

### ME & THIS

- ▶ Jeff Kramer: <u>github.com/jeffk</u> & <u>twitter.com/jeffk</u>
- Work: Data Engineering at Vox Media
- Play: Late 80's to Early 00's Mac, SGI & NeXT Machines
- Life: {"better\_half": "Irma Kramer", "kids": ["Clara", "Arthur"]}

▶ Presentation: <u>GitHub.com/jeffk/PyTexas-BigQuery</u>

### WHAT IS BIG DATA?

- Big Data is about Time, not about Space
- Problems bigger than a single core can handle must be distributed, it's the new Data Warehouse
- It's really hard to do this well, and you have to make compromises
- BigQuery solves this problem well for data archiving and analytics tasks

### WHAT IS BIGQUERY?

- ▶ Google Cloud's Semi-Structured Database: SQL at Scale
- \$5 per TB of data processed
- Converts your data into Google's custom Capacitor file format
- ▶ Fast at scale: SELECT ... GROUP BY ... across 100 Billion rows (3.6 TB) of Wikipedia logs in ~36 sec
- Loads/exports data through HTTP and bulk jobs

#### SUPPORTED DATA TYPES, QUERIES & ORGANIZATION

- SQL 2011: INT64, FLOAT64, BOOL, STRING, BYTES, DATE, DATETIME, TIME, TIMESTAMP, ARRAY, STRUCT
- STRUCTS == Nested | Repeated
- Standard SQL 2011 & Legacy BigQuery Extensions
- JavaScript UDFs
- Project -> Dataset -> Table

visit_id	INTEGER	REQUIRED	•
visit_time	TIMESTAMP	REQUIRED	•
payload	RECORD	REQUIRED	•
payload.visit_location	STRING	REQUIRED	•
payload.metadata	RECORD	REPEATED	
payload.metadata.key	STRING	REQUIRED	•
payload.metadata.value	STRING	NULLABLE	
payload.metrics	RECORD	REPEATED	
payload.metrics.key	STRING	REQUIRED	•
payload.metrics.value	FLOAT	NULLABLE	

# CREATING DATASETS

```
client = bigquery.Client()
dataset_ref = client.dataset(name)
dataset = bigquery.Dataset(dataset_ref)
dataset.description = description
dataset = client.create_dataset(dataset)
```

# CREATING SCHEMAS

```
SCHEMA = [
  SchemaField('visit id', 'INT64',
   mode='required', description="Visit ID"),
  SchemaField('visit time', 'TIMESTAMP',
   mode='required', description="Visit Time"),
 SchemaField('payload', 'STRUCT', mode='REQUIRED', fields = [
    SchemaField('visit location', 'STRING',
     mode='required', description="Visit Location"),
    SchemaField('metadata', 'STRUCT', mode='REPEATED',fields = [
      SchemaField('key', 'STRING', mode='REQUIRED'),
      SchemaField('value', 'STRING')
   ]),
    SchemaField('metrics', 'STRUCT', mode='REPEATED', fields = [
      SchemaField('key', 'STRING', mode='REQUIRED'),
      SchemaField('value', 'FLOAT64')
```

### CREATING TABLES

```
client = bigquery.Client()
table_ref = dataset.table(name)
table = bigquery.Table(table_ref, schema=SCHEMA)
table.description = description
created_table = client.create_table(table)
```

# LOADING DATA VIA HTTP STREAMS

```
ROWS TO INSERT = [
    { 'visit id': 1,
    'visit time': '2017-04-01T12:21:32',
    'payload':
        { 'visit location': 'NORTH',
        'metadata': [
            {'key':'first name', 'value':'Alice'},
            {'key':'favorite color', 'value':'red'},
            {'key':'last purchase id', 'value':'1243'},
            {'key':'last purchase total', 'value':'34.53'},
        ], 'metrics': [
            {'key': 'checkout time', 'value': 82.4},
            {'key': 'net promoter', 'value': 5},
            {'key':'visit count', 'value': 12},
client = bigquery.Client()
errors = client.create_rows(table, ROWS_TO_INSERT)
```

### QUERYING STRUCT DATA

```
client = bigquery.Client()
QUERY = """
SELECT visit id, visit time, payload.visit location,
(SELECT value FROM UNNEST(payload.metadata) WHERE key = "first name")
AS first name,
(SELECT value FROM UNNEST(payload.metrics) WHERE key = "net promoter")
AS net promoter
FROM `%s.%s. ORDER BY visit id LIMIT 100
""" % (client.project, dataset name, table_name)
rows = list(client.query rows(QUERY, timeout=30))
for row in rows:
    print("%s\t%s\t%s\t%s\t%s" % (row[0], row[1], row[2], row[3], row[4]))
```

#### LIMITING COST BY QUERYING PARTITIONED TABLES

11 11 11

```
SELECT visit_id, visit_time
FROM `%s.%s.%s`
WHERE _PARTITIONTIME BETWEEN
   TIMESTAMP('2017-11-01') AND TIMESTAMP('2017-11-15');
""" % (client.project, dataset_name, table_name)
```

# QUERYING DATA WITH A UDF

```
QUERY = """
CREATE TEMPORARY FUNCTION rot13(x STRING)
RETURNS STRING
LANGUAGE js
AS \"\"\"
x = x.replace(/[a-zA-Z]/g,function(c){}
   return String.fromCharCode((c<='Z'?90:122)\>=(c=c.charCodeAt(0)+13)?c:c-26);
});
return x;
\"\"\";
SELECT visit id, visit time, payload.visit location,
  (SELECT rot13(value) FROM UNNEST(payload.metadata) WHERE key = "first name")
    AS first name,
  (SELECT value FROM UNNEST(payload.metrics) WHERE key = "net promoter")
     AS net promoter
FROM `%s.%s. ORDER BY visit id LIMIT 100
""" % (client.project, dataset_name, table_name)
```

### QUERYING INTO A TABLE

```
client = bigquery.Client()
    OUERY = """
SELECT visit id, visit time, payload.visit location,
  (SELECT value FROM UNNEST(payload.metadata) WHERE key = "first_name") AS first_name,
  (SELECT value FROM UNNEST(payload.metrics) WHERE key = "net promoter")AS net promoter
FROM `%s.%s.%s`
""" % (client.project, dataset_name, source_table)
    dataset = client.dataset(dataset name)
    job config = QueryJobConfig()
    job config.destination = dataset.table(dest table)
    job config.write disposition = 'WRITE TRUNCATE'
    query job = QueryJob(str(uuid.uuid4()), QUERY, client=client, job config=job config)
    query job. begin()
    while not query job.done():
       time.sleep(5)
    print("%s bytes processed." % query_job.total_bytes_billed)
```

### **EXPORTING A TABLE**

```
client = bigquery.Client()
dataset = client.dataset(dataset_name)
table ref = dataset.table(table)
job config = ExtractJobConfig()
job config.destination format = 'AVRO'
dest = ['gs://%s/complex query output-*.avro' % bucket name]
query job = ExtractJob(str(uuid.uuid4()),
    table ref, dest, client, job config=job config)
query job. begin()
while not query job.done():
    time.sleep(5)
if query job.errors:
    print(query job.errors)
print("%s file(s) created." %
  query_job._job_statistics().get('destinationUriFileCounts')[0])
```

# LOADING DATA WITH A LOAD JOB: GOOGLE CLOUD STORAGE

```
client = bigquery.Client()
dataset ref = client.dataset(dataset_name)
table ref = dataset ref.table(table+'$20171118')
job config = bigquery.LoadJobConfig()
job config.source format = 'AVRO'
job config. properties['timePartitioning'] = {'type': 'DAY'}
job config.write disposition = 'WRITE TRUNCATE'
GS URL = 'gs://{}/{}'.format(bucket name, blob name)
job = client.load table from uri(
    GS URL, table ref, job config=job config)
while not job.done():
    time.sleep(5)
if job.errors:
   print(job.errors)
else:
    print('Loaded %s rows.' % job.output rows)
```

# WHY BIGQUERY

- Super easy to run. It just works.
- > Storage costs are low: \$.02 per gig per month, \$.01 if not edited in last 90 days
- Pay on demand: \$5 per TB of data processed, \$.05 per GB streamed in
- Easy connections to other Google products (Google Analytics Data, YouTube data, Data Studio, Cloud Dataflow, Google Sheets, etc)

# WHY NOT BIGQUERY

- Application-style workloads: Just use PostgreSQL
- Vendor lock-in: You may be in AWS already, you may be on-premises and need something like Apache Impala
- Limits on table UPDATE/DELETEs: 100 per table per day