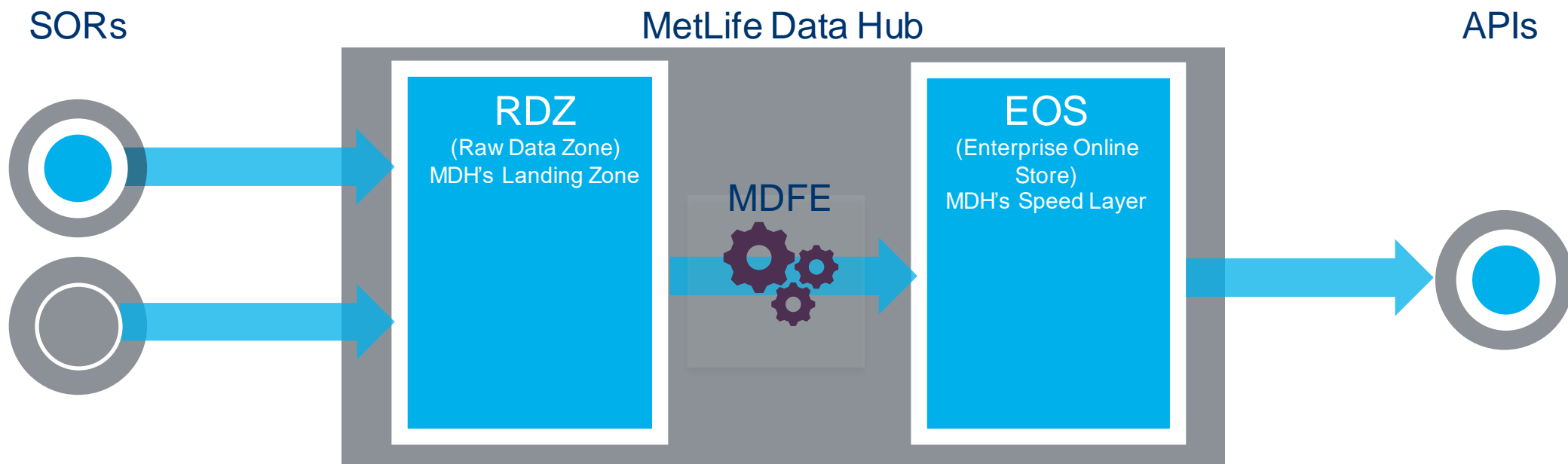# MDFE

MetLife Data Flow Engine

## Engine that moves data within MetLife Data Hub

# MDFE Overview

MetLife Data Flow Engine is a configuration and metadata driven framework that moves data within the MetLife Data Hub environment. Current version of MDFE is built using custom spark code.

# What does MDFE solve for?

## Key Constraints

- SORs can communicate data async
- Changes can come in out of sequence
- Latency can change in future
- SOR object hierarchy and EOS object hierarchy can be different

## MDFE Solves for

- Data move from RDZ (file, messages) to EOS
- Validation rules
- Transformation to conform to target model
    - Attribute level transformation
    - Relation between objects / entities
    - Aggregate multiple components into an entity
- Relationship management
    - Direct relationships (foreign keys)
    - Inferred – rule based (post processing logic)
    - Orphaned relationships (dirty records)
- Data flow tracking
    - Meta data about source at record / vertex level
    - Process Management
        - Run logs, metrics, rows / vertexes tied to the process

# MDFE Coding vs. Config
## *Common operations through config*

**Configuration driven**

- Attribute level Validations
- Attribute mapping
- Attribute Transformations
- Key generation
- Relationships
- Composite objects

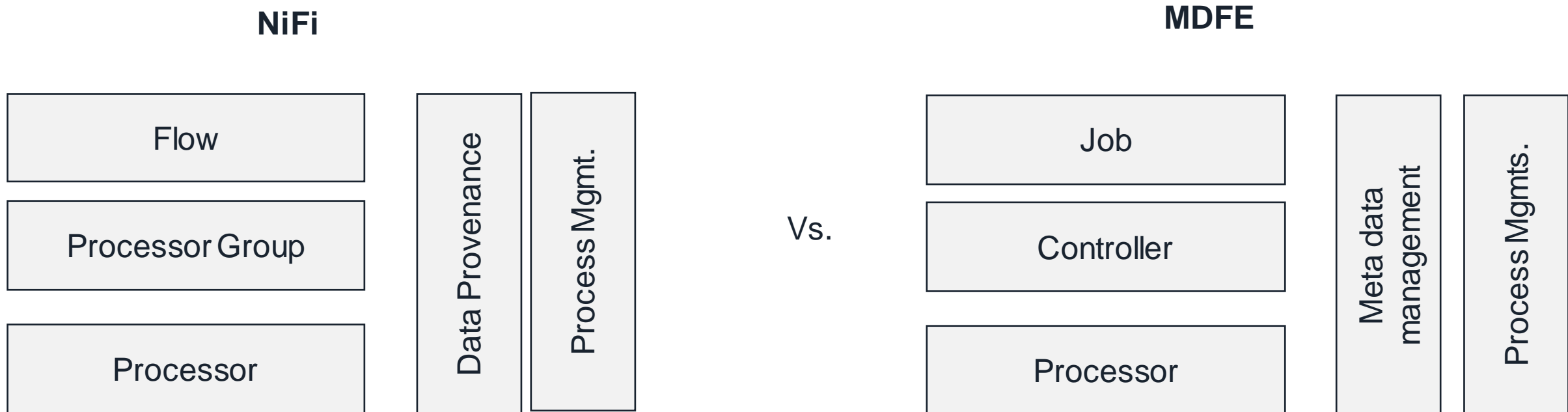Externalization to feed into data lineage analysis
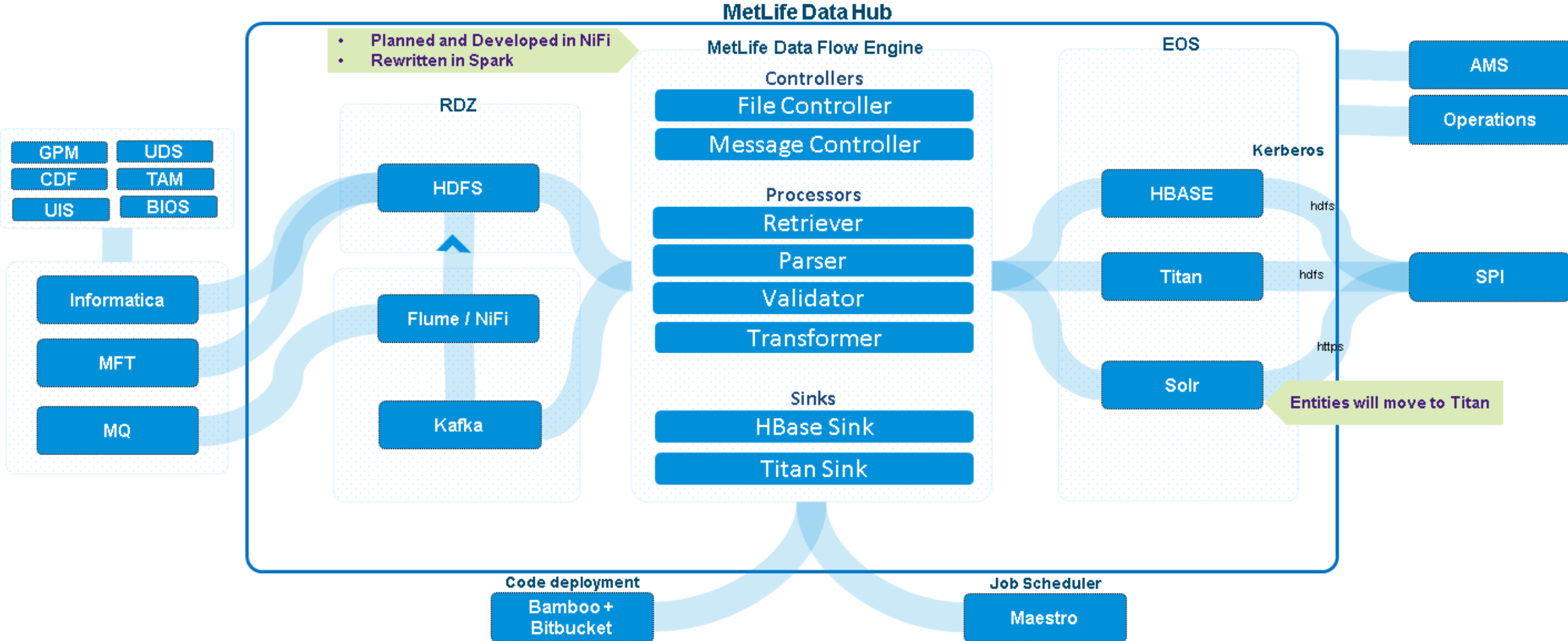
Vs.

**Code Driven**

- Complex aggregations
  - Pluggable rules / post processing

- SOR file level nuances
  - UIS delta file vs. one time differences
  - Generating inferred plan versions for UDS

Note: Graph schema is configured in MDFE, but Titan indexes are created and managed outside MDFE

# Key MDFE concepts are borrowed from NiFi World …

**NiFi**

| Flow |
| --- |

| Processor Group |
| --- |

| Processor |
| --- |

| Data Provenance | Process Mgmt. |
| --- | --- |

Vs.

**MDFE**

| Job |
| --- |

| Controller |
| --- |

| Processor |
| --- |

| Meta data management | Process Mgmts. |
| --- | --- |

MDFE was originally created for NiFi … and then refactored and optimized for Spark
*(even the naming conventions of the processors follow NiFi conventions)*
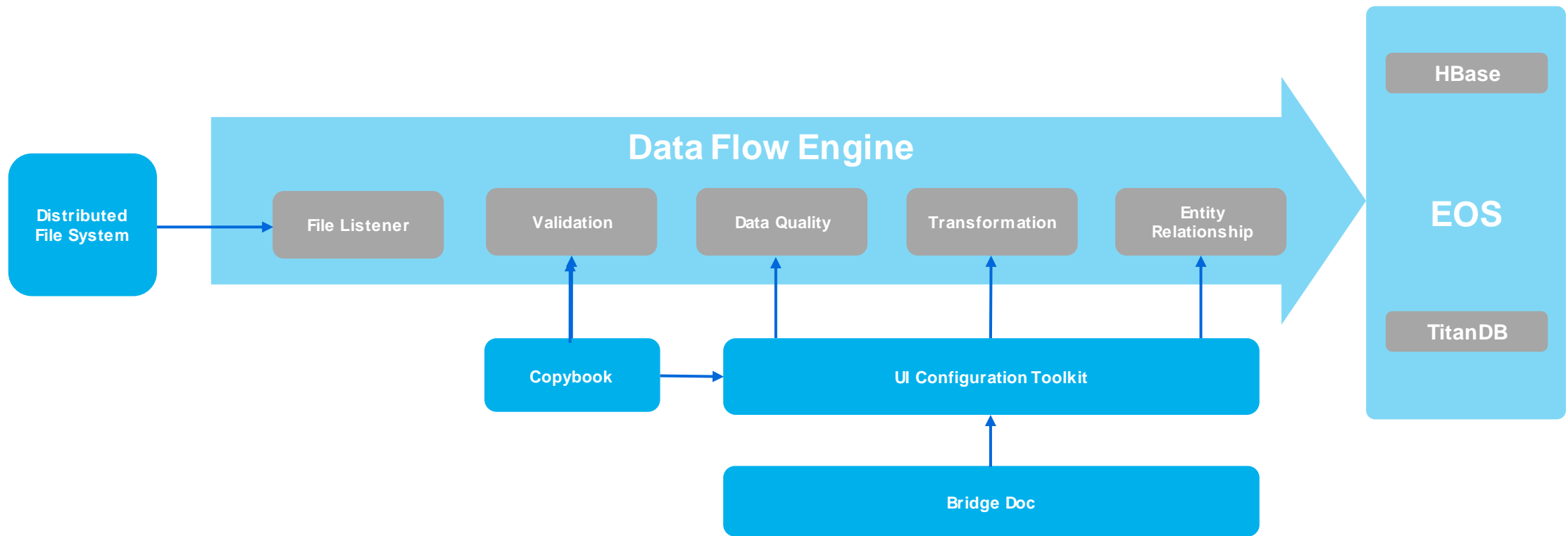
MDH Data Flow – A high level view

| Latency / Data Delivery mechanism | Land data in RDZ (HDFS) | Qualify File Structure | Demarcation Point | Minimum Viable Record Check | Map from SOR to EOS format | Create Relationships | Orphan record monitoring |
|---|---|---|---|---|---|---|---|
| | Port Adapter | Entity Port | | D² | E³ | R² | WD |
| Near Real Time - MQ | Flume | Spark | | Spark | Spark | Spark | Spark |
| Batch (micro / macro) - sFTP | Informatica MFT / ETL | Spark | | Spark | Spark | Spark | Spark |
| CDC — Change stream | Informatica CDC | Spark | | Spark | Spark | Spark | Spark |
| Exception | Operational Fail File | Operational Fail File | | Fail Record (With Flag) | Fail Record (With Flag) | Exception Queue | Exception Queue |
| Restart / Rerun | Operational Notifications | | | Functional Notification | | | |

# Why MDFE?

To make data coming from disparate SOR systems in various formats and latencies available on MDH's speed layer, along with Metadata, Metaprocess and relationships attached to it.
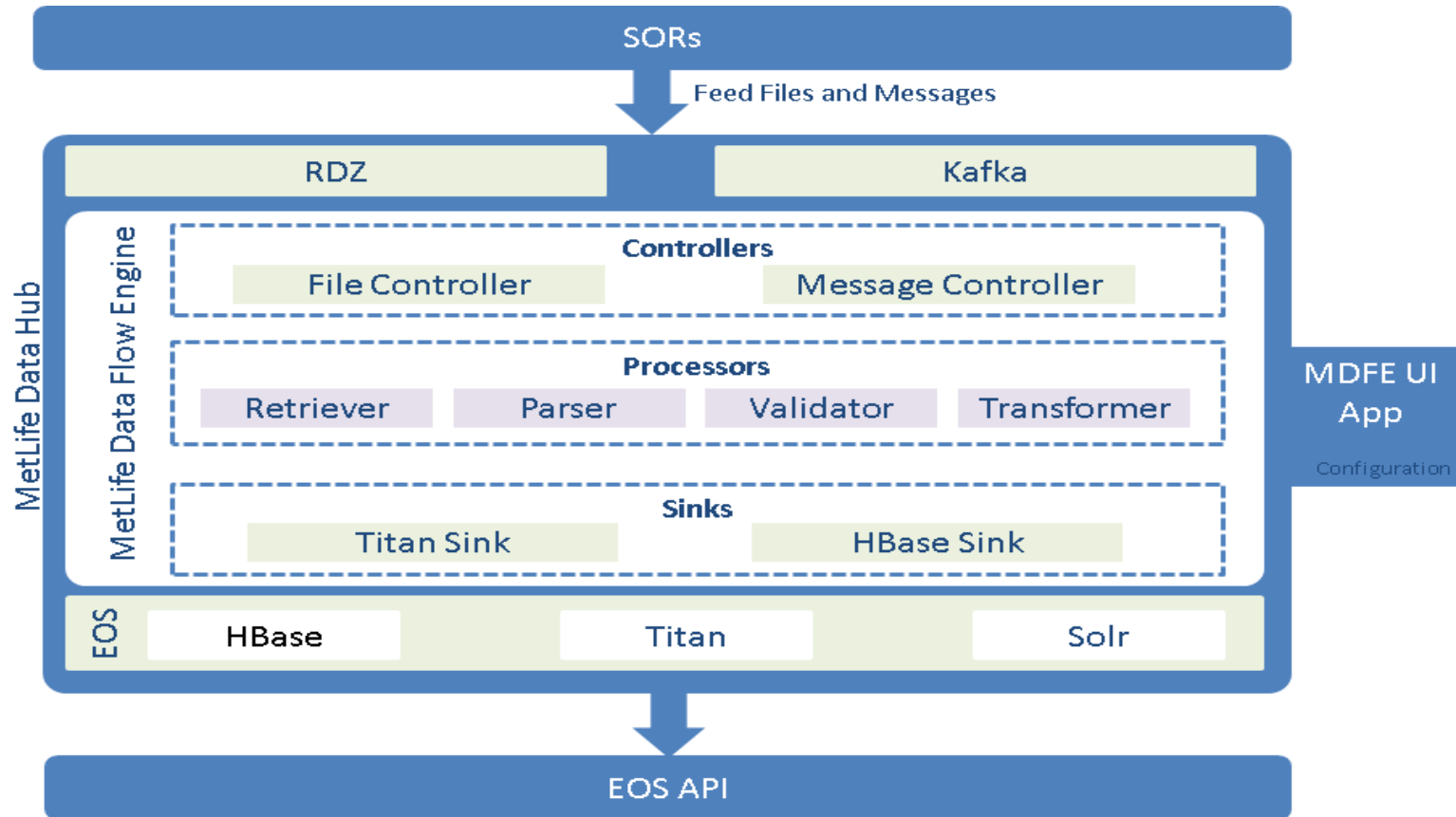
Latencies - Near real time: <5 min, Micro batch: 5-60 min, Macro batch: >60 min and < 24hrs

# What does it do ?

MDFE Ingesting data into EOS, a simplified view

# MDFE
# Code
# Structure

# MDFE Component Structure

| Components | Versions |
|---|---|
| Hadoop | 2.7.3 |
| hBase | 1.2.4 |
| titan | 1.0.0 |
| hive | 1.2.1 |
| solr | 6.3.0 |
| kafka | 0.10.1.0 |
| flume | 1.7.0 |
| spark | 2.1.0 |

# Tech Stack

MDFE is a spark based engine, which runs on IBM Big Insights Platform.

# MDFEController.scala

START

If argument = null

Initialize the Spark Context

Check trigger folder, Get the RDZ path from the Configuration Table

If trigger folder present

TRUE

Convert feed file to Dataset

TRUE

STOP

# MDFEPutTitanBulk.scala

Create Graphson for the Harmonized JSON

Put the data in the Watchdog Table

Create Vertex and Edges in Titan

# MDFEPutJsonHbase.scala

Put the checksum in HBase Recon Table

Put the data in HBase Entity Table

# DStreamController.scala

Check/Parse & Iterate the Job Configuration Json

Orchestrate processor Instantiation & Initialization

# GetData.scala

Call the GetData processor to pass the Dataset to the next processor

**THINGS TO DO BEFORE STARTING THE DATA PIPE:**

1. Configuration table '**fileregister**' should be populated with all the configurations assigned to a particular application ID we are running.
2. Trigger file should be present in the trigger folder of the RDZ path.
3. Trigger file should be of the same name with the data feed file.

# MDFEProcessor.scala

Fingerprinting the

**1**

**Build & Package : Bamboo**
- **Build Jar**
- **Create tar file containing Jar and script files**
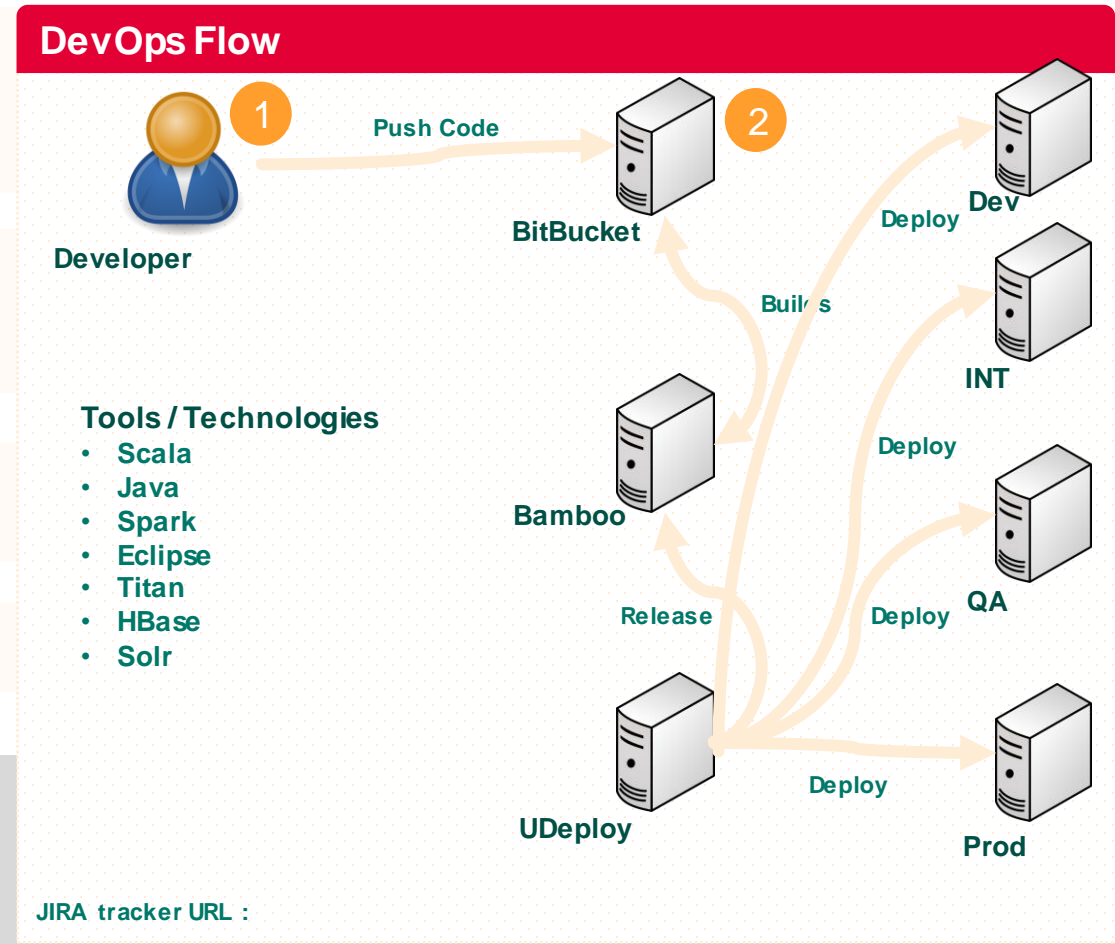- **Create a release version**
- **URL:**

**2**

**Repository : Bitbucket**
Development Branch: depends on the relase
Baselined Version          : ***
URL:

**3**

**Deploy : Udeploy [Not in Prod]**
- **Create a Snapshot of the baselined version - Eg: 10856-MDFE-Rel2018.1.0-<relDate>**
- **Open tickets and get necessary approvals to run Prod deployment procedure**

**4**

**Run: Maestro**

## DevOps Flow



**Developer**

**1** Push Code → BitBucket

**2** → Dev, INT, QA, Prod

Deploy

Builds

**Tools / Technologies**
- **Scala**
- **Java**
- **Spark**
- **Eclipse**
- **Titan**
- **HBase**
- **Solr**

**Bamboo**

**Release**

**Deploy**

**UDeploy**

**Deploy**

**Prod**

**JIRA tracker URL :**

# How do we build and Deploy

MetLife Data Flow Engine is a custom built engine that moves data within the MetLife Data Hub environment.

Typical reviews include

- High level and detailed design document reviews and manual and automated code reviews and security reviews (veracode scans)
- Architecture reviews to confirm that the development confirms to the architecture

Current status

- MetLife and Cognizant Architecture and Development teams have completed few rounds of review
- Knowledgent has completed design and code reviews
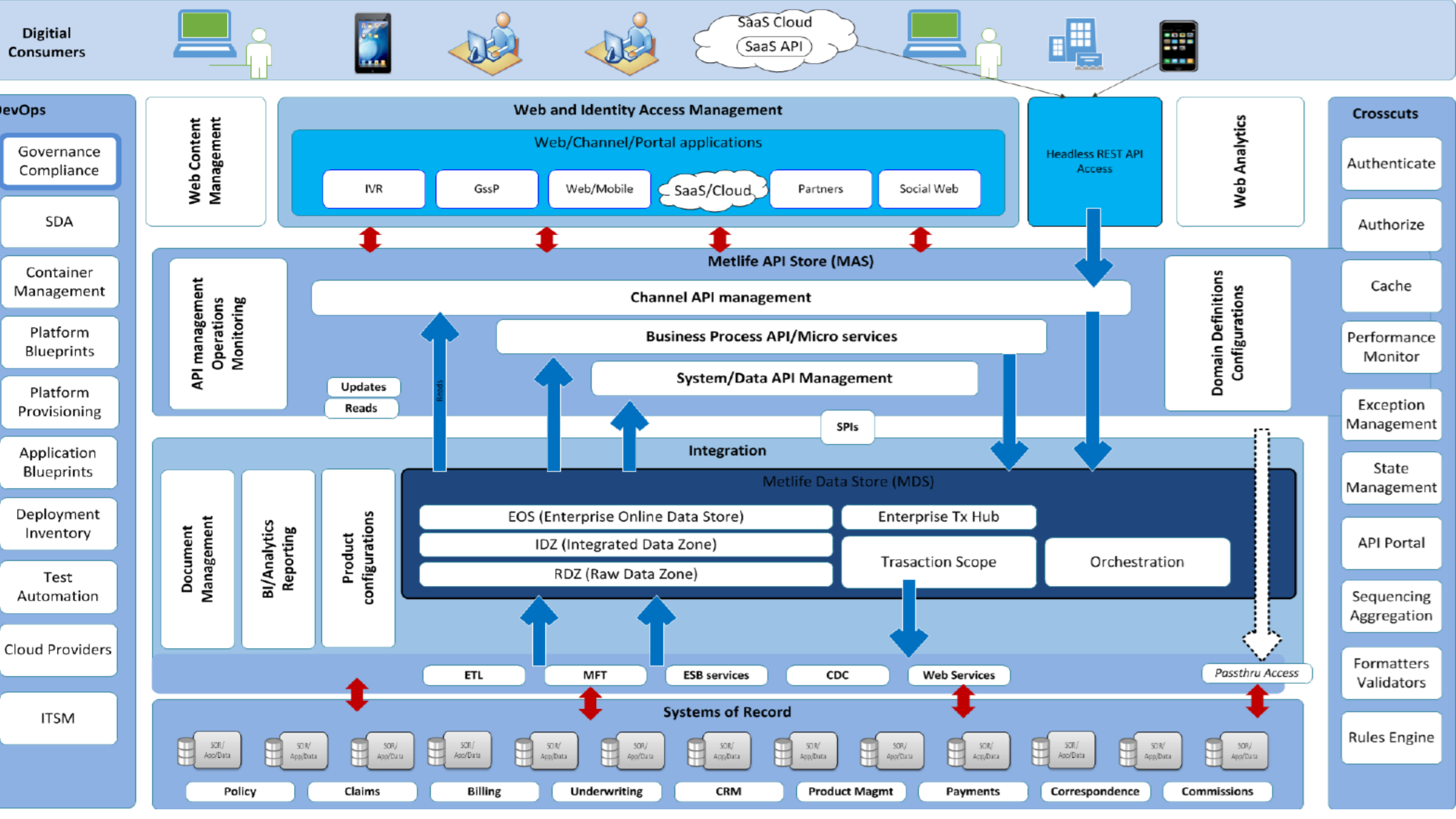
Design Document Reviews

Sharepoint links ->

Code Reviews

Bitbucket location ->

# Review Process

# Appendix

# Data Curation