
    message_date TIMESTAMP(3),  
    message_id STRING,  
    message_in_reply_to STRING,  
    message_from_name STRING,  
    message_from_email STRING,  
    message_subject STRING,  
    message_body_html STRING,  
    message_body_plain STRING,  
    PRIMARY KEY(message_id) NOT ENFORCED  
)  
WITH (  
    'connector' = 'jdbc',  
    'url' = 'jdbc:postgresql://postgres:5432/postgres',  
    'table-name' = 'perceval.stg_flink_user_ml',  
    'username' = 'postgres',  
    'password' = 'postgres'  
)"""  
  
create_jdbc_source = t_env.execute_sql(ddl_jdbc_source)
```

Connector Properties

```
ddl_jdbc_sink = """CREATE TABLE flink_user_ml_topics (  
    message_id STRING,  
    message_date TIMESTAMP(3),  
    message_from_name STRING,  
    topic VARCHAR(2)  
)  
WITH (  
    'connector' = 'jdbc',  
    'url' = 'jdbc:postgresql://postgres:5432/postgres',  
    'table-name' = 'perceval.flink_user_ml_topics',  
    'username' = 'postgres',  
    'password' = 'postgres'  
)"""  
  
create_jdbc_sink = t_env.execute_sql(ddl_jdbc_sink)
```

DDL Statement Execution



Step 2. Write and register a UDF to clean and classify the messages.

```
@udf(input_types=DataTypes.STRING(), result_type=DataTypes.STRING())
def class_model(m):

    lda = LdaModel.load("model/lda_model/lda_model_user_ml")

    v = tokenizer.pre_process(m)

    vector = lda[v]

    topics = sorted(vector, key=lambda x: x[1], reverse=True)

    return str(topics[0][0])

t_env.register_function("CLASS_TOPIC", class_model)
```

Text pre-processing and classification logic

UDF Registration

➤ How would this perform if it were defined as a [Pandas UDF](#)?



Step 3. Build your query, that will insert your results into the sink table.

```
t_env.from_path("user_ml_messages") \
    .filter("message_subject.isNotNull") \
    .select("message_id, message_date, message_from_name, CLASS_TOPIC(message_subject)")

t_env.execute_insert("flink_user_ml_topics")
```

OR

```
t_env.execute_sql("INSERT INTO flink_user_ml_topics \
    SELECT message_id, message_date, message_from_name, CLASS_TOPIC(message_subject) \
    FROM user_ml_messages WHERE message_subject IS NOT NULL")
```



Step 4. Submit the job (and dependencies) to the cluster.

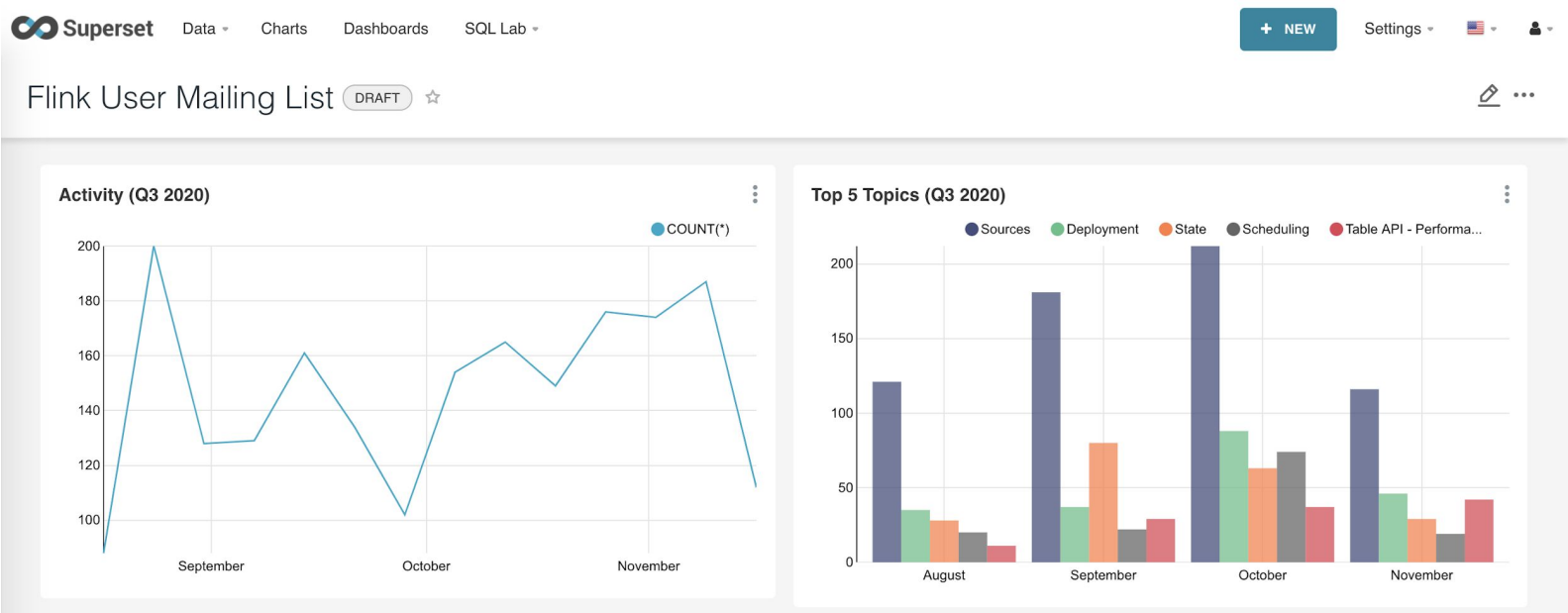
```
docker-compose exec jobmanager ./bin/flink run -py /opt/pyflink-ff2020/pipeline.py \  
--pyArchives /opt/pyflink-ff2020/lda_model.zip#model \ ] LDA Model + Dictionary  
--pyFiles /opt/pyflink-ff2020/tokenizer.py -d ] Pre-processing Class
```

Flink Web UI

Running Jobs						
Job Name	Start Time	Duration	End Time	Tasks	Status	
insert-into_default_catalog.default_database.flink_user_ml_topics	2020-10-22 18:09:14	1m 6s	-	1 1	RUNNING	



Step 5. Visualize in Superset!



PyFlink in a Nutshell*

- Native SQL integration
- Unified APIs for batch and streaming
- Support for a large set of operations (incl. complex joins, windowing, pattern matching/CEP)



PyFlink in a Nutshell*

- Native SQL integration
- Unified APIs for batch and streaming
- Support for a large set of operations (incl. complex joins, windowing, pattern matching/CEP)

Execution



UDF Support



+ UDTF, UDAF (1.12) + UDTF, UDAF (1.12)



PyFlink in a Nutshell*

- Native SQL integration
- Unified APIs for batch and streaming
- Support for a large set of operations (incl. complex joins, windowing, pattern matching/CEP)

Execution

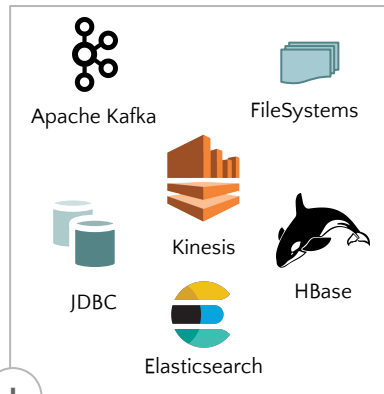


UDF Support

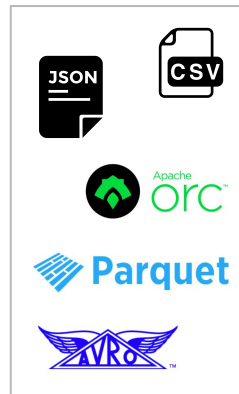


+ UDTF, UDAF (1.12) + UDTF, UDAF (1.12)

Native Connectors



Formats



PyFlink in a Nutshell*

- Native SQL integration
- Unified APIs for batch and streaming
- Support for a large set of operations (incl. complex joins, windowing, pattern matching/CEP)

Execution

Streaming

Batch

UDF Support

Python UDF

Pandas UDF

+ UDTF, UDAF (1.12) + UDTF, UDAF (1.12)

Native Connectors



Apache Kafka



FileSystems



Kinesis



HBase



JDBC



Elasticsearch

+

Formats



JSON



CSV



Apache
orc

Parquet



Avro

ML Library (WIP)

[FLIP-39](#)

Notebooks



Apache Zeppelin





Thank you, Data Science UA!

Follow me on Twitter: @morsapaes

Learn more about Flink: <https://flink.apache.org/>