

[Return to Table of Contents](#)

Choose a Lesson

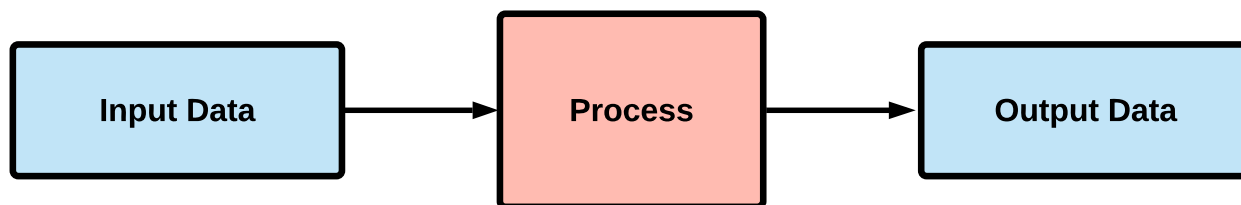
[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Data Processing Challenges

[Next](#)

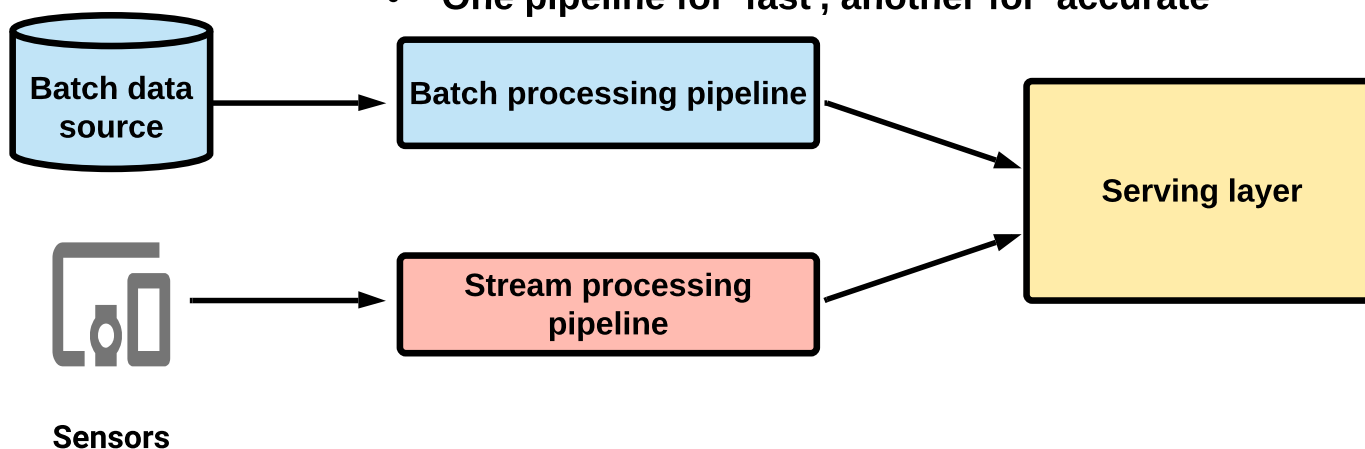
What is Data Processing?

- Read Data (Input)
 - Transform it to be relevant - Extract, Transform, and Load (ETL)
 - Create output
- to our particular interest*



Challenge: Streaming and Batch data pipelines:

- Until recently, separate pipelines are required for each
- Difficult to compare recent and historical data
- One pipeline for 'fast', another for 'accurate'



Why both?

- Credit card monitoring
- Compare streaming transactions to historical batch data to detect fraud



[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Data Processing Challenges

[Previous](#)

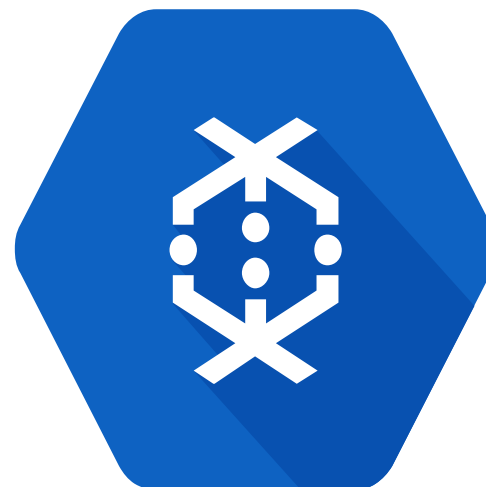
Challenge: Complex element processing:

- Element = single data input
- One at a time element ingest from single source = easy
- Combining elements (aggregation) = hard
- Processing data from different sources, streaming, and out of order (composite) = REALLY hard

Solution: Apache Beam + Cloud Dataflow



beam



Cloud Dataflow



create batch and
data processing pipeline

execute pipelines



[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Cloud Dataflow Overview

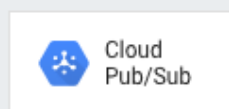
What is it?

- Auto-scaling, **No-Ops**, **Stream, and Batch Processing**
- Built on Apache Beam:
 - Documentation refers to Apache Beam site
 - Configuration is 100% code-based (Java / Python)
- Integrates with other tools (GCP and external):
 - Natively - Pub/Sub, BigQuery, Cloud ML Engine
 - Connectors - Bigtable, Apache Kafka
- Pipelines are regional-based

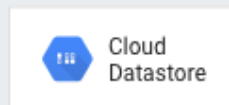
[Next](#)

Big Picture - Data Transformation

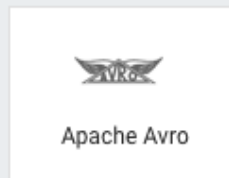
Ingest



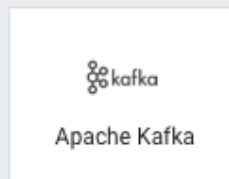
Cloud Pub/Sub



Cloud Datastore



Apache Avro



Apache Kafka

Stream

Batch

Process



Cloud Dataflow

Analyze



Data Studio



Third-Party Tools



Cloud BigQuery



Cloud Machine Learning



Cloud Bigtable



Data Warehouse



Predictive Analytics



Caching & Serving



[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Cloud Dataflow Overview

[Previous](#)

[Next](#)

IAM:

- **Project-level only** - all pipelines in the project (or none)
- Pipeline data access separate from **pipeline access**
- Dataflow Admin - Full **pipeline access** plus machine type/storage bucket config access
- Dataflow Developer - Full **pipeline access**, no machine type/storage bucket access
- Dataflow Viewer - view permissions only
- Dataflow Worker - Specifically for service accounts

Dataflow vs **Dataproc**?

Beam vs. **Hadoop/Spark**?

Dataproc:

- Familiar tools/packages
- Employee skill sets
- Existing pipelines

Dataflow:

- Less Overhead
- Unified batch and stream processing
- **Pipeline portability across Dataflow, Spark, and Flink as runtimes**

↑
Dataproc can not do!

WORKLOADS	CLOUD DATAPROC	CLOUD DATAFLOW
Stream processing (ETL)		X
Batch processing (ETL)	X	X
Iterative processing and notebooks	X	
Machine learning with Spark ML	X	
Preprocessing for machine learning		X (with Cloud ML Engine)



[Return to Table of Contents](#)

Cloud Dataflow Overview

[Previous](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Dataflow vs. Dataproc decision tree

Team topic: case study

Do you have existing system?

Do you have dependencies on specific tools/packages in the Apache Hadoop/Spark ecosystem?

YES

NO

-Starting fresh

Do you favor a hands-on/DevOps approach to operations, or a hands-off/serverless one?

DevOps

Serverless



Cloud Dataproc



Cloud Dataflow

[Return to Table of Contents](#)

Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Key Concepts

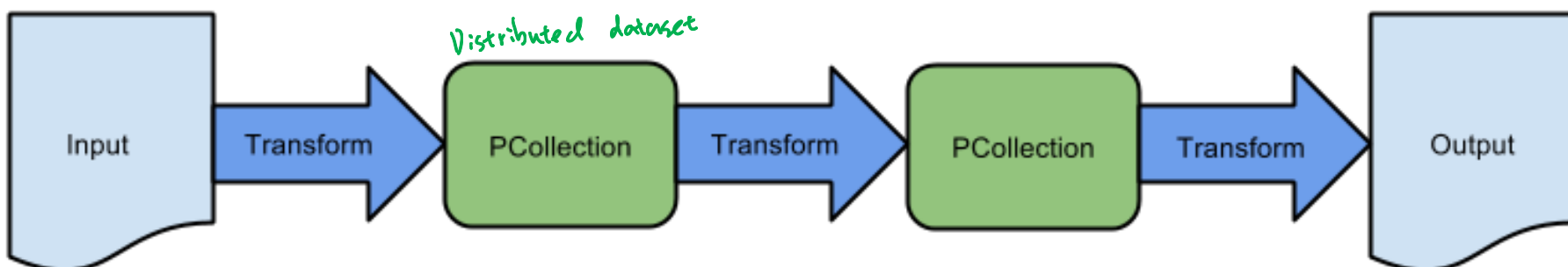
Course/exam perspective:

[Next](#)

- Dataflow is very code-heavy
- Exam does not go deep into coding questions
- **Some key concepts/terminology** will be tested

Key terms: (exam)

- **Element** - single entry of data (e.g., table row)
- **PCollection** - Distributed data set, data input and output
- **Transform** - Data processing operation (or step) in pipeline:
 - Uses programming conditionals (for/while loops, etc.)
- **ParDo** - Type of transform applied to individual elements:
 - Filter out/extract elements from a large group of data



PCollection and ParDo in example Java code.

One step in a multi-step transformation process.

```
PCollection<LaneInfo> currentConditions = p //
    .apply("GetMessages", PubsubIO.readStrings().fromTopic(topic)) //
    .apply("ExtractData", ParDo.of(new DoFn<String, LaneInfo>() {
        @ProcessElement
        public void processElement(ProcessContext c) throws Exception {
            String line = c.element();
            c.output(LaneInfo.newLaneInfo(line));
        }
    }));
```

[Return to Table of Contents](#)

Choose a Lesson

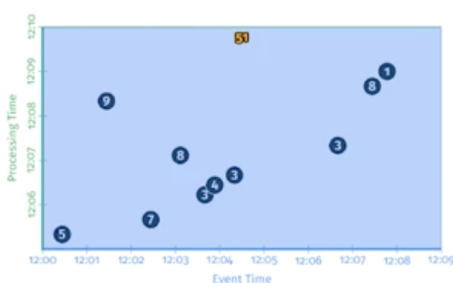
[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Cloud Dataflow Overview

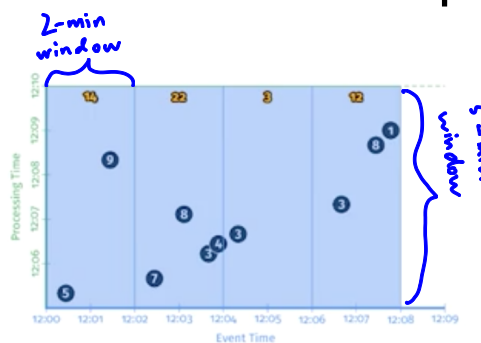
[Previous](#)*(Streaming data)*

Dealing with late/out of order data:

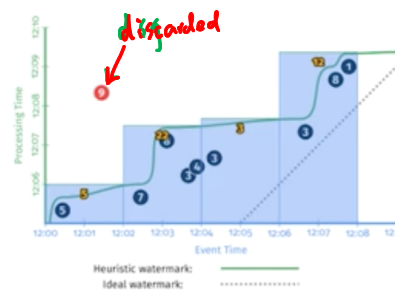
- Latency is to be expected (network latency, processing time, etc.)
- **Pub/Sub does not care about late data, that is resolved in Dataflow**
- Resolved with Windows, Watermarks, and Triggers
- **Windows** = logically divides element groups by time span
- **Watermarks** = 'timestamp':
 - Event time = when data was generated
 - Processing time = when data processed anywhere in the processing pipeline
 - Can use Pub/Sub-provided watermark or source-generated
- **Trigger** = determine when results in window are emitted (submitted as complete):
 - Allow late-arriving data in allowed time window to re-aggregate previously submitted results
 - Timestamps, element count, combinations of both

Watermark types

1

**Classic
Batch**

2

**Windowed
Batch**

3

**Streaming
w/Discarding**

4

**Streaming
+ Accumulation**



[Return to Table of Contents](#)

Template Hands On

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

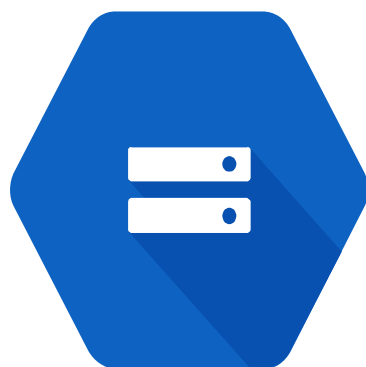
Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

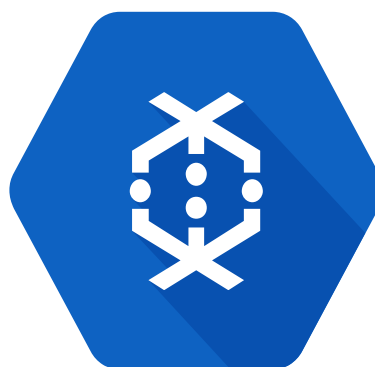
Additional Best Practices

- Google-provided templates
- Simple word count extraction



Cloud Storage

romeoandjuliet.txt



Cloud Dataflow
Read lines

output.txt



Cloud Storage

you can view
how many element
in/out

(transformation)

Read Lines



Word Count
template



Map
elements



Writes
Count



[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Streaming Ingest Pipeline Hands On

[Next](#)

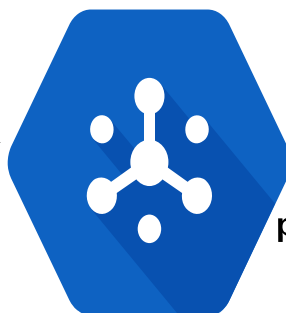
- Take San Diego traffic data
- Ingest through Pub/Sub
- Process with Dataflow
- Analyze results with BigQuery
- First: Enable Dataflow API from API's and Services



Traffic data

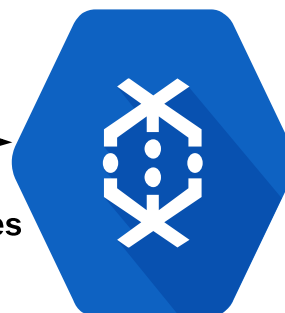
Data ingest

Published
streaming
sensor data



Topic: sandiego

Subscription
pulls messages



Cloud Dataflow
Transform data
to calculate
average speed.
Output to
BigQuery.



BigQuery

[Return to Table of Contents](#)

Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Streaming Ingest Pipeline Hands On

[Previous](#)[Next](#)

Quick command line setup (Cloud Shell)

- **Create BigQuery dataset** for processing pipeline output:
 - `bq mk --dataset $DEVSHHELL PROJECT_ID:demos`
- **Create Cloud Storage bucket** for Dataflow staging:
 - `gsutil mb gs://$DEVSHHELL PROJECT_ID` *will resolve the current project name*
- **Create Pub/Sub topic** and stream data: *to the topic*
 - `cd ~/googledataengineer/courses/streaming/publish`
 - `gcloud pubsub topics create sandiego` *Sandiego*
 - `./download_data.sh`
 - `sudo pip install -U google-cloud-pubsub`
 - `./send_sensor_data.py --speedFactor=60 --project=$DEVSHHELL PROJECT_ID`

Open a new Cloud Shell tab:

- **Execute Dataflow pipeline** for calculating average speed:
 - `cd ~/googledataengineer/courses/streaming/process/sandiego`
 - `./run_oncloud.sh $DEVSHHELL PROJECT_ID $DEVSHHELL PROJECT_ID AverageSpeeds`
- Error resolution: *dataflow pipeline creation process*
 - Pub/Sub permission denied, re-authenticate
 - `gcloud auth application-default login`
 - Dataflow workflow failed - enable Dataflow API

AverageSpeeds.java

./run_oncloud.sh project Id storage bucket java file we're printing

[Return to Table of Contents](#)

Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Streaming Ingest Pipeline Hands On

[Previous](#)

View results in BigQuery:

- List first 100 rows:
 - `SELECT * FROM [<PROJECTID>:demos.average_speeds] ORDER BY timestamp DESC LIMIT 100`
- Show last update to table:
 - `SELECT MAX(timestamp) FROM [<PROJECTID>:demos.average_speeds]`
- Look at results from the last minute:
 - `SELECT * FROM [<PROJECTID>:demos.average_speeds@-60000] ORDER BY timestamp DESC`

\downarrow
ms = | min

Shut down pipeline:

- **Drain** - finishing processing buffered jobs before shutting down
 - **Cancel** - **full stop**, cancels existing buffered jobs
- (stop accepting new data, and shut down gracefully)

[Return to Table of Contents](#)

Additional Best Practices

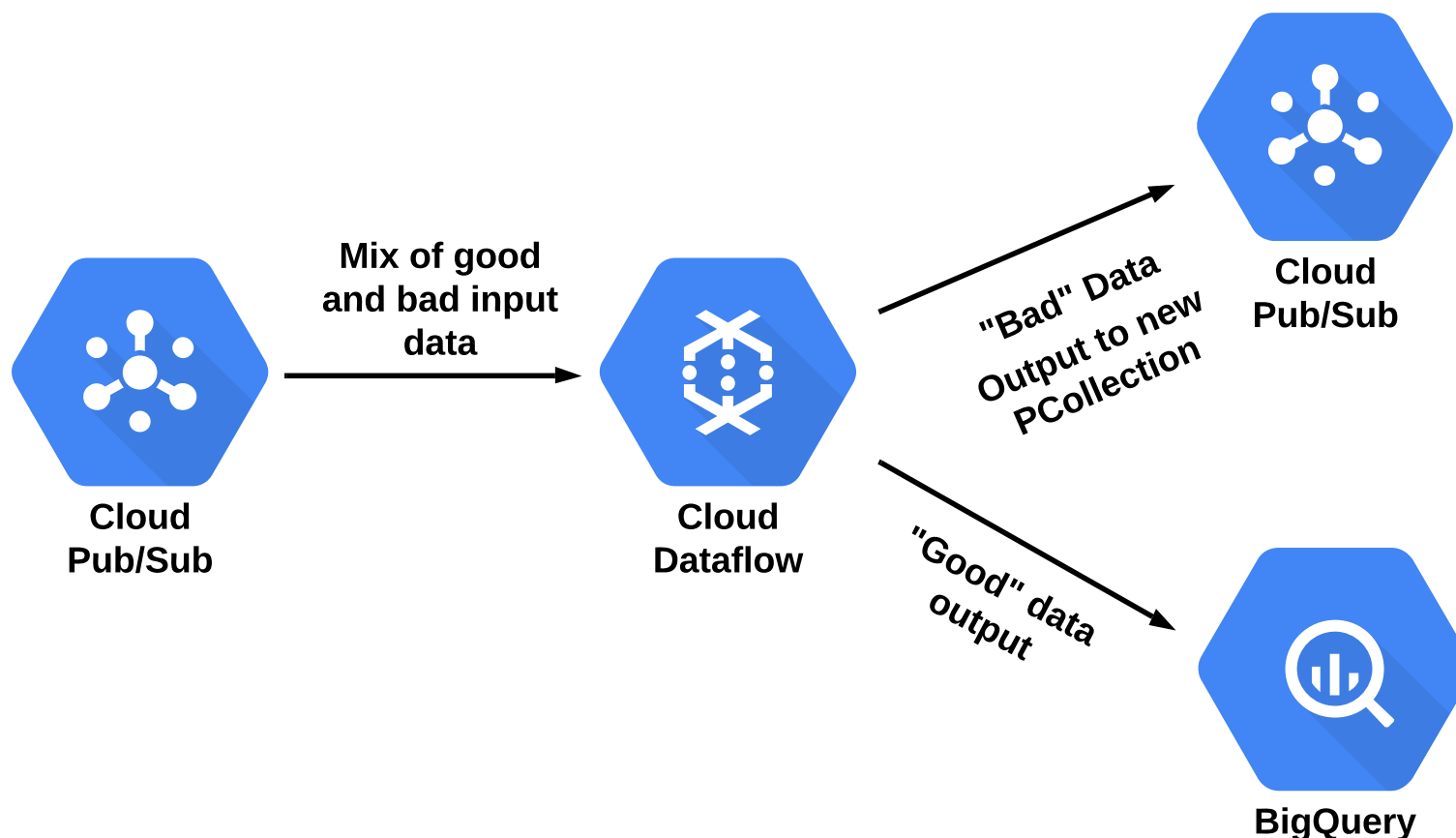
[Next](#)

Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

① Handling Pipeline Errors

- If you do not have a mechanism in place to handle input data errors in your pipeline, the job can fail. How can we account for this?
- Gracefully catch errors:
 - Create separate output: (in Apache)
 - **Try-catch** block handles errors
 - Output errors to new **PCollection** - Send to **collector** for later analysis (Pub/Sub is a good target)
 - Think of it as *recycling* the *bad* data
- Technique is also valid for troubleshooting missing messages:
 - Scenario: Streaming pipeline missing some messages
 - Solution: Run a batch of the streaming data, and check output:
 - Create additional output to capturing and processing error data.



[Return to Table of Contents](#)

Additional Best Practices

[Previous](#)[Next](#)

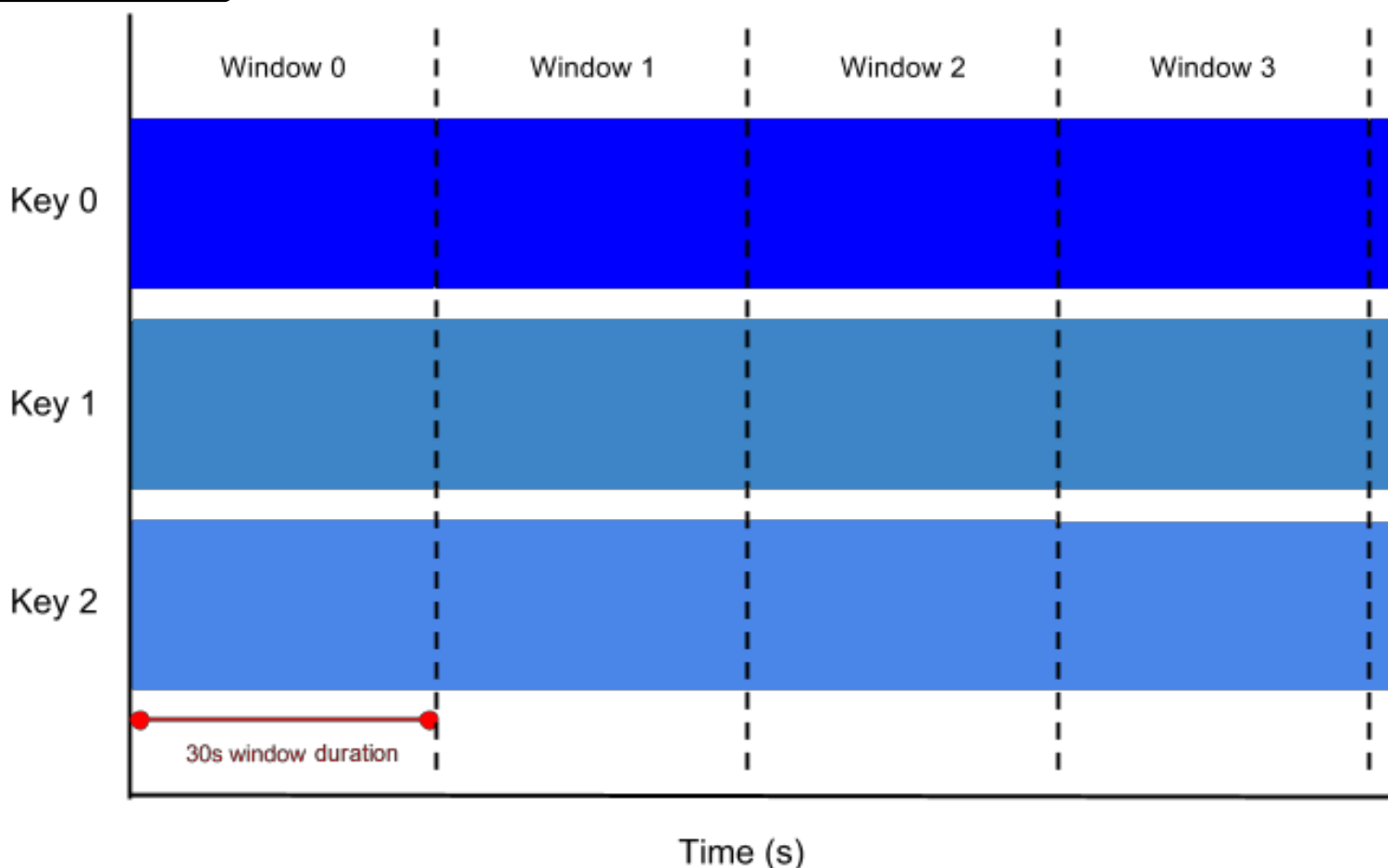
Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

2 Know your window types

- **Global, Fixed, Sliding, Session** (4 types)
- **Global** - The default, uses a single window for entire pipeline
- **Fixed time** - Every (x) period of time
 - Every 5 seconds, 10 minutes, etc.
- **Sliding time** - Overlapping time windows
- **Session** - Within certain time of certain elements:
 - For example, *Time since last user/mouse activity*

Fixed time Window

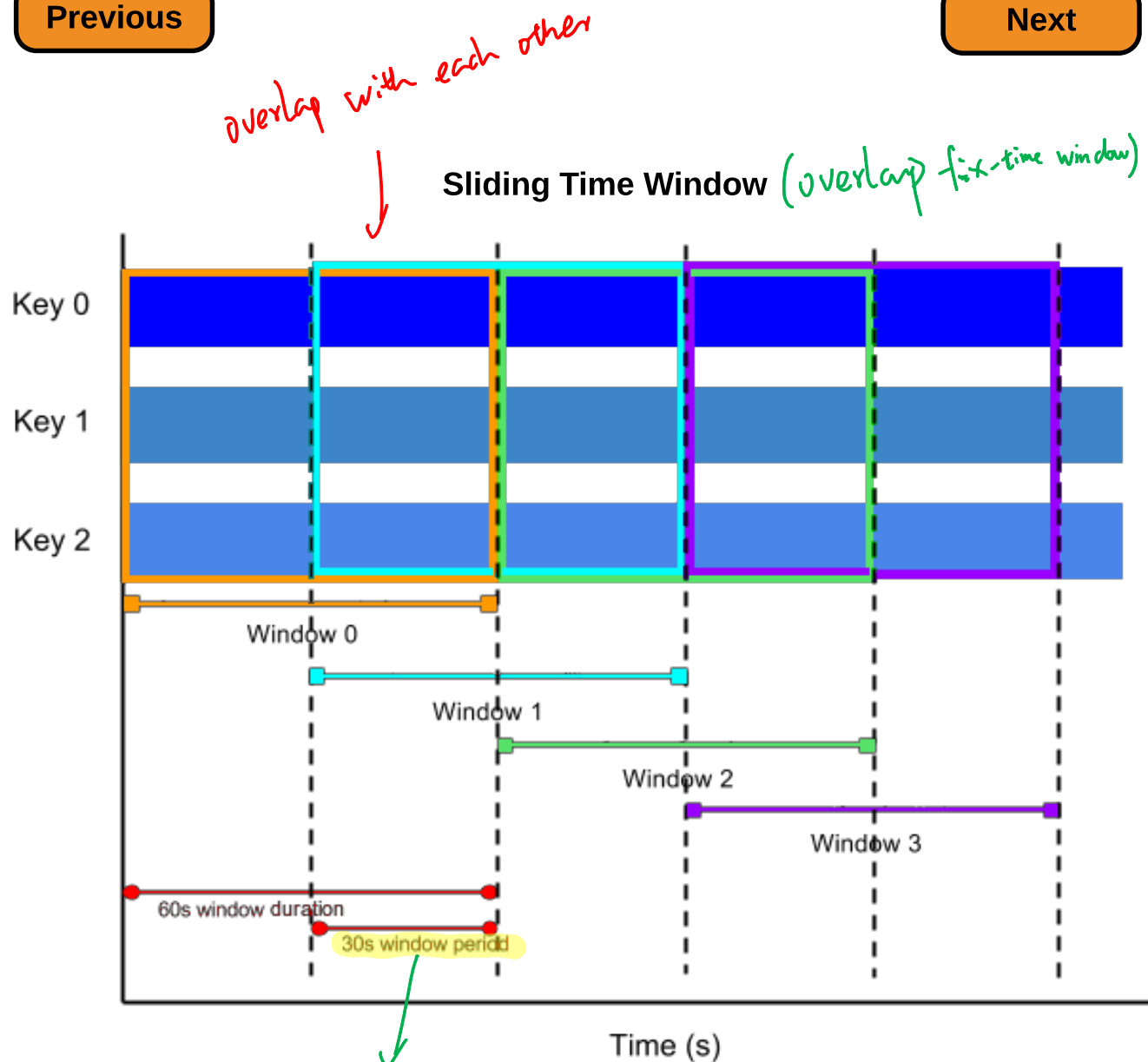


[Return to Table of Contents](#)

Choose a Lesson

[Data Processing Challenges](#)[Cloud Dataflow Overview](#)[Key Concepts](#)[Template Hands On](#)[Streaming Ingest Pipeline Hands On](#)[Additional Best Practices](#)

Additional Best Practices

[Previous](#)[Next](#)



[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

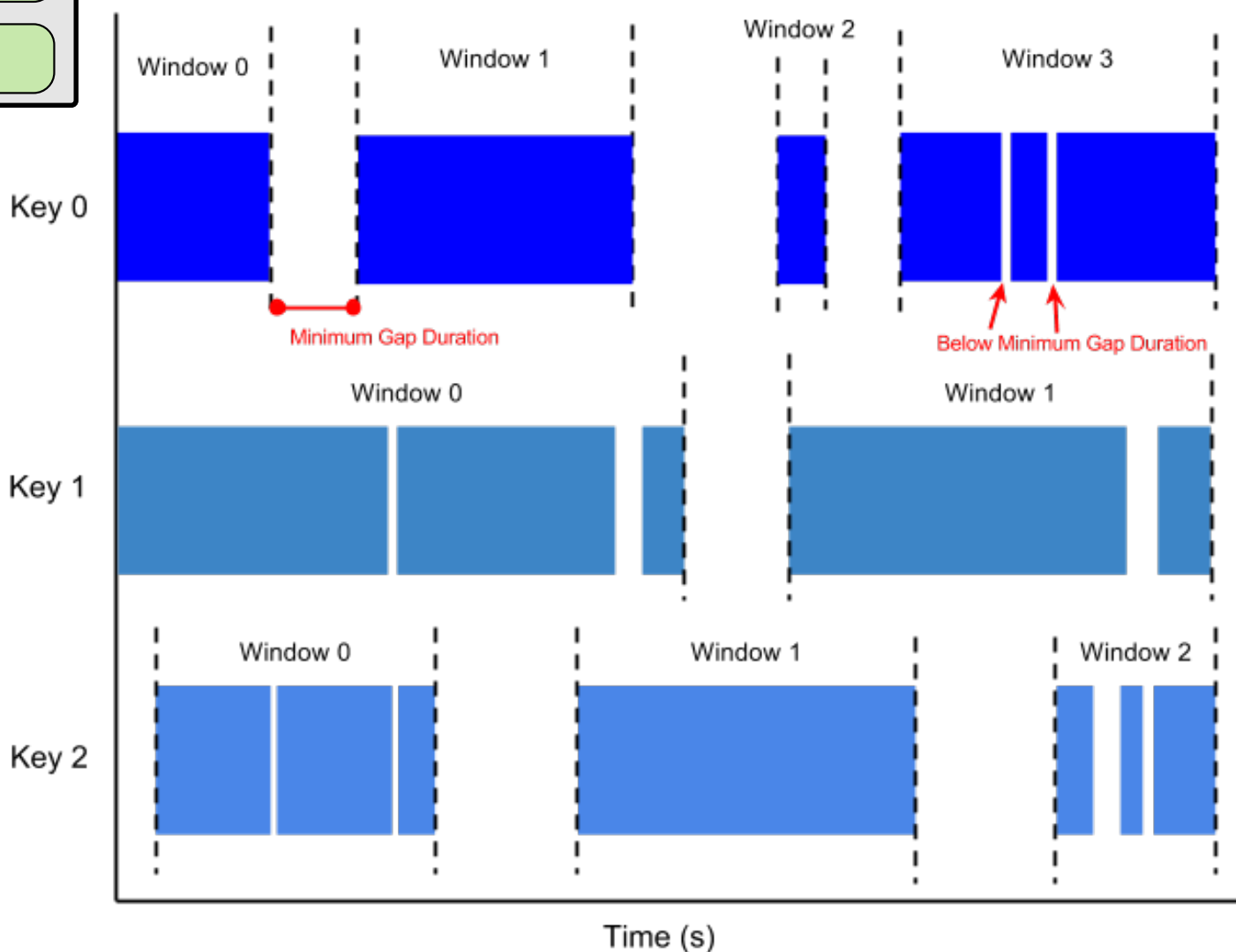
Additional Best Practices

[Previous](#)

[Next](#)

Q - Commerce website
reword "how long for each activity"

Session Window (based on incoming elements)





[Return to Table of Contents](#)

Choose a Lesson

Data Processing Challenges

Cloud Dataflow Overview

Key Concepts

Template Hands On

Streaming Ingest Pipeline Hands On

Additional Best Practices

Additional Best Practices

[Previous](#)

Updating Dataflow Pipelines *from a pub/sub topic*

- **Scenario:** Update streaming Dataflow pipeline with **new code**:
 - New code = new pipeline **not compatible** with current version
 - Need data to **switch over to new job/pipeline** without losing anything in the process
- **Solution:** Update job:
 - Creates **new job with same name**/new *jobID*
- Compatibility between old/new jobs:
 - Map old to new job transforms with **transform mapping**
 - "Bridge" between old and new **code base** *
 - After compatibility check:
 - Buffered data transferred to new job, using transform mapping to translate changes