

The **Data Lakehouse** for Simpler, Faster and More Integrated Data Analysis and Data Science

Carlos Costa

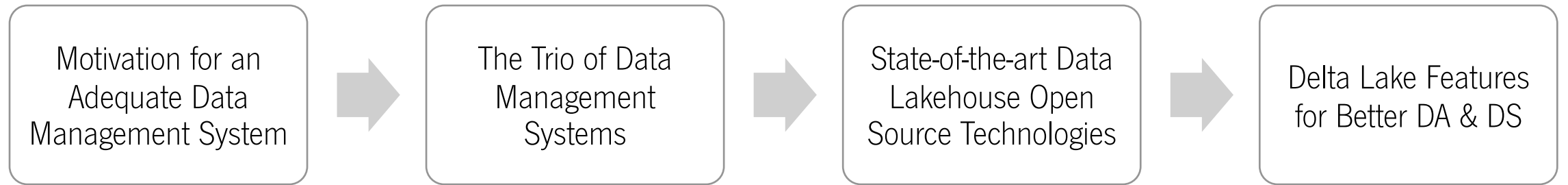
DSPT Webinar

April, 2022



DATA SCIENCE PORTUGAL

Agenda



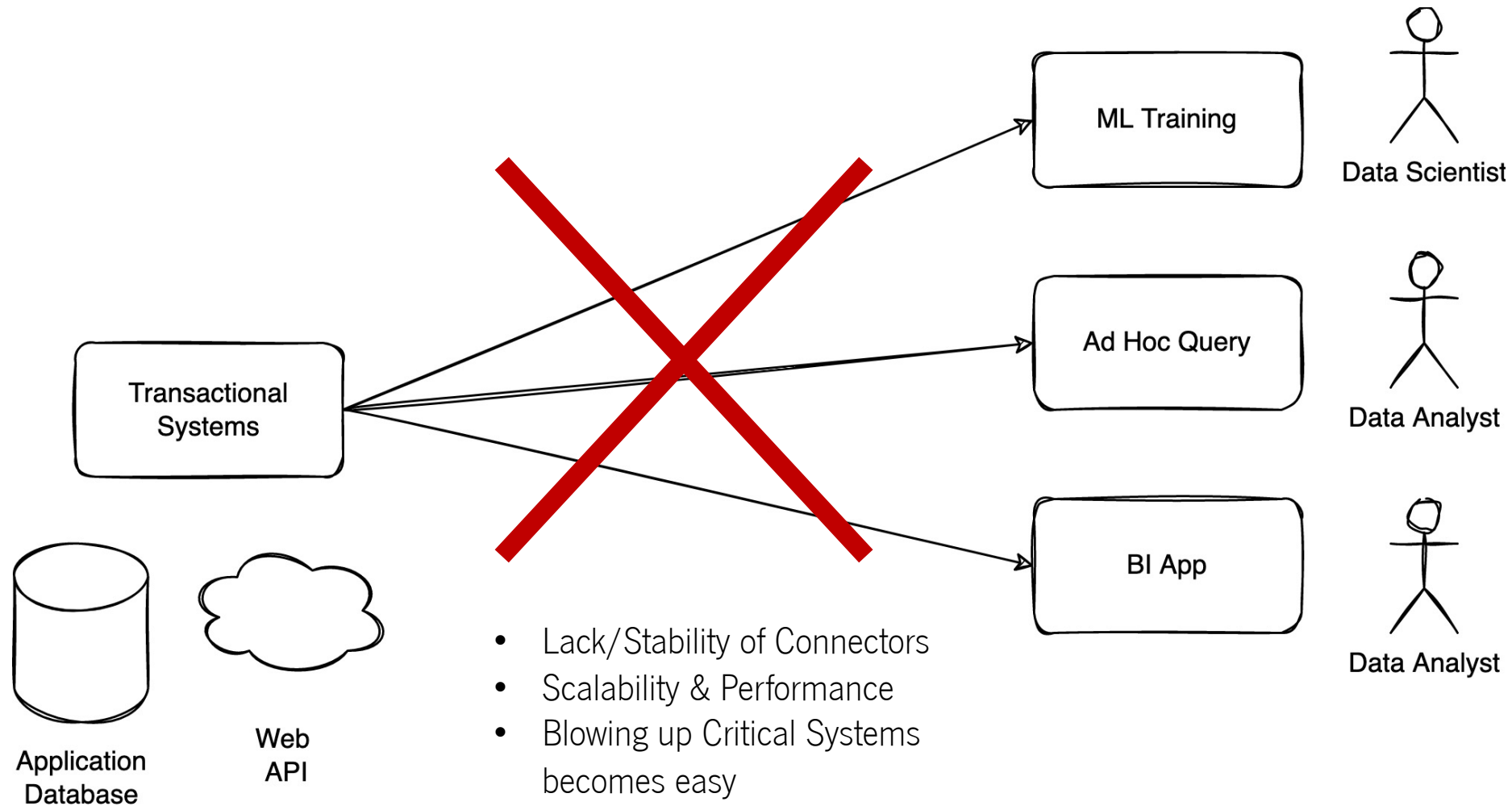
Disclaimer

These are my opinions and interpretations of the state of the art... does not mean I'm wrong or right! Use it as food for thought... or ammo for discussion :P

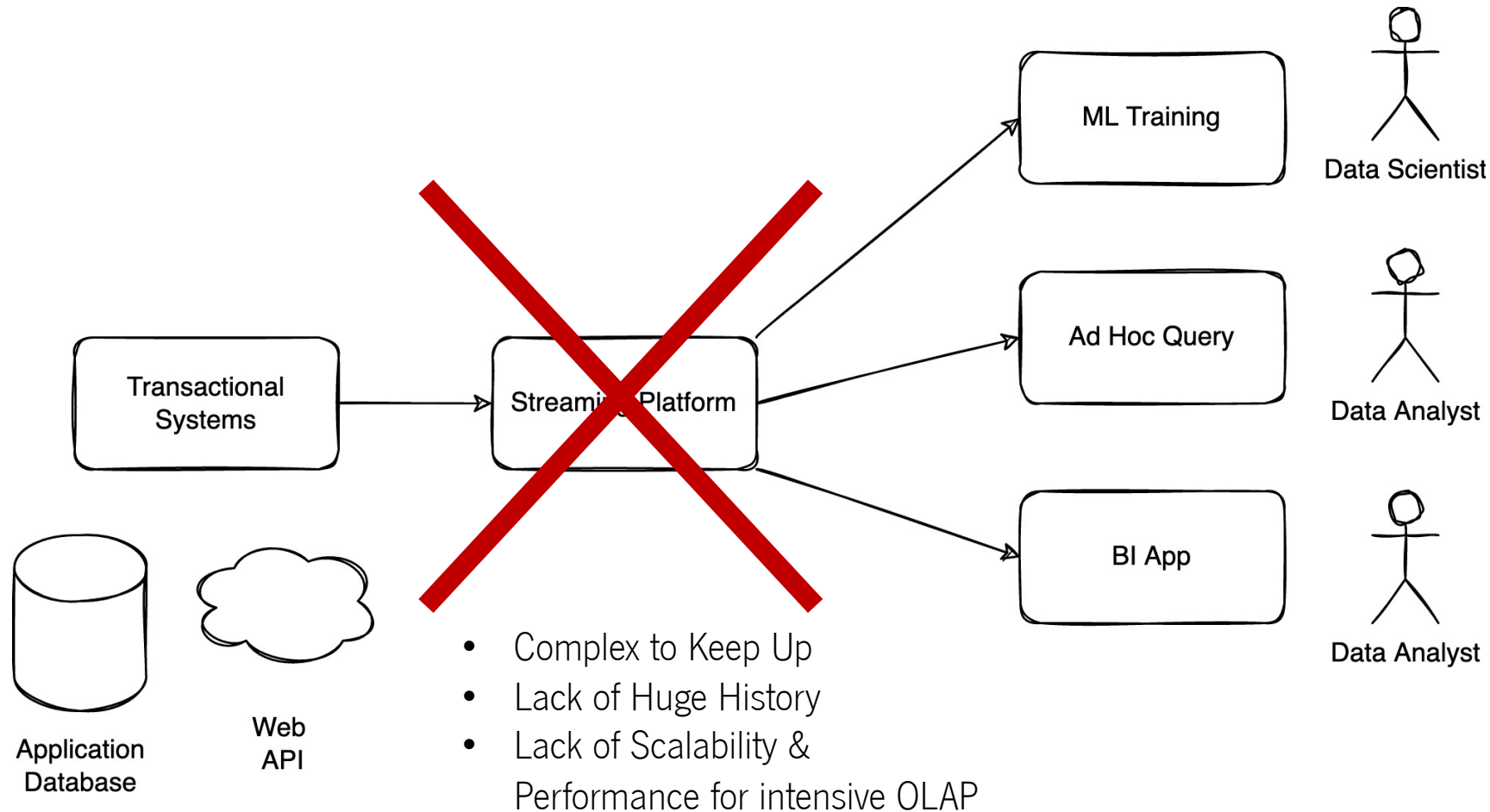


Motivation for an Adequate Data Management System

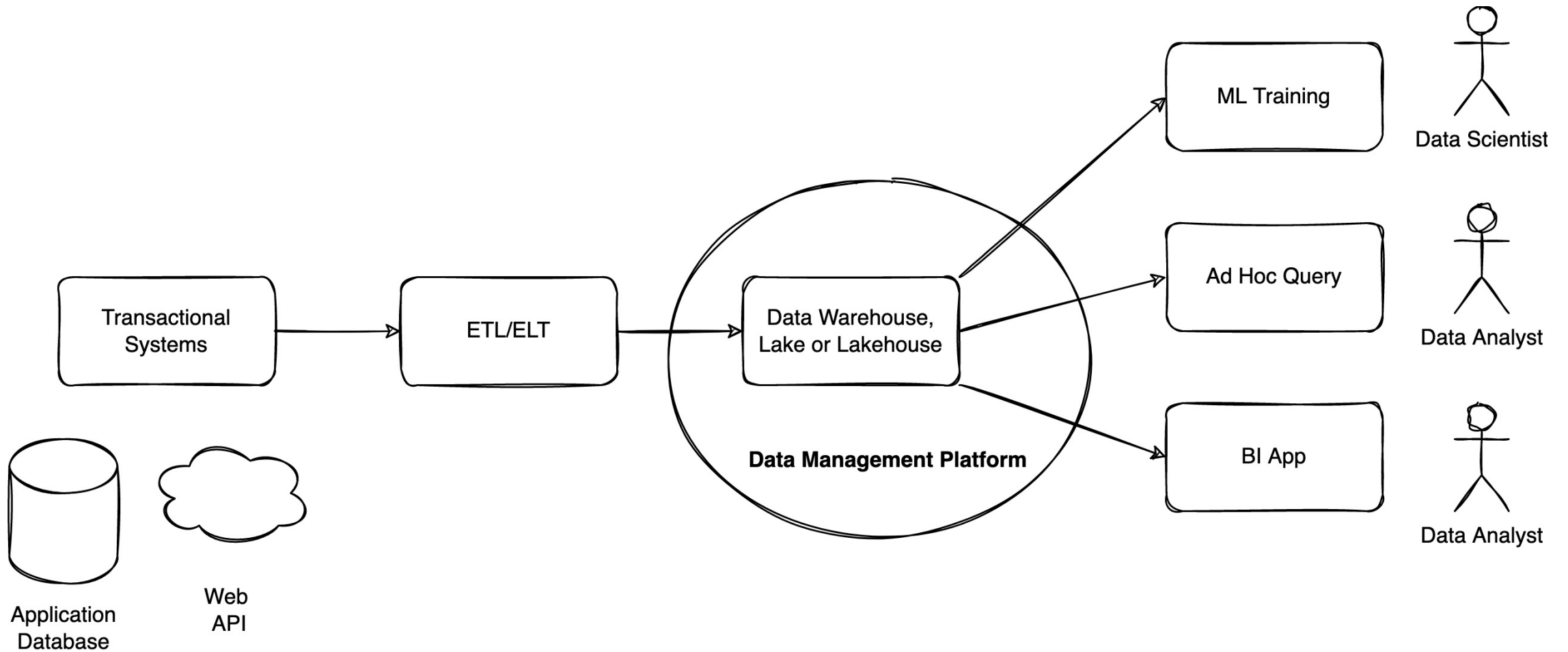
"Give me access to the Source" Pattern

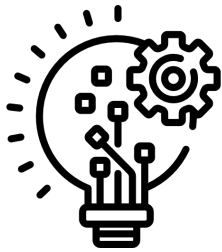


The “Streaming, Streaming, Streaming!!!” Pattern



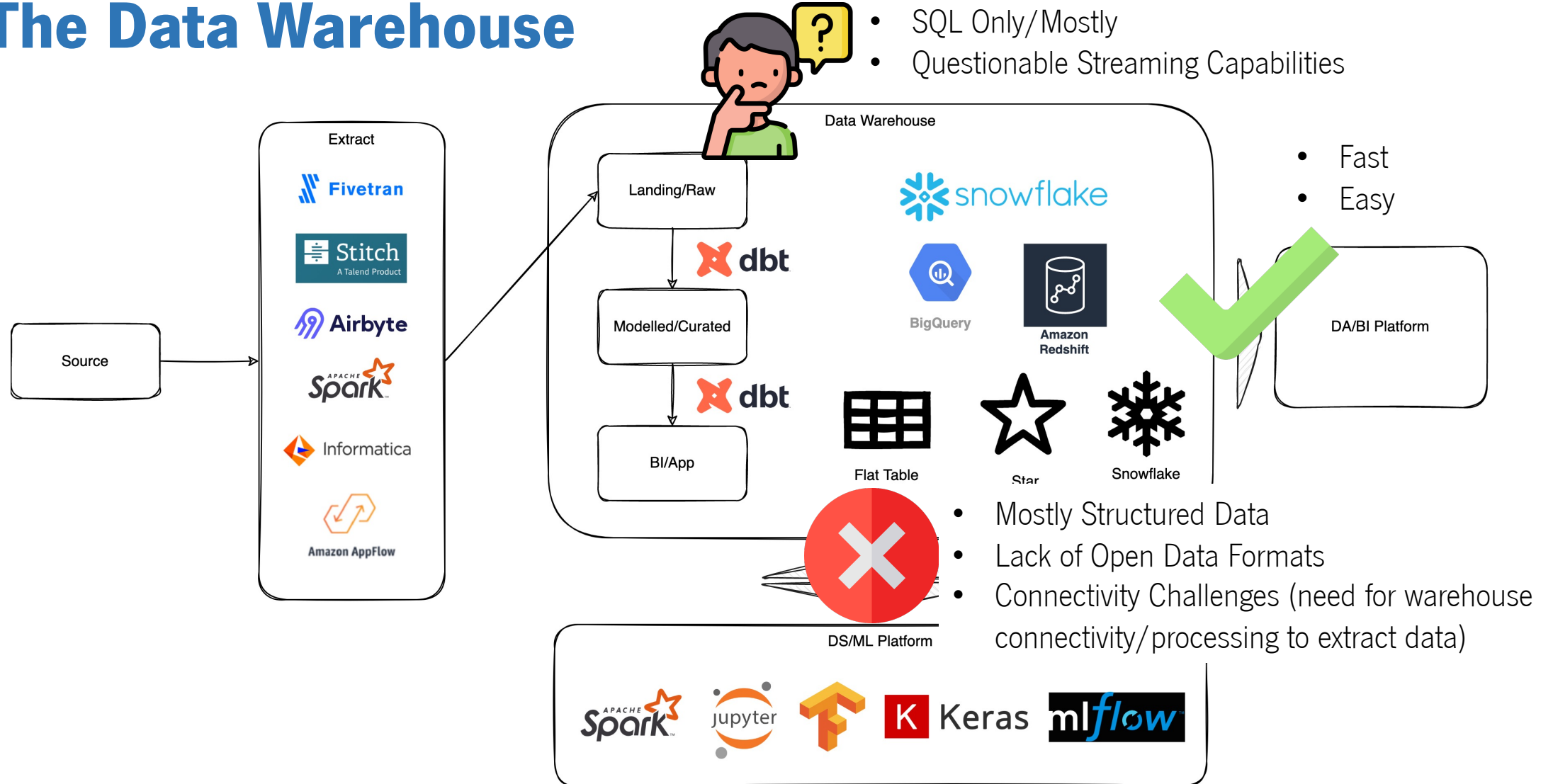
The Data Management Platform



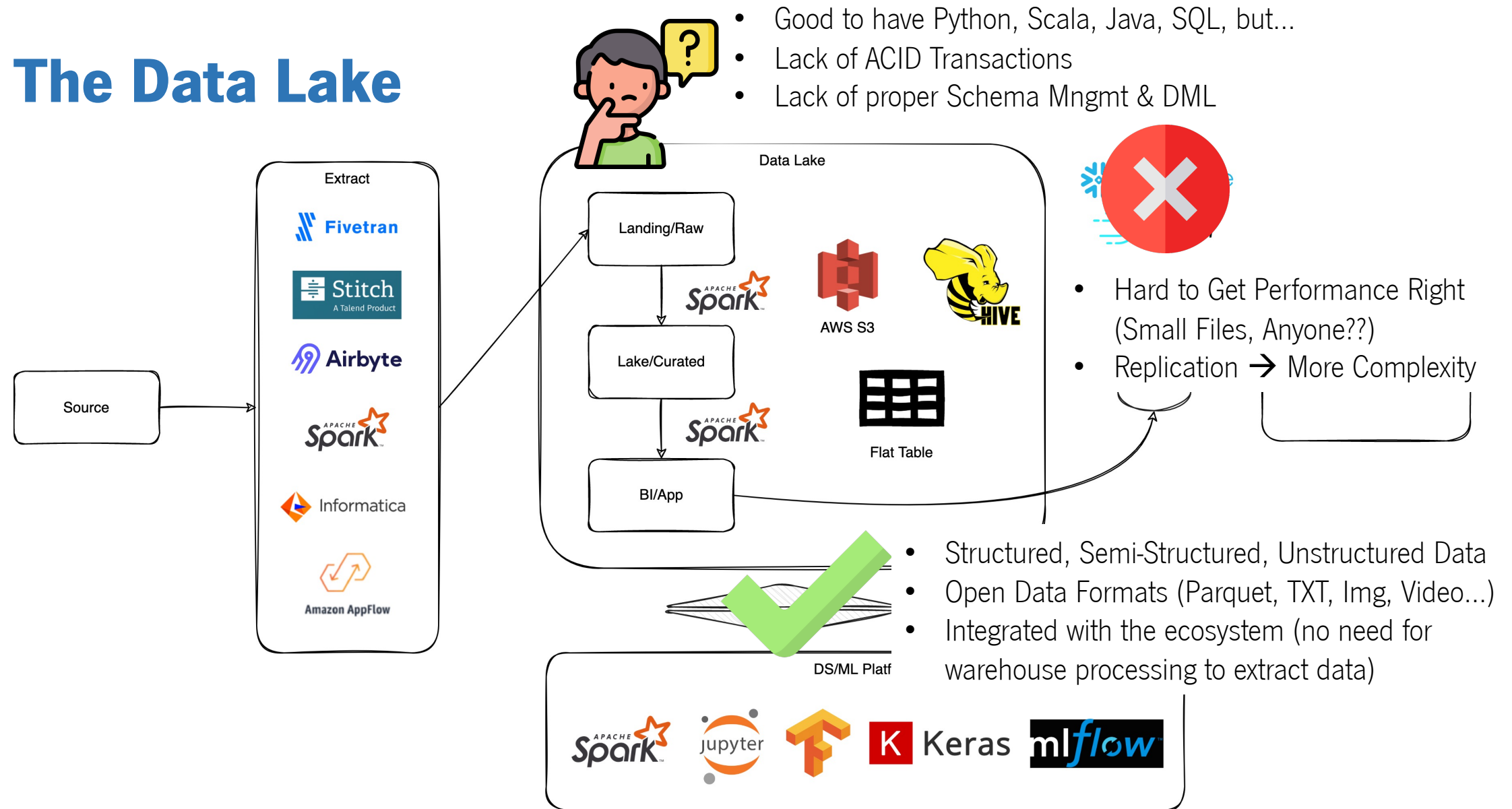


The Trio of Data Management Platforms

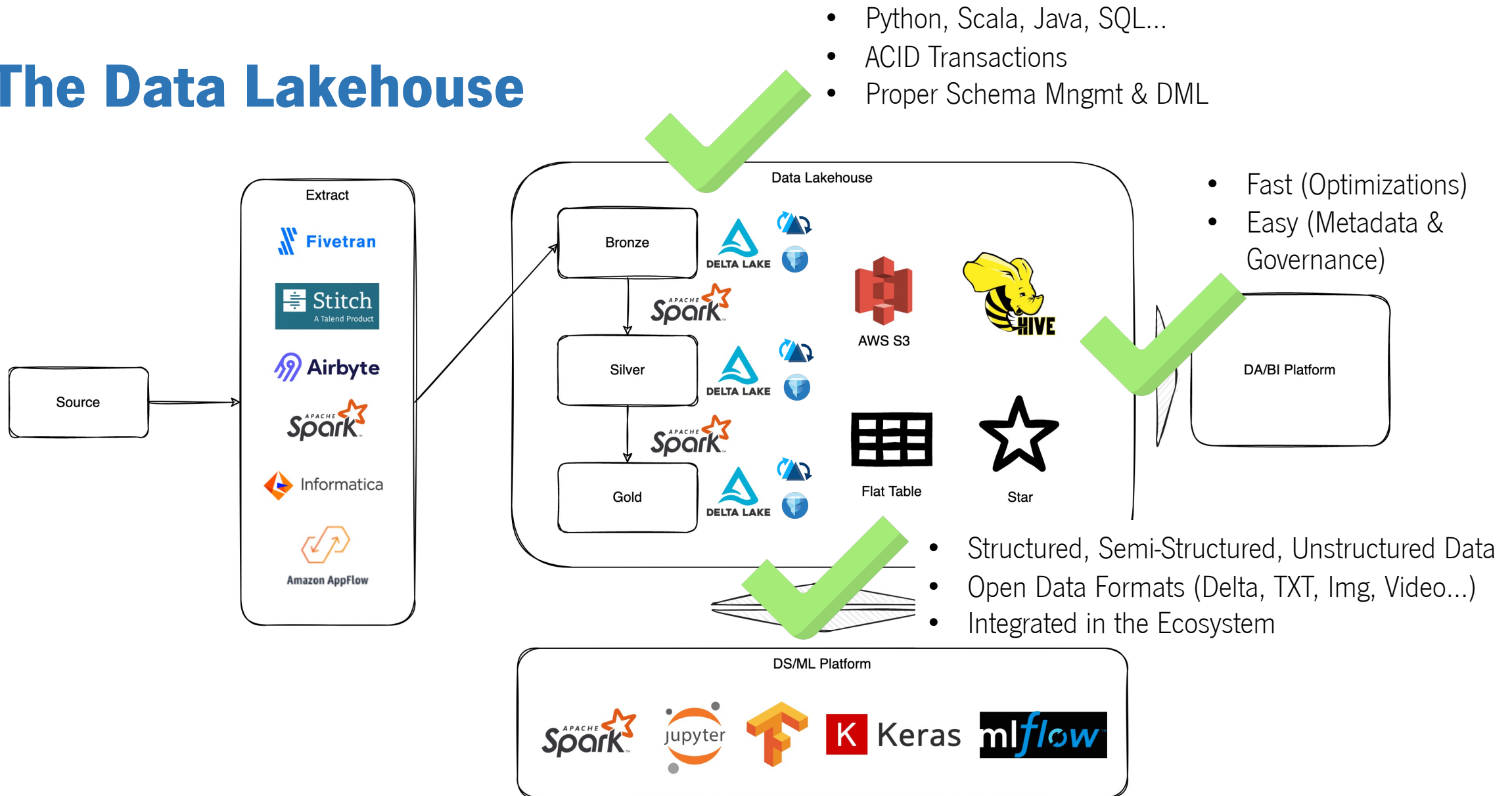
The Data Warehouse



The Data Lake

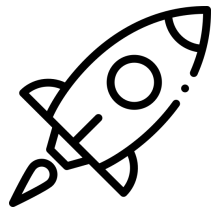


The Data Lakehouse



Wrapping it Up

	<i>Data Warehouse</i>	<i>Data Lake</i>	<i>Data Lakehouse</i>
The Good	<ul style="list-style-type: none">• BI/OLAP Speed• Easiness of Management and Access	<ul style="list-style-type: none">• Handling All Volumes, Varieties and Velocities of Data• Several Programming Languages• Good Data Science Integration	<ul style="list-style-type: none">• All the Benefits of the Data Warehouse and Data Lakehouse
The Bad	<ul style="list-style-type: none">• Unstructured Data Handling• Streaming Capabilities are Questionable• SQL Only/Mostly	<ul style="list-style-type: none">• Lack of ACID• Lack of proper DML• Lack of proper Schema Mngmt.	<ul style="list-style-type: none">• <i>(Not so bad)</i> Performance is influenced by the querying processing technology (e.g., Spark, Presto, etc....)
The Ugly	<ul style="list-style-type: none">• An “Alien” to the Data Science Ecosystem	<ul style="list-style-type: none">• Replication into BI/OLAP oriented Stores (e.g., DW)• Complexity of the Replication	<ul style="list-style-type: none">• ???



State-of-the-art Data Lakehouse Open-Source Technologies



- Open-Sourced by **Uber**
- *“Hadoop Upserts Deletes and Incrementals”*
- “Hudi brings transactions, record-level updates/deletes and change streams to data lakes”
- *Can use Parquet underneath*
- <https://hudi.apache.org>



- Open-Sourced by **Netflix**
- *“The open table format for analytic datasets”*
- “Iceberg is a high-performance format for huge analytic tables. Iceberg brings the reliability and simplicity of SQL tables to big data...”
- *Can use Parquet underneath*
- <https://iceberg.apache.org>



- Open-Sourced by **Databricks**
- *“Build Lakehouses with Delta Lake”*
- *“Delta Lake is an open-source storage framework that enables building a Lakehouse architecture with compute engines...”*
- *Parquet only underneath*
- <https://delta.io>



Delta Lake Features For Better DA & DS

ACID Transactions

- Readers see a consistent snapshot view of the data since the start of the query...
 - Even if the table was modified since then
- Multiple Writers can write to the same table partition concurrently and:
 - Always see a consistent snapshot view of the data
 - A serial order for these writes is ensured (isolation)
- General Steps:
 - Read latest state of the table to identify the files that need to be modified
 - Stage the new files to write
 - Validate if there are no conflicts and commit.



my_table/

_delta_log/

000000.json *Add 2.parquet*
000001.json *Remove 1.parquet*
...
000010.checkpoint.parquet

sales_date=2022-04-18/

sales_date=2022-04-19/

sales_date=2022-04-20/

Sources:

<https://docs.delta.io/latest/concurrency-control.html>

<https://databricks.com/blog/2019/08/21/diving-into-delta-lake-unpacking-the-transaction-log.html>

Comprehensive Data Manipulation Language (DML)

```
DELETE FROM dummy_table WHERE date < '2022-04-25'  
DELETE FROM delta.`/foo/bar` WHERE date < '2022-04-25'
```

```
UPDATE dummy_table SET month = 'January' WHERE month = 'Jan';
```

```
MERGE INTO dummy_table  
USING new_data  
ON dummy_table.id = new_data.id  
WHEN MATCHED THEN  
    UPDATE SET dummy_table.id = new_data.id, dummy_table.firstName = new_data.firstName, ...  
WHEN NOT MATCHED THEN  
    INSERT (...) VALUES (...)
```

Source: <https://docs.delta.io/latest/>

Schema Enforcement & Evolution

Fail if schema does not match ... **OR**

```
> df.write...option("mergeSchema", "true")
```

OR (to set it across the Spark Session)

```
> spark.conf.set("spark.databricks.delta.schema.autoMerge.enabled", "true")
```

OR (for overwrites)

```
> df.write.option("overwriteSchema", "true")
```

Source: <https://docs.delta.io/latest/>

Constraints

NOT NULL

```
> CREATE TABLE dummy_table (  
    dummy_id INT NOT NULL,  
    ...  
)  
> ALTER TABLE table CHANGE COLUMN dummy_name DROP NOT NULL;
```

CHECK CONSTRAINTS

```
> ALTER TABLE dummy_table ADD CONSTRAINT c CHECK (date > '2022-04-25')
```

Source: <https://docs.delta.io/latest/>

Unified Batch & Streaming

Read the Stream

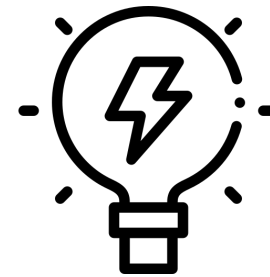
```
events = spark.readStream.format("delta").load("/foo/bar")
```

Transform

...

Write the Stream

```
events
  .writeStream
  .format("delta")
  .outputMode("append")
  .option("checkpointLocation", ".../_checkpoints/...")
  .toTable("table")
```



Tip of the day!

Batch = Streaming

```
.trigger(Trigger.Once())
```

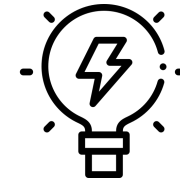
Source: <https://docs.delta.io/latest/>

Audit History

DESCRIBE HISTORY delta.`/foo/bar/` -- full history
DESCRIBE HISTORY `/data/events/` **LIMIT** 1 -- last version only
DESCRIBE HISTORY eventsTable -- use table name

OR

```
from delta.tables import *  
delta_table = DeltaTable.forPath(spark, '/foo/bar')  
all_history_df = deltaTable.history()  
last_version_df = deltaTable.history(1)
```



Git Log for your Data!

version	timestamp	userId	userName	operation	operationParameters	job	clusterId	operationMetrics
5	2019-07-29 14:07:47	null	null	DELETE	[predicate -> ["(... null		null	[numTotalRows -> ...
4	2019-07-29 14:07:41	null	null	UPDATE	[predicate -> (id... null		null	[numTotalRows -> ...
3	2019-07-29 14:07:29	null	null	DELETE	[predicate -> ["(... null		null	[numTotalRows -> ...
2	2019-07-29 14:06:56	null	null	UPDATE	[predicate -> (id... null		null	[numTotalRows -> ...
1	2019-07-29 14:04:31	null	null	DELETE	[predicate -> ["(... null		null	[numTotalRows -> ...

Source: <https://docs.delta.io/latest/>

Time Travel

```
df1 = spark.read.format("delta").option("timestampAsOf", timestamp_string)...\ndf2 = spark.read.format("delta").option("versionAsOf", version)...
```

Use it for Fixing Problems (re-insert deleted user events)

```
df = spark.read.format("delta") \> .option("timestampAsOf", yesterday) \> .load("/foo/bar") \> df.where("user = 'dummy_user'") \> .write.format("delta") \> .mode("append") \> .save("/foo/bar")
```



Git Reset for your Data!

Source: <https://docs.delta.io/latest/>

Time Travel for Analytics

```
df = spark.read.format("delta") \
    .option("timestampAsOf", last_week) \
    .load("/foo/bar")
```

```
last_week_count = df.select("user") \
    .distinct() \
    .count()
```

```
count = spark.read.format("delta") \
    .load("/foo/bar") \
    .select("user") \
    .distinct() \
    .count()
```

```
new_customers_count = count - last_week_count
```



Feature Store

(Customer, Product,
Product Images Metadata,
Stores...)



Reproducible Experiments



Read Changes

between 2 versions
with *Except*

Source: <https://docs.delta.io/latest/>

Compactation & Other Optimizations

Optimize the layout of data (compact small files into larger ones)

```
OPTIMIZE table
```

```
OPTIMIZE delta.`/foo/bar`
```

```
-- only optimize a subset of the table
```

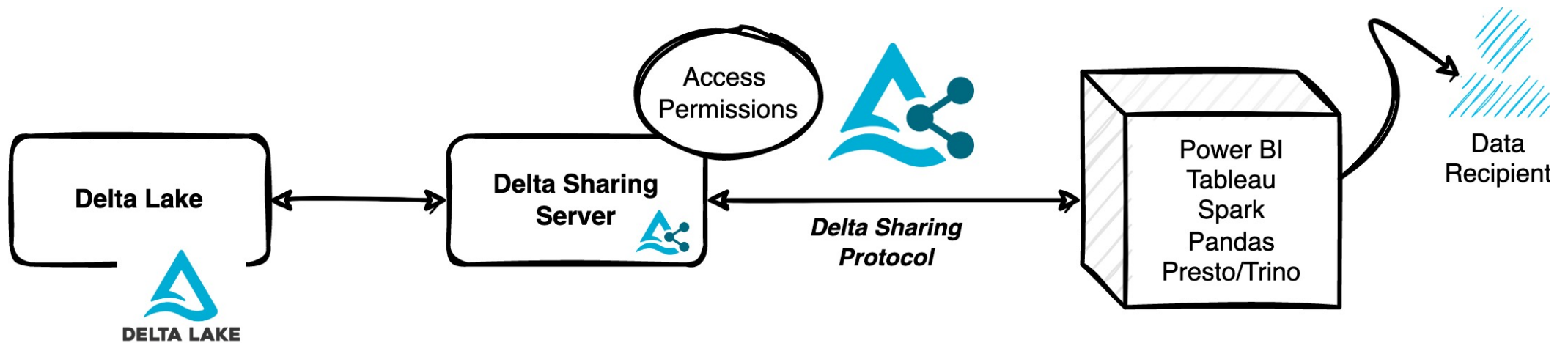
```
OPTIMIZE table WHERE date >= '2022-04-25'
```

Data Skipping

- Used at query time to **skip unnecessary data**
- **Min** and **Max** values for each column **collected automatically** when writing
 - for x configurable number of columns
- Best Performance when used in combination with “**ordering before writing**” pattern

Source: <https://docs.delta.io/latest/>

Delta Sharing



```
> spark.read.format("deltaSharing").load(table_path)
```

```
> delta_sharing.load_as_pandas(...)
```

Source: <https://delta.io/sharing/>

Thank You!
Ask me anything 😊

Carlos Costa

DSPT Webinar

April, 2022



DATA SCIENCE PORTUGAL