

“What is the Role of a Data Engineering Team?”

Jorge Figueiredo



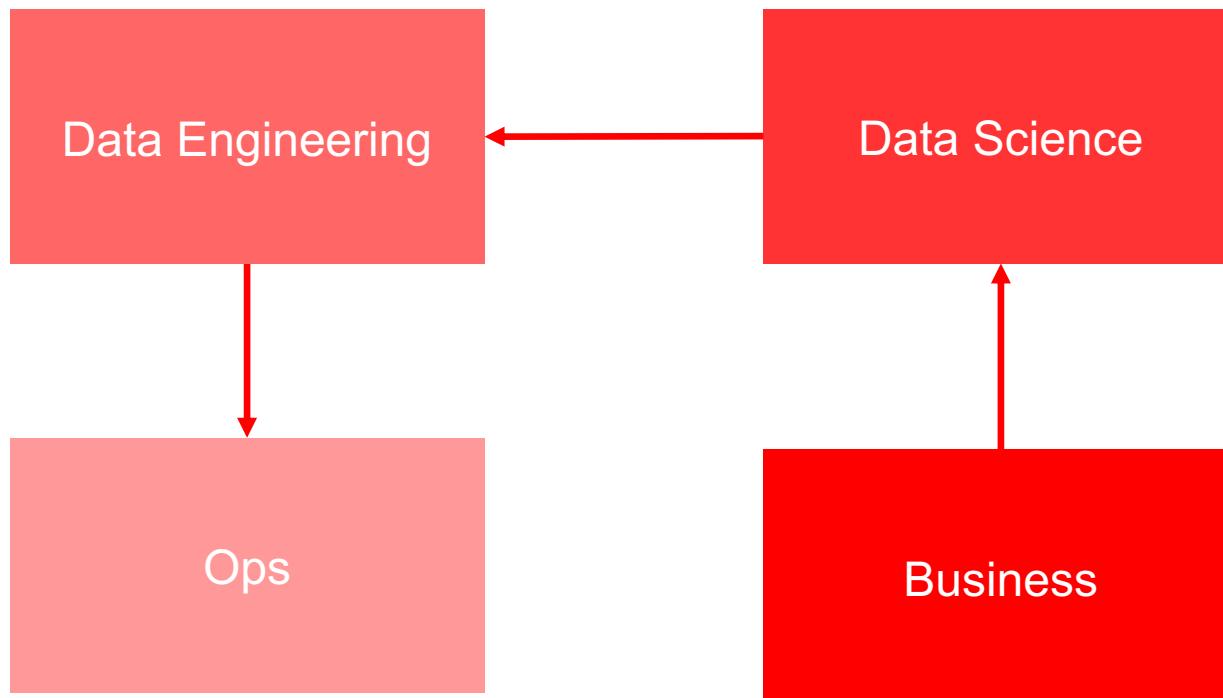
What is the problem?



2



What organizations do in this data age?



Data Scientist vs Data Engineer

Data Scientist

- Advanced Math/Statistics
- AI/ML
- Advanced Analytics

Overlapping Skills

- Programming
- Analysis
- Big Data

Data Engineer

- Advanced Programming
- Distributed Systems
- Data Pipelines



What have we used for ETL in Data Engineering?

- JVM driven ETL
- SQL driven ETL



What are the options for Data Engineering?

- JVM (Java/Scala) ETL pipelines

```
package com.twitter.scalding.examples

import com.twitter.scalding._
import com.twitter.scalding.source.TypedText

class WordCountJob(args: Args) extends Job(args) {
    TypedPipe.from(TextLine(args("input")))
        .flatMap { line => tokenize(line) }
        .groupBy { word => word } // use each word for a key
        .size // in each group, get the size
        .write(TypedText.tsv[(String, Long)](args("output")))

    // Split a piece of text into individual words.
    def tokenize(text: String): Array[String] = {
        // Lowercase each word and remove punctuation.
        text.toLowerCase.replaceAll("[^a-zA-Z0-9\\s]", "").split("\\s+")
    }
}
```

What are the options for Data Engineering?

- SQL (Hive) ETL pipelines

```
1   -- Dropping Database
2   DROP DATABASE IF EXISTS ipl_data CASCADE;
3   DROP DATABASE IF EXISTS ipl_stats CASCADE;
4
5   -- Creating a Database
6   CREATE DATABASE ipl_data ; -- Database for Loading CSV Input File
7   CREATE DATABASE ipl_stats ; -- Database for Loading stats tables
8
9   -- Creating two tables to load input files ( matches and deliveries )
10  CREATE TABLE ipl_data.matches(
11      `id` int,
12      `season` int,
13      `city` string,
14      `date` date,
15      `team1` string,
16      `team2` string,
17      `toss_winner` string,
18      `toss_decision` string,
19      `result` string,
20      `dl_applied` int,
21      `winner` string,
22      `win_by_runs` int,
23      `win_by_wickets` int,
24      `player_of_match` string,
25      `venue` string,
26      `umpire1` string,
27      `umpire2` string,
28      `umpire3` string)
29      ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
30      WITH SERDEPROPERTIES (
31          "separatorChar" = ",",
32          "quoteChar" = "\""
33      )
34      STORED AS TEXTFILE
35      TBLPROPERTIES (
36          'serialization.null.format' = '',
37          'skip.header.line.count' = '1');
```



What are the options for Data Engineering?

■ Bash "pipelines"

```
#!/bin/bash
# Counting the number of lines in a list of files
# function version

# function storing list of all files in variable files
get_files () {
    files="`ls *.[ch]`"
}

# function counting the number of lines in a file
count_lines () {
    local f=$1 # 1st argument is filename
    l=`wc -l $f | sed 's/^([0-9]*\).*$/\1/'` # number of lines
}

# the script should be called without arguments
if [ $# -ge 1 ]
then
    echo "Usage: $0 "
    exit 1
fi

# split by newline
IFS=$'\012'

echo "$0 counts the lines of code"
# don't forget to initialise!
```



What is the other way?



What is the other way?

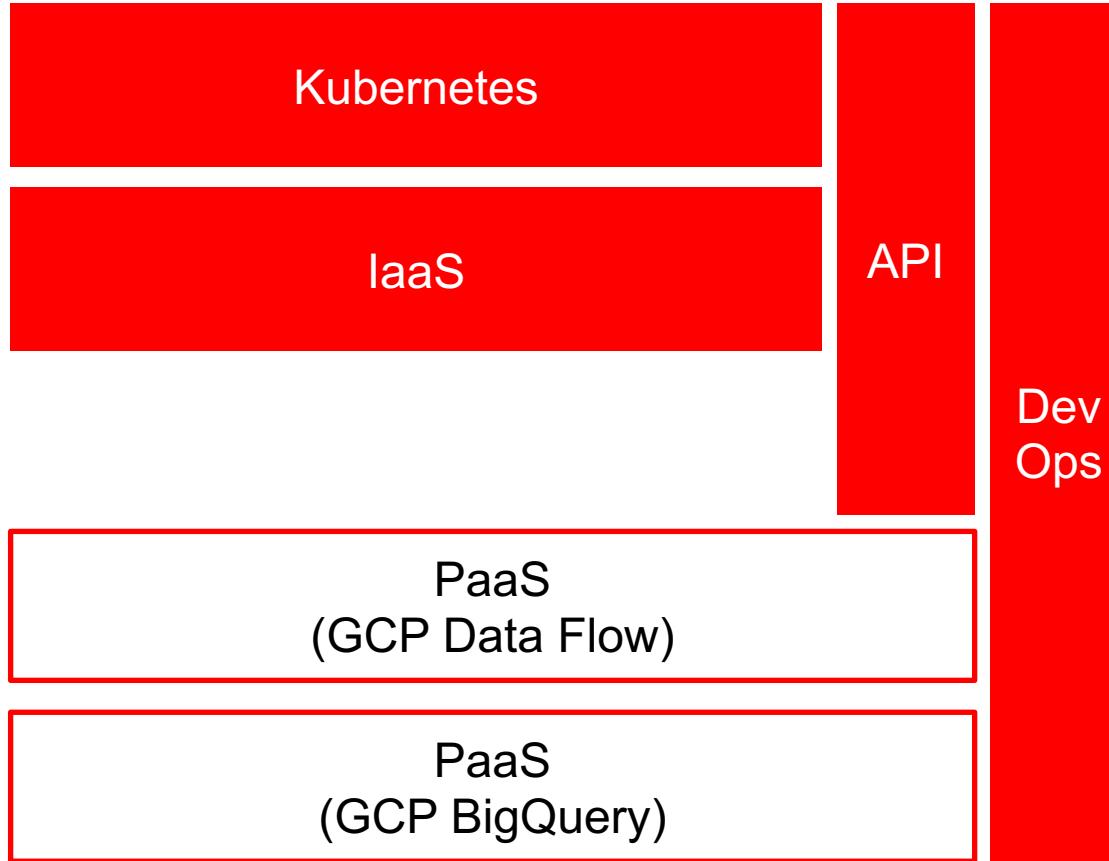
- APIs
- No-Ops
- Self-service



APIs



DE



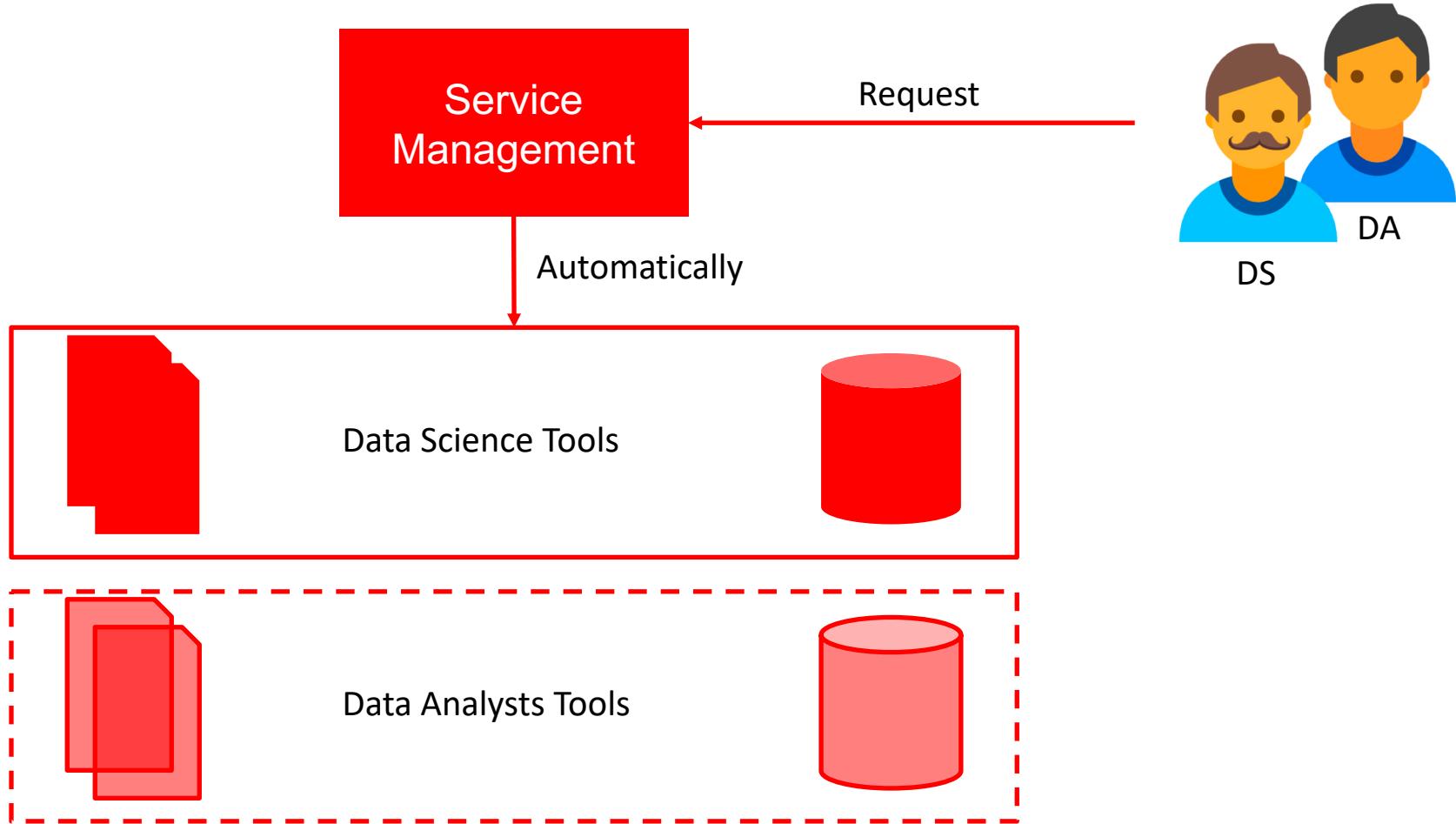
DS

NoOps

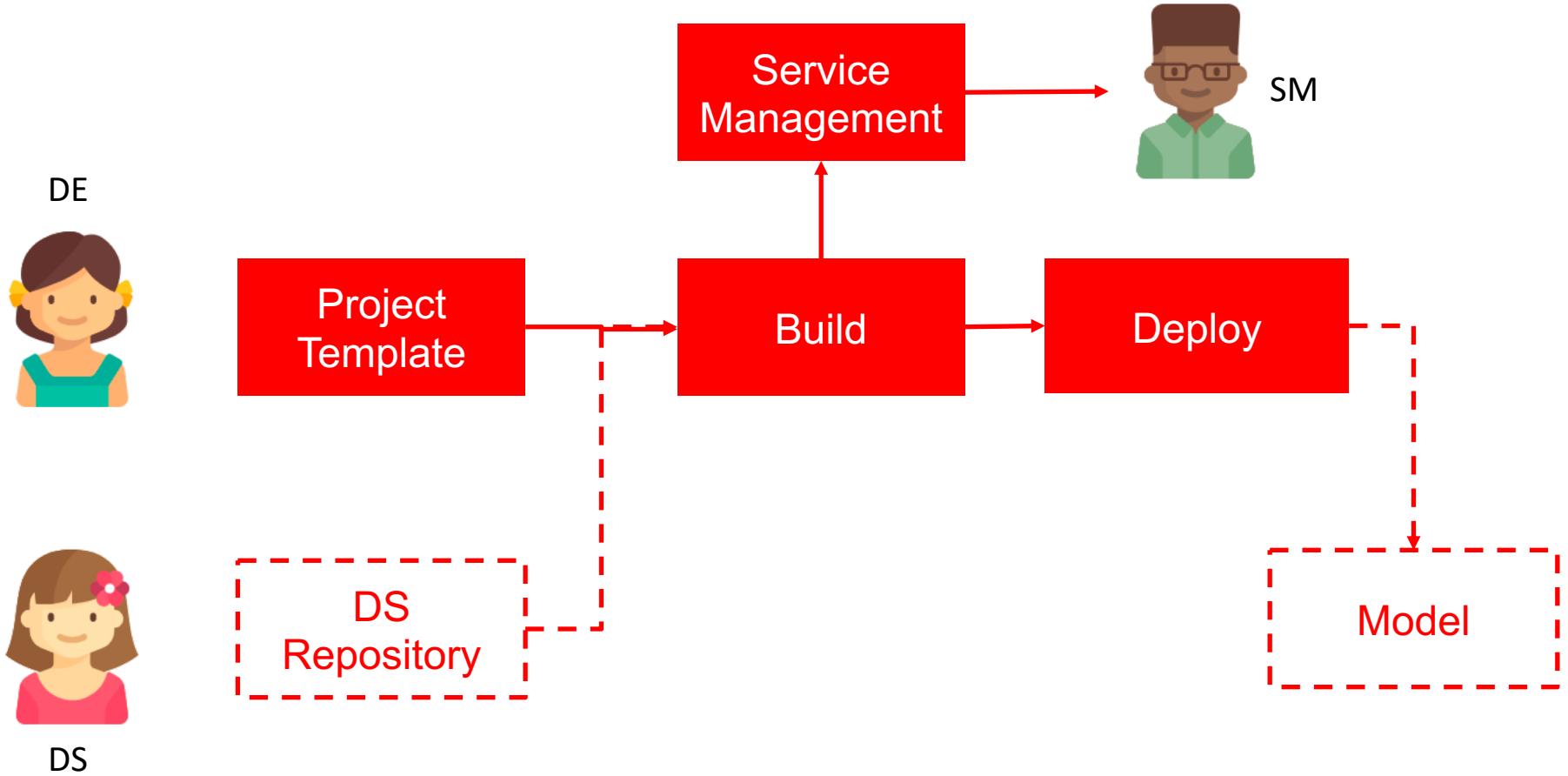
DevOps	NoOps
Development and operations work together.	Development and operations never need to interact.
Dev and Ops work together to select, build, monitor and maintain self-service solutions.	Self-service infrastructure, continuous integration and deployment solutions.
Based on principles and practices.	Only possible through use of specific software solutions.



Technical patterns – Self-service Infrastructure toolkit



Technical patterns – Project Template



Technical patterns – Project template

Branch feature/GBIDE-524

Full project name: scala-template/feature%2FGBIDE-524

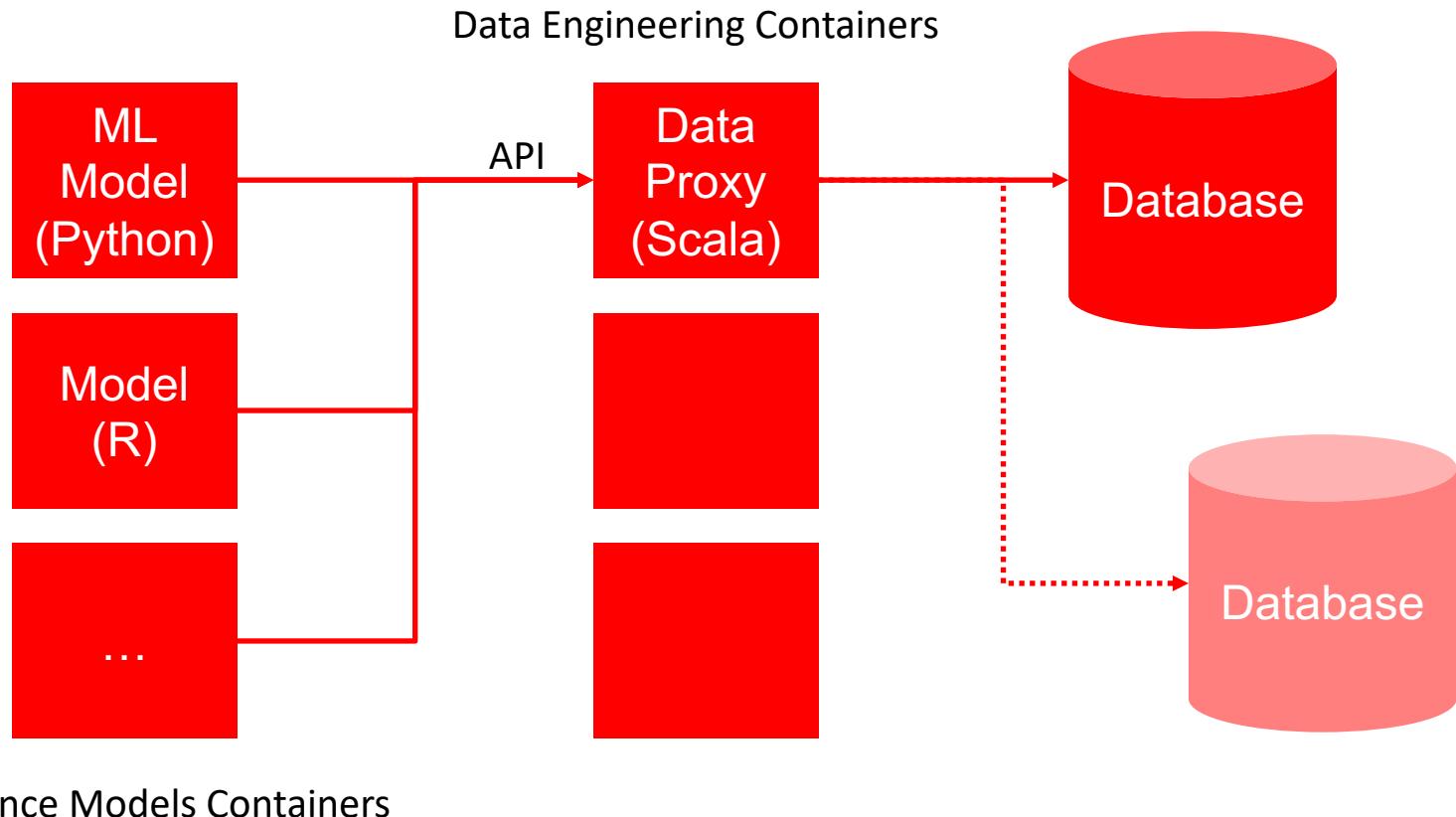


[Recent Changes](#)

Stage View



Technical patterns – Containers



In the end what should be the focus of a Data Engineer?

1. Create building blocks or abstractions that the Data Scientists can use autonomously.
2. Looking after the end to end development cycle of data applications.
3. Consulting to Data Scientists in case of more extreme cases.



Anti Patterns

1. Ratios of Data Engineers compared with Data Scientists.
2. Are your Data Scientist doing work that is Data Engineering or Operations?
3. Are you sure that you have Big Data? Scalability problems?



Questions

