# The Dataflow Job Builder
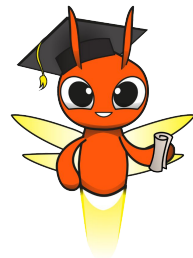
Ryan Madden

# Agenda

- The Job Builder & Beam YAML
- Why use the Job Builder?
- Tour of the Job Builder
- Working with YAML

# What is the Dataflow Job Builder?

Create custom batch and streaming jobs using the Dataflow job builder and YAML editor. You can transform, filter, and combine data using Python and SQL. Use the builder form to compose your job steps or edit the YAML specification directly. Learn more ⎘

**BUILDER FORM** | YAML EDITOR | LOAD BLUEPRINTS ▾

Job name *

Job names can contain lowercase letters, numbers, and dashes

**Job type**

⦿ Batch

◯ Streaming

## Sources

Read data from BigQuery, Pub/Sub, or Cloud Storage
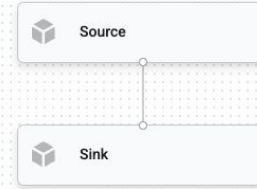
⌃ New source 🗑

Source name *
Source

Source type * ▾

DONE

ADD A SOURCE

## Transforms

Optionally manipulate, aggregate, and join data from sources and transforms

ADD A TRANSFORM



Source

Sink

# What is Beam YAML?

Write:

```yaml
pipeline:
  transforms:
    - type: Create
      config:
        elements: [1, 2, 3]
    - type: LogForTesting
        input: Create
```

Run:

```
python -m apache_beam.yaml.main --yaml_pipeline_file=pipeline.yaml
```

# Why use the Dataflow Job Builder?

## New source

**Source name ***

Source

**Source type ***

| Pub/Sub Topic |
| --- |
| Pub/Sub Subscription |
| BigQuery Table |
| CSV from Cloud Storage |
| JSON from Cloud Storage |
| Text files from Cloud Storage |
| YAML Source |
| JDBC |

Tra

Optic

ADD A TRANSFORM

## New source

**Source name ***

Source

**Source type ***

CSV from Cloud Storage ▼

📄 gs:// CSV location *                                             BROWSE

The location of your CSV file(s) in Cloud Storage. You can use wildcards to specify multiple files. Learn more ↗

**PREVIEW SOURCE DATA**

**CSV delimiter character**

Comma ▼

Character used to parse CSV fields

**DONE**

∧ **New source**                                                    🗑

Source name *
Source

Source type *
BigQuery Table                                                      ▾

BigQuery table *
☑ bigquery-public-data.samples.shakespeare          BROWSE

PREVIEW SOURCE DATA

**Table fields ***

◯ Read all table fields

◉ Specify table fields

Table fields *

⊟ 1 of 4 selected

☑ **word**
Type: STRING

☐ **word_count**
Type: INTEGER

☐ **corpus**
Type: STRING

☐ **corpus_date**
Type: INTEGER

CANCEL       OK

ADD A TRANSFORM

## New transform

**Transform name ***

Transform

**Transform type ***

| Filter (Python) |
| Filter records with a Python expression |

SQL Transform

Manipulate records or join multiple inputs with a SQL statement

Map Fields (Python)

Add new fields or re-map entire records with Python expressions and functions

Map Fields (SQL)

Add or map record fields with SQL expressions

YAML transform

Use any transform from the Beam YAML SDK

## New transform

Transform name *

Transform

Transform type *

Filter (Python)

Python filter expression *                                    ?

Records that do not match the expression will be filtered out

☐ Handle errors

Route records that throw exceptions when processed to a separate step output. This output m
transform or sink.

Input step for the transform *

Python filter expressions can filter records by                 ✕
performing checks against the record's fields. For
example:

```
myField == 'value'
myField > 0 and myField < 10
myField != False
```

DONE

## New transform

**Transform name ***

Transform

**Transform type ***

Map Fields (Python) ▼

☑ **Preserve existing fields**
Append new fields to the record. If unchecked, existing fields are dropped.

## Mapped fields

ADD A FIELD

## Dropped fields

+ ADD A FIELD

## Python dependencies

Add PyPI packages as dependencies to make them available to import in callables.

+ ADD A DEPENDENCY

☐ Handle errors
Route records that throw exceptions when processed to a separate step output. This output must be consumed by a downstream transform or sink.

**Input step for the transform ***

Source ▼

DONE

## New field

**Field name \***

year

☑ Callable
If checked, the field value must be a Python function rather than a simple expression

**Python callable** ❓
Python function that defines the field value

Press Option+F1 for Accessibility Options.

```
1  import roman
2  def convert(row):
3    return roman.toRoman(row.year)
```

DONE

ADD A FIELD

## Python dependencies

Add PyPI packages as dependencies to make them available to import in callables.

**Dependency 1 \***

roman>=4.2

The requirement specifier with an optional version specifier

➕ ADD A DEPENDENCY

## Edit transform

**Transform name ***

Join

**Transform type ***

SQL Transform ▼

**Input steps for the transform ***

Taxi ride data and Zone lookup data ▼

### Input aliases

If a single input is selected, it can be referenced in the query as PCOLLECTION. If multiple inputs are selected, aliases for each input will appear above the query editor.

Taxi ride data: input0    Zone lookup data: input1

Press Option+F1 for Accessibility Options.

```
1    SELECT * FROM input0
2    JOIN input1 ON input0.PULocationID = input1.LocationID
```

A Beam Calcite SQL ↗ query to be run over one or more inputs

DONE

## New sink

Sink name *

Sink

Sink type *

| Pub/Sub Topic |
| --- |
| BigQuery Table |
| CSV files on Cloud Storage |
| JSON files on Cloud Storage |
| Text files on Cloud Storage |
| YAML Sink |
| Log Sink |

# Working with YAML

# Working with YAML

Create custom batch and streaming jobs using the Dataflow job builder and YAML editor. You can transform, filter, and combine data using Python and SQL. Use the builder form to compose your job steps or edit the YAML specification directly. Learn more ☑

**BUILDER FORM** | **YAML EDITOR** | LOAD BLUEPRINTS ▾

Job name *
word-count

Job names can contain lowercase letters, numbers, and dashes

Press Option+F1 for Accessibility Options.

```
 1  pipeline:
 2    transforms:
 3      - type: ReadFromText
 4        name: Read from GCS
 5        config:
 6          path: gs://dataflow-samples/shakespeare/kinglear.txt
 7      - type: MapToFields
 8        name: Split words
 9        config:
10          language: python
11          fields:
12            word:
13              callable: |
14                import re
15                def my_mapping(row):
16                  return re.findall(r'[A-Za-z\']+', row.line.lower())
17            value: "1"
18        input: Read from GCS
19      - type: Explode
20        name: Explode word arrays
21        config:
22          fields: [word]
23        input: Split words
24      - type: Combine
```

**Read from GCS**
Text files from Cloud Storage

**Split words**
Map Fields (Python)

**Explode word arrays**
Explode

**Count words**
Group by

**Format output**
Map Fields (Python)

Beam College 2025

# YAML Providers

# YAML Providers

# Getting Started



Create custom batch and streaming jobs using the Dataflow job builder and YAML editor. You can transform, filter, and combine data using Python and SQL. Use the builder form to compose your job steps or edit the YAML specification directly. Learn more ⧉

**BUILDER FORM**　**YAML EDITOR**　**LOAD BLUEPRINTS** ▾

Job name *

Job names can contain lowercase letters, nu

**Job type**
◉ Batch
◯ Streaming

**Sources**

Read data from BigQuery, Pub/Sub, or Cloud St

⌃　New source　🗑

Source name *
Source

Source type *　▾

DONE

ADD A SOURCE

Load Blueprints dropdown:
- Text files on Cloud Storage to BigQuery
- CSV files on Cloud Storage to BigQuery
- Pub/Sub topic to BigQuery
- Pub/Sub subscription to BigQuery
- Pub/Sub subscription to Pub/Sub topic
- SQL Server to BigQuery
- PostgreSQL to BigQuery
- MySQL to BigQuery
- Oracle to BigQuery
- Word Count

Patterns
- Stream Pub/Sub messages to BigQuery
- Filter data with Python
- Map data with Python dependencies
- Join sources with SQL
- Divert data with error handling
- Log data to worker logs

# Demo!

# Learn more

- Dataflow Job Builder

  - http://console.cloud.google.com/dataflow/createjob;createMode=builder

- Job Builder documentation

  - https://cloud.google.com/dataflow/docs/guides/job-builder

- Beam YAML documentation

  - https://beam.apache.org/documentation/sdks/yaml/