# Google Cloud for Scientific Infrastructure

Karan Bhatia, PhD

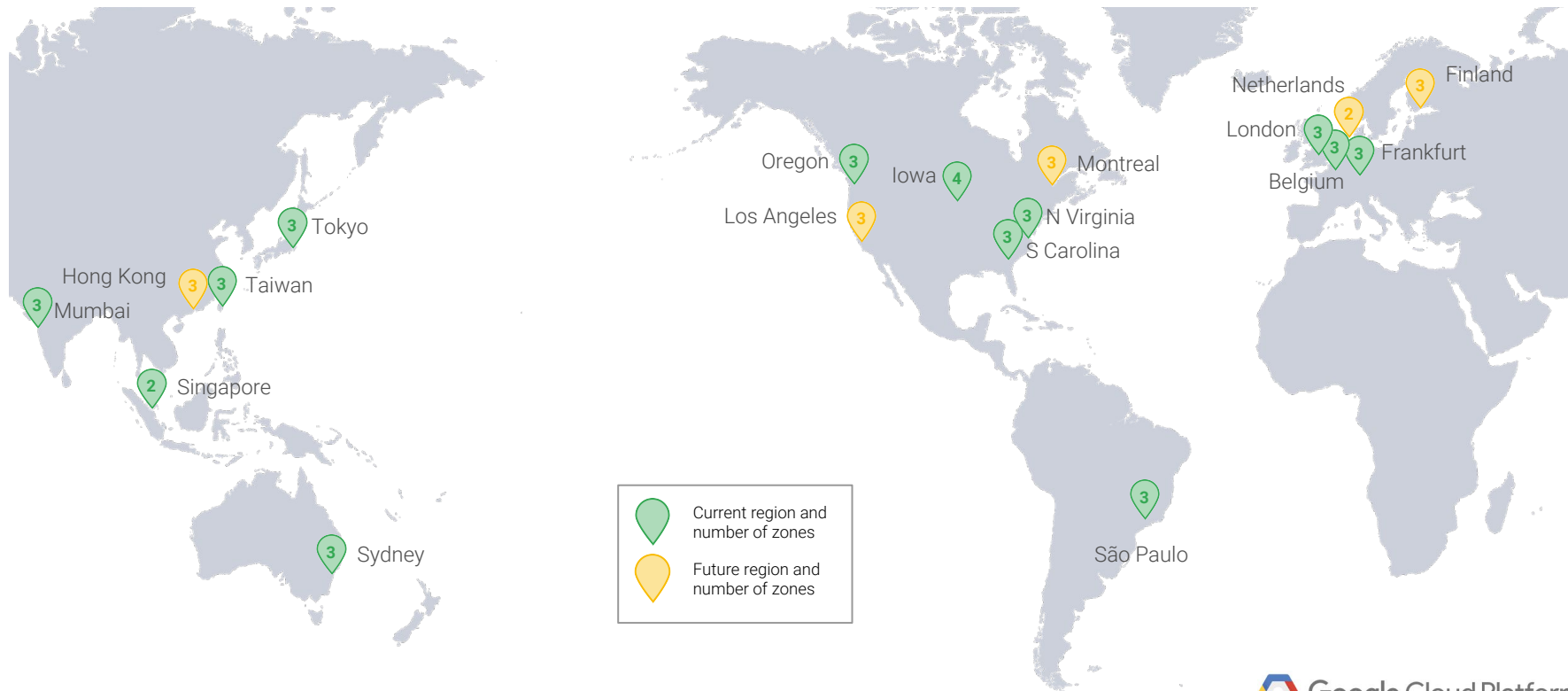ESGF Annual Meeting, Dec 2017

Google Cloud

Google Cloud Platform

# $29.4 Billion

3 year trailing Google
CAPEX investment

# Google Cloud Platform Regions
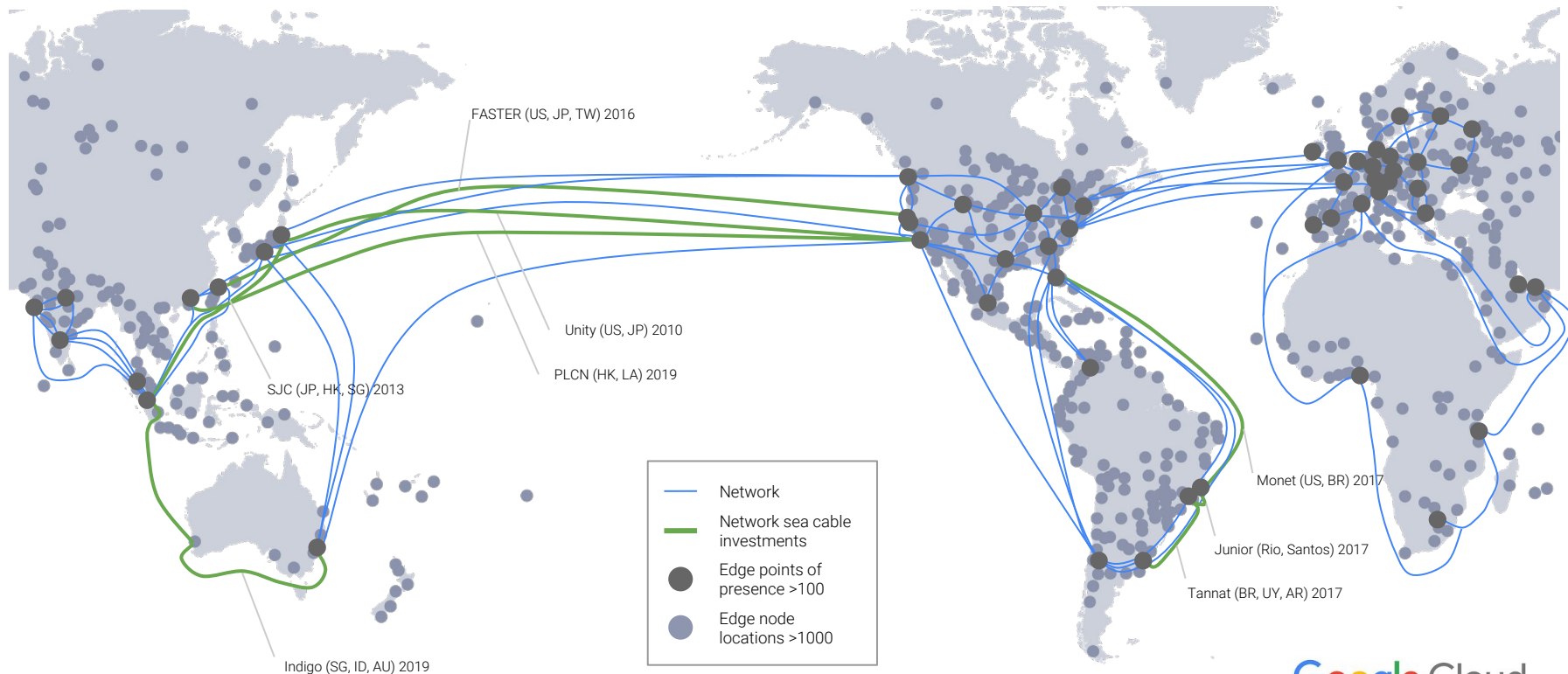Select from 13 regions. 5 new regions coming in 2018.

Netherlands
Finland 3
London 3
Belgium 3 3 Frankfurt

Oregon 3 Iowa 4 3 Montreal
Los Angeles 3 3 N Virginia
3 S Carolina

Tokyo 3
Hong Kong 3 3 Taiwan
3 Mumbai

Singapore 2

São Paulo 3

Sydney 3

Current region and number of zones

Future region and number of zones

Google Cloud Platform

# Google Cloud Network

The largest cloud network, comprised of >100 points of presence



FASTER (US, JP, TW) 2016

Unity (US, JP) 2010

PLCN (HK, LA) 2019

SJC (JP, HK, SG) 2013

Monet (US, BR) 2017

Junior (Rio, Santos) 2017

Tannat (BR, UY, AR) 2017

Indigo (SG, ID, AU) 2019

**Legend:**
— Network
— Network sea cable investments
● Edge points of presence >100
● Edge node locations >1000

Google Cloud

# Google Cloud Network

The largest cloud network, comprised of >100 points of presence
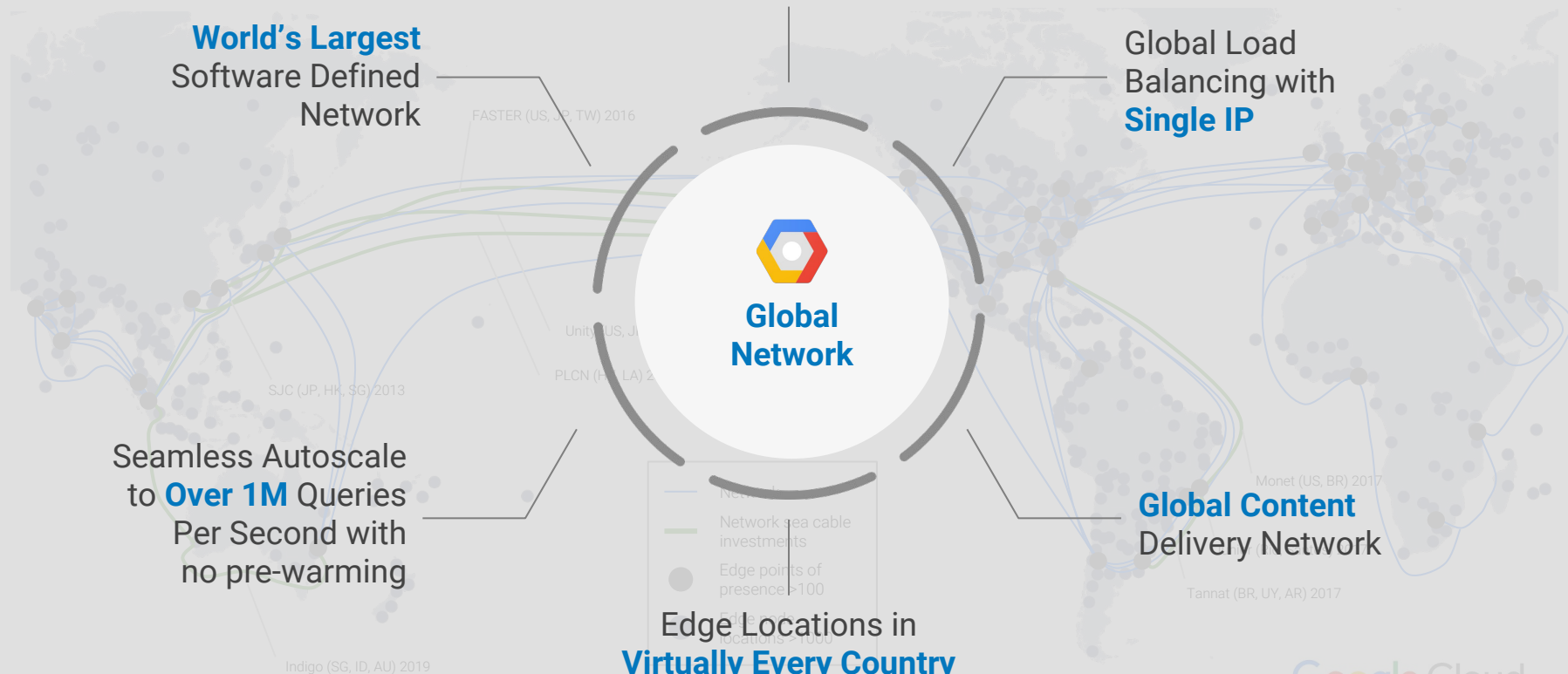
**More than 100**
Peering Locations

**World's Largest**
Software Defined
Network

Global Load
Balancing with
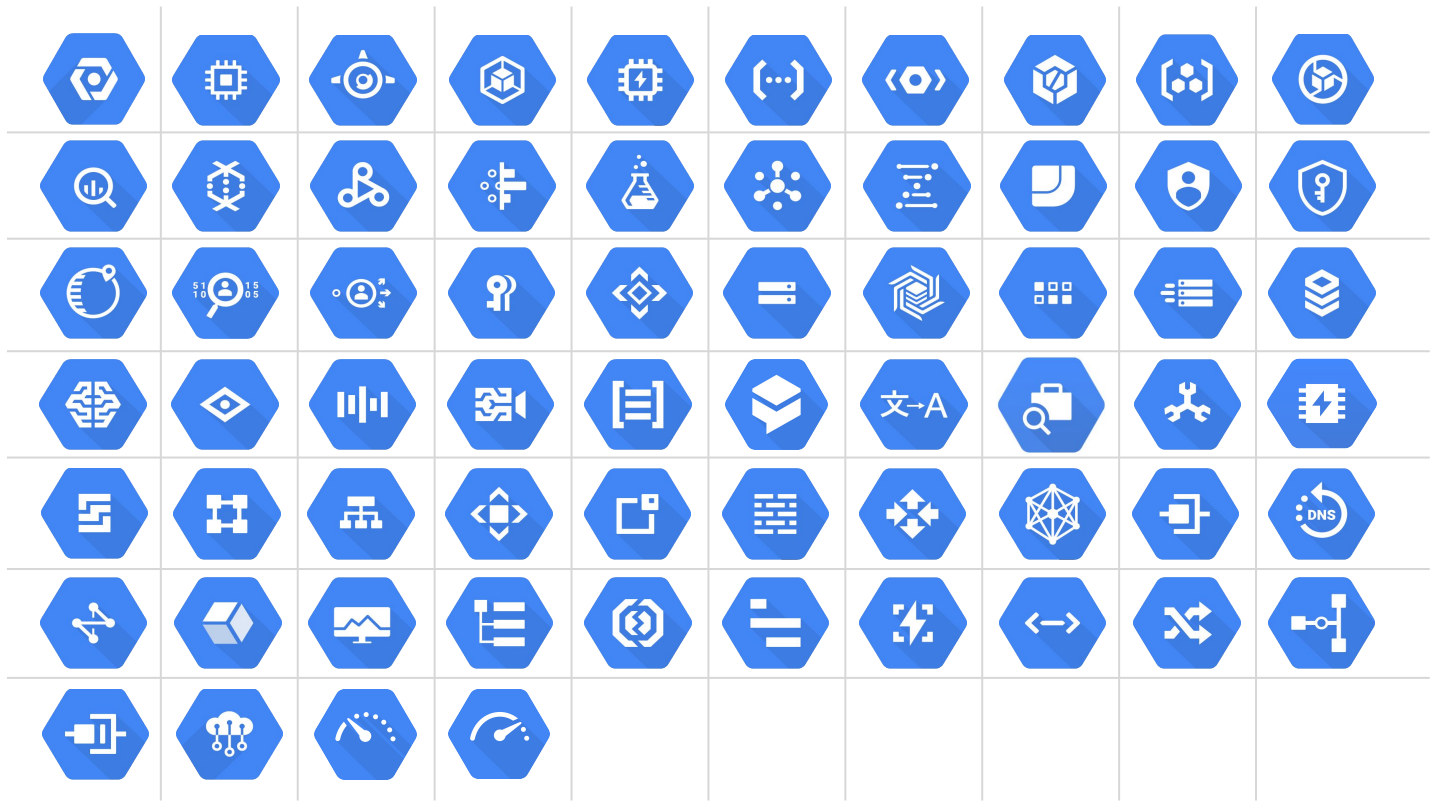**Single IP**

**Global Network**

Seamless Autoscale
to **Over 1M** Queries
Per Second with
no pre-warming

**Global Content**
Delivery Network

Edge Locations in
**Virtually Every Country**

FASTER (US, JP, TW) 2016

Unity (US, J

SJC (JP, HK, SG) 2013

PLCN (H, LA) 2

Indigo (SG, ID, AU) 2019

Monet (US, BR) 2017

Tannat (BR, UY, AR) 2017

Netw
Network sea cable
investments
Edge points of
presence >100

Google Cloud

# Google Cloud Platform

# Agenda

Compute

Data

Machine Learning

Academic / Research Programs

# Compute

# Infrastructure

—

**Lightning fast & scalable:** Fast VM startup time, millisecond access for all storage classes, high IOPS for VCPUs,  high bandwidth global networking

**Reliable**: Built-in redundancy and scale, live Migration, Google Site Reliability Engineering for your workload.

**Customer friendly pricing:** simple and efficient: Pay-per-second, custom VMs, automatic discounts, flexible buy-in-bulk discounts

**Geographic coverage:**  11 new regions in 2017-18 for a total of 17, HA in each region

- Significant "per core" performance improvements
- Intel® Advanced Vector Extension 512 (Intel® AVX-512)
  - 2x flops/second
- Accelerated IO with Intel® Omni-Path Architecture (Fabric)
- Integrated Intel® QuickAssist Technology (crypto & compression offload)
- Intel® Resource Director Technology (Intel® RDT) for Efficiency & TCO

Google Cloud Platform Blog

Product updates, customer stories, and tips and tricks on Google Cloud Platform

Google Cloud Platform is the first cloud provider to offer Intel Skylake

Friday, February 24, 2017

By Urs Hölzle, Senior Vice President, Google Cloud Infrastructure

I'm excited to announce that Google Cloud Platform (GCP) is the first cloud provider to offer the next generation Intel Xeon processor, codenamed Skylake.

# Hardware Accelerated



- Available Today: **NVIDIA K80 GPU, P100s**
- Coming Soon: **Tensor Processing Unit (TPU)**
- Custom ASIC built and optimized for TensorFlow
- Used in production at Google for over 16 months
- 7 years ahead of GPU performance per watt
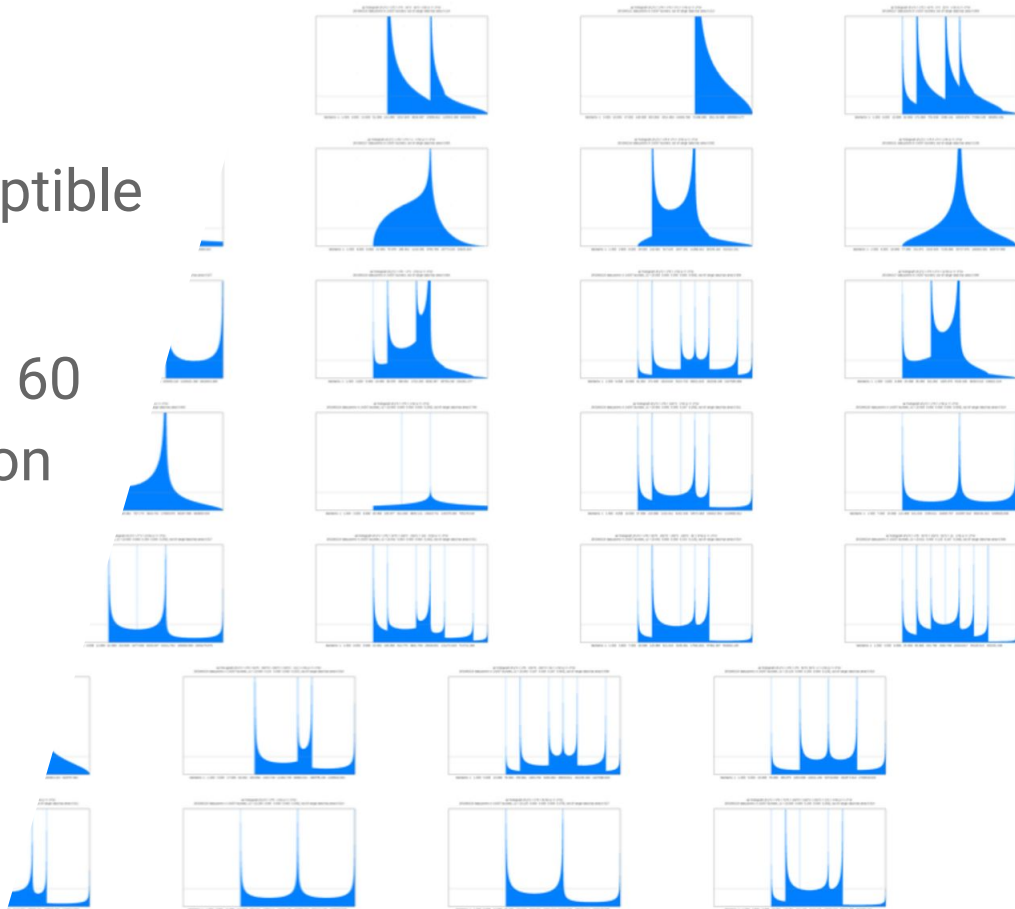
MIT Research w/ VMs

580,000 cores
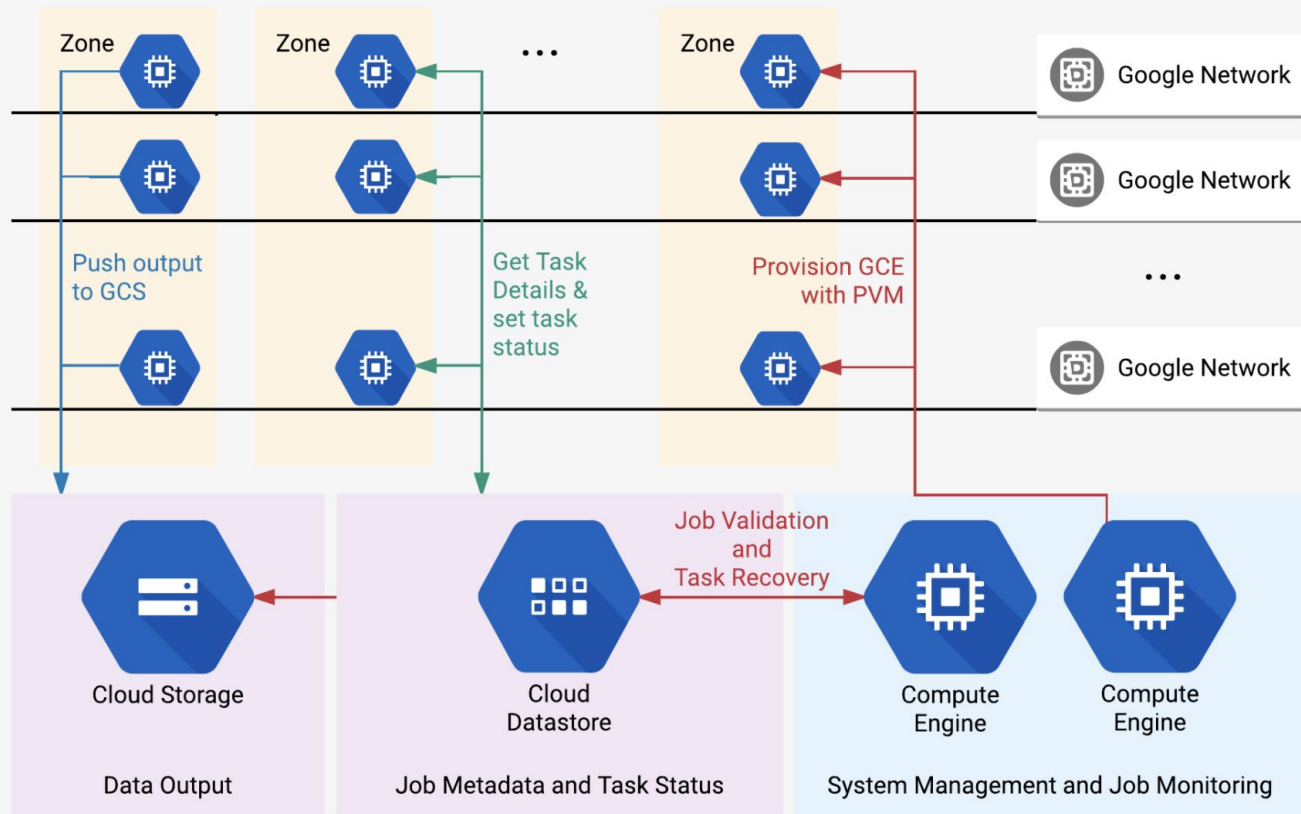
~~220,000~~ cores on preemptible VMs

2,250 32-core instances, 60 CPU-years of computation in a single afternoon
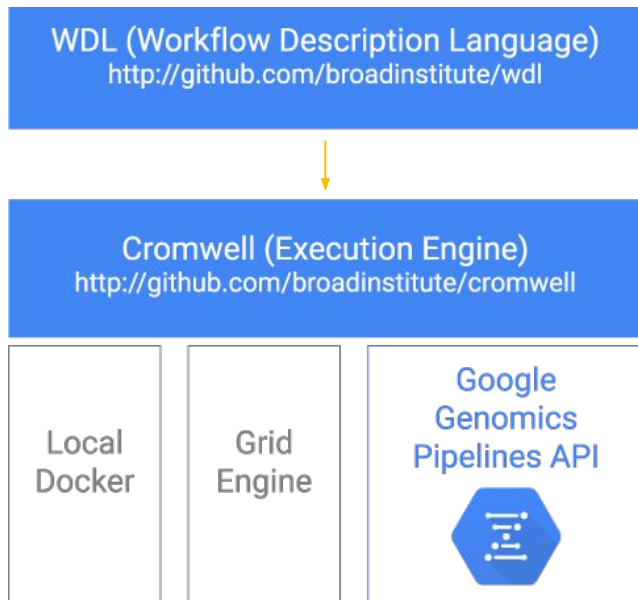
Answers in hours v. months

Products used: Google Compute Engine, Cloud Storage, DataStore

Google Cloud

Broad Firecloud:
WDL, Cromwell and Google Genomics



WDL (Workflow Description Language)
http://github.com/broadinstitute/wdl

Cromwell (Execution Engine)
http://github.com/broadinstitute/cromwell

Local Docker

Grid Engine

Google Genomics Pipelines API

A full stack for use by the community!
See software.broadinstitute.org/wdl

**WDL:** an external DSL used by computational biologists to express the analytical pipelines

**Cromwell:** a scalable, robust engine for executing WDL against pluggable backends including local, Docker, Grid Engine or …

**Google Genomics Pipelines API:** co-developed by Broad and Google Genomics, a scalable Docker-as-a-Service with data scheduling

Google Cloud

## Pipeline definition

```json
{
  "name": "samtools index",
  "description": "Run samtools index to generate a BAM index file",
  "inputParameters": [
    {"name": "inputFile",
      "localCopy": {
        "disk": "data",
        "path": "input.bam"
      }
    },
    {"name": "outputFile",
      "localCopy": {
        "disk": "data",
        "path": "output.bam.bai"
      }
    },
  ],
  "resources": {
    "minimumCpuCores": 1,
    "minimumRamGb": 1,
    "disks": [{
      "name": "data",
      "type": "PERSISTENT_HDD"
      "sizeGb": 200,
      "mountPoint": "/mnt/data",
    }]
  },
  "docker": {
    "imageName": "quay.io/cancercollaboratory/dockstore-tool-samtools-index",
    "cmd": "samtools index /mnt/data/input.bam /mnt/data/output.bam.bai"
  }
}
```

# Create, run, monitor, and kill pipelines

## Create

```
$ gcloud alpha genomics pipelines create --pipeline-json-file PIPELINE-FILE.json --pipeline-json-file samtools_index.json
Created samtools index, id: PIPELINE-ID
```

## Run

```
$ gcloud alpha genomics pipelines run --pipeline_id PIPELINE-ID \
--logging gs://YOUR-BUCKET/YOUR-DIRECTORY/logs \
--inputs inputFile=gs://genomics-public-data/gatk-examples/example1/NA12878_chr22.bam \
--outputs outputFile=gs://YOUR-BUCKET/YOUR-DIRECTORY/output/NA12878_chr22.bam.bai
Running: operations/OPERATION-ID
```

## Status

```
$ gcloud alpha genomics operations describe OPERATION-ID
```

## Kill

```
$ gcloud alpha genomics operations cancel OPERATION-ID
```

# DSUB (google genomics pipelines)

googlegenomics / dsub

Watch 14   ★ Star 16   ⑂ Fork 4

<> Code   ⓘ Issues 11   Pull requests 0   Projects 0   Pulse   Graphs

Branch: master ▾   dsub / README.md   Find file   Copy path

```
./dsub \
    --project my-cloud-project \
    --zones "us-*" \
    --logging gs://my-bucket/logs \
    --input BAM=gs://genomics-public-data/1000-genomes/bam/HG00114.mapped.ILLUMINA.bwa.GBR.low_covera
    --output BAI=gs://my-bucket/HG00114.mapped.ILLUMINA.bwa.GBR.low_coverage.20120522.bam.bai \
    --image quay.io/cancercollaboratory/dockstore-tool-samtools-index \
    --command 'samtools index ${BAM} ${BAI}' \
    --wait
```
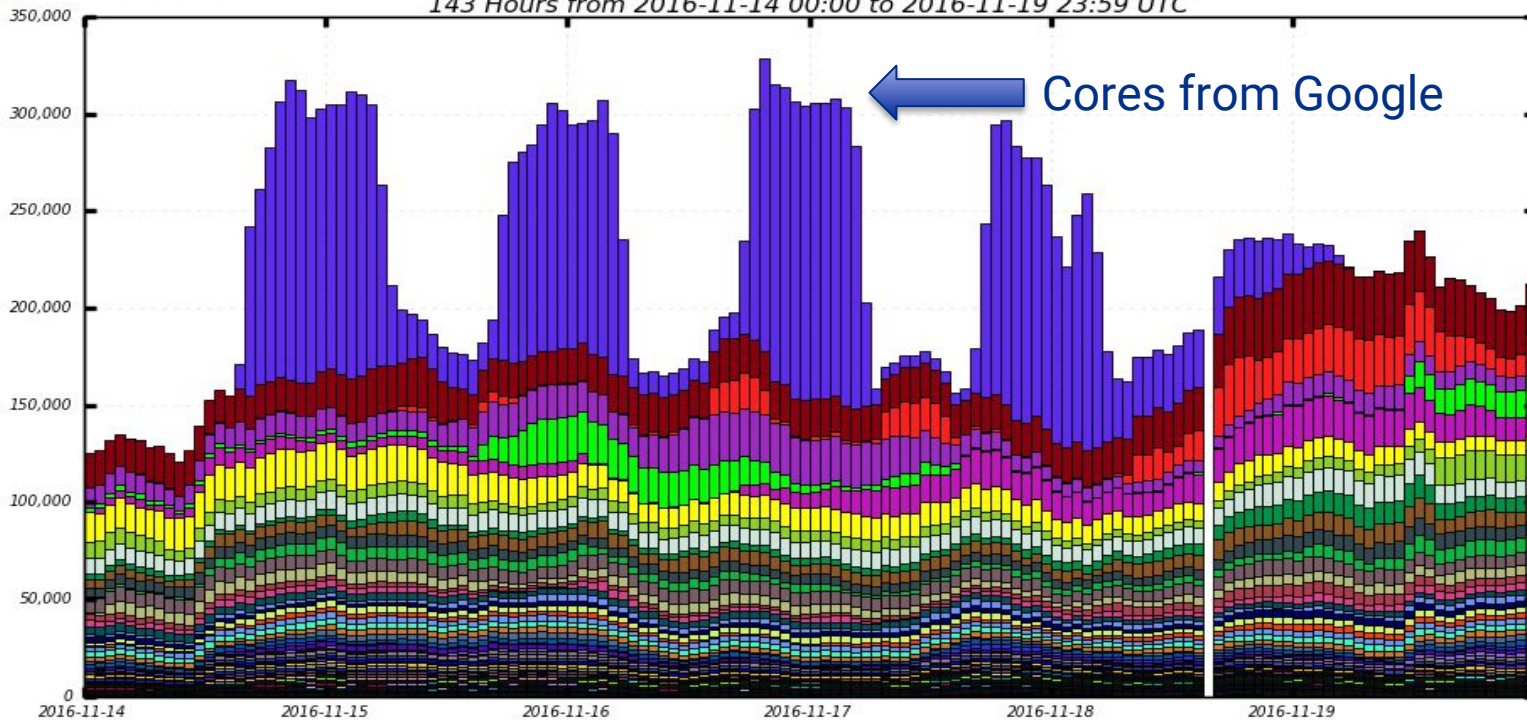
# SC16 CMS Demonstrator

Target: generate 1 Billion events in 48 hours during Supercomputing 2016 on Google Cloud via HEPCloud

35% filter efficiency = stage out 380 million events → 150 TB output

Double the size of global CMS computing resources



CMS Higgs Event - credit: CERN
https://commons.wikimedia.org/wiki/File:CMS_Higgs-event.jpg

# On-prem vs. Cloud

Average cost per core-hour (~25% error)
- On-premises Fermilab:
  0.9 cents per core-hour
  (assumes 100% utilization)
- Google Cloud:
  1.6 cents per core-hour
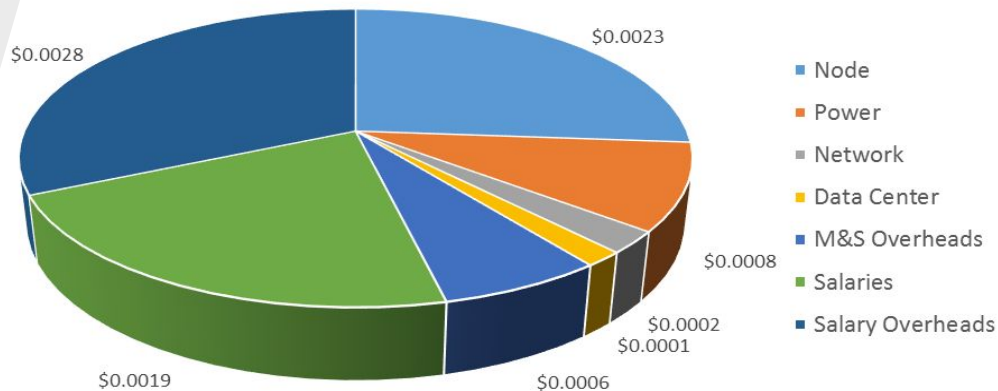  (comparable to other vendors)

Fermilab has years of experience in optimizing its facility

Cloud costs larger, but approaching equivalence

Considered well worth the cost of adding 160,000 core in a few hours
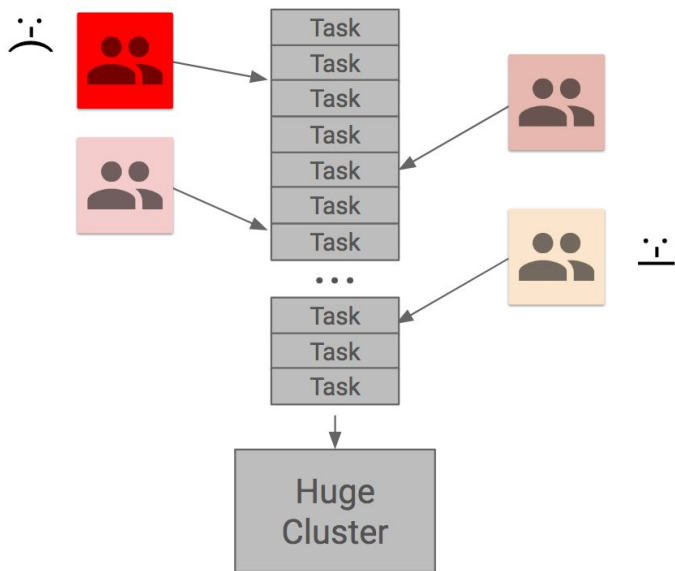
## Fermilab CMS Tier1 Costs

Cost per core-hour
Total = $0.0088



$0.0023
$0.0028
$0.0008
$0.0002
$0.0001
$0.0006
$0.0019

- Node
- Power
- Network
- Data Center
- M&S Overheads
- Salaries
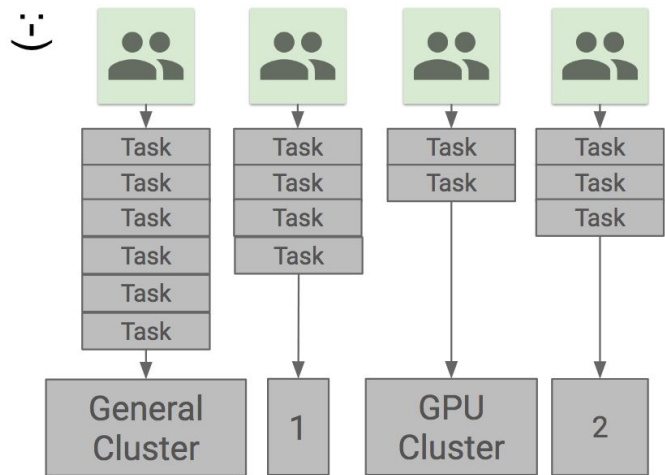- Salary Overheads

Google Cloud

# Task Tailored Resources

**On-prem**: One cluster, one queue, one size fits all hardware, angry users.

**In the cloud**: Tailored clusters, less queue sharing, happy users.



1 - n1-highcpu-16, Preemptible VM Cluster
2 - n1-highmem-32 Cluster

# Preemptible VM Instances

- **What Preemptible VMs are**
    - Up to 80% cheaper than regular VMs. (~$0.01 per core hour)
    - Very easy to use -- just flip one switch in the UI, API or command line
    - Many of our biggest customers run huge clusters (10k+ cores) with great success and savings.

- **Things to keep in mind**
    - Same great disk, OS images and network
    - Google Compute Engine can *preempt* (i.e. shutdown/take-away) the VM with 30 seconds of notice
    - Maximum 24 hours of uptime
    - No SLAs or guarantees of any kind but we historically see preemption rates of 5-15%

# Data

# Fully Managed Storage & Database Services

| Object | Key-value | Non-relational | | Relational | | Warehouse |
|--------|-----------|----------------|---|-----------|---|-----------|
| **Cloud Storage** | **App Engine Memcache** | **Cloud Datastore** | **Cloud Bigtable** | **Cloud SQL** | **Cloud Spanner** | **BigQuery** |
| Binary or object data | Web/mobile applications, gaming | Hierarchical, mobile, web | Heavy read + write, events | Web frameworks | RDBMS+scale, HA, HTAP | Enterprise Data Warehouse |
| Images, Media serving, backups | Game state, user sessions | User profiles, Game State | AdTech, Financial, IoT | CMS, eCommerce | Transactions, Ad/Fin/MarTech | Analytics, Dashboards |

# Block storage

Reliable, high-performance block storage for any GCE VM instance
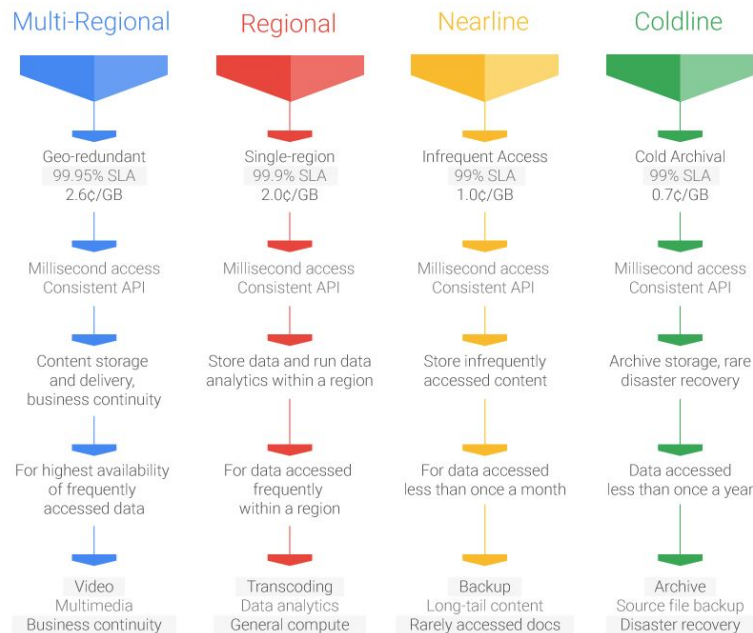
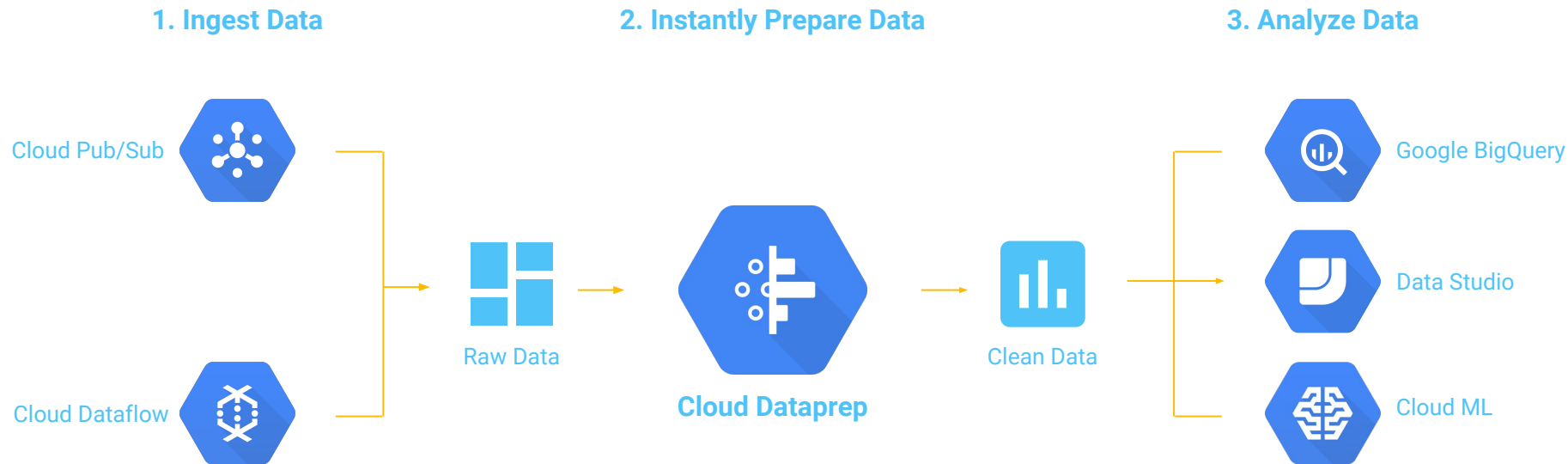| | Local SSD<br>Fastest, Attached, Ephemeral | Persistent Disk: SSD<br>Fast, Persistent, Durable, Remote | Persistent Disk: HDD<br>Cheapest, Persistent, Durable, Remote |
|---|---|---|---|
| **Target scenarios** | - High-performance scratch space. Frequently accessed data.<br>- Excellent for scientific workloads, especially when combined with fast compute VMs like GPU instances | - Latency sensitive applications and files.<br>- High performance database and enterprise applications<br>- Databases | - Large data processing workloads<br>- Latency incentive tasks with lots of data: Genomics processing, video transcoding in GCE |
| **Features** | - Ephemeral storage<br>- Highest-performance ($0.218 GB)<br>- **IOPS:** 680k read / 360k write | - Persistent storage<br>- Performance sensitive ($0.17GB)<br>- **IOPS:** up to 40k read / 30k write | - Persistent storage<br>- Cost sensitive ($.04 GB)<br>- **IOPS:** 3k read / 15k write |
| | Encryption<br>3TB - 375 GB per partition, up to 8 partitions | Encryption, Snapshots<br>64 TB, Disk Size sets performance<br>(Attach larger VMs for max SSD performance) | |

# GCS: Object/Blob store

- [Google Cloud Storage](#) is a scalable object storage service suitable for all kinds of unstructured data

- **Cloud Storage vs Perst. Disk:**
  - Scales to exabytes
  - Accessible from anywhere; REST interface
  - Higher latency than PD
  - Write semantics include insert and overwrite file only
  - Offers versioning
  - Cheaper - put your data here until you need it

- Lots of guidelines on picking storage on our [site](#)

## Google Cloud Storage Classes

| Multi-Regional | Regional | Nearline | Coldline |
|---|---|---|---|
| Geo-redundant 99.95% SLA 2.6¢/GB | Single-region 99.9% SLA 2.0¢/GB | Infrequent Access 99% SLA 1.0¢/GB | Cold Archival 99% SLA 0.7¢/GB |
| Millisecond access Consistent API | Millisecond access Consistent API | Millisecond access Consistent API | Millisecond access Consistent API |
| Content storage and delivery, business continuity | Store data and run data analytics within a region | Store infrequently accessed content | Archive storage, rare disaster recovery |
| For highest availability of frequently accessed data | For data accessed frequently within a region | For data accessed less than once a month | Data accessed less than once a year |
| Video Multimedia Business continuity | Transcoding Data analytics General compute | Backup Long-tail content Rarely accessed docs | Archive Source file backup Disaster recovery |

Next

Google Cloud

# Data Prep

**1. Ingest Data**

**2. Instantly Prepare Data**

**3. Analyze Data**

Cloud Pub/Sub

Cloud Dataflow

Raw Data

**Cloud Dataprep**

Clean Data

Google BigQuery

Data Studio

Cloud ML

Google Cloud

# Cloud Dataprep

**Instant Data Exploration**
Visually explore and interact with data in seconds. Instantly understand data distribution and patterns. There is no need for one to write code. You can prepare data with a few clicks.

**Intelligent Data Cleansing**
Cloud Dataprep automatically identifies data anomalies and helps you to take corrective actions fast. Get data transformation suggestions based on your usage pattern. Standardize, structure, and join datasets easily with a guided approach.

**Serverless**
Cloud Dataprep is a serverless service, so you do not need to create or manage infrastructure.

**Seriously Powerful**
Cloud Dataprep is built on top of powerful Google Cloud Dataflow service. Cloud Dataprep is auto-scalable and can easily handle processing massive data sets.

**Supports Common Data Sources of Any Size**
Process diverse datasets - structured and unstructured. Transform data stored in CSV, JSON, or relational Table formats. Prepare datasets of any size, megabytes to terabytes, with equal ease.

Google Cloud

Grid Columns | Untitled Flow > package_log_3 – 2

Full Dataset - 785.29kB | 5 Columns | 20,000 Rows | 3 Data Types

Preview

| # column2 | 🏴 column3 | # column4 | # column5 |
|---|---|---|---|
| 0 - 1 | 13 Categories | | |
| 0 | AZ | | |
| 0 | CA | | |
| 0 | AK | | |
| 0 | NM | | |
| 0 | WA | | |
| 0 | KS | | |
| 0 | OK | | |
| 0 | OK | | |
| 0 | MN | | |
| 0 | MN | | |
| 1 | MN | | |
| 1 | MN | | |
| 1 | OK | | |
| 1 | OK | | |
| 1 | KS | | |

SUGGESTIONS

Delete rows where sourcerownumber() == 4

| # column2 | colum |
|---|---|
| 0 | NM |
| 0 | AZ |
| 0 | CA |

Affects all columns, 1 row

---

Untitled Flow > package_log_3 – 2

Full Dataset - 785.29kB | 7 Columns | 20,000 Rows | 3 Data Types | Columns: ✓ All | Transformed - 3 Columns

Source | to be dropped | Preview

| # column4 | # column5 | ABC column6 | ABC column1 | # column7 |
|---|---|---|---|---|
| 1.42B - 1.45B | 0 - 499 | 10,000 Categories | 1,000 Categories | 152 - 9M |
| 1424860400 | 37 | 1LOCWA4000790 | 1LOC | 4000790 |
| 1430020400 | 495 | 1LOCKS5000950 | 1LOC | 5000950 |
| 1438298400 | 395 | 1LOCOK6000228 | 1LOC | 6000228 |
| 1438380400 | 72 | 1LOCOK7000310 | 1LOC | 7000310 |
| 1436414400 | 320 | 1LOCMN8000344 | 1LOC | 8000344 |
| 1438167400 | 266 | 1LOCMN900097 | 1LOC | 900097 |
| 1438271020 | 195 | 1LOCMN900097 | 1LOC | 900097 |
| 1437008760 | 351 | 1LOCMN8000344 | 1LOC | 8000344 |
| 1438623160 | 429 | 1LOCOK7000310 | 1LOC | 7000310 |
| 1438557900 | 96 | 1LOCOK6000228 | 1LOC | 6000228 |
| 1430173580 | 164 | 1LOCKS5000950 | 1LOC | 5000950 |

Switch to editor | Cancel | Add to Recipe

Transformation | Column required | Between two positions ⌄ ? | required | Number

column6 | starting from 5 | Number

column6 | ending at 7

---

Untitled Flow > package_log_3 – 2

Grid Columns | Full Dataset - 785.29kB | 5 Columns | 20,000 Rows | 3 Data Types

# column5

Overview | Patterns

SUMMARY | TOP VALUES | MISMAT
1 → → 0 | Run Job | None

57 56 56 55 54 54 54 54 54 54 53 53 53

OUTLI
None

Untitled Flow > package_log_3 ⌄

Untitled Flow > package_log_3 ⌄                                   1 → ⊞ → 0    Run Job

Grid    Columns    Full Dataset - 785.29kB ⌄    5 Columns    21,000 Rows    3 Data Types

# column5                                                              ✕

Overview    Patterns

**SUMMARY**

| | | |
|---|---|---|
| Valid ● | 20,000 | 95.2% |
| Unique | 500 | 2.4% |
| Outliers | 0 | 0.0% |
| Mismatched ● | 0 | 0.0% |
| Missing ● | 1,000 | 4.8% |

**STATISTICS**

| | |
|---|---|
| Minimum | 0.00 |
| Lower Quartile | 124.00 |
| Median | 249.00 |
| Upper Quartile | 373.00 |
| Maximum | 499.00 |
| Average | 249.05 |
| Standard Deviation | 144.17 |

**TOP VALUES**

| | |
|---|---|
| 369 | 57 |
| 111 | 56 |
| 37 | 56 |
| 386 | 55 |
| 11 | 54 |
| 149 | 54 |
| 273 | 54 |
| 280 | 54 |
| 368 | 54 |
| 379 | 54 |
| 452 | 54 |
| 172 | 53 |
| 237 | 53 |
| 300 | 53 |

**MISMATCHED VALUES**

None

**OUTLIERS**

None

**VALUE HISTOGRAM**

💡 SUGGESTIONS                                    Cancel    Modify    Add to Recipe

| Keep rows where ismismatched(column3, ['State']) | Delete rows where ismismatched(column3, ['State']) | Create a new column from ismismatched(column3, ['State']) | Set column3 to IFMISMATCHED(column3, ['State |
|---|---|---|---|
| **column3** | **column3** | **column3** | **column1** | **column3** | col |
| · 401 · ERROR | · 401 · ERROR | · AZ | false | · AZ | · AZ |
| · 401 · ERROR | · 401 · ERROR | · CA | false | · CA | · CA |
| · 401 · ERROR | · 401 · ERROR | · AK | false | · AK | · AK |
| Affects all columns, 1000 rows | Affects all columns, 1000 rows | Affects 1 column, all rows | Creates 1 column | Affects 1 column, all rows | Changes 1 column |

Google Cloud

# Machine Learning

# Two flavors of machine learning

**API**

Vision    Translation    Speech [BETA]    Natural Language

## Pre-Trained Models

Storage    BigQuery    Datalab    Tensor Flow    Pipelines    Model Management

## Build Your Own Model

Google Cloud

# Google Cloud Machine Learning Services

cloud.google.com/translate/

Enter a word or phrase:

Hello world

Translate from:

English

Translate to:

Albanian

TRANSLATE

cloud.google.com/natural-language/

**Try the API**

Google, headquartered in Mountain View, unveiled the new Android phone at the Consumer Electronic Show.  Sundar Pichai said in his keynote that users love their new Android phones.

ANALYZE

Enter text in English, Spanish or Japanese

cloud.google.com/vision/

Try the API

Drag image file here or
Browse from your computer

cloud.google.com/speech/

Convert your voice to text right now

Click on the microphone icon to start recording

English (United States)

Google Cloud

# Cloud ML Engine

- **PaaS** for Tensorflow

- **Scale** your training up to 100 workers

- Automatic **monitoring** and **logging**

- Easy transition from training to **prediction**

- Built in model **version management**

- **No lock-in.** Option to download your trained models for on-premise or mobile deployment

# CloudML is part of a bigger picture



| Capture | Store | Process | Analyze | Insight |
|---------|-------|---------|---------|---------|
| Pub/Sub | Cloud Storage | Dataflow | BigQuery | Cloud ML Engine |
| | BigQuery | Dataproc | Dataflow | |
| | Cloud SQL | | Datalab | |
| | Datastore | | | |
| | BigTable | | | |

Google Cloud

# TensorFlow



- World's most popular ML framework
- Developer friendly yet performance optimized
- **Powers over 100 Google services**
- Managed infrastructure with Cloud ML
- Tutorials at https://www.tensorflow.org

# Linear Regression VS Neural Network

```
1  import tensorflow as tf
2
3  #Define input feature columns
4  sq_footage = tf.contrib.layers.real_valued_column("sq_footage")
5  feature_columns = [sq_footage]
6
7  #Define input function
8  def input_fn(feature_data,label_data=None):
9    return {"sq_footage":feature_data}, label_data
10
11 #Instantiate Linear Regression Model
12 estimator = tf.contrib.learn.LinearRegressor(
13   feature_columns=feature_columns,
14   optimizer=tf.train.FtrlOptimizer(learning_rate=100))
15
16 #Train
17 estimator.fit(
18   input_fn=lambda:input_fn(tf.constant([1000,2000]),
19                     tf.constant([100000,200000])),
20   steps=100)
21
22 #Predict
23 estimator.predict(input_fn=lambda: input_fn(tf.constant([3000])))
```
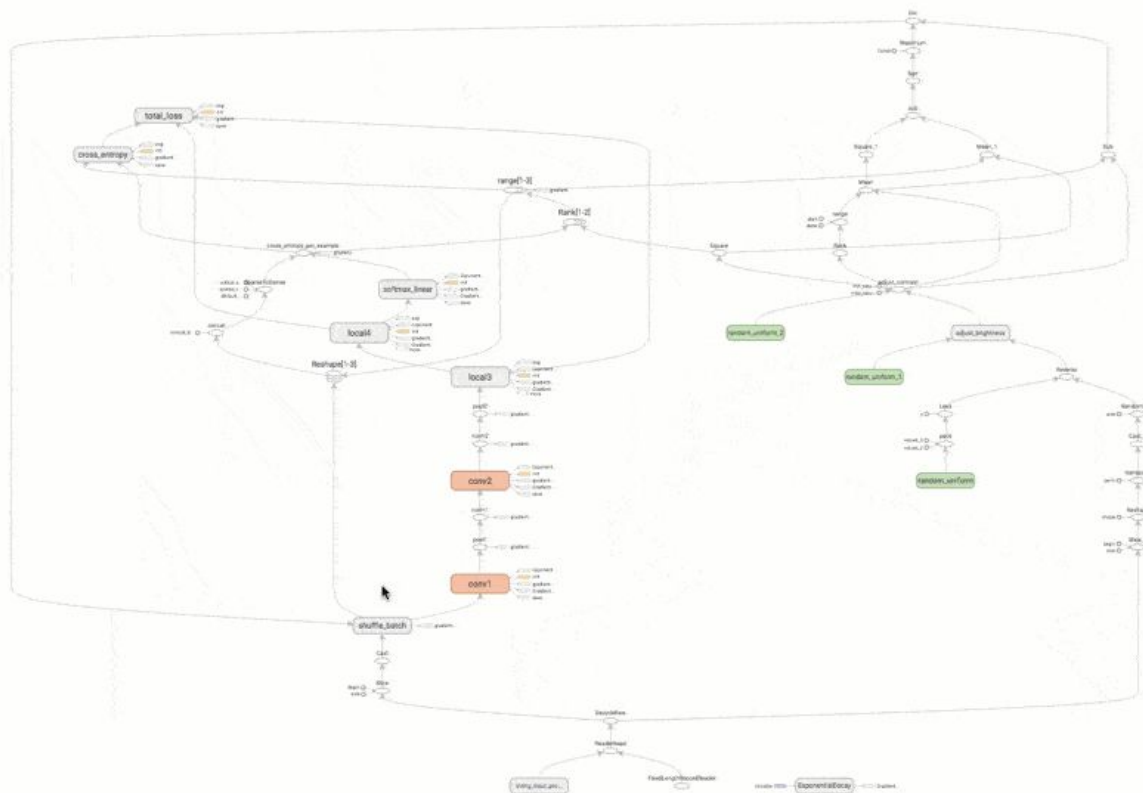
```
1  import tensorflow as tf
2
3  #Define input feature columns
4  sq_footage = tf.contrib.layers.real_valued_column("sq_footage")
5  feature_columns = [sq_footage]
6
7  #Define input function
8  def input_fn(feature_data,label_data=None):
9    return {"sq_footage":feature_data}, label_data
10
11 #Instantiate Neural Network Model
12 estimator = tf.contrib.learn.DNNRegressor(
13   feature_columns=feature_columns, hidden_units=[10,10])
14
15
16 #Train
17 estimator.fit(
18   input_fn=lambda:input_fn(tf.constant([1000,2000]),
19                     tf.constant([100000,200000])),
20   steps=100)
21
22 #Predict
23 estimator.predict(input_fn=lambda: input_fn(tf.constant([3000])))
```

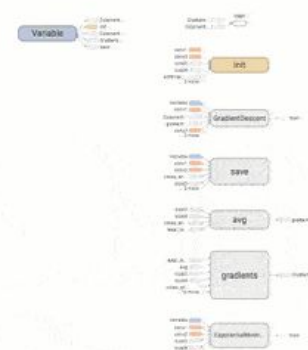Google Cloud

Fit to screen

Run    cifar-train

Upload    Choose File

Color    Structure

color: same substructure
gray: unique substructure

Main Graph

Auxiliary nodes

Variable

init

GradientDescent

save

avg

gradients

ExponentialMovin...

total_loss

cross_entropy

range[1-2]

Rank[1-2]

cross_entropy_per_example

softmax_linear

local4

local3

conv2

conv1

shuffle_batch

random_uniform_2

random_uniform_1

adjust_brightness

random_uniform

Reshape[1-2]

Graph    (* = expandable)

Namespace*

OpNode

Unconnected series*

Connected series*

Constant

Summary

Dataflow edge

Control dependency edge

Reference edge

# Copy Number se

The goal of this notebook is to in

This table contains all available T
Genome Wide SNP6 array, as of
recent archives (egbroad.mit.
types was downloaded from the
Each of these segmentation files
During ETL the sample identifer
the SDRF file in the associated m

In order to work with BigQuery,
the name(s) of the table(s) you a

```
import gcp.bigquery a
cn_BQtable = bq.Table
```

From now on, we will refer to thi
table name each time.

Let's start by taking a look at the

```
%bigquery schema --ta
```

| name | type |
|------|------|
| ParticipantBarcode | STRIN |
| SampleBarcode | STRIN |
| SampleTypeLetterCode | STRIN |
| AliquotBarcode | STRIN |
| Study | STRIN |
| Platform | STRIN |
| Chromosome | STRIN |
| Start | INTE |
| End | INTE |
| Num_Probes | INTE |
| Segment_Mean | FLOA |

Unlike most other molecular dat
microRNAs, this data is produce
sizes and positions of these segn

---

Now we'll use matplotlib to create some simple visual

```
import numpy as np
import matplotlib.pyplot as plt
```

For the segment means, let's invert the log-transform

```
%%sql --module getCNhist

SELECT
    lin_bin,
    COUNT(*) AS n
FROM (
    SELECT
        Segment_Mean,
        (2.*POW(2,Segment_Mean)) AS lin_
        INTEGER(((2.*POW(2,Segment_Mean)
    FROM
        $t
    WHERE
        ( (End-Start+1)>1000 AND SampleT
GROUP BY
    lin_bin
HAVING
    ( n > 2000 )
ORDER BY
    lin_bin ASC
```

```
CNhist = bq.Query(getCNhist,t=cn_BQt
bar_width=0.80
plt.bar(CNhist['lin_bin']+0.1,CNhist
plt.xticks(CNhist['lin_bin']+0.5,CN
plt.title('Histogram of Average Copy
plt.ylabel('# of segments');
plt.xlabel('integer copy-number');
```

Histogram of Average Co



The histogram illustrates that the vast majority of the
either side representing deletions (left) and amplific
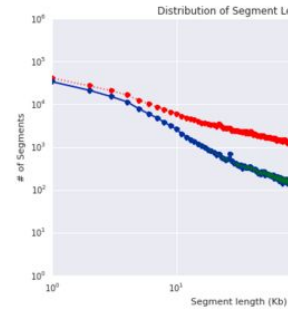
---

```
    bin
ORDER BY
    bin ASC
```

```
%%sql --module getSLhist_1k_amp

SELECT
    bin,
    COUNT(*) AS n
FROM (
    SELECT
        (END-Start+1) AS segLength,
        INTEGER((END-Start+1)/1000) AS b
    FROM
        $t
    WHERE
        (END-Start+1)<1000000 AND Sample
GROUP BY
    bin
ORDER BY
    bin ASC
```

```
SLhistDel = bq.Query(getSLhist_1k_de
SLhistAmp = bq.Query(getSLhist_1k_am
```
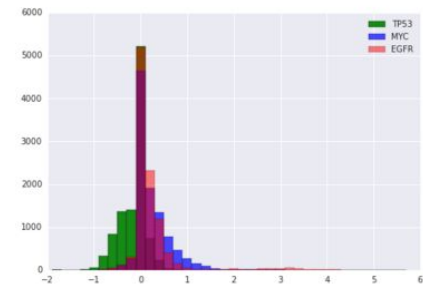
```
plt.plot(SLhist_1k['bin'],SLhist_1k[
plt.plot(SLhistDel['bin'],SLhistDel[
plt.plot(SLhistAmp['bin'],SLhistDel[
plt.xscale('log');
plt.yscale('log');
plt.xlabel('Segment length (Kb)');
plt.ylabel('# of Segments');
plt.title('Distribution of Segment L
```
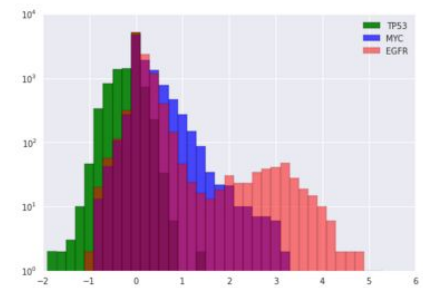
Distribution of Segment Le



The amplification and deletion distributions are nearly
from this graph that a majority of the segments less th
lengths >100Kb are copy-number neutral.

---

And now we'll take a look at histograms of the average copy-number for these three genes. TP53 (in green) shows a significant number of partial deletions (CN<0), while MYC (in blue) shows some partial amplifications -- more frequently than EGFR, while EGFR (pale red) shows a few extreme amplifications (log2(CN/2) > 2). The final figure shows the same histograms on a semi-log plot to bring up the rarer events.
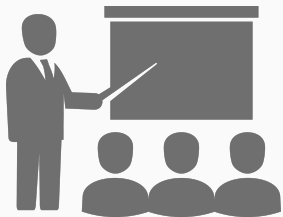
```
binwidth = 0.2
binvals = np.arange(-2+(binwidth/2.), 6-(binwidth/2.), binwidth)
plt.hist(tp53CN['avgCN'],bins=binvals,normed=False,color='green',alpha=0.9,label='TP53');
plt.hist(mycCN ['avgCN'],bins=binvals,normed=False,color='blue',alpha=0.7,label='MYC');
plt.hist(egfrCN['avgCN'],bins=binvals,normed=False,color='red',alpha=0.5,label='EGFR');
plt.legend(loc='upper right');
```



```
plt.hist(tp53CN['avgCN'],bins=binvals,normed=False,color='green',alpha=0.9,label='TP53');
plt.hist(mycCN ['avgCN'],bins=binvals,normed=False,color='blue',alpha=0.7,label='MYC');
plt.hist(egfrCN['avgCN'],bins=binvals,normed=False,color='red',alpha=0.5,label='EGFR');
plt.yscale('log');
plt.legend(loc='upper right');
```
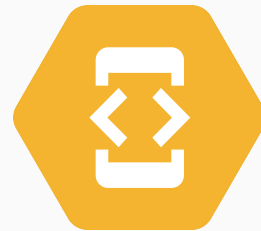
# Programs

Teaching

Faculty in select countries + Teaching university courses + In computer science or related fields

Google Cloud

# Funding Agency Partnerships

- ## National Science Foundation
  - BIGDATA

- ## National Institutes of Health
  - Data Commons

Google Cloud Public Datasets Program

**Mission:**
Facilitate the onboarding of datasets into Google Cloud products

# Thank you

KaranBhatia@google.com
@sdksb

Google Cloud