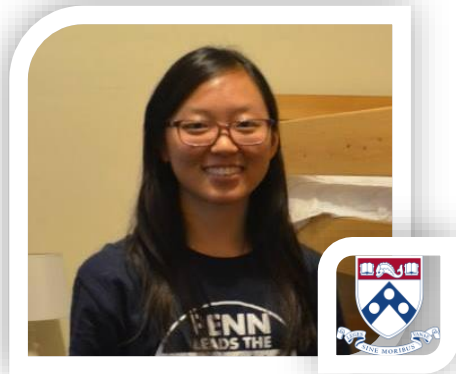# DataSpark

clarity redefined

# SparkGen: Spark Job Configurator

Data Engineering Team

# The Team



**Karen Her**
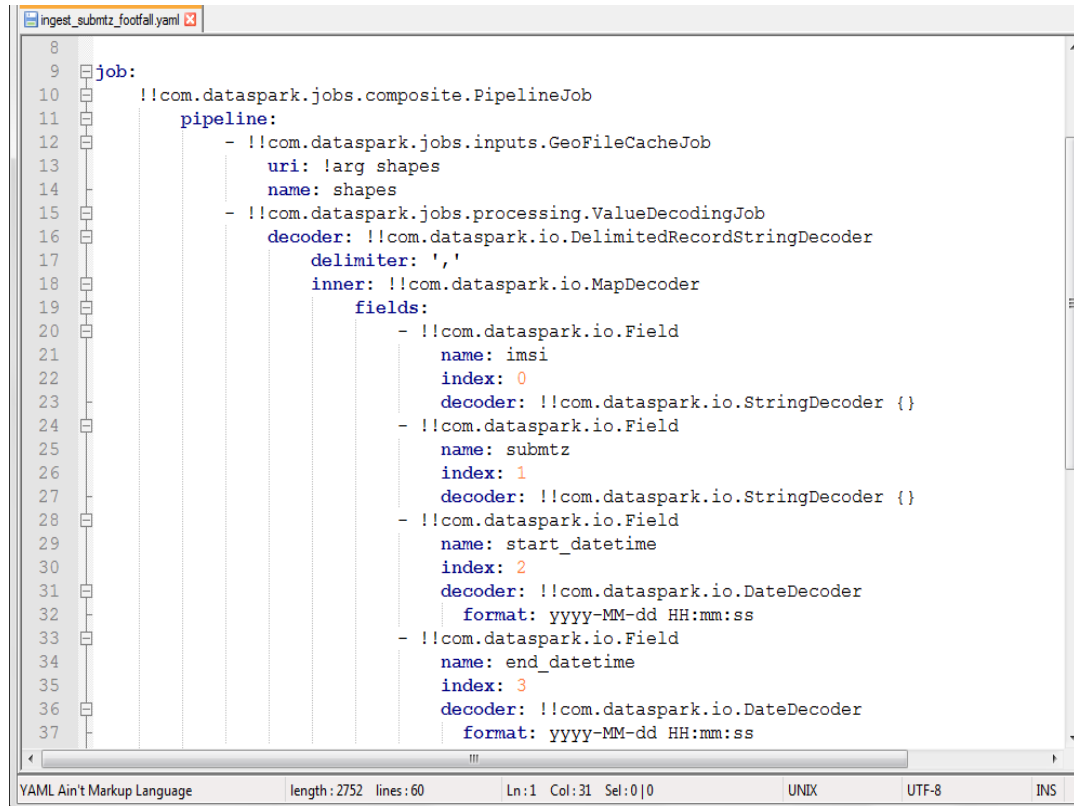University of Pennsylvania
Computer Science

**Skyler Sin**
Stanford University
Symbolic Systems





**Karthik Balasubramanian**
Carnegie Mellon University
Information Systems Management

DataSpark

# Executive Summary

- SparkGen is a Spark job configurator



```
ingest_submtz_footfall.yaml
  8
  9  job:
 10      !!com.dataspark.jobs.composite.PipelineJob
 11          pipeline:
 12              - !!com.dataspark.jobs.inputs.GeoFileCacheJob
 13                  uri: !arg shapes
 14                  name: shapes
 15              - !!com.dataspark.jobs.processing.ValueDecodingJob
 16                  decoder: !!com.dataspark.io.DelimitedRecordStringDecoder
 17                      delimiter: ','
 18                      inner: !!com.dataspark.io.MapDecoder
 19                          fields:
 20                              - !!com.dataspark.io.Field
 21                                  name: imsi
 22                                  index: 0
 23                                  decoder: !!com.dataspark.io.StringDecoder {}
 24                              - !!com.dataspark.io.Field
 25                                  name: submtz
 26                                  index: 1
 27                                  decoder: !!com.dataspark.io.StringDecoder {}
 28                              - !!com.dataspark.io.Field
 29                                  name: start_datetime
 30                                  index: 2
 31                                  decoder: !!com.dataspark.io.DateDecoder
 32                                    format: yyyy-MM-dd HH:mm:ss
 33                              - !!com.dataspark.io.Field
 34                                  name: end_datetime
 35                                  index: 3
 36                                  decoder: !!com.dataspark.io.DateDecoder
 37                                    format: yyyy-MM-dd HH:mm:ss

YAML Ain't Markup Language    length : 2752   lines : 60    Ln : 1   Col : 31   Sel : 0 | 0    UNIX    UTF-8    INS
```
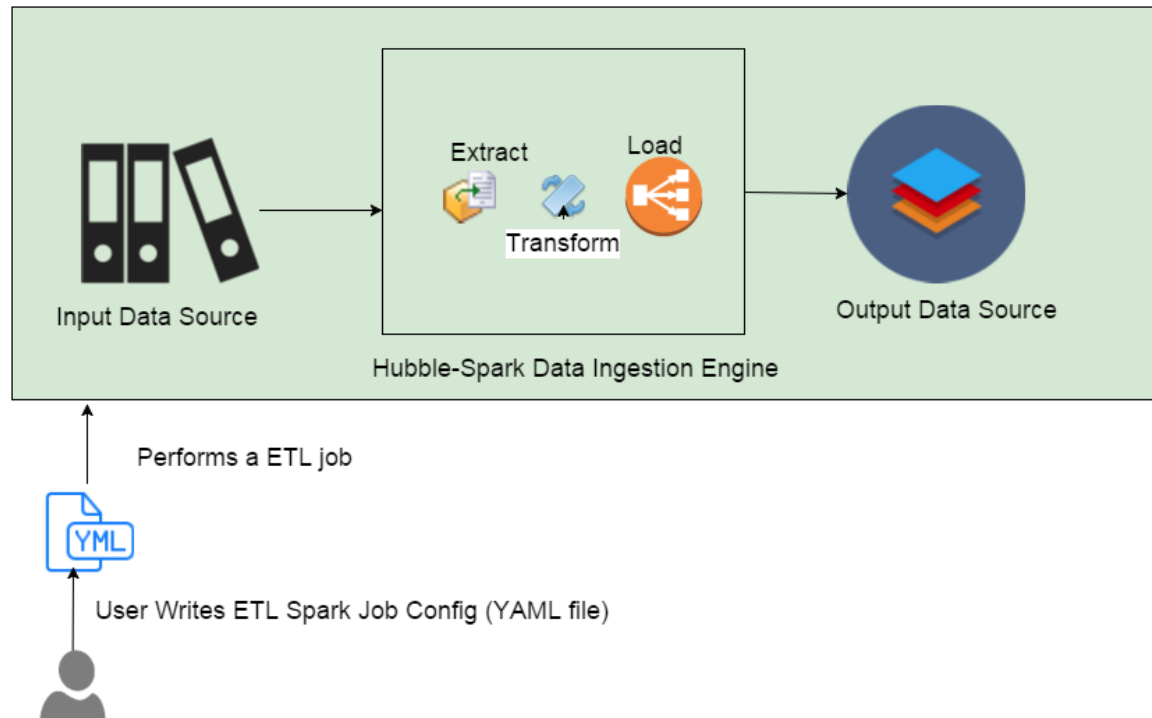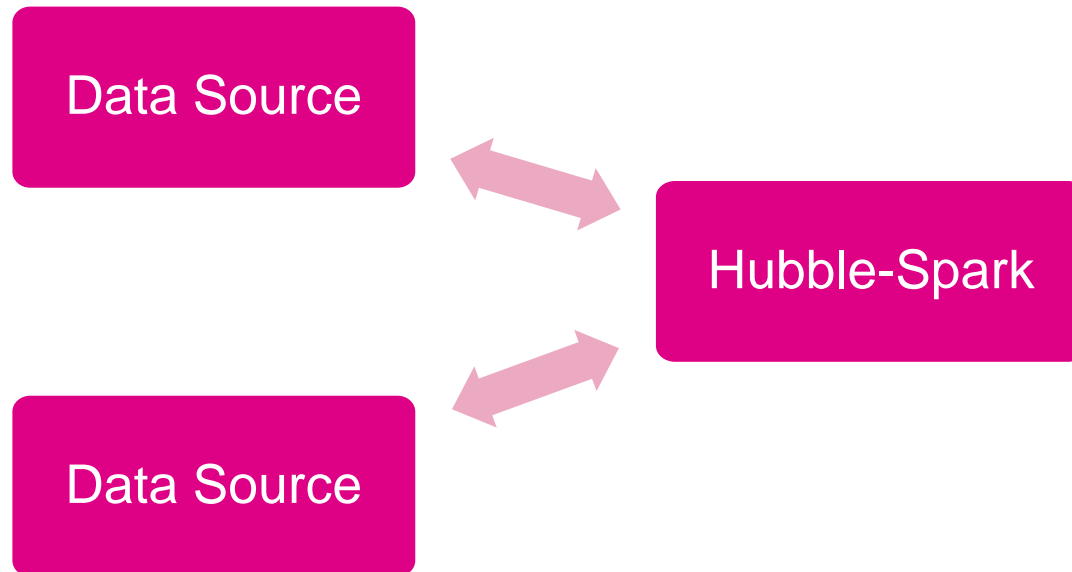
Job Configuration File

**DataSpark**

3

# Executive Summary

- SparkGen is a Spark job configurator

- Hubble-Spark – A Data Ingestion Engine in DataSpark environment. It is written in Java.
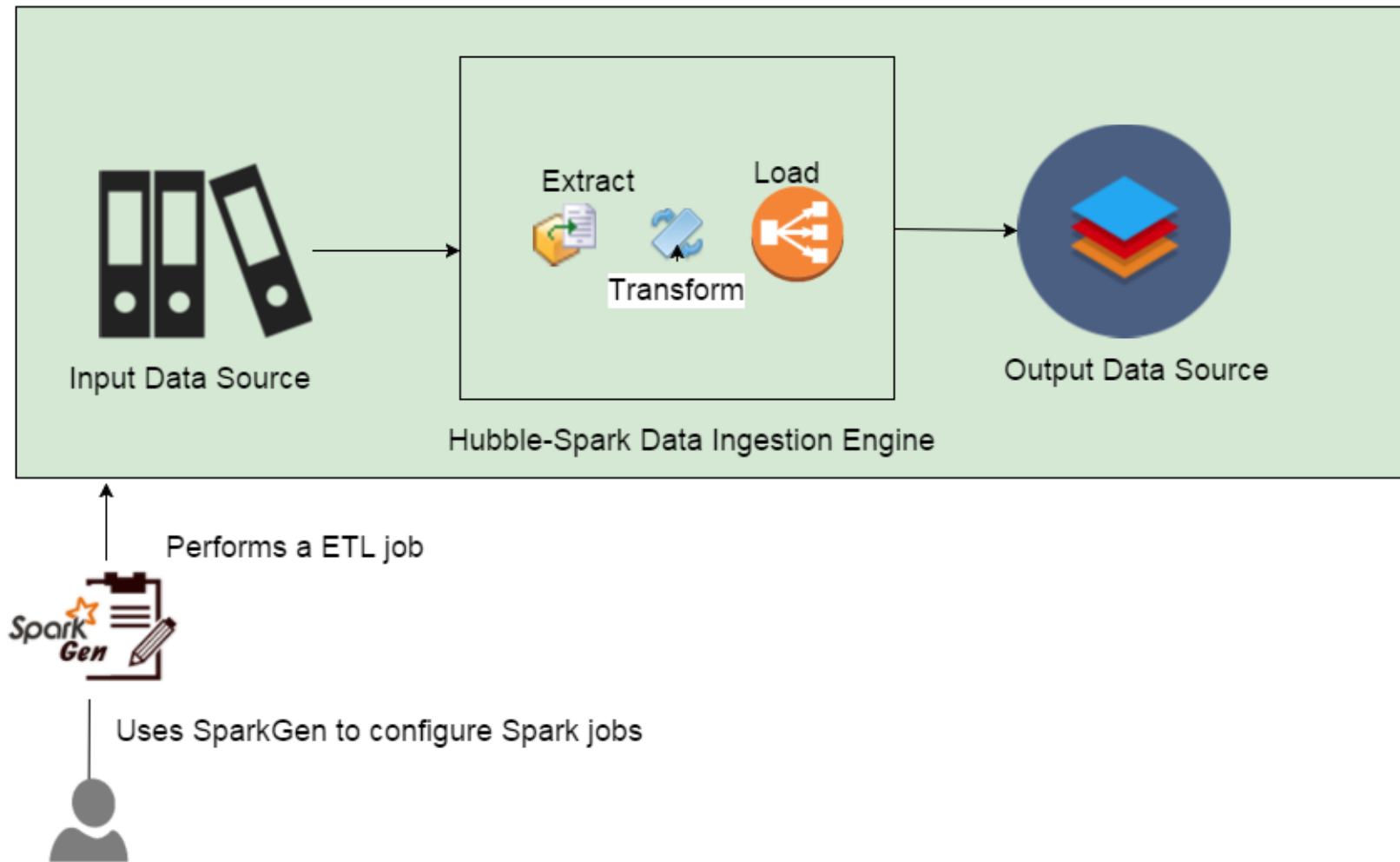


**DataSpark**

# Executive Summary

- SparkGen is a Spark job configurator

- Hubble-Spark – A Data Ingestion Engine in DataSpark environment. It is written in Java.

- Has connectors to multiple varieties of data sources to read, write and perform Extract, Transform and Load (ETL) jobs.

| Data Source | | Hubble-Spark |
|:---:|:---:|:---:|
| | ⬌ | |
| Data Source | ⬌ | |

DataSpark

# Executive Summary

- SparkGen is a Spark job configurator

- Hubble-Spark – A Data Ingestion Engine in DataSpark environment. It is written in Java.

- Has connectors to multiple varieties of data sources to read, write and perform Extract, Transform and Load (ETL) jobs.

- Uses markup language called YAML to configure ETL Spark Jobs

- We developed a SparkGen API component to expose internal components of Hubble-Spark ETL engine and  a UI component to generate a Spark Job Configuration file (YAML file)

**DataSpark**

# Executive Summary



Input Data Source

Extract

Transform

Load

Hubble-Spark Data Ingestion Engine

Output Data Source

Performs a ETL job

SparkGen

Uses SparkGen to configure Spark jobs

DataSpark

# Problem Statement

*How do we simplify the robust configuration driven*

*Hubble-Spark to a team member with limited expertise*

*in Spark and achieve seamless ETL activities?*

DataSpark

# Problem Background

- Hubble Spark is an execution framework for any ETL Spark Job
- Has hundreds of components to extract, transform, and load terabytes of data which facilitate DataSpark's business
- Ideology: One code base – Million jobs
- All you have to do is write your ETL configurations using the existing templates

DataSpark

# Problem Background

However,

- Hubble-Spark requires deep understanding of Java-Spark to configure Spark ETL Jobs
- YAML file used to configure complex Spark Jobs is a new configuration markup, is created manually and error prone

DataSpark

# Our Idea

*What if?*

- *We provide an abstraction to Hubble-Spark as a Graphical User Interface and give users the limited information they need to know to run any ETL Jobs using Hubble-Spark*

**DataSpark**

# Goal

- Create a user-friendly interface that will produce a YAML file to run Spark Jobs.



Internal Components
of Hubble-Spark

# Goal

- Create a user-friendly interface that will produce a YAML file to run Spark Jobs.



Produce YAML
File

# Solution Description

- Three components working together
  - Spark Jar File
    - An execution framework for all ETL Spark Jobs
  - SparkGen API
    - Exposes Hubble-Spark internal components to user interface
  - SparkGen UI
    - Creates ETL job configuration file (YAML)
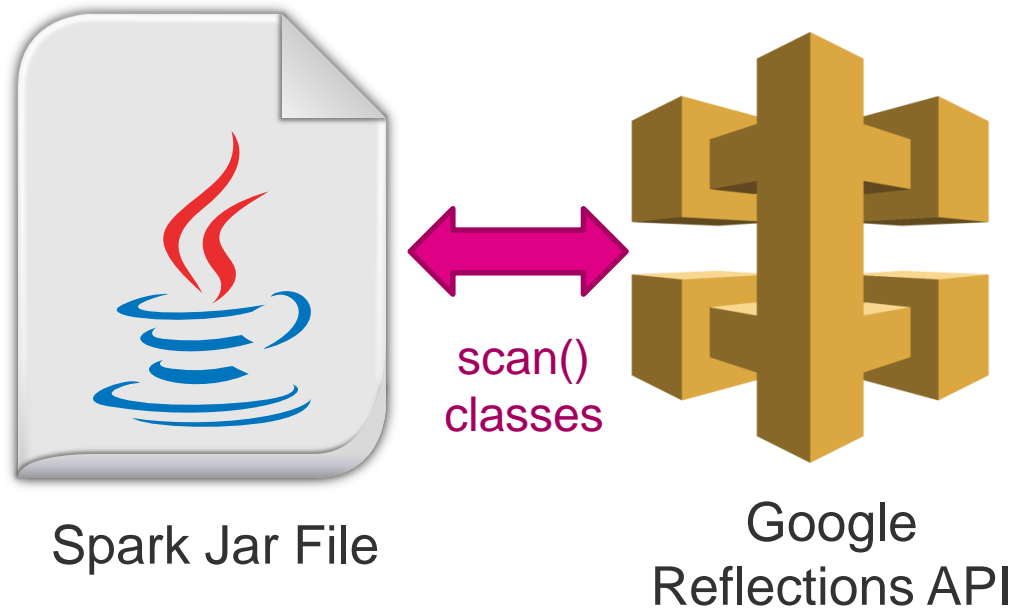
DataSpark
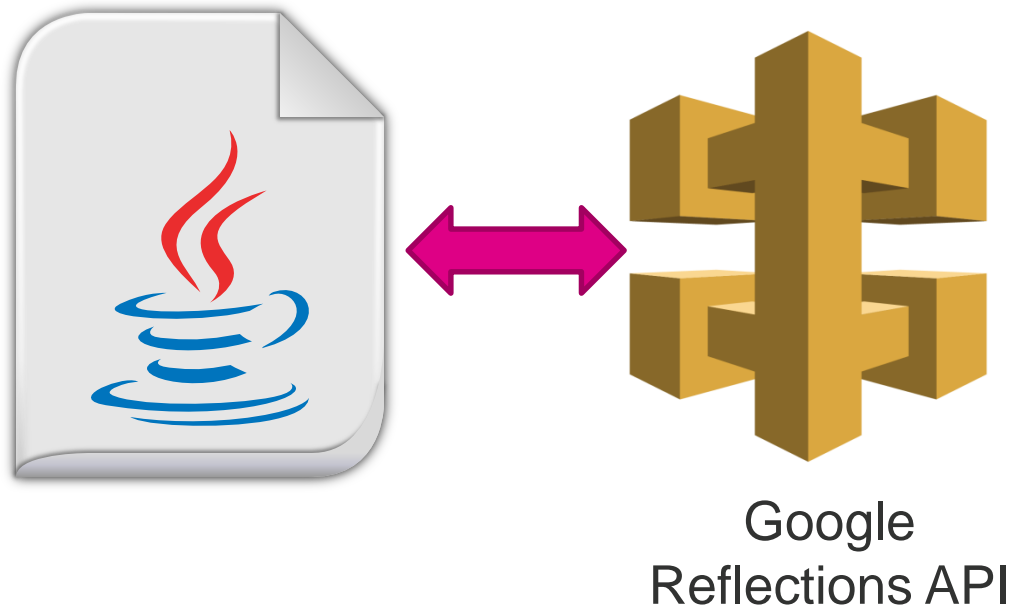
# Solution Steps


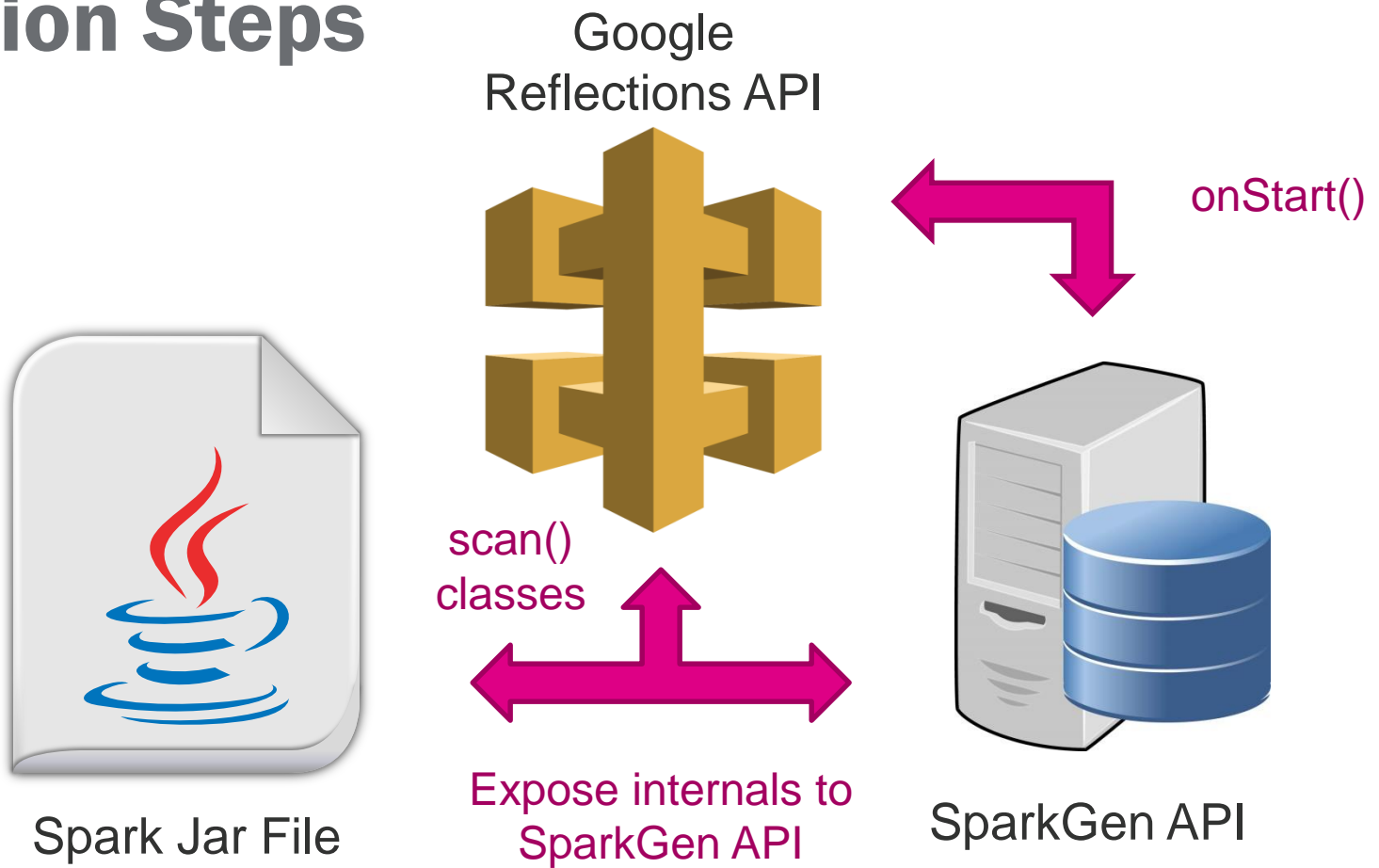
Spark Jar File

# Solution Steps
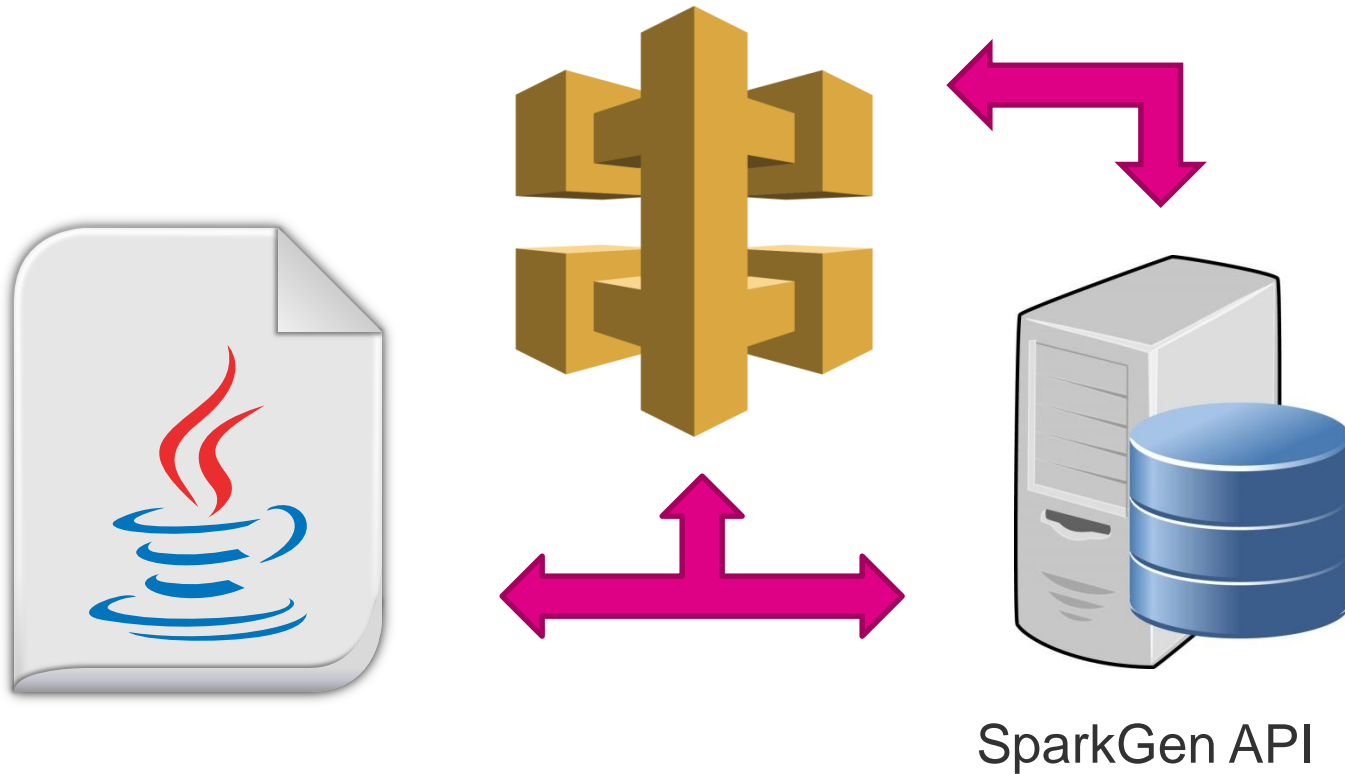


Spark Jar File

# Solution Steps



scan() classes

Spark Jar File

Google Reflections API

# Solution Steps



Google
Reflections API

# Solution Steps

Google
Reflections API

onStart()

Spark Jar File

scan()
classes

Expose internals to
SparkGen API

SparkGen API

# Solution Steps



SparkGen API

DataSpark

# Solution Steps

# Solution Steps

Google
Reflections API

onStart()

scan()
classes

Spark Jar File

Expose internals to
SparkGen API

SparkGen API

Spark
Gen

DataSpark

# Solution Steps

SparkGen API
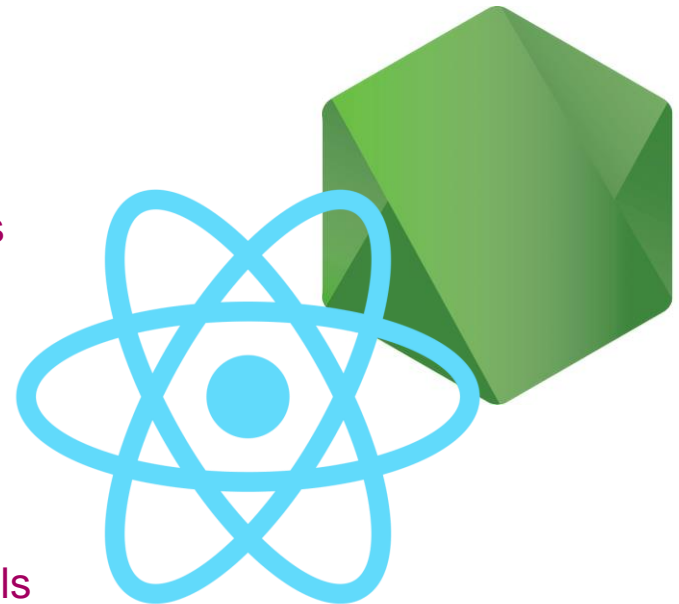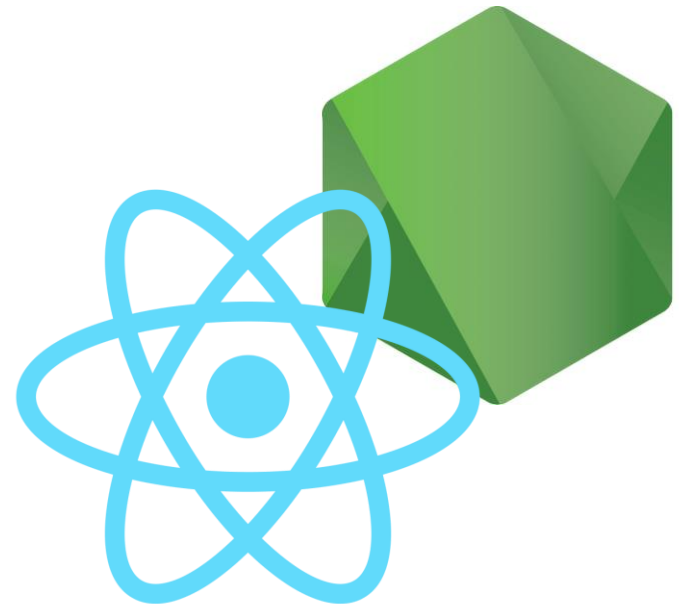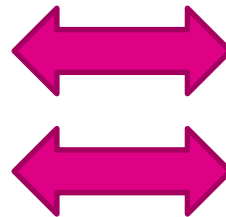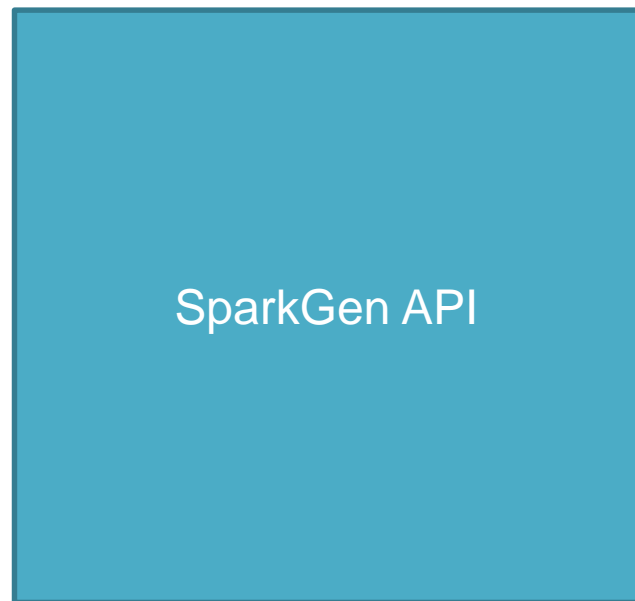
getAllClasses

getClassDetails

React Form with
Node Server

# Solution Steps

SparkGen API

React Form with
Node Server

# Solution Steps



SparkGen API
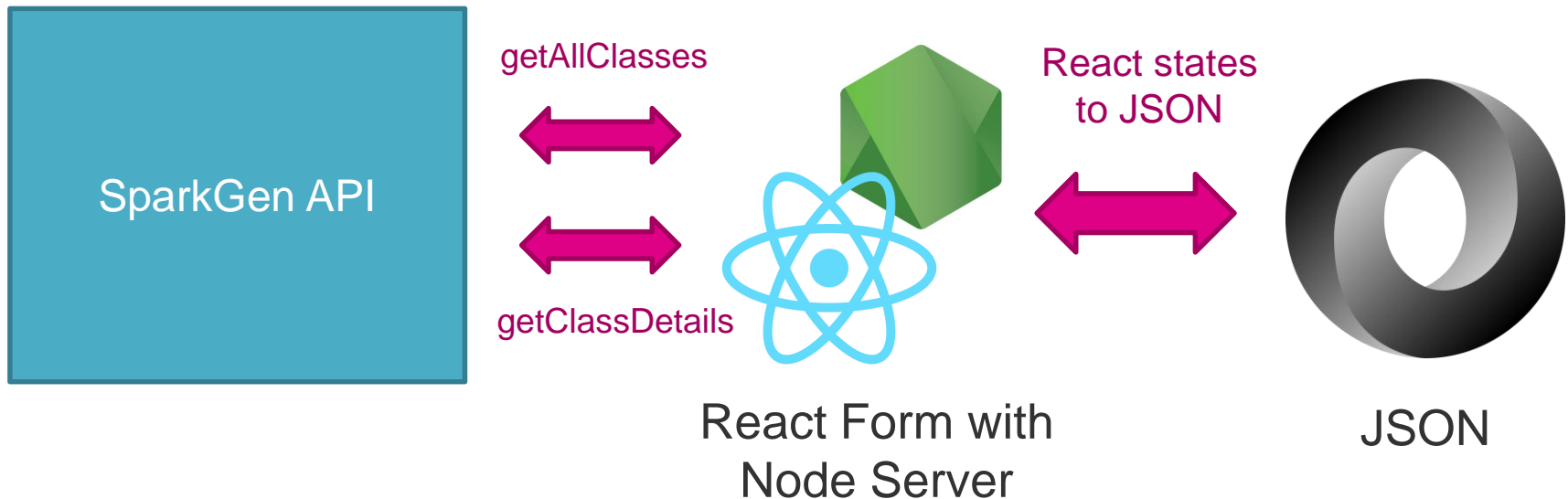
getAllClasses

getClassDetails

React Form with
Node Server

React states
to JSON

JSON

DataSpark

# Solution Steps



SparkGen API

JSON

# Solution Steps



getAllClasses

SparkGen API

React states to JSON

SnakeYaml

getClassDetails

React Form with Node Server

JSON

YAML

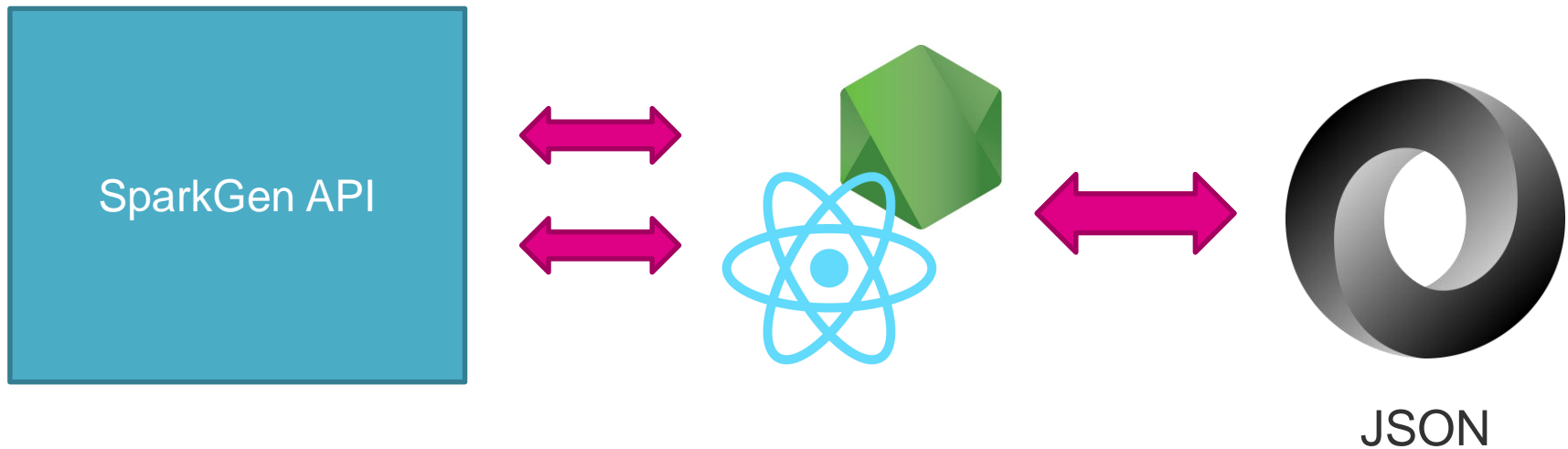DataSpark

# Solution Overview

Flow Chart

scan() classes — Google Reflections API — onStart() — getAllClasses — Node Server — React states to JSON — SnakeYaml — getClassDetails — Expose internals to SparkGen API

Spark Jar File — SparkGen API — React Form — JSON — YAML

# SparkGen

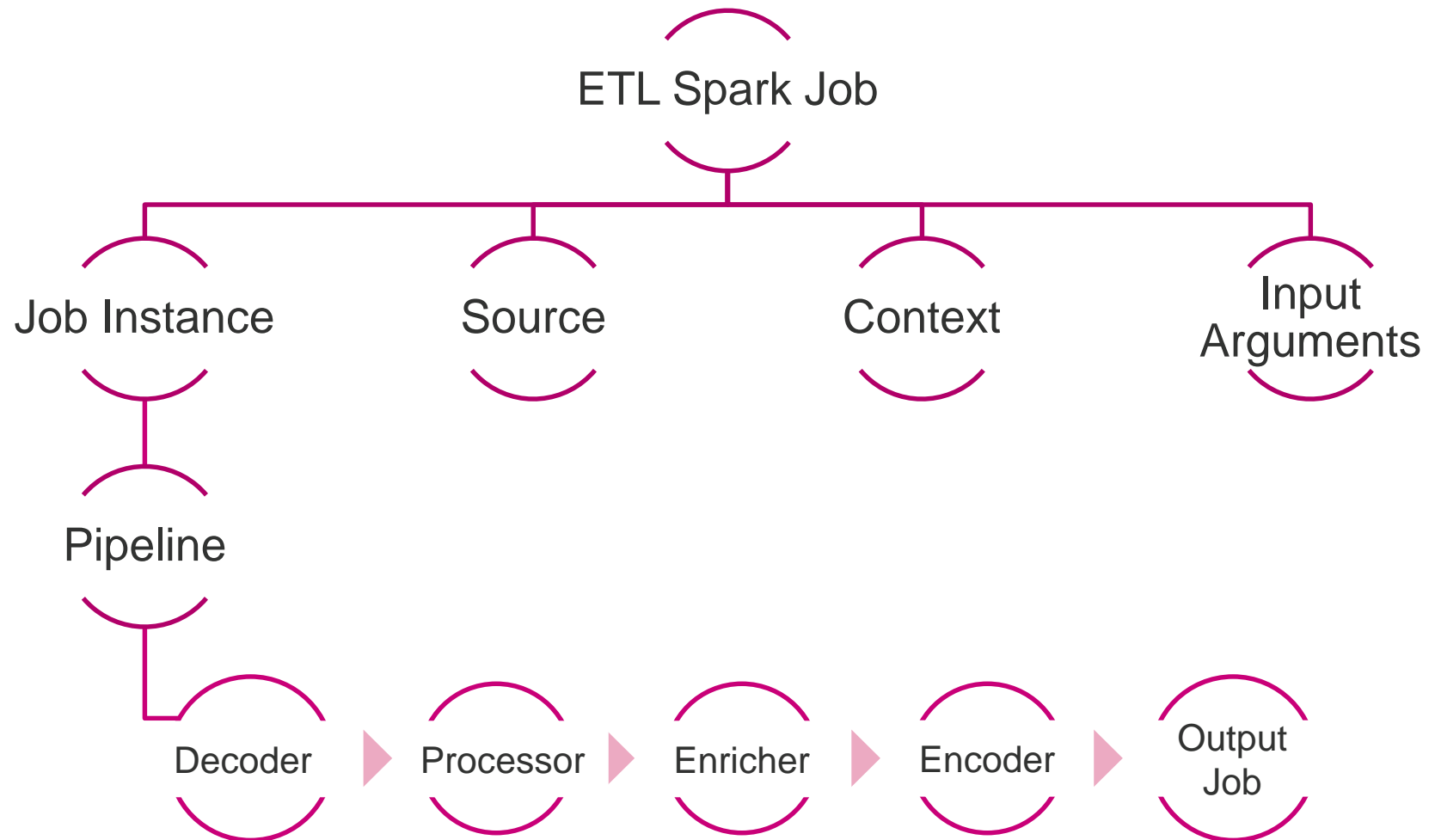- SparkGen is a dynamic form
  - Grows based on the user's inputs
  - User's inputs are recorded in the browser
  - On submit, the form is created and shows a YAML file in the CodeMirror palette of the UI
  - Users can copy to clipboard or edit the YAML directly in the palette

**DataSpark**

# Components of a YAML File

ETL Spark Job

Job Instance    Source    Context    Input Arguments

Pipeline

Decoder ▶ Processor ▶ Enricher ▶ Encoder ▶ Output Job

**DataSpark**

# Next Steps

- Refining YAML file
- Work with Product Team to validate YAML generation process

**DataSpark**

# DEMO

DataSpark

# DataSpark

Thank you