

Data Engineering and the Death of Visual ETL



Wednesday | April 1, 2020
1:00 PM

Stewart Bryson, Red Pill Analytics
@stewartbryson



Founder & CEO
Red Pill Analytics

twitter: @stewartbryson

medium: @stewartbryson

linkedin: stewartbryson



Your Data Warehouse is elastic.
Shouldn't your team be?

ELASTIC DELIVERY

Elastic Crew



MULTIFUNCTIONAL TEAM

Each Crew has the necessary players to move from question to insight in a single sprint.

AUTONOMOUS

Whether you have an existing team or not, a Crew works autonomously to solve whatever data problems you assign to it.

SCALABLE

Scale to answer all of your data requirements by adding additional Crews.

Elastic Surge



SUPPLEMENT YOUR TEAM

Increase your delivery capacity without adding fulltime resources.

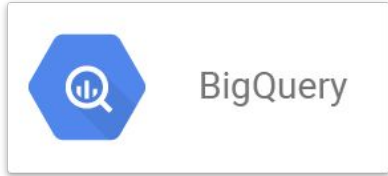
SEMI-PERMANENT

Designed as a long-term enhancement to organizations with functioning teams.

RESPOND TO DEMAND

Be more responsive to user demand: increase in size, swap out functional roles, or simply add a new technical skill for a few sprints.

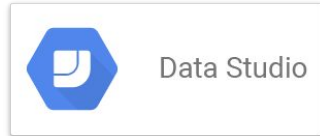
Data Warehouse



Amazon Athena



Analytics



Data Engineering / ETL



databricks




AWS Glue





The Rise of the Data Engineer

 Maxime Beauchemin January 20, 2017 · 12 min read





The Rise of the Data Engineer

 Maxime Beauchemin January 20, 2017 · 12 min read

<https://medium.com/free-code-camp/the-rise-of-the-data-engineer-91be18f1e603>



“There’s a multitude of reasons why complex pieces of software are not developed using drag and drop tools: it’s that ultimately *code is the best abstraction* there is for software.”





What's so
different
about
ETL?

Code *versus* Clicks

Source control
CI/CD
Collaboration
Code anywhere

Faster
Easier to start
Any discipline
Business friendly

A top-down photograph of a silver laptop on the left and a dark blue vintage typewriter on the right, both resting on a light-colored wooden desk. A semi-transparent dark grey horizontal band is overlaid across the center of the image, containing text.

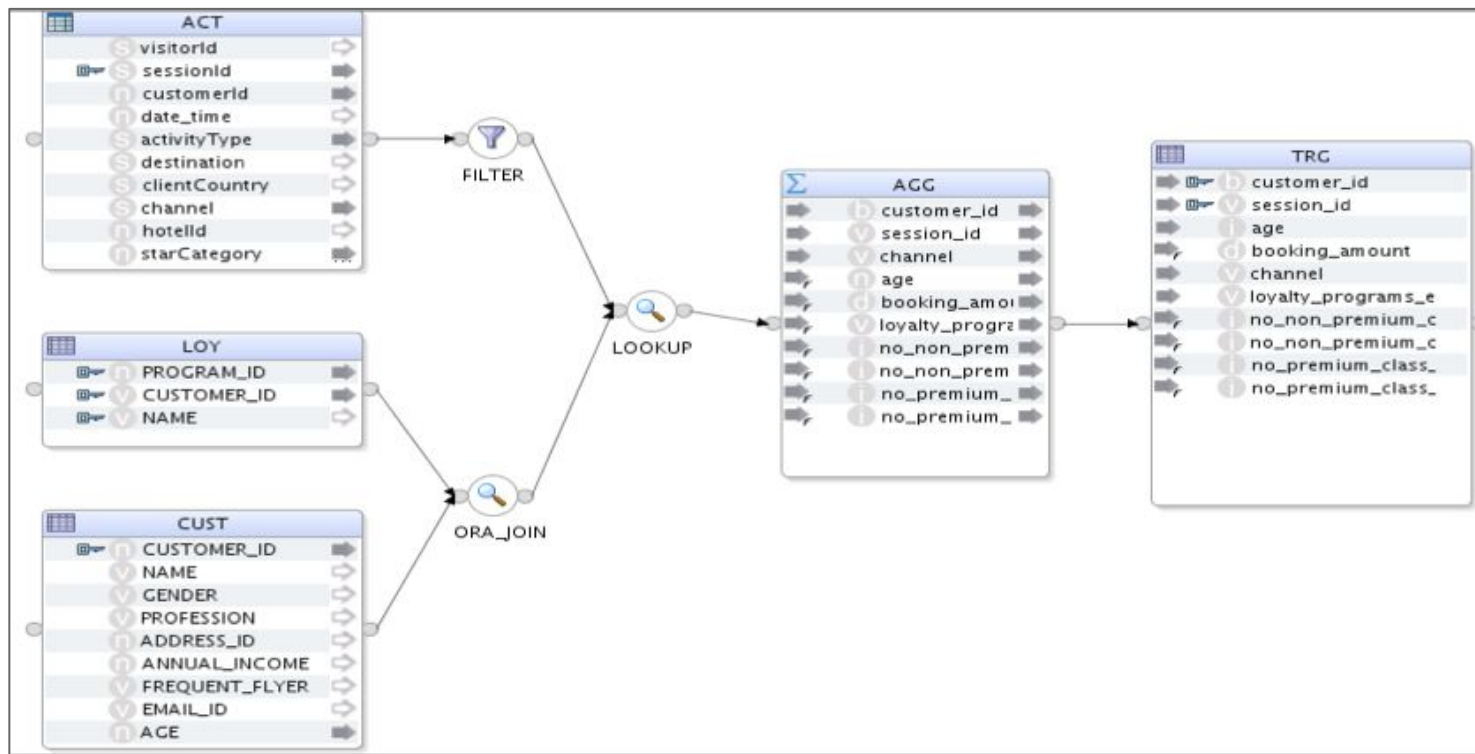
Assumption:
Code is better, but harder.

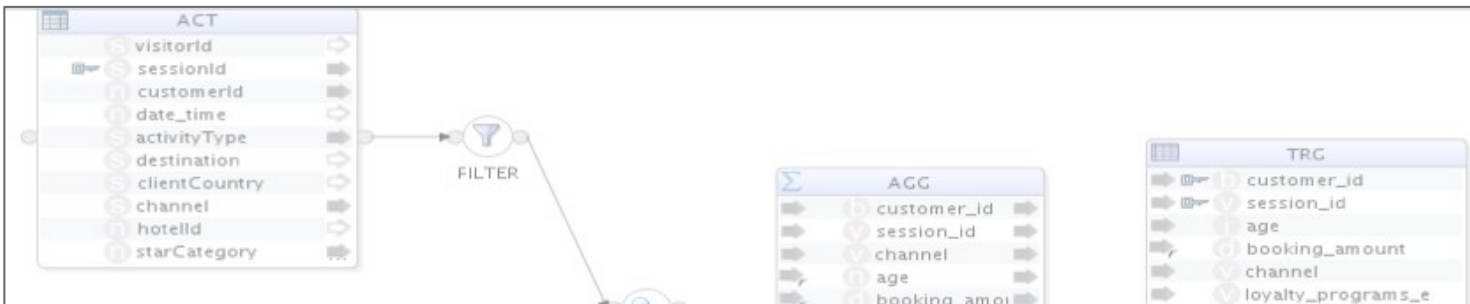
Data Warehouses

SQL >> GUI (ETL) >> GUI (ELT)

Modern Data Platforms

SQL >> GUI (ETL) >> GUI (ELT) >> Distributed >> SQL





snowflake | Databases | Shares | Warehouses | Worksheets | History | Partner Connect | Help | Stewart Bryson SYSADMIN

Worksheet 1 | Find database objects | | **Run** | All Queries | Saved a few seconds ago | Context: SYSADMIN | DATALOAD (S) | BROKER_DATA_D... | TRT | ...

```

1 select *,
2     rank() over (partition by customer_id order by action_ts) customer_rank
3 from CUSTOMER

```

Results | Data Preview | Open History

Query ID | SQL | 149ms | 50,000 rows

Filter result... | Download | Copy

Row	ACTION_TS	TIER	ACCOUNT_ID	STATE	LAST_NAME	FIRST_NAME	ACCOUNT_...	COUNTRY	ACCOU
1	2015-07-20 1...	NULL	24476	ME	Bezanson	Van	JXHMCAS...	United States	
2	2017-02-15 ...	1	29282	UT	Baudais	Kat	Bkx zMTYH...	United States	
3	2009-06-13 ...	3	5919	SK	Remers	Idus	xvOMLRKLS...	Canada	
4	2010-06-05 1...	NULL	8889	NULL	NULL	NULL	JjmVkihBjFQ...	NULL	
5	2010-12-27 ...	3	NULL	MT	NULL	NULL	NULL	United States	

Name	Datatype
PLAN_ACCO...	varchar2
PRODUCT	varchar2
PERIOD	num ber
FORECAST	num ber

Name	Datayn
FSC_LYEAR...	num b
START_DAT...	num b
END_DATE...	num b
NEXTYR_ST...	num b
NEXTYR_EN...	num b

Name	Datayn
PLANTO_WID	num b
PLANTO	varchl
PLANTO_KEY	varchl
PLANTO_NA...	varchl
PLANTO_CH...	varchl

Name	Datayn
LVIS_KEY	varchl
ITEMGRP	varchl
PROMOGRP	varchl
TRADEGRP	varchl
PROMOCAT	varchl

Edit Transformations

Transformation | Ports | Properties

Select transformation: **SQ SQ_HD**

Transformation type: Source Qualifier

Transformation Attributes

- Sql Query
- User Defined Join
- Source Filter
- Number Of Sorted Ports
- Tracing Level
- Select Distinct
- Pre SQL
- Post SQL
- Output is deterministic
- Output is repeatable

Sql Query

User-defined SQL statement

SQL Editor - SQ_HCGS_JOTS_DELL_FORECAST_STG (Expression)

Ports | Variables

- HCGS_JOTS_DELL_FORECAST_STG
- PLAN_ACCOUNT
- PRODUCT
- PERIOD
- FORECAST
- HMRM_FSC_L_WK_D
- HMRM_PLANTO_DH
- HMRM_SKU_DH

Instance Name: **HCGS_JOTS_DELL_FORECAST**

Transformation Type: **Source Definition**

SQL:

```
with week as
(select
count(fsc_l_yearweek) week_count
_fsc_l_yearper
from hmr.hmr_fsc_l_wk_d
where fsc_l_year > 2018
group by fsc_l_yearper
)
select
d.fsc_l_yearweek
.c.kev_wid
```

Connect to database:

ODBC data source: **DSDW (Oracle in OraClient12H)**

Use Kerberos Authentication

Username:

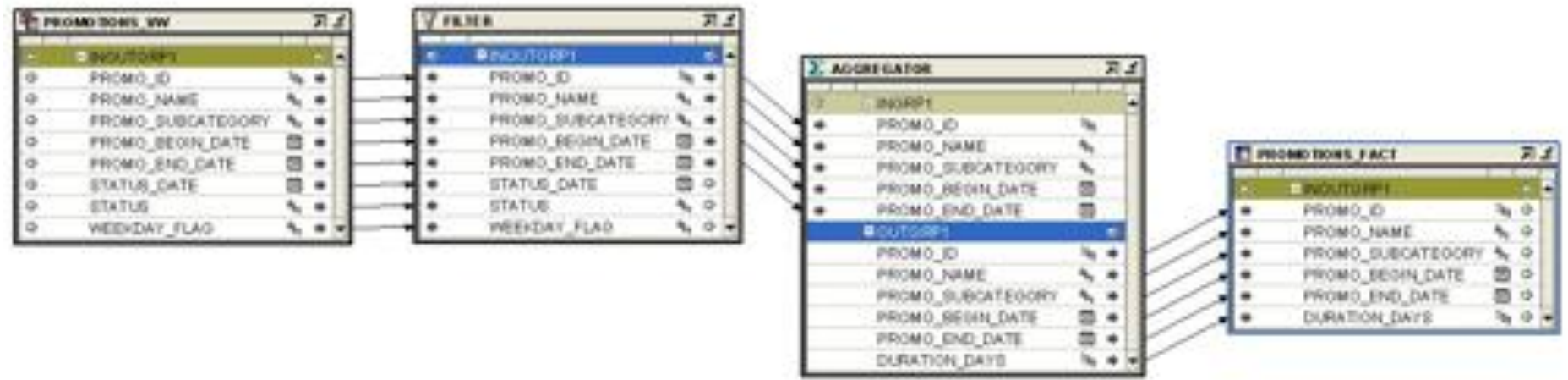
Password:

OK Cancel Generate SQL Validate Help

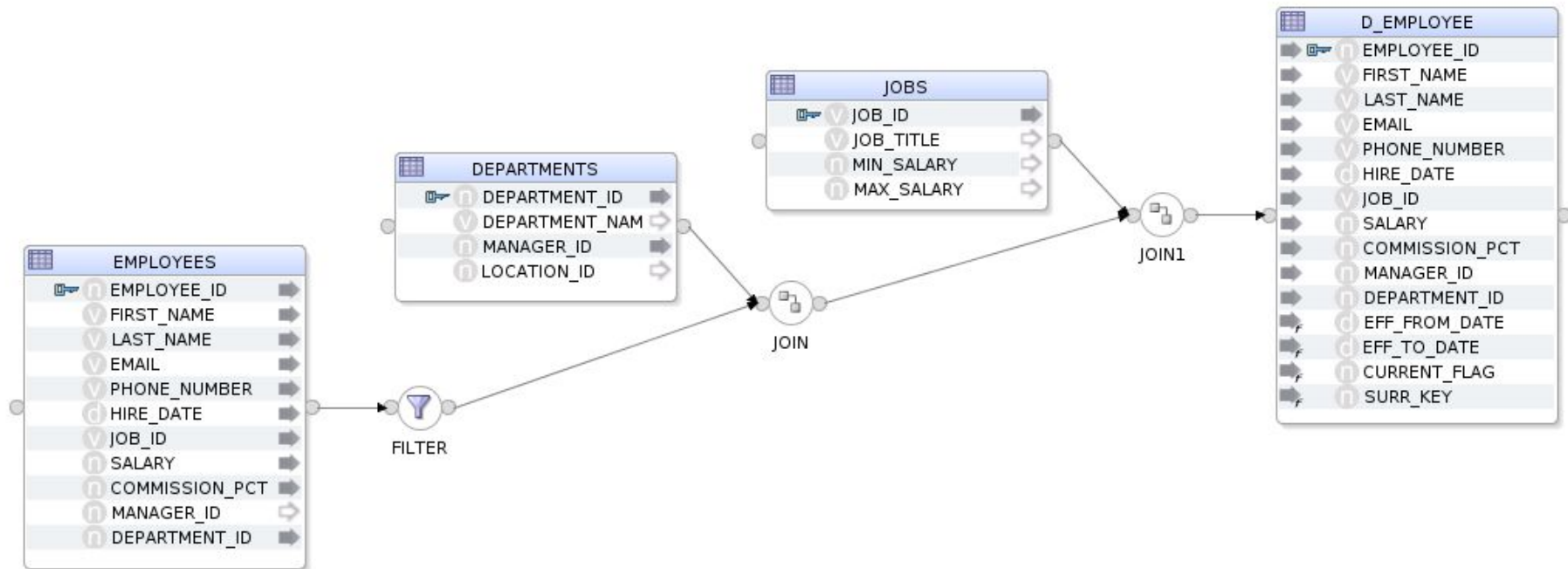
Port	Datatype	Length
FSC_LYEAR...	double	1...
F_WID	double	1...
PLANTO_WID	double	1...
FORECAST	double	1...
PLAN_ACCO...	string	2...
PRODUCT	string	1...
PERIOD	double	1...
FSC_LYEAR...	double	1...
FSC_LYEAR...	double	1...
PLANTO_KEY	string	3...
F_KEY	string	3...

Name	Datatype
FSC_LYEAR...	num ber
SKU_WID	num ber
PLANTO_WID	num ber
FORECAST	num ber

```
CREATE OR REPLACE VIEW stage.promotions_vw
AS
SELECT promo_id,
       promo_name,
       promo_subcategory,
       promo_begin_date,
       promo_end_date,
       time_id status_date,
       status,
       weekday_flag
FROM ( SELECT time_id,
             CASE
             WHEN day_name IN ('Saturday','Sunday') THEN 'N'
             ELSE 'Y'
             END weekday_flag,
             promo_id,
             MAX(promo_name) OVER (partition BY promo_id) promo_name,
             MAX(promo_subcategory) OVER (partition BY promo_id) promo_subcategory,
             MAX(promo_begin_date) OVER (partition BY promo_id) promo_begin_date,
             MAX(promo_end_date) OVER (partition BY promo_id) promo_end_date,
             nvl(status, 'ongoing') status
       FROM stage.promotions_stg partition BY (promo_id)
       right outer JOIN sh.times
       ON time_id=status_date)
WHERE time_id BETWEEN promo_begin_date AND promo_end_date
```



```
mapping.create('MY_PROJECT', 'DEMO_FOLDER', 'EMPLOYEE_DIM_LOAD')
    .datastores([
        [name: "HR.EMPLOYEES"],
        [name: "HR.DEPARTMENTS"],
        [name: "HR.JOBS"],
        [name: "PERF.D_EMPLOYEE", integration_type: "SCD"],
    ])
    .select("EMPLOYEES")
    .filter('NAME_FILTER', [filter_condition: "EMPLOYEES.FIRST_NAME LIKE 'D%' " ])
    .join('EMP_DEPT', ['DEPARTMENTS'],
        [join_condition: "EMPLOYEES.DEPARTMENT_ID = DEPARTMENTS.DEPARTMENT_ID" ])
    .join('DEPT_JOBS', ['JOBS'],
        [join_condition: "EMPLOYEES.JOB_ID = JOBS.JOB_ID" ])
    .connect("D_EMPLOYEE",
        [[ attr: "employee_id", key_indicator: true ],
        [ attr: "eff_from_date", expression: "sysdate", execute_on_hint: "TARGET" ],
        [ attr: "eff_to_date", expression: "sysdate", execute_on_hint: "TARGET" ],
        [ attr: "current_flag", expression: 1, execute_on_hint: "TARGET" ],
        [ attr: "surr_key",
            expression: ":RM_PROJECT.D_EMPLOYEE_SEQ_NEXTVAL",
            execute_on_hint: "TARGET" ],])
    .commit()
    .validate()
```



```
SELECT first_name,  
       last_name,  
       email,  
       phone_number,  
       hire_date,  
       job_id,  
       manager_id,  
       department_id,  
       salary,  
       commission_pct,  
       eff_from_date,  
       eff_to_date,  
       current_flag  
FROM EMPLOYEES  
JOIN DEPARTMENTS using (DEPARTMENT_ID)  
JOIN JOBS using (JOB_ID)
```

Code >> GUI >> Code

I would *rather* write this.

```
SELECT first_name,  
       last_name,  
       email,  
       phone_number,  
       hire_date,  
       job_id,  
       manager_id,  
       department_id,  
       salary,  
       commission_pct,  
       eff_from_date,  
       eff_to_date,  
       current_flag  
FROM EMPLOYEES  
JOIN DEPARTMENTS using (DEPARTMENT_ID)  
JOIN JOBS using (JOB_ID)
```




Navigation

[Installation](#)

[Usage instructions](#)

[API Reference](#)

[History](#)

Quick search

StreamSets SDK for Python

The StreamSets SDK for Python enables users to interact with StreamSets products programmatically using Python 3.4+. As an example, with a running instance of StreamSets Data Collector, you can create and import a functional pipeline in less than 10 lines of code:

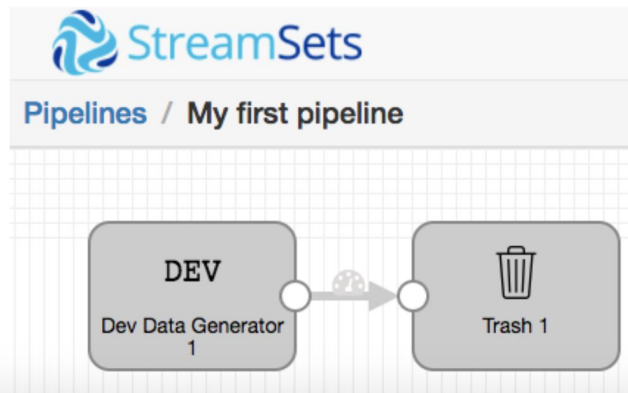
```
from streamsets.sdk import DataCollector
server_url = 'http://localhost:18630'
data_collector = DataCollector(server_url)

builder = data_collector.get_pipeline_builder()
dev_data_generator = builder.add_stage('Dev Data Generator')
trash = builder.add_stage('Trash')

dev_data_generator >> trash # connect the Dev Data Generator origin

pipeline = builder.build('My first pipeline')
data_collector.add_pipeline(pipeline)
```

The resulting pipeline can be examined by opening the StreamSets Data Collector user interface and selecting the My first pipeline pipeline:





StreamSets SDK for Python

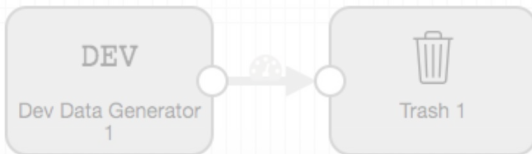
The StreamSets SDK for Python enables users to interact with StreamSets products programmatically using Python 3.4+. As an example, with a running instance of StreamSets Data Collector, you can create and import a functional pipeline in less than 10 lines.

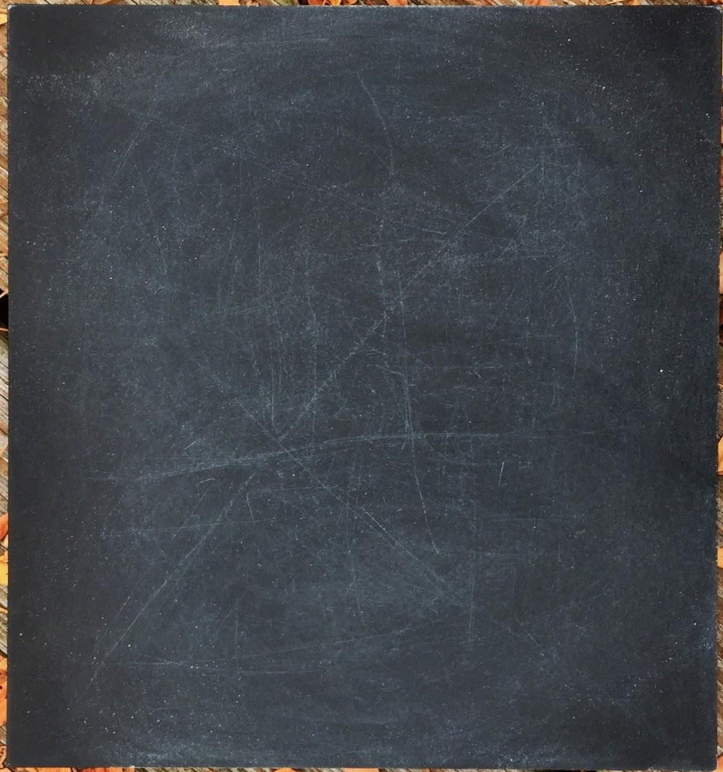
```
from streamsets.sdk import DataCollector
server_url = 'http://localhost:18630'
data_collector = DataCollector(server_url)

builder = data_collector.get_pipeline_builder()
dev_data_generator = builder.add_stage('Dev Data Generator')
trash = builder.add_stage('Trash')

dev_data_generator >> trash # connect the Dev Data Generator origin

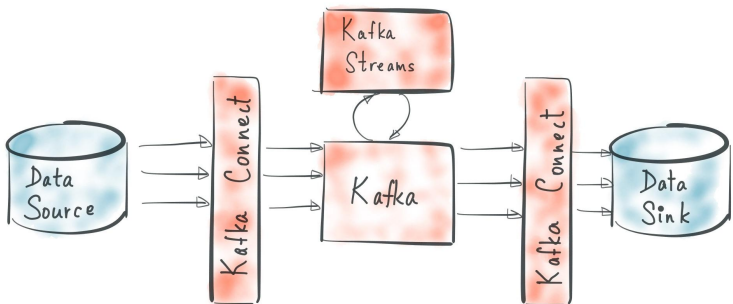
pipeline = builder.build('My first pipeline')
data_collector.add_pipeline(pipeline)
```



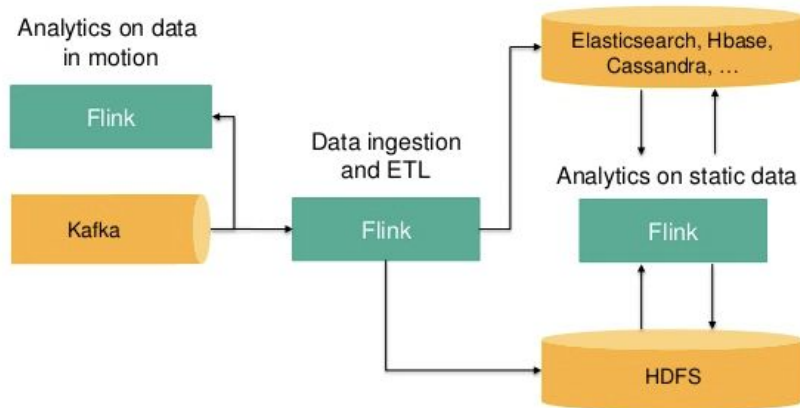
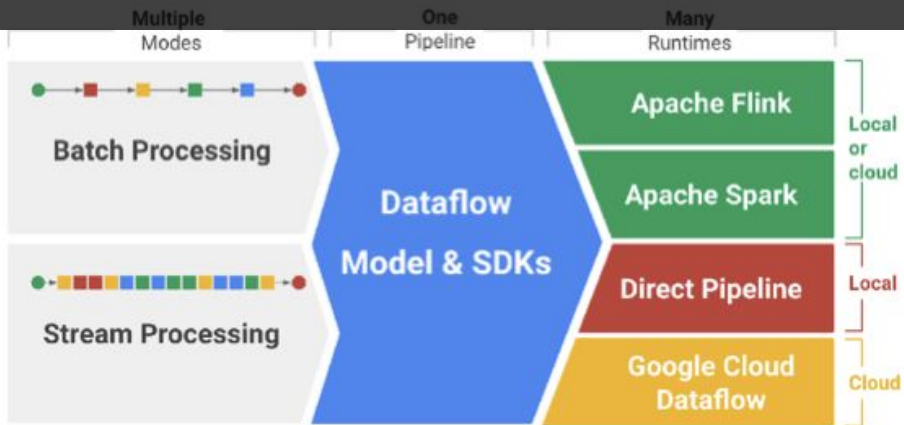


We aren't starting with a *blank* slate.

KAFKA CONNECT + STREAMS



An *embarrassment* of riches.

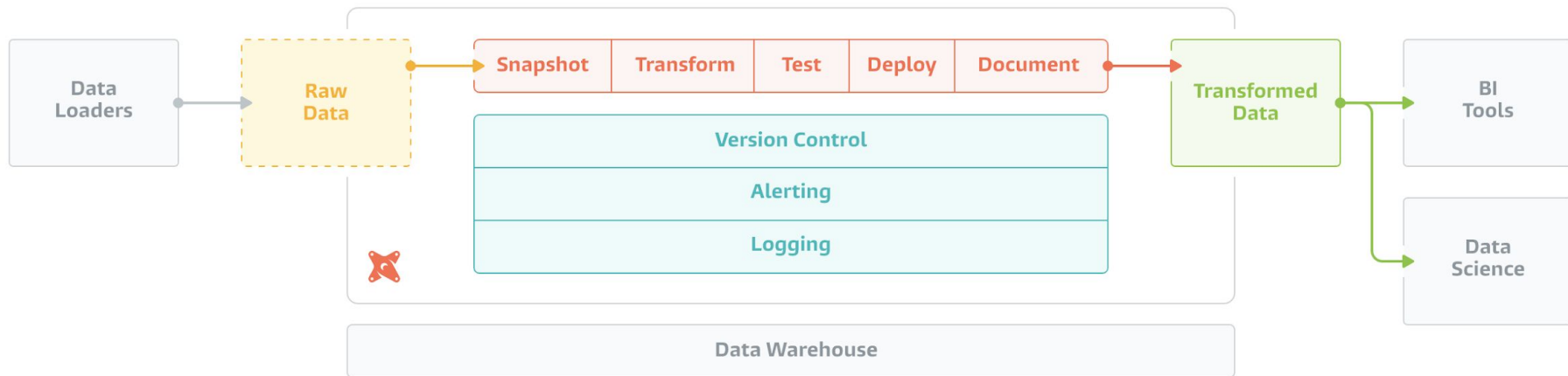


Shouldn't my *pipelines* just
live in my data warehouse?

Data Build Tool (*dbt*)

Analytics engineering tools designed for analysts

dbt is a command line tool that speaks the preferred language of data analysts everywhere—SQL. With dbt, analysts take ownership of the entire analytics engineering workflow, from writing data transformation code to deployment and documentation.



Dataform

Dataform SDK

Develop reliable and scalable SQL pipelines

Open source. Run Dataform locally or anywhere.

Manage dependencies between your tables.

Reuse code across all your scripts.

Write tests to assert your data quality.

```
# Install Dataform CLI
npm i -g @dataform/cli

# Create new project
dataform init snowflake my_project

# Create new table
echo 'config { type: "table" } SELECT 1 as one'
> my_table.sqlx

# Run project
dataform run my_project
```

Read the [docs](#) or view on [GitHub](#)

Dataform web

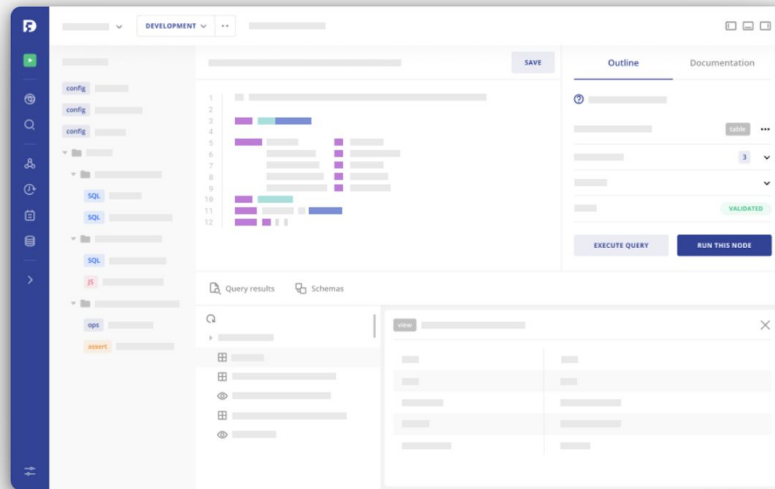
A complete solution for data warehouse management

Develop as a team in a collaborative web environment.

Version control all your code with a native GitHub integration.

Orchestrate your pipelines and get alerted if anything goes wrong.

Share a data catalog with all your data definitions.



[Learn more about Dataform web](#) →

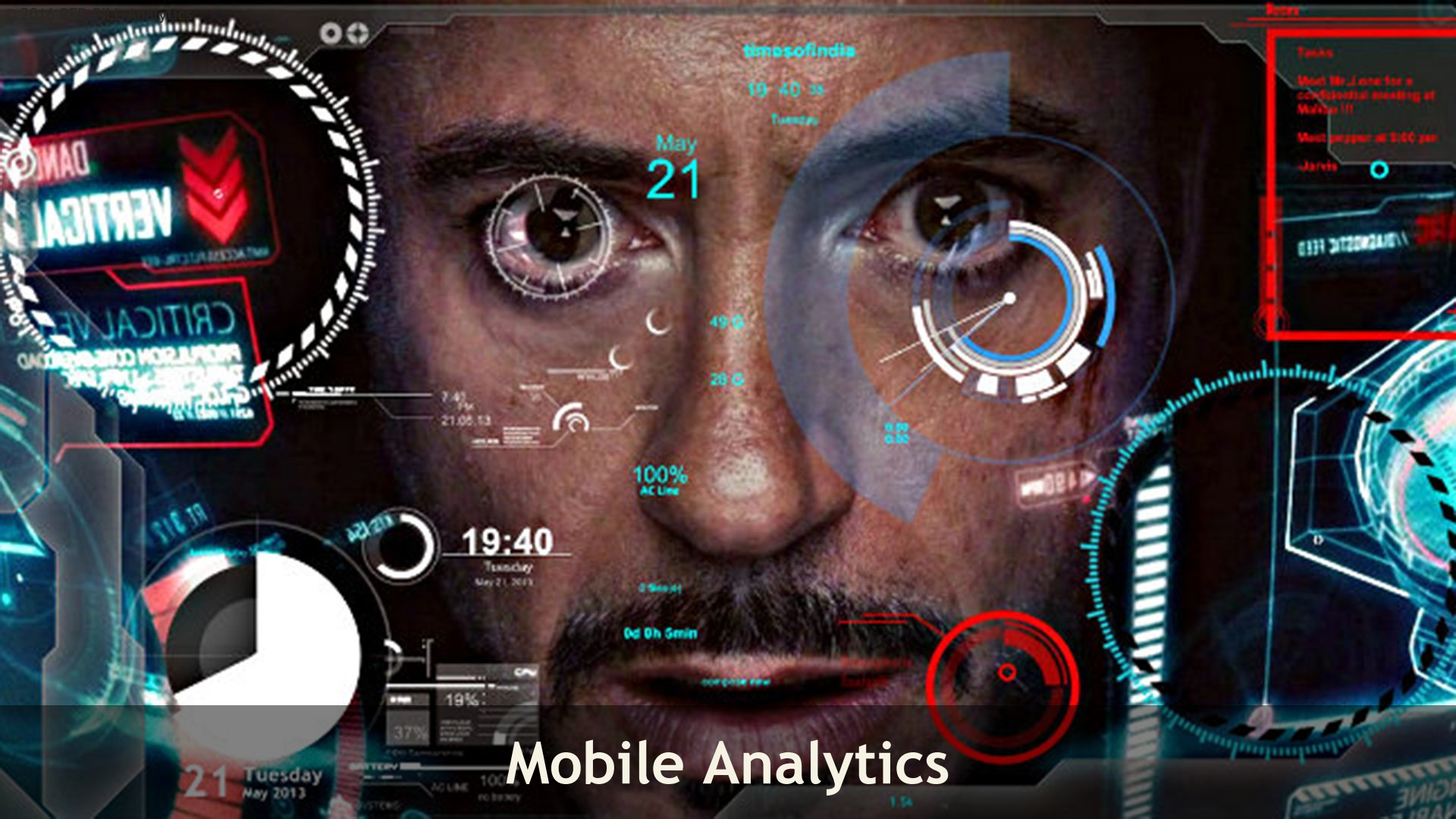
Am I feeding other
downstream systems?

Data warehouses *typically*
aren't very good at that.

What are some other *use cases* we might want to consider?

Realtime Events





timesofindia

19:40:35

Tuesday

May
21

49.0

28.0

100%
AC Line

0.50

19:40

Tuesday
May 21, 2013

0.50

0d 0h 5min

0.50

Tasks

Meet Mr. Jones for a
confidential meeting at
Mulberry III

Meet supplier at 3:00 pm

Jarvis

VERTICAL FEED

21
Tuesday
May 2013

Mobile Analytics

1.54

Search





Machine *Learning*

Why *Spark*?

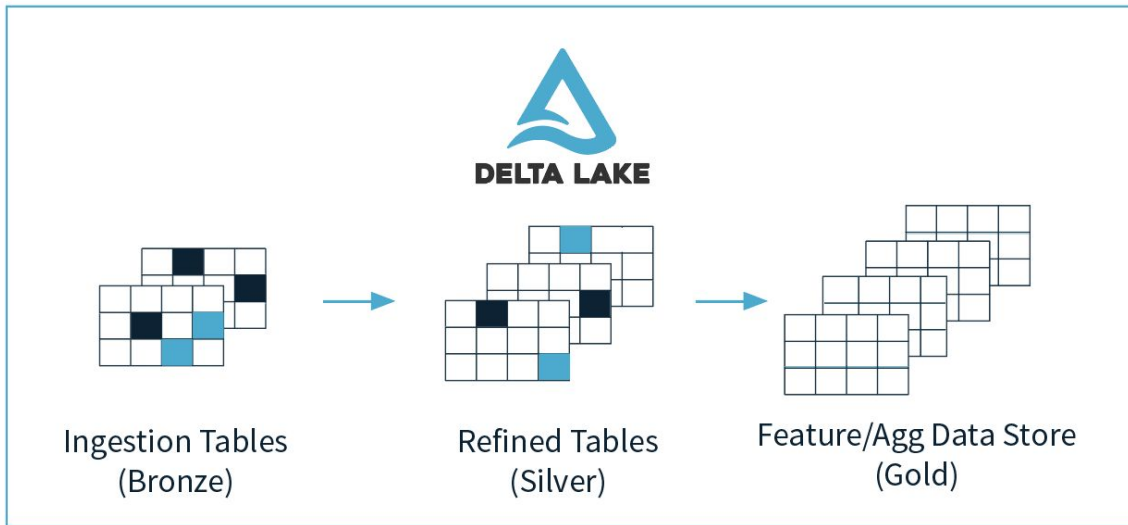
Why *Databricks*?



Streaming →



Batch →



Ingestion Tables
(Bronze)

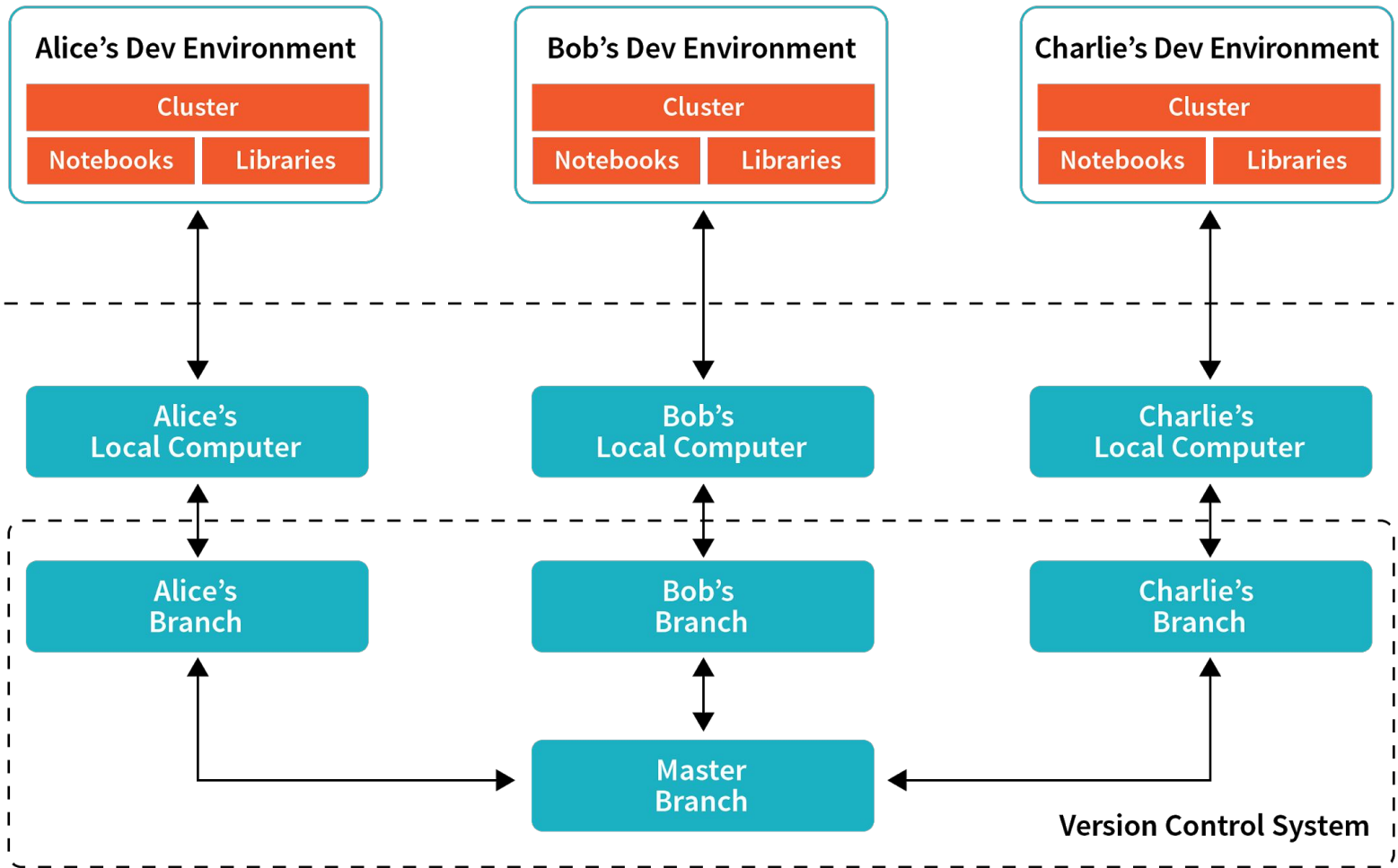
Refined Tables
(Silver)

Feature/Agg Data Store
(Gold)

→ Analytics
and Machine
Learning

Your Existing Data Lake







Demonstration

YOU CAN
CHOOSE
TO SEE
DATA
DIFFERENTLY.



RED
PILL
ANALYTICS

RedPillAnalytics.com

[@RedPillA](https://twitter.com/RedPillA)