# Beam Notebooks

Reza Rokni

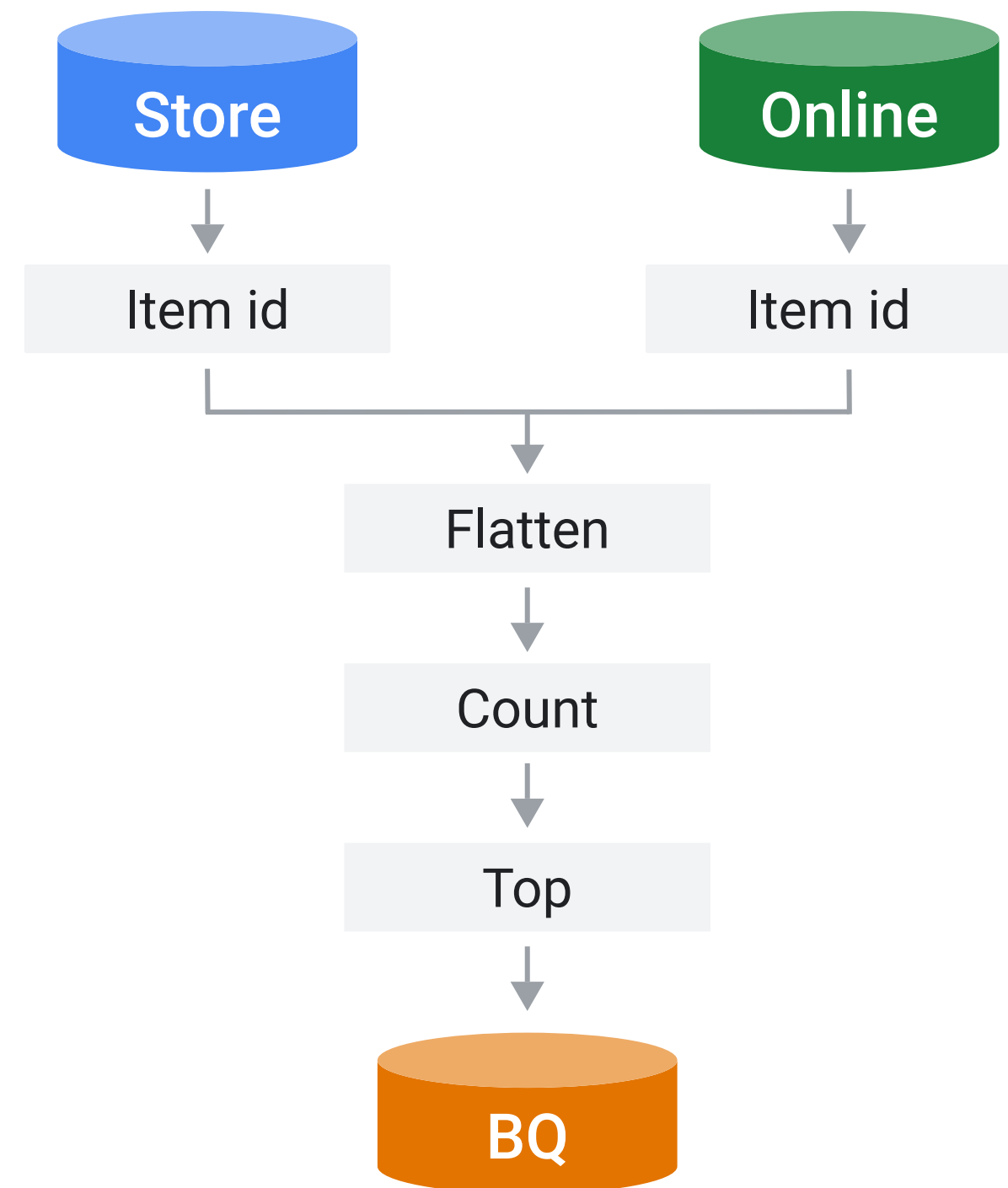Developer Advocate,
Google Cloud

# Agenda

# Apache Beam and interactive development

```
storeSales = p  |
beam.io.ReadFromText("purchases-store")

              | beam.Map(lambda s: ...)

onlineSales = p |
beam.io.ReadFromText("purchases-online")

              | beam.Map(lambda s: ...)

topSales = (storeSales, onlineSales)

              | beam.Flatten()

              | beam.Combiners.Count.perKey()

              | beam.Combiners.Top.of(10, key = lambda
x: x[1])

topSales        | beam.io.WriteToBigQuery(topSales)
```

# Apache Beam interactive runner

The interactive runner module allows

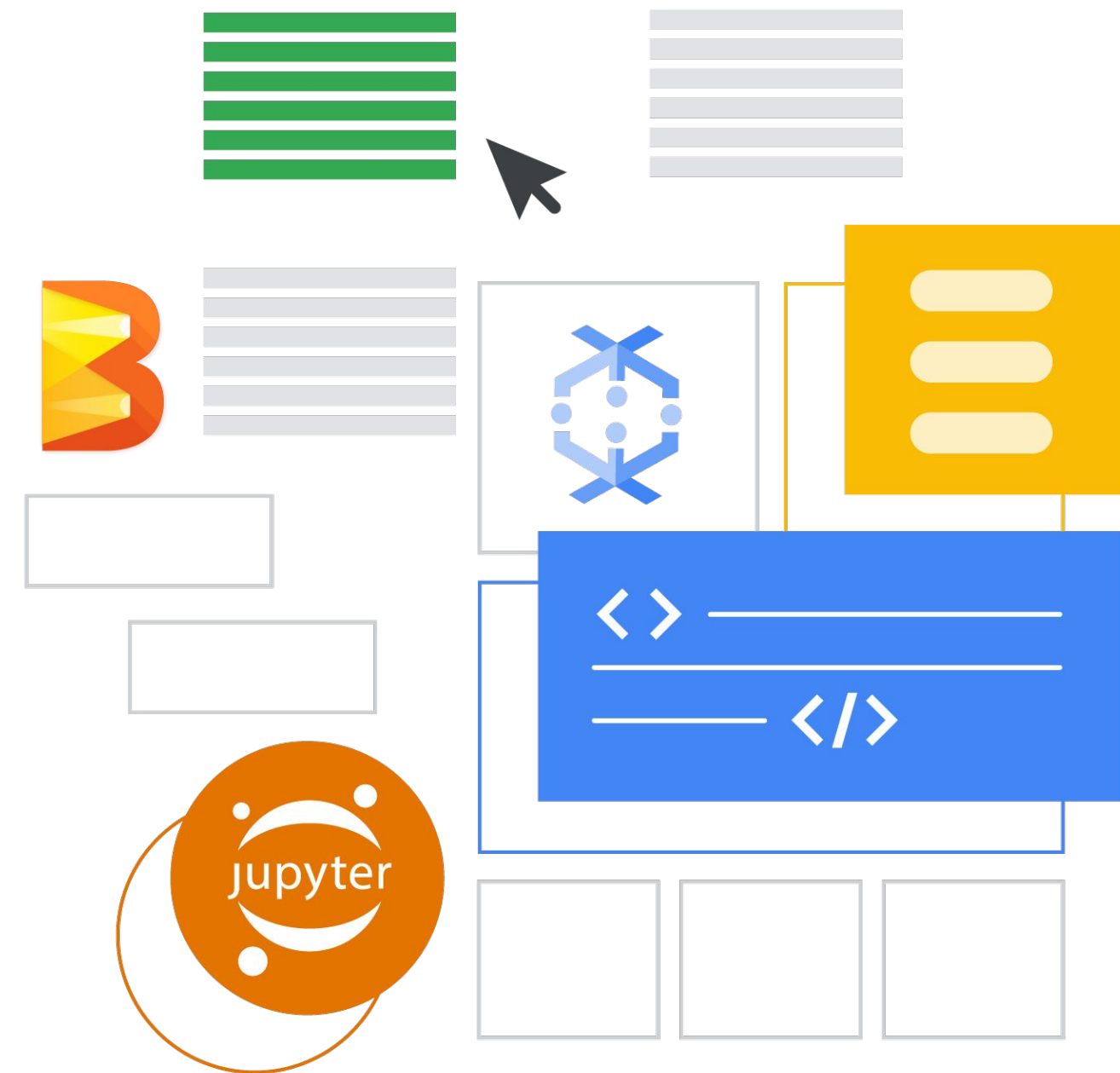**1**  Interactive development

# Apache Beam interactive runner

The interactive runner module allows

**1**  Interactive development

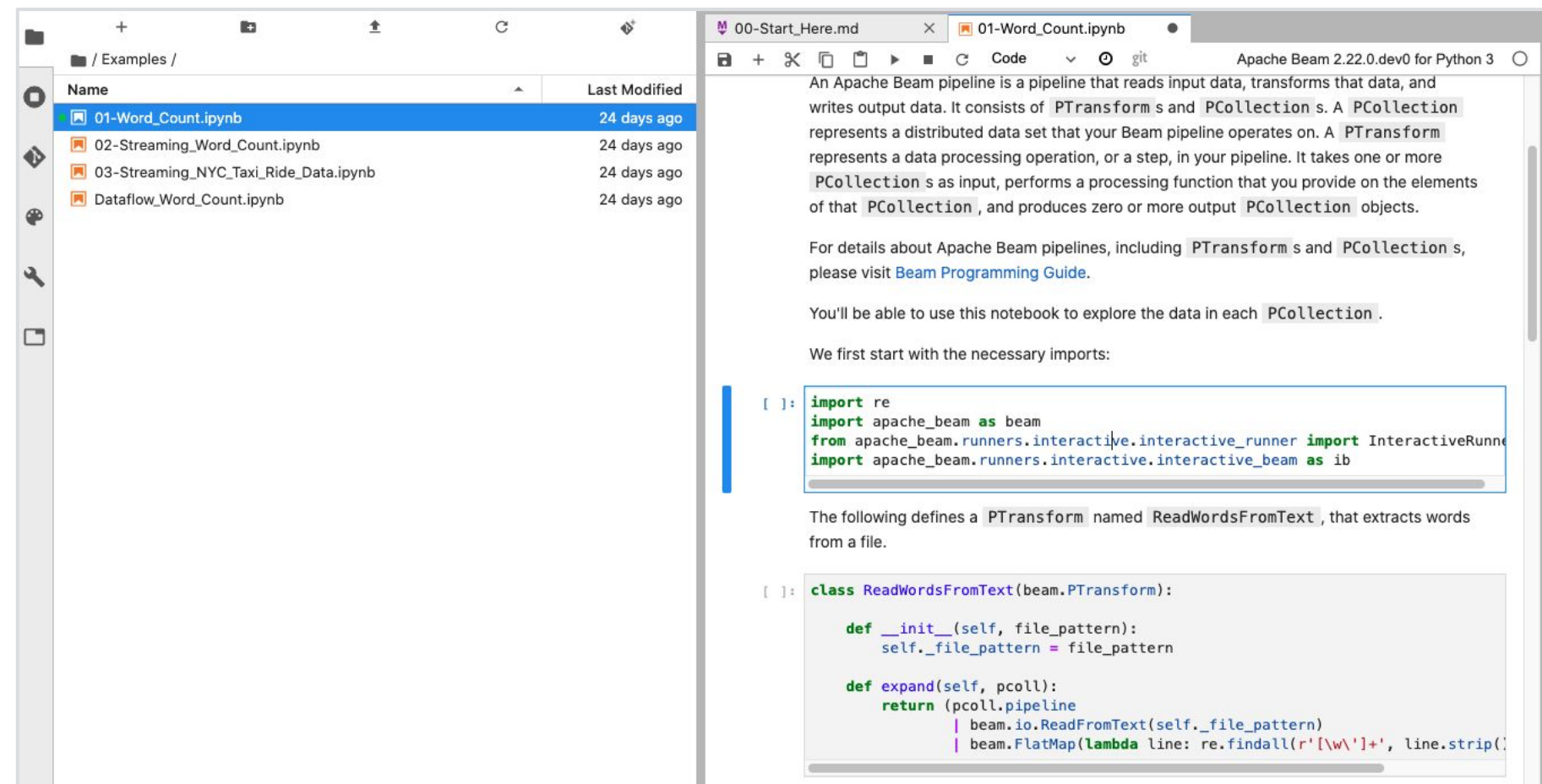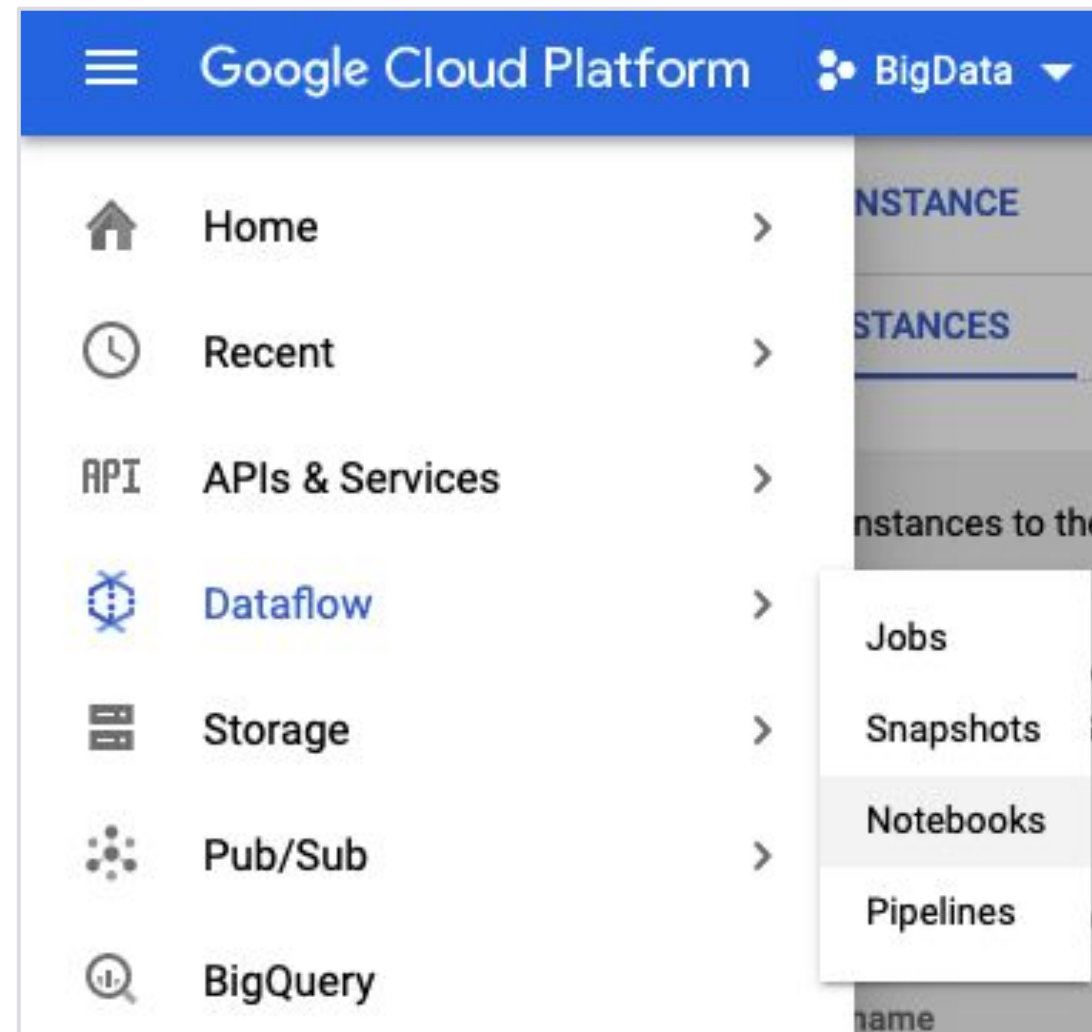**2**  Access to intermediate results

# Apache Beam interactive runner

The interactive runner module allows

**1** Interactive development

**2** Access to intermediate results

**3** Stream or batch sources

# Beam Notebooks

# Add a transform

```python
words = p | "read" >> beam.io.ReadFromPubSub(topic=topic)

windowed_words = (words
    | "window" >>
beam.WindowInto(beam.window.FixedWindows(10)))

windowed_words_counts = (windowed_words
    | "count" >> beam.combiners.Count.PerElement())
```

# Set interactivity options before we run the cell

### ib.options.recording_duration

Sets the amount of time the InteractiveRunner records data from an unbounded source

```
# Set the recording duration to 10 min
ib.options.recording_duration = '10m'
```

### ib.options.recording_size_limit

Sets the amount of data the InteractiveRunner records (in bytes) from an unbounded source

```
# Set the recording size limit to 1 GB
ib.options.recording_size_limit = 1e9
```

# Access Transform output



```
# Materializes the resulting PCollection in a table
ib.show(windowed_word_counts, include_window_info=True)

# Load the output in a Pandas DataFrame
ib.collect(windowed_word_counts, include_window_info=True)
```

# Visualize Transform output



```python
# Visualize the data in the notebook
ib.show(windowed_word_counts, include_window_info=True, visualize_data=True)
```

# Going from development to production

```python
# Import the production Dataflow runner
from apache_beam.runners import DataflowRunner

# Set up Apache Beam pipeline options
options = pipeline_options.PipelineOptions()

# Run the pipeline
runner = DataflowRunner()
runner.run_pipeline(p, options=options)
```