

CS 327E Class 7

November 5, 2018

Check your GCP Credits :)



Hello,

We are emailing you because we noticed some of your students have used over 60% of their credit. This is great; we love to see students use our platform! We also want to make sure your students don't run out of credit.

If you are still using GCP, we can provide more credit for just these students or your whole class if you think others will also be using as much. If you do need more credit, please reply back to this message with how much more you will need and we will review your request.

Thanks,

Google Cloud Platform Education Grants Team



Google Cloud

Learn more about the [GCP Grants Programs](#)

[Higher Ed Learning Center](#)

iClicker Question

Are you running low on GCP credits?

- A. Yes
- B. No

No Quiz Today

Dataflow Concepts

- A system for processing arbitrary computations on large amounts of data
- Can process batch data and streaming data using the same code
- Uses Apache Beam, an open-source programming model
- Designed to be very scalable, millions of QPS

Apache Beam Concepts

- A model for describing data and data processing operations:
 - `Pipeline`: a data processing task from start to finish
 - `PCollection`: a collection of data elements
 - `Transform`: a data transformation operation
- SDKs for Java, Python and Go
- Executed in the cloud on Dataflow, Spark, Flink, etc.
- Executed locally with Direct Runner for dev/testing

Beam Pipeline

- Pipeline = A directed acyclic graph where the nodes are the Transforms and the edges are the PCollections
- General Structure of a Pipeline:
 - Reads one or more data sources as input PCollections
 - Applies one or more Transforms on the PCollections
 - Outputs resulting PCollection as one or more data sinks
- Executed as a single unit
- Run in batch or streaming mode

PCollection

- `PCollection` = A collection of data elements
- Elements can be of any type (String, Int, Array, etc.)
- `PCollections` are distributed across machines
- `PCollections` are immutable
- Created from a data source or a `Transform`
- Written to a data sink or passed to another `Transform`

Transform Types

- Element-wise:
 - maps 1 input to (1, 0, many) outputs
 - Examples: `ParDo`, `Map`, `FlatMap`
- Aggregation:
 - reduces many inputs to (1, fewer) outputs
 - Examples: `GroupByKey`, `CoGroupByKey`
- Composite: combines element-wise and aggregation
 - `GroupByKey` \rightarrow `ParDo`

Transform Properties

- Serializable
- Parallelizable
- Idempotent

ParDo

- ParDo = “Parallel Do”
- Maps 1 input to (1, 0, many) outputs
- Takes as input a `PCollection`
- Applies the user-defined `ParDo` to the input `PCollection`
- Outputs results as a new `PCollection`
- Typical usage: filtering, formatting, extracting parts of data, performing computations on data elements

ParDo Example

```
1 import apache_beam as beam
2 from apache_beam.io import ReadFromText
3 from apache_beam.io import WriteToText
4
5 # DoFn to perform on each element in the input PCollection.
6 class ComputeWordLengthFn(beam.DoFn):
7     def process(self, element):
8         words = element.strip().split(' ')
9         result_list = []
10        for word in words:
11            result_list.append((word, len(word)))
12        return result_list
13
14 # Create a Pipeline using a local runner for execution.
15 with beam.Pipeline('DirectRunner') as p:
16
17     # create a PCollection from the file contents.
18     in_pcoll = p | 'Read' >> ReadFromText('input.txt')
19
20     # apply a ParDo to the PCollection
21     out_pcoll = in_pcoll | beam.ParDo(ComputeWordLengthFn())
22
23     # write PCollection to a file
24     out_pcoll | 'Write' >> WriteToText('output.txt')
```

Aggregation Example

```
1 import apache_beam as beam
2 from apache_beam.io import ReadFromText
3 from apache_beam.io import WriteToText
4
5 # DoFn to perform on each element in the input PCollection.
6 class ComputeWordLengthFn(beam.DoFn):
7     def process(self, element):
8         words = element.strip().split(' ')
9         result_list = []
10        for word in words:
11            result_list.append((len(word), word))
12        return result_list
13
14 # Create a Pipeline using a local runner for execution.
15 with beam.Pipeline('DirectRunner') as p:
16
17     # create a PCollection from the file contents.
18     in_pcoll = p | 'Read' >> ReadFromText('input.txt')
19
20     # apply a ParDo to the PCollection
21     word_pcoll = in_pcoll | 'ParDo' >> beam.ParDo(ComputeWordLengthFn())
22
23     # apply GroupByKey to the PCollection
24     out_pcoll = word_pcoll | 'GroupByKey' >> beam.GroupByKey()
25
26     # write PCollection to a file
27     out_pcoll | 'Write' >> WriteToText('output.txt')
```

Source File: https://github.com/cs327e-fall2018/snippets/blob/master/group_words_by_length.py

BigQuery Data Sink Example

```
1 import logging
2 import apache_beam as beam
3 from apache_beam.io import ReadFromText
4 from apache_beam.io import WriteToText
5
6 # DoFn to perform on each element in the input PCollection.
7 class ComputeWordLengthFn(beam.DoFn):
8     def process(self, element):
9         words = element.strip().split(' ')
10        result_list = []
11        for word in words:
12            result = {'word' : word, 'length' : len(word)}
13            result_list.append(result)
14        return result_list
15
16 # Create a Pipeline using a local runner for execution.
17 with beam.Pipeline('DirectRunner') as p:
18
19     # create a PCollection from the file contents.
20     in_pcoll = p | 'Read from File' >> ReadFromText('input.txt')
21
22     # apply a ParDo to the PCollection
23     out_pcoll = in_pcoll | beam.ParDo(ComputeWordLengthFn())
24
25     # write PCollection to a file
26     out_pcoll | 'Write to File' >> WriteToText('output.txt')
27
28     # write PCollection to a BQ table
29     qualified_table_name = 'cs327e-fa2018:beam.Words'
30     table_schema = 'word:STRING,length:INTEGER'
31     out_pcoll | 'Write to BigQuery' >> beam.io.Write(beam.io.BigQuerySink(qualified_table_name,
32                                                                    schema=table_schema,
33                                                                    create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED,
34                                                                    write_disposition=beam.io.BigQueryDisposition.WRITE_TRUNCATE))
```

Source File: https://github.com/cs327e-fall2018/snippets/blob/master/word_length_bq_out.py

How to “Dataflow”

1. Start with some initial working code.
2. Test and debug **each** new line of code.
3. Write temporary and final `PCollections` to log files.
4. Test and debug **end-to-end** pipeline locally before running on Dataflow.
5. If you get stuck, make a Piazza post that has **enough** details for the instructors to reproduce the error and help you troubleshoot.
6. Start assignments **early**. The Beam Python documentation is sparse and learning Beam requires patience and experimentation.

Dataflow Set Up

<https://github.com/cs327e-fall2018/snippets/wiki/Dataflow-Setup-Guide>

Milestone 7

<http://www.cs.utexas.edu/~scohen/milestones/Milestone7.pdf>