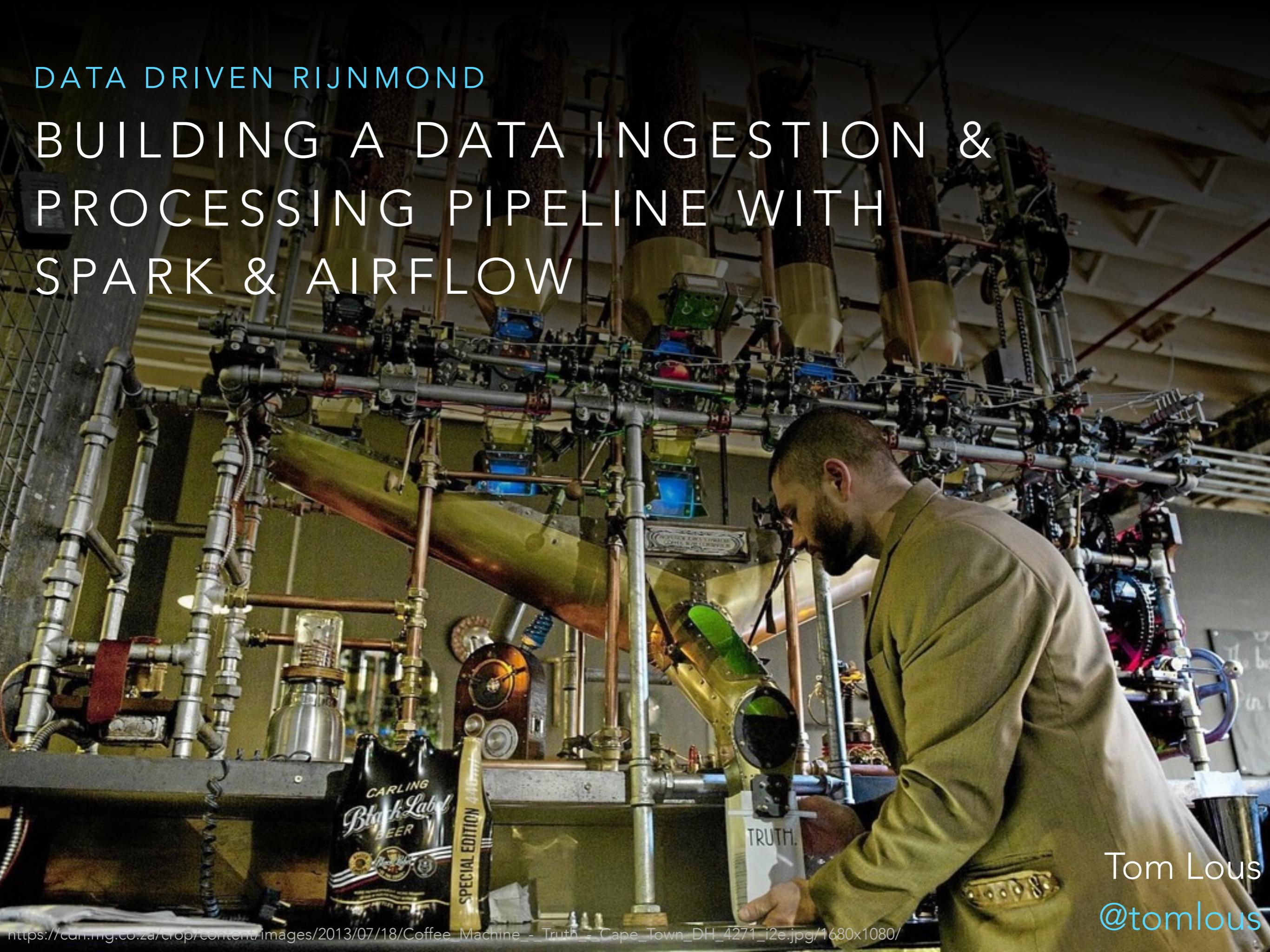


DATA DRIVEN RIJNMOND

# BUILDING A DATA INGESTION & PROCESSING PIPELINE WITH SPARK & AIRFLOW



Tom Lous  
@tomlous

WHO ARE WE?

# DATLINQ BY NUMBERS

- 14 countries

Netherlands, Belgium, France, UK, Germany, Spain, Italy, Luxembourg

- 100 customers across Europe

Coca-Cola, Unilever, FrieslandCampina, Red Bull, Spa, AB InBev & Heineken

- 2000 active users

Working with our SaaS, insights & mobile solutions

- 1.2M outlets

Restaurants, bars, hotels, bakeries, supermarkets, kiosks, hospitals

- 5000 mutations per day

Out of business, starters, ownership change, name change, etc

- 30 students

Manually checking many changes

WHAT DO WE DO?

# DATLINQ COLLECTS DATA

- Outlet features

Address, segmentation, opening hours, contact info, lat/long, terrace, parking, pricing

- Outlet surroundings

Nearby public transportation, school, going out area, shopping, high traffic

- Demographics

Median income, age, density, marital status

- Social Media

Presence, traffic, likes, posts



WHY DO WE DO IT?

# DATLINO'S MISSION

# FRIETWINKEL

by dapp  bio food

- Food service market transparent
- Effective marketing & sales
- Which food products are relevant where?

## SOLUTION

# WHY SPARK?

- Scaling is hard

Mo' outlets, mo' students

- Data != information

One-of analysis projects should be done continuously and immediately

- Work smarter not harder

Have students train AI instead of repetitive work

- Faster results

Distribute computing allows for increased processing speeds

## SOLUTION?



SOLUTION

# HOW TO START

- Focus on MVP

Minimal Viable Product that acts like a Proof of Concept

- Declare riskiest assumption

We can automatically match data from different sources

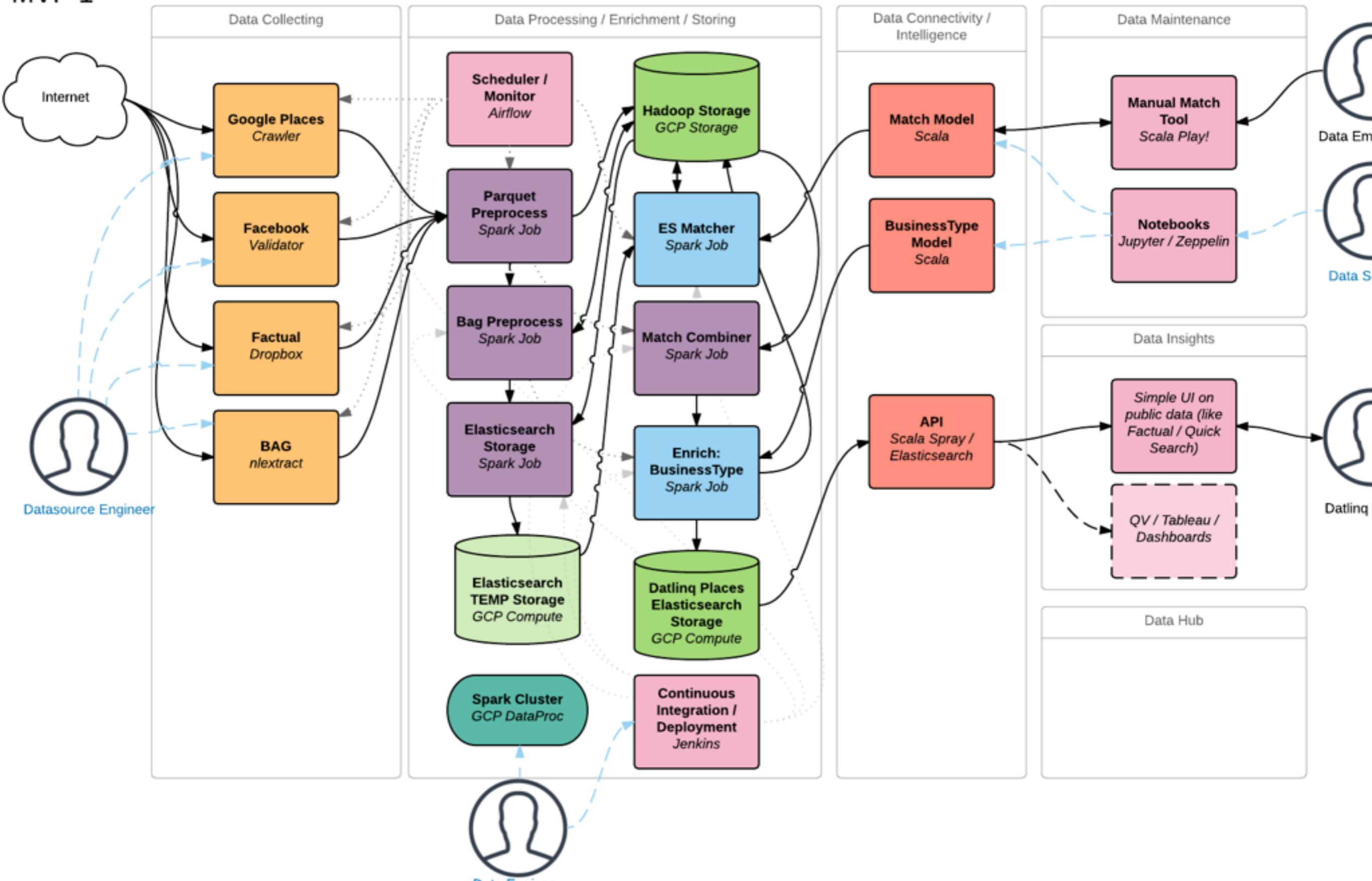
- Prove it

Small scale ML models



# BUILDING THE PIPELINE

MVP 1



# USE THE CLOUD

- Flexible usage

Pay for what you use

- Scaling out of the box

The only limit seems to be your credit card

- CLI > UI

CLI tools are your friend. UI is limited

```
tomlous@~/Development/Scala/datalabs-job/gcloud> gcloud
ERROR: (gcloud) too few arguments
Usage: gcloud [optional flags] <group | command>
  group may be      app | auth | components | compute | config |
                    container | dataflow | dataproc | datastore | debug |
                    deployment-manager | dns | iam | organizations |
                    projects | service-management | source | sql | topic
  command may be    docker | feedback | help | info | init | version
```

For detailed information on this command and its flags, run:  
gcloud --help

accept-cluster-sw-g8WX Preemptible Worker

Equivalent REST

## BUILDING THE PIPELINE

# DEVELOP LOCALLY

- Check versions

Your cloud provider sets the max version (e.g. Spark 2.0.2 instead of 2.1.0)

- Same code smaller data

Use samples of your data to test on

- Use at least 2 threads

Otherwise parallel execution cannot be truly tested

- Exclude libraries from deployment

Some libraries are provided on the cloud, don't add them to your final build

```
implicit val spark = SparkSession
  .builder()
  .master("local[2]")
  .appName(this.getClass.getSimpleName)
  .config("spark.ui.enabled", "false")
//    .config("spark.eventLog.dir", "/tmp/spark-events")
//    .config("spark.eventLog.enabled", "true")
  .getOrCreate()
```

```
val sparkDependencyScope = sys.props.getOrElse("sparkDependencyScope", default = "compile")

// spark
libraryDependencies += "org.apache.spark" % "spark-core" % "2.0.2" % sparkDependencyScope
libraryDependencies += "org.apache.spark" % "spark-sql" % "2.0.2" % sparkDependencyScope
libraryDependencies += "org.apache.spark" % "spark-streaming" % "2.0.2" % sparkDependencyScope
```

The screenshot shows a Java IDE interface with multiple tabs open. The main focus is a Scala code editor containing a build.sbt file. The file defines dependencies for Spark, Elasticsearch, and logging libraries, specifying versions and dependency scopes. A callout box highlights the sparkDependencyScope variable and its value of "2.0.2". Below the code editor is a terminal window showing the output of a build process, indicating 14 of 15 tests have been completed. The overall theme is local development and testing of a data pipeline application.

## BUILDING THE PIPELINE

# CONTINUOUS INTEGRATION

- Deploy code

Manual uploads or deployments is asking for trouble.

Use polling instead of pushing for security

Create pipeline script to stage builds

Automated tests

### Checkout

10s

### Build

5min 39s

Deploy  
accept

3s

```
1 node {
2     try {
3         stage("Checkout") {
4             git url: 'git@bitbucket.org:datling-datalabs/job.git'
5         }
6
7         stage("Build") {
8             wrap([$class: 'AnsiColorBuildWrapper', 'colorMapName': 'XTerm']) {
9                 sh "sbt -DsparkDependencyScope=provided -Dspark.es.cluster.name=eshub -Dspark.es.nodes=192.168.0.2 -Dspark.es.port=9200 -Dspark.es.cluster.port=9300 clean
10            }
11        }
12
13        stage("Deploy to accept") {
14            sh "git rev-parse HEAD > target/_GIT_REVISION"
15            sh "gsutil cp target/scala-2.11/job.jar gs://datling/accept/jobs/"
16            sh "gsutil cp target/_GIT_REVISION gs://datling/accept/jobs/"
17            slackSend (color: '#4F8A10', message: "SUCCESSFUL: Job '${env.JOB_NAME} [${env.BUILD_NUMBER}]' (${env.BUILD_URL})")
18        }
19
20    } catch (e) {
21        currentBuild.result = "FAILED"
22        slackSend (color: '#D8000C', message: "FAILED: Job '${env.JOB_NAME} [${env.BUILD_NUMBER}]' (${env.BUILD_URL})")
23        throw e
24    }
25 }
26 }
```

- Commit config

Store CI integration

## BUILDING THE PIPELINE

# SPARK JOBS

- All jobs in one project

Unless there is a very good reason, keep all your Jobs in one file

- Use Datasets as much as possible

Avoid RDD's and DataFrames

- Use Case Class (encoders) as much as possible

Generic Row object will only get you so far

- Verbosity only in main class

Log messages from nodes won't make it to the master

- Look at the execution plan

Remove temp files, clean the slate

```
= Optimized Logical Plan =
project [address_components#95, adr_address#141, formatted_address#166, formatted_phone_number#191, geometry#99, icon#216, id#101, int
- Project [address_components#95, UDF(adr_address#96) AS adr_address#141, UDF(formatted_address#97) AS formatted_address#166, UDF(f
 + LogicalRDD [_id#94, address_components#95, adr_address#96, formatted_address#97, formatted_phone_number#98, geometry#99, icon#10
= Physical Plan =
Project [address_components#95, adr_address#141, formatted_address#166, formatted_phone_number#191, geometry#99, icon#216, id#101, int
- *Project [address_components#95, UDF(adr_address#96) AS adr_address#141, UDF(formatted_address#97) AS formatted_address#166, UDF(f
 + Scan ExistingRDD[_id#94, address_components#95, adr_address#96, formatted_address#97, formatted_phone_number#98, geometry#99, icon#10
```

## BUILDING THE PIPELINE

## PARQUET

Buckets / daliq / accept / data / parquet / facebook\_places\_detail / PARTITION\_KEY=netherlands

- What is it

Apache Parquet is a columnar storage format available to any project in the Hadoop ecosystem

- Compressed

Small footprint.

- Columnar

Easier to drill down into data based on fields

- Metadata

Stores meta information easy to use by spark (counts, etc)

- Partitioning

Allows for partitioning on values for faster access

	Name	Size	Type
<input type="checkbox"/>	part-r-00000-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	51.79 MB	application/octet-stream
<input type="checkbox"/>	part-r-00001-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	52.19 MB	application/octet-stream
<input type="checkbox"/>	part-r-00002-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	50.36 MB	application/octet-stream
<input type="checkbox"/>	part-r-00003-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	50.36 MB	application/octet-stream
<input type="checkbox"/>	part-r-00004-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	49.83 MB	application/octet-stream
<input type="checkbox"/>	part-r-00005-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	47.77 MB	application/octet-stream
<input type="checkbox"/>	part-r-00006-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	14.87 MB	application/octet-stream
<input type="checkbox"/>	part-r-00007-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	27.67 MB	application/octet-stream
<input type="checkbox"/>	part-r-00008-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	37.4 MB	application/octet-stream
<input type="checkbox"/>	part-r-00009-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	33.15 MB	application/octet-stream
<input type="checkbox"/>	part-r-00010-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	39.91 MB	application/octet-stream
<input type="checkbox"/>	part-r-00011-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	35.5 MB	application/octet-stream
<input type="checkbox"/>	part-r-00012-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	10.27 MB	application/octet-stream
<input type="checkbox"/>	part-r-00013-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	30.56 MB	application/octet-stream
<input type="checkbox"/>	part-r-00014-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	30.85 MB	application/octet-stream
<input type="checkbox"/>	part-r-00015-967fb101-db9c-4efa-90cd-465b6af0b082.sn...	34.41 MB	application/octet-stream

## BUILDING THE PIPELINE

# ELASTICSEARCH

- NxN matching is computationally expensive

Querying Elasticsearch / Lucene / Solr less

Fuzzy searching using n-grams and other tokenizers

- Cluster out of the box

name\_whitespace: "Theater Aan Het Spui",

Auto replication

address\_locality: "DEN HAAG",

name\_ngram: "Theater Aan Het Spui"

- Config differs from normal use

It's not a search feature for a website.

Don't care about update performance

Don't load balance queries (\_local)

- Discovery Plugins for GC

So nodes can find each other on local network without Zen

Create VM's with --scopes=compute-rw

address\_street: "SPUI",

address\_locality: "DEN HAAG",

id\_uuid: "111c0065-2171-44b2-b55d-784eee7a5874",

name\_ngram: "Studio's Spui"

}

,

\_index: "unittest",

\_type: "google",

\_id: "76ae37dc-800b-4040-9cle-c2d2fd6ff6ff",

score: 5.179006

## BUILDING THE PIPELINE

# ANSIBLE

- Automate VM's / agent less  
Elasticsearch cluster
- Dependant on Ansible GC Modules

gcloud keeps updating. Ansible GC doesn't directly follow suite

```
23
24      name: Elasticsearch Link Cluster
25      hosts: localhost
26
27      vars:
28
29      tasks:
30          - name: create boot disks
31            gce_pd:
32                disk_type: pd-ssd
33                image: "tt image"
34                name: "{{ item.node }}-disk"
35                size_gb: 100
36                state: present
37                zone: "us-central1-a"
38            service_account_email: "{{ service_account_email }}"
39            credentials_file: "{{ credentials_file }}"
40            project_id: "{{ project_id }}"
41            with_items: "{{nodes}}"
42
43          - name: create instances
44            gce:
45                instance_names: "{{item.node}}"
46                zone: "{{ zone }}"
47                machine_type: custom-16-30720
48                preemptible: true
49                disk_auto_delete: true
50                disks:
51                    - name: "{{ item.node }}-disk"
52                      mode: READ_WRITE
53                state: present
54                service_account_email: "{{ service_account_email }}"
55                service_account_permissions: "compute-rw"
56                credentials_file: "{{ credentials_file }}"
57                project_id: "{{ project_id }}"
58                tags: "elasticsearch,es5,{{cluster_name}},link"
59                register: gce_raw_results
60                with_items: "{{nodes}}"
```

## BUILDING THE PIPELINE

# AIRFLOW

- Create workflows

Define all tasks and dependencies

- Cron on steroids

Starts every day

- Dependencies

Tasks are dependent on tasks upstream

- Backfill / rerun

Run as if in the past

facebook\_details\_to\_parquet

facebook\_to\_parquet

facebook\_parquet

google\_parquet

create\_spark\_cluster



AirFlow

DAGs

Tools ▾

Browse ▾

Admin ▾

Docs ▾

Configuration

Connections

Users

Reload DAGs

List (4)

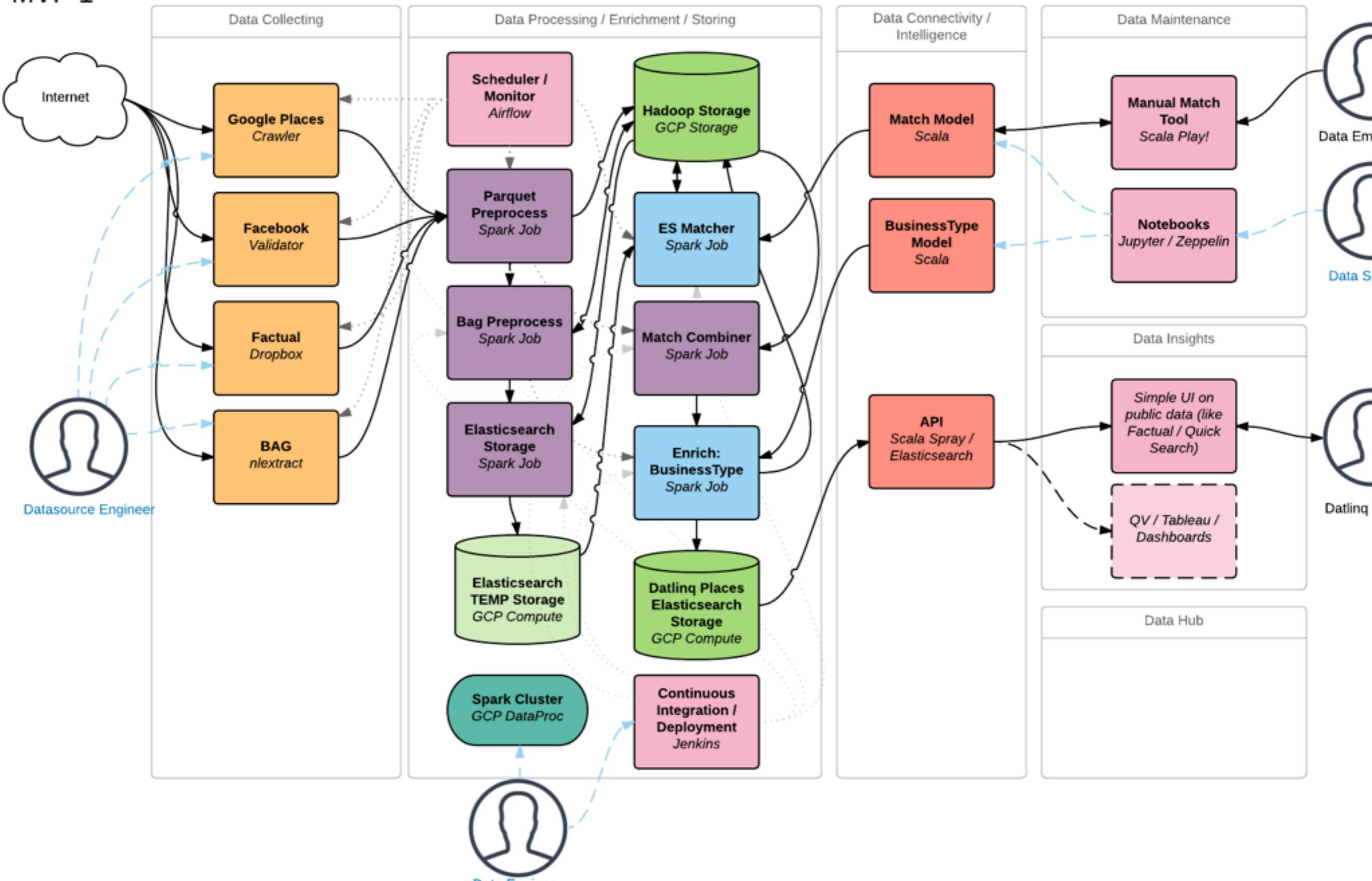
Create

With selected ▾

		Conn Id	Type
<input type="checkbox"/>		local_mysql	mysql
<input type="checkbox"/>		mysql_default	mysql
<input type="checkbox"/>		presto_default	presto
<input type="checkbox"/>		hive_default	hive

# BUILDING THE PIPELINE

MVP 1



## CONCLUSION

# FINALLY

- What's next?

- Improving MVP

- Adding more machine learning

- Adding more datasources

- Adding more countries

- Improving architecture

- Production

- Create an API for use in our products

- Create UI for exploring the data

- Exporting to standard databases to further use

- Team

- Hire Data Engineer(s)

- Hire Data Scientist(s)



## CONCLUSION

# THANKS FOR WATCHING

- Questions?

AMA

- Vacancies!

Ask me or check our site <http://www.datlinq.com/en/vacancies>

- Break

Grab a drink and let's reconvene after a 15 min break

