# Coral & Transport UDFs

## Building Blocks of a Postmodern Data Warehouse
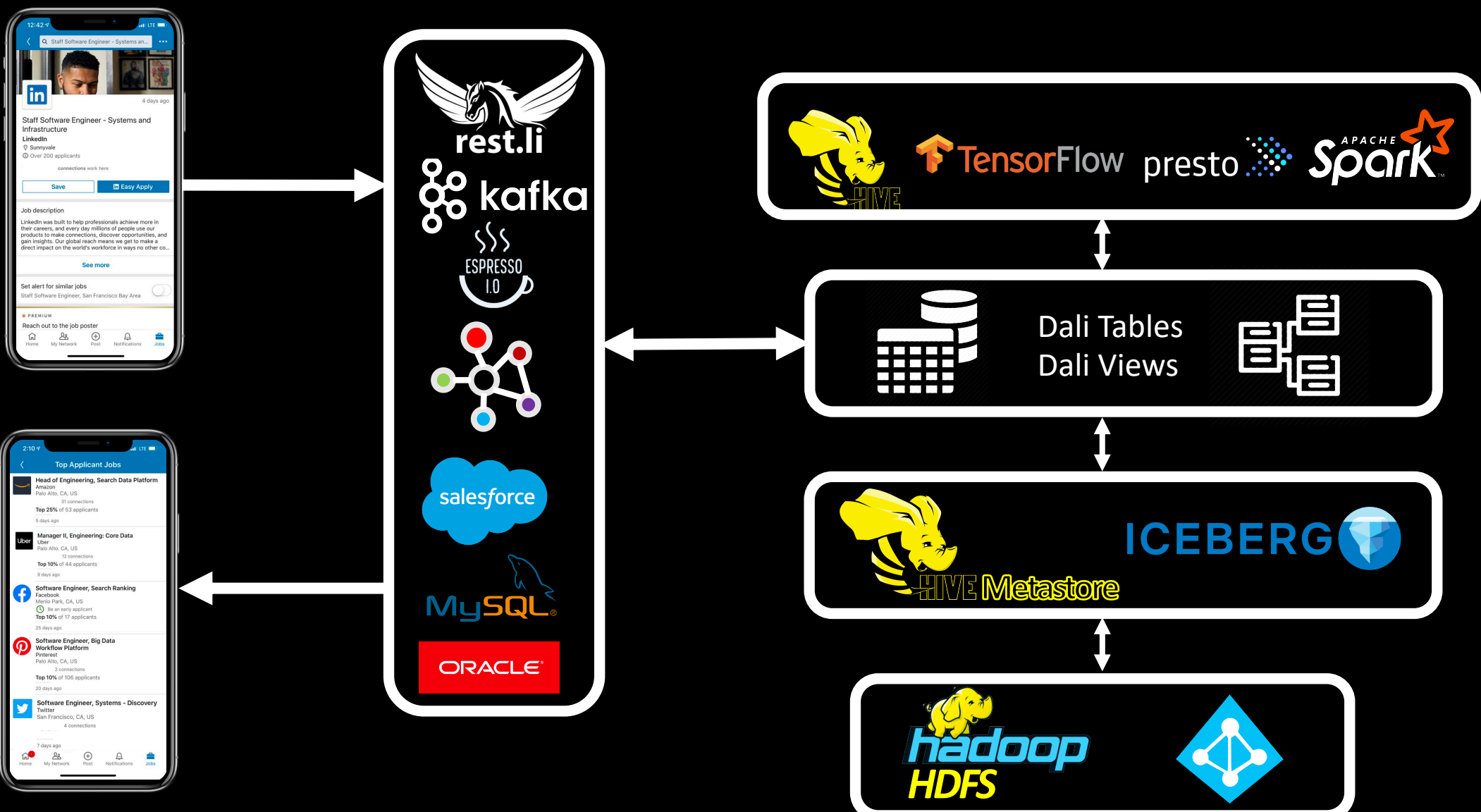
Walaa Eldin Moustafa
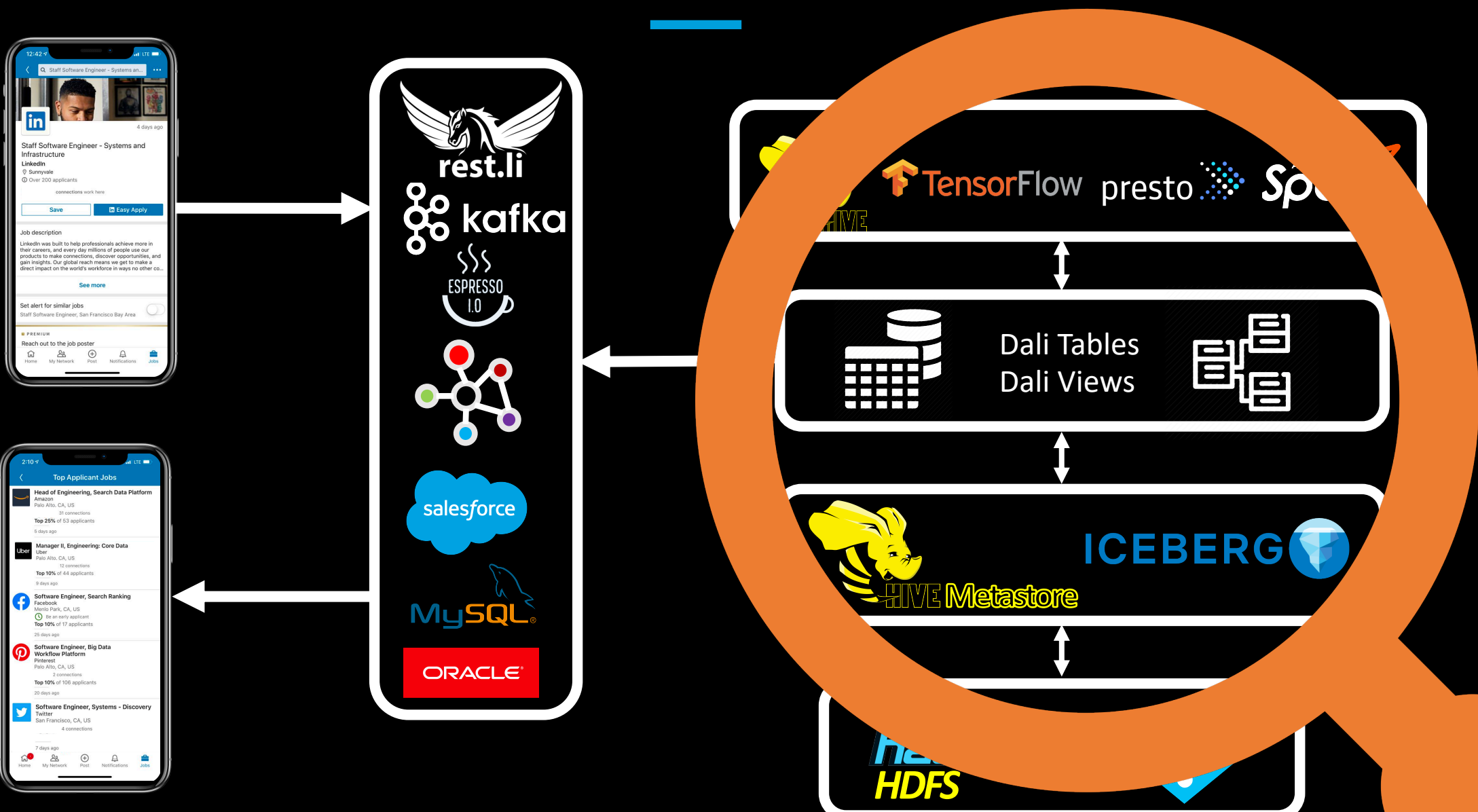
Staff Software Engineer @ **Linked in**
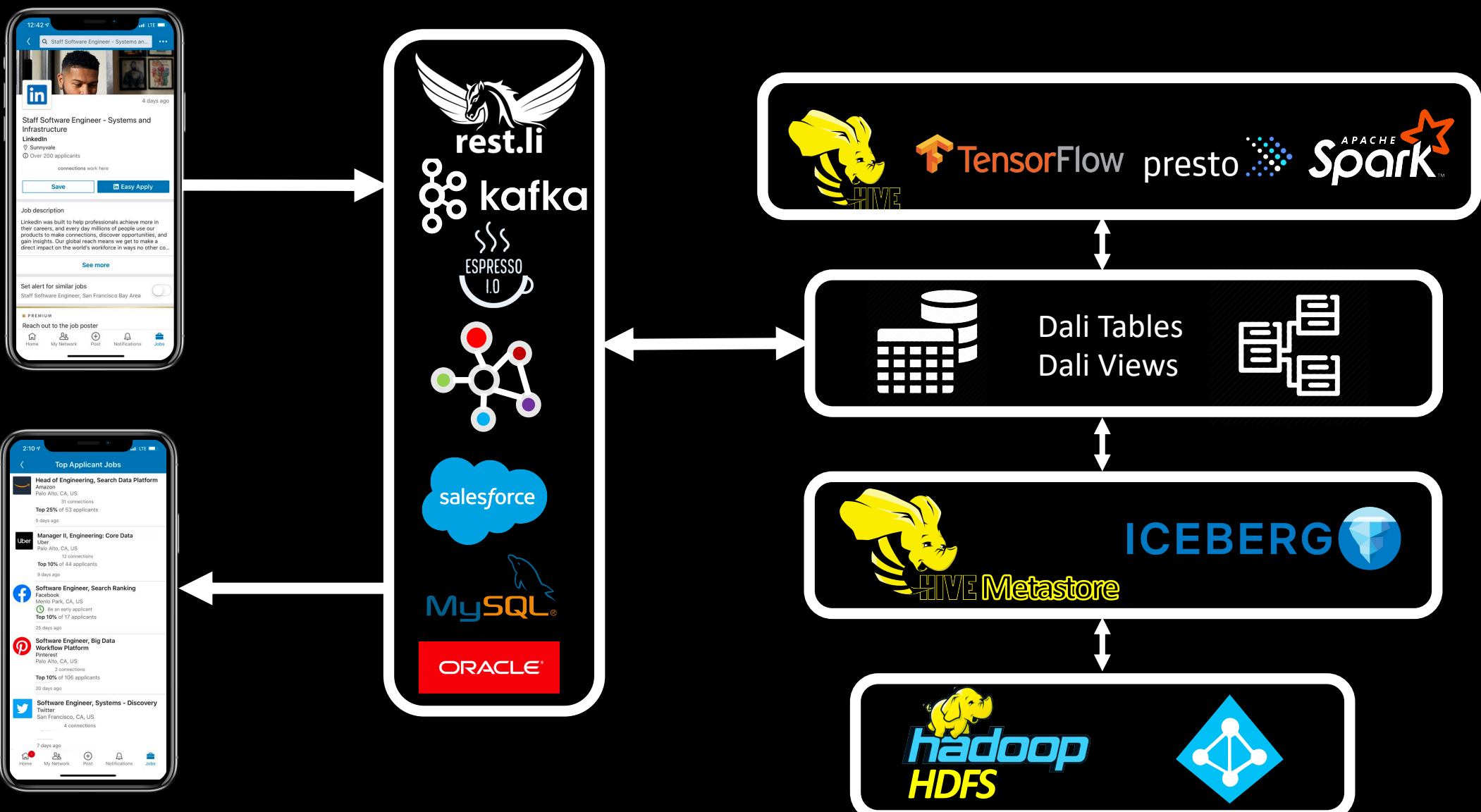
# The lifecycle of a data application at LinkedIn [*Simplified*]

# The lifecycle of a data application at LinkedIn [*Simplified*]
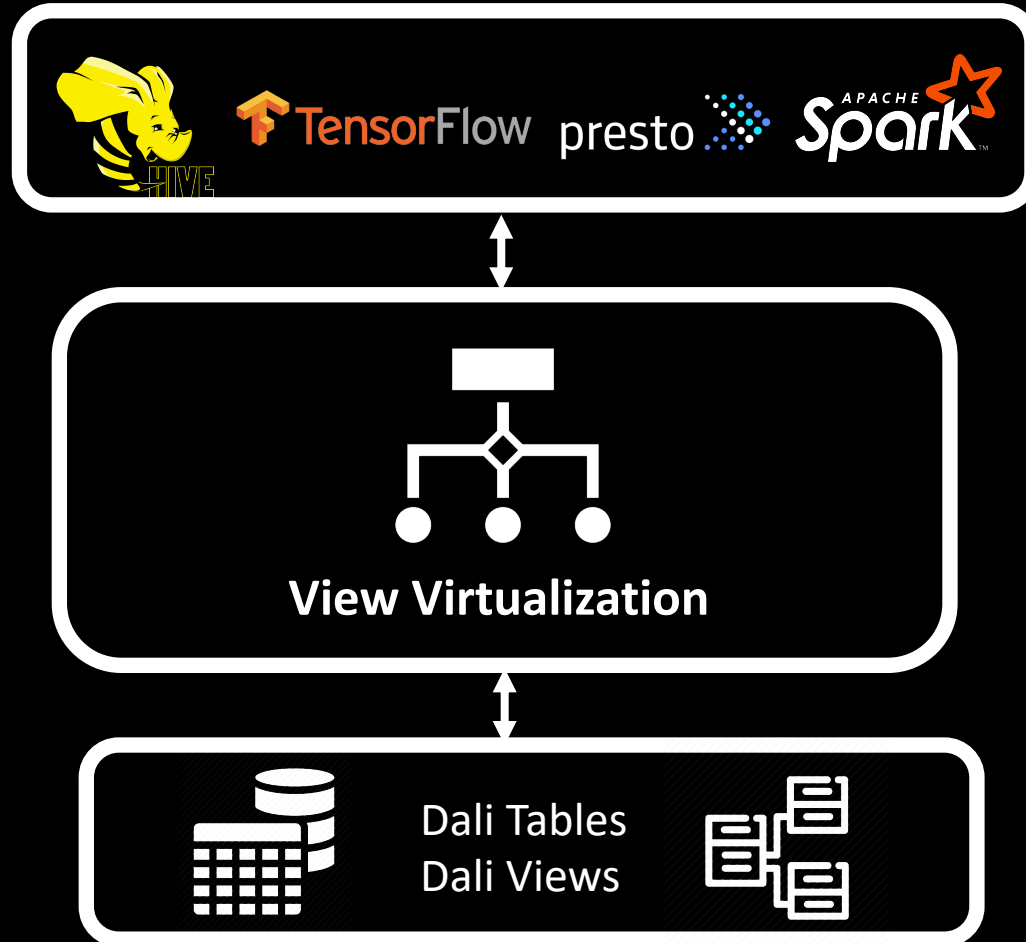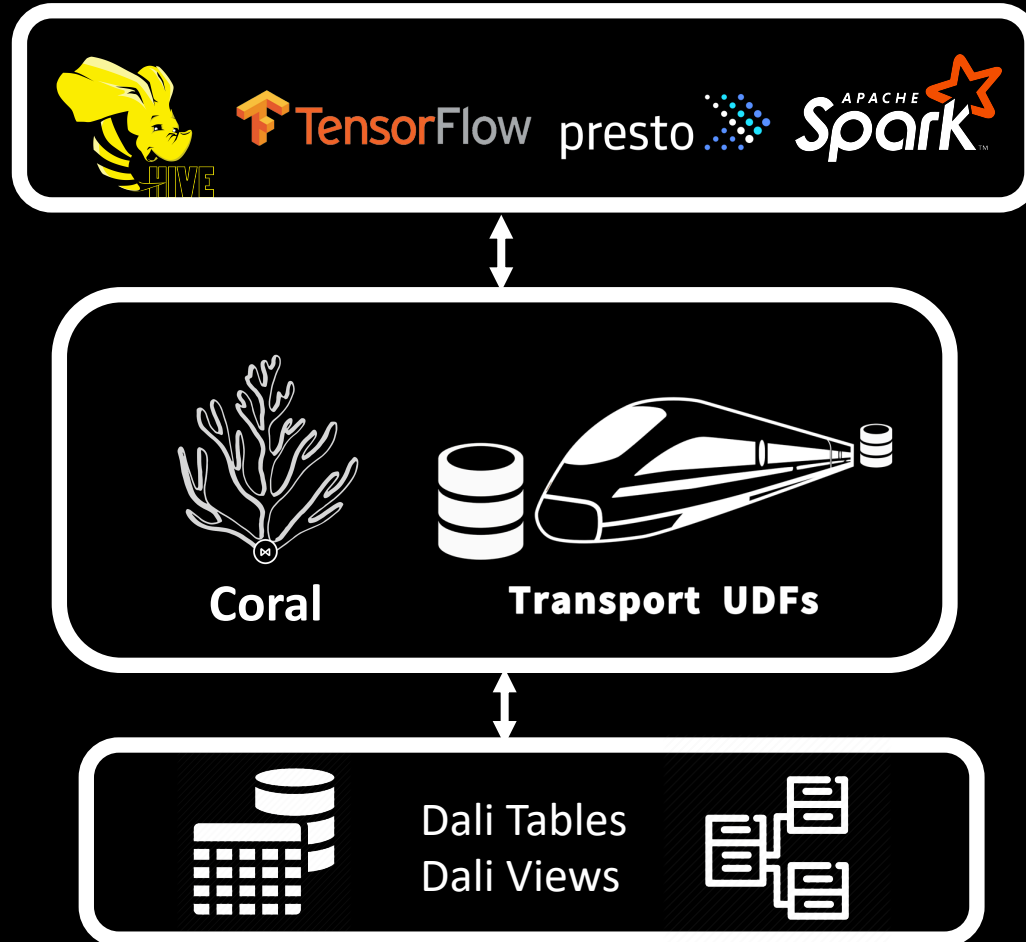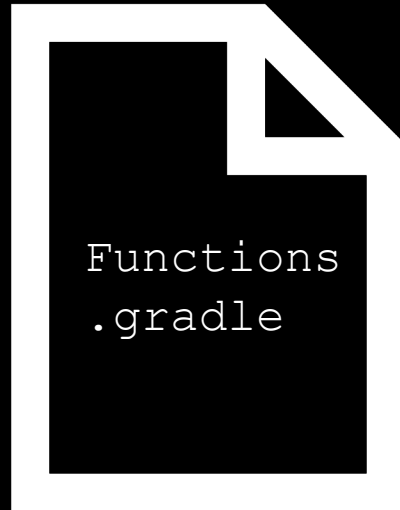
# The lifecycle of a data application at LinkedIn *[Simplified]*

# View Virtualization

# Dali Views



git

HIVE Metastore

- ✔ Code Review
- ✔ Test
- ✔ Publish
- ✔ Deploy

View.SQL

Functions
.gradle

```
SELECT my_udf(c)
FROM R JOIN S JOIN T
WHERE date > today() - 5
AND date <= today()

TBLPROPERTIES(
'functions' = 'my_udf:
com.linkedin.MyUDF',
'dependencies' =
'group:artifact:0.0.1')
```
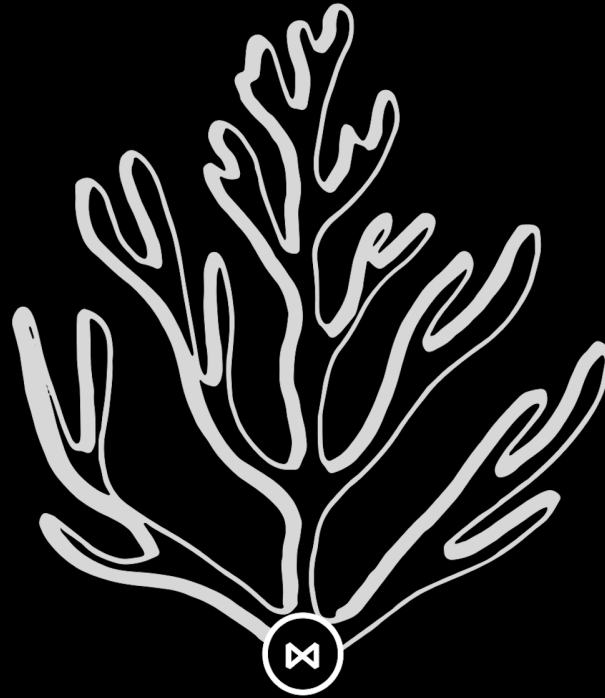
**SQL + UDFs**

# Coral

# Coral

```
SELECT my_udf(c)
FROM R JOIN S JOIN T
WHERE date > today() - 5
AND date <= today()

TBLPROPERTIES(
'functions' = 'my_udf',
'dependencies' =
'group:artifact:0.0.1')
```

coral-hive →



Coral IR

# Coral

# Coral-Presto



Coral IR

Presto Rel

```
SELECT presto_udf(c + 1)
FROM R JOIN S JOIN T
WHERE date > today() - 5
AND date <= today()
```

# Coral-Pig



Coral IR

Pig Rel

```
R = LOAD "R"
S = LOAD "S"
T = LOAD "T"
RS = R JOIN S
RSF = FILTER RS BY ...
RSFT = RSF JOIN T
```

# What does the future look like?

# Flexible Language Interface

user {
  name
  address{}
}

**GraphQL**

A = LOAD
B = LOAD
C = A
JOIN B

**Pig Latin**

d = load()
.filter()
.map()
.count()

**Dataflow**

TwoHop
(X,Y) :-
E(X,Z),
E(Z,Y).

**Datalog**

# Workload Analytics

Materialized View Selection

# Data Governance

# PII Lineage and Derivation

# Format-specific View Schema Derivation

# Transport UDFs: *Translatable Portable* UDFs

# IR

- **SQL has pretty well-understood IR: Relational Algebra**

- **Standard Operators**

  - **Scan, Filter, Project, Join, Group By, etc**

- **UDFs**

  - **Opaque**

  - **Use imperative language**

  - **Not portable or translatable**

# UDF Denormalization

## Duplication

Multiple versions of the same UDF. Not clear which is the source of truth.

## Inconsistency

Duplicate implementations can diverge causing data inconsistency

## Low Productivity

Developers need to learn multiple APIs, implement same logic multiple times.

## Low Performance

In some cases, use tuple conversion adapters to enable portability.

# UDF APIs

- **API Complexity**
  - APIs expose low-level details of engines
  - Data types may not intuitvely map to SQL type-system

- **API Disparity**
  - APIs differ in what to expect from developer
  - APIs differ in features they can provide

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }

  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }

  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }
    return map;
  }
}
```

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }

  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }

  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }
    return map;
  }
}
```

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }

  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }

  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }
    return map;
  }
}
```

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }

  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }

  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }
    return map;
  }
}
```

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }

  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }

  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }
    return map;
  }
}
```

# Transport UDFs

```java
public class MapFromTwoArrays
    extends StdUDF2<StdArray, StdArray, StdMap> {

  @Override
  public List<String> getInputParameterSignatures() {
    return ImmutableList.of(
        "array(K)",
        "array(V)"
    );
  }


  @Override
  public String getOutputParameterSignature() {
    return "map(K,V)";
  }


  @Override
  public StdMap eval(StdArray a1, StdArray a2) {
    StdMap map = getStdFactory().createMap(
        getOutputParameterSignature());
    for (int i = 0; i < a1.size(); i++) {
      map.put(a1.get(i), a2.get(i));
    }

    return map;
  }
}
```
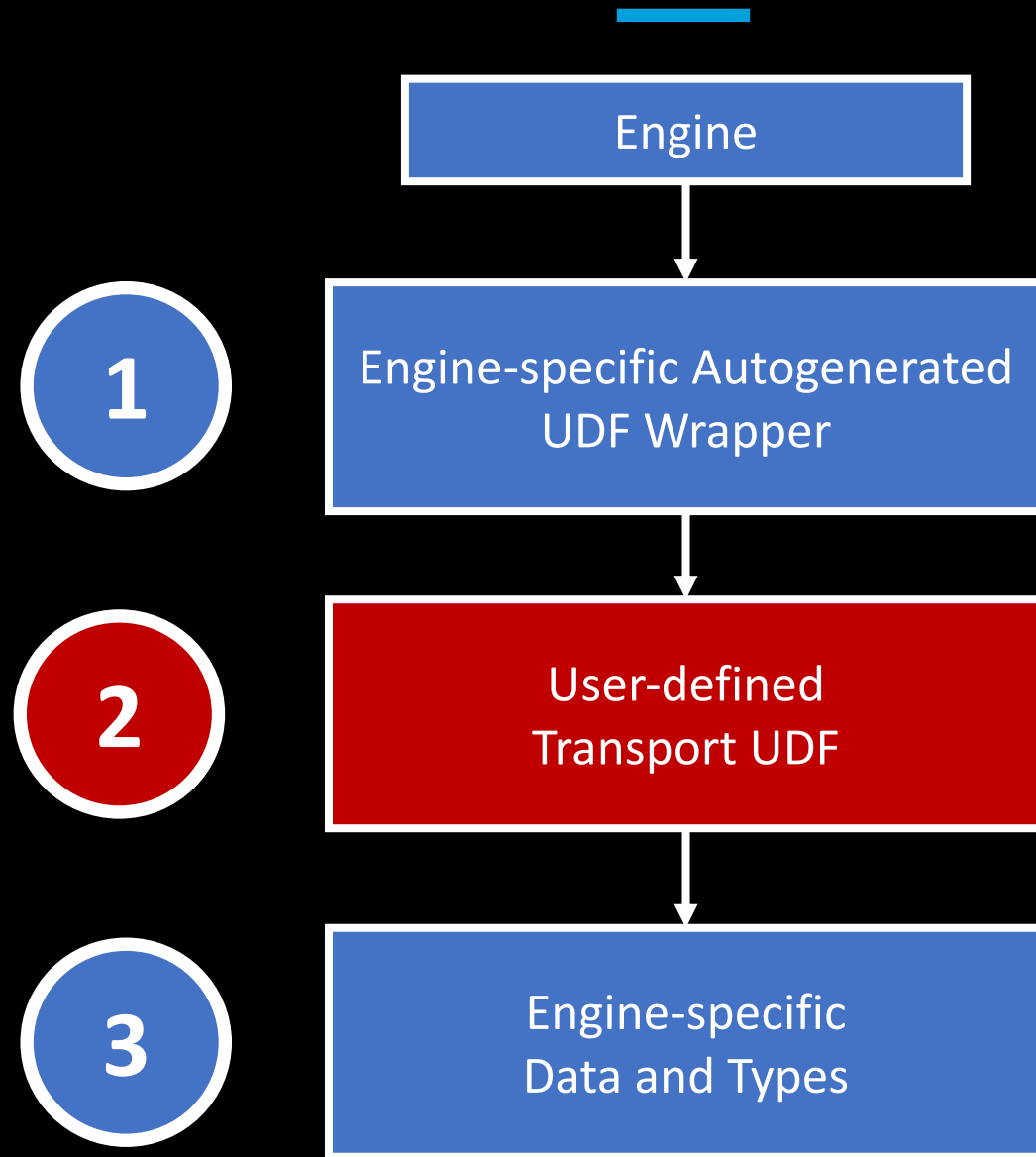
# Then What?

---

```
> gradle build


> ls build/my-udfs/libs

    my-udfs-presto.jar

    my-udfs-hive.jar

    my-udfs-spark.jar
```

# Architecture

Engine

**1** Engine-specific Autogenerated
UDF Wrapper

**2** User-defined
Transport UDF

**3** Engine-specific
Data and Types

# Contributors