# Google Cloud HPC & AI Capabilities

## Overview

Leonid Kuligin
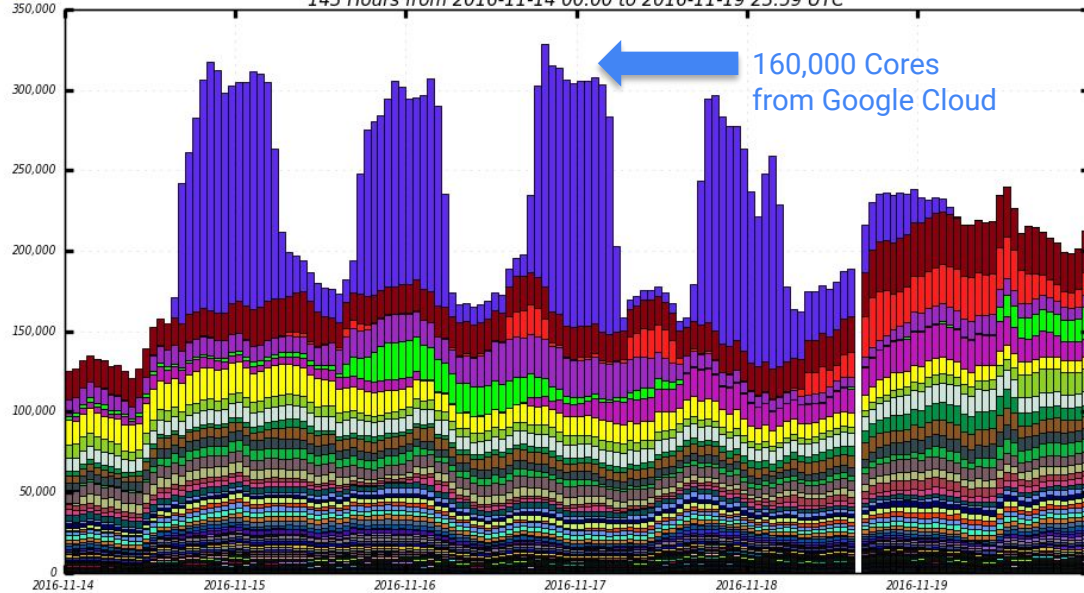
Karan Bhatia

Dec 6 2019

# Doubling LHC CMS compute capacity.

# Combating Traffic Congestion using Massive HPC Analytics in the Google Cloud Platform



**Follow up 2,138,000 vCPU Cluster**

GCP CPU Core Ramp and Count

- Sustained over 2 million vCPU for over an hour
- 2,138,000 vCPUs and 133,573 instances at peak
- Average of $0.008 USD per vCPU hour
- Only took 1.5 hours to hit our previous 1.5 million vCPUs

CLEMSON UNIVERSITY

> 1Tbps throughput

Cloud Storage Throughput - Send

- 128 GiB/s peak in Cloud Storage
- 59 GiB/s peak in one bucket

Data Intensive Computing Ecosystems (DICE)
School of Computing, Clemson University

TrafficVision

# Rediscovering Higgs boson on GCP

- https://github.com/mmm/higgs-tutorial
- Kubecon Barcelona 2019 keynote
  https://www.youtube.com/watch?v=CTfp2woVEkA

Google Cloud

# Google Cloud Background

# HPC Infrastructure

Compute

Storage

Networking

Google Cloud

# TPUs are ASICs focused on Machine Learning



**TPU v1**
**(2015)**
92 teraops
First Generation

**TPU v2**
**(2017)**
180 teraflops
Available via Google Cloud

**TPU v3**
**(2018)**
420 teraflops
Available via Google Cloud
~2.3x the power of v2

**Edge TPU**
**(2018 EAP)**
Inference Accelerator

Google Cloud

# Cloud TPU Pods - Product Offerings



### Cloud TPU v2 Pod[Beta]
11,500 teraflops
Up to 256 chips
4,000 GB HBM
2-D toroidal mesh network



### Cloud TPU v3 Pod[Beta]
100,000+ teraflops
Up to 1,024 chips
32,000 GB HBM
2-D toroidal mesh network

Cloud TPU Configurations

- <u>TPU</u>: The Tensor Processing Unit (TPU) is a custom-design chip, built from the ground up by Google for machine learning workloads.
- <u>Cloud TPU</u>: a device containing four TPU chips along with a fraction of a CPU host.
- <u>Cloud TPU pods</u>: Cloud TPUs are connected via a high-speed 2D toroidal mesh network to form Cloud TPU Pods.
- <u>Cloud TPU slices</u>: Slices, or smaller sections of pods, are scalable to address as much performance is needed for the workload. Slices are internal allocations consisting of different numbers of TPU cores. Pod slices come in 32, 128, 256, 512, 1024, and 2048 core-count configurations.



Hosts    Hosts

TPU v2-32
(32 cores, 4x4 slice)

TPU v2-128
(128 cores, 8x8 slice)

TPU v2-256
(256 cores, 8x16 slice)

# The network matters

134 points of presence and 13 subsea cable investments around the globe



Current regions and number of zones

Future regions and number of zones

Edge points of presence

CDN nodes

Network

Dedicated Interconnect

Seattle
San Francisco
Los Angeles
Denver
Dallas
Oregon
Salt Lake City
Los Angeles
Chicago
Montréal
Toronto
Iowa
N. Virginia
S. Carolina
New York
Washington DC
Atlanta
Miami
Montréal

Rio de Janeiro
São Paulo
São Paulo
Buenos Aires

Amsterdam
London
Stockholm
Finland
Netherlands
London
Belgium
Frankfurt
Hamburg
Frankfurt
Paris
Madrid
Zurich
Zurich
Milan
Marseille
Munich

Seoul
Tokyo
Osaka
Tokyo
Hong Kong
Hong Kong
Taiwan
Taipei
Osaka
Mumbai
Mumbai
Chennai
Kuala Lumpur
Singapore
Singapore
Jakarta
Sydney
Sydney

JUC (JP, HK, SG) 2

# HPC → Machine Learning

Google

# ML improves with data size



The unreasonable effectiveness of data
https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/35179.pdf

Deep Learning scaling is predictable, empirically
https://arxiv.org/abs/1712.00409

Google Cloud

# Increases in accuracy require much more compute.



Learning Transferable Architectures for Scalable Image Recognition
Barret Zoph, Vijay Vasudevan, Jonathon Shlens, Quoc V. Le
https://arxiv.org/abs/1707.07012

AI and Compute

MAY 16, 2018

OpenAI

https://blog.openai.com/ai-and-compute/

Google Cloud

### AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



Petaflop/s-day (Training) vs Year

- AlphaGo Zero
- AlphaZero
- Neural Machine Translation
- Neural Architecture Search
- TI7 Dota 1v1
- Xception
- DeepSpeech2
- VGG
- Seq2Seq
- ResNets
- GoogleNet
- AlexNet
- Visualizing and Understanding Conv Nets
- Dropout
- DQN

LOG SCALE   LINEAR SCALE

## AI and Compute

MAY 16, 2018

OpenAI

https://blog.openai.com/ai-and-compute/

Google Cloud

**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute**



New capabilities

Compute

Time

- AlphaGo Zero
- AlphaZero
- Neural Machine Translation
- Neural Architecture Search
- TI7 Dota 1v1
- Xception
- DeepSpeech2
- VGG
- ResNets
- Seq2Seq
- GoogleNet
- Visualizing and Understanding Conv Nets
- AlexNet
- Dropout
- DQN

10,000
1,000
100
10
1
.1
.01
.001
.0001
.00001

2013 2014 2015 2016 2017 2018 2019

LOG SCALE  LINEAR SCALE

**AI and Compute**

MAY 16, 2018

"

... since 2012, the amount of compute used in the largest AI training runs **has been increasing exponentially with a 3.5-month doubling time** (by comparison, Moore's Law had an 18-month doubling period).

Since 2012, **this metric has grown by more than 300,000x** (an 18-month doubling period would yield only a 12x increase).



**OpenAI**

... within many current domains, **more compute seems to lead predictably to better performance**, and is often complementary to algorithmic advances.

... we believe **the relevant number** is not the speed of a single GPU, nor the capacity of the biggest datacenter, but **the amount of compute that is used to train a single model** – this is the number most likely to correlate to how powerful our best models are.

https://blog.openai.com/ai-and-compute/

# Leading to increased model quality

% ranking on Kaggle private leaderboard



■ AutoML Tables

For each product:
● Relevant tables joined by given IDs
● Some minimal preprocessing done to match input requirements
● Run until converge
● Benchmarks run between H2 2018 to today (as they became available)

# AutoML + experts = even better!



Competition Submissions vs AutoML Benchmark

- AutoML Tables Benchmark
- 90-Percentile of Submissions

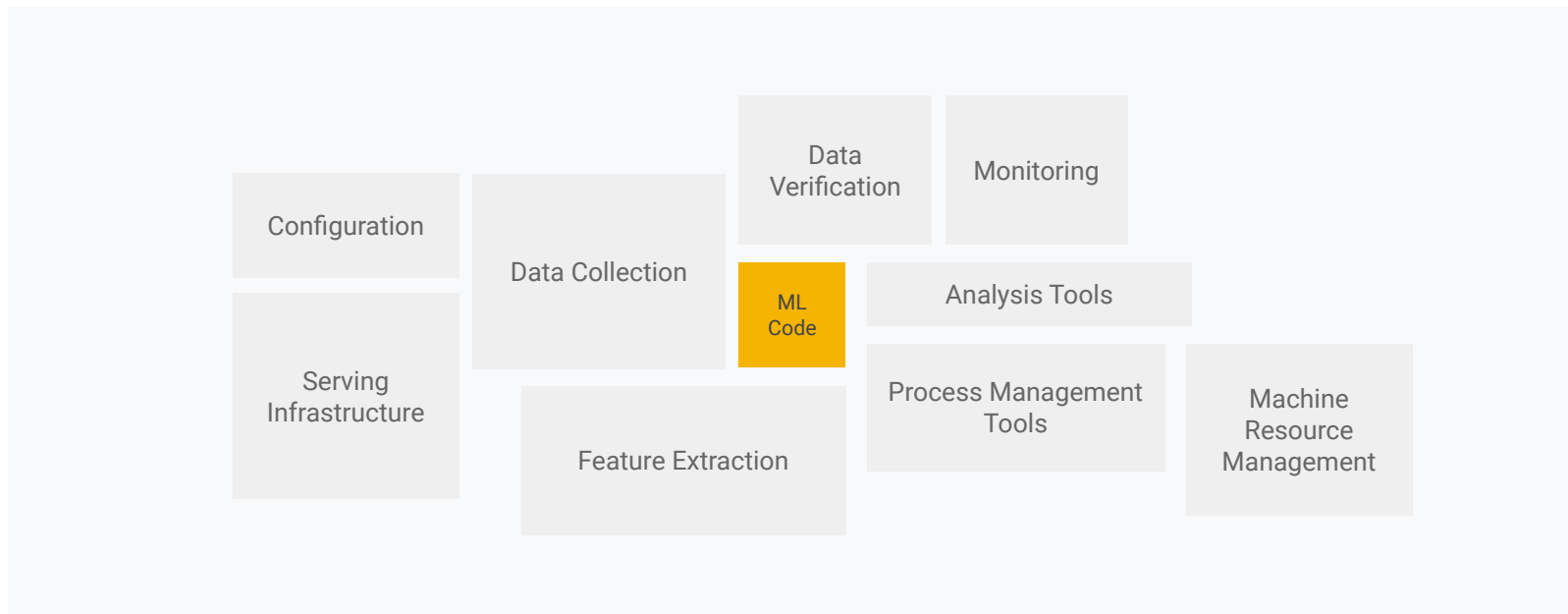| team | score |
|---|---|
| Erkut & Mark,Google AutoML | 0.618492 |
| Erkut & Mark | 0.616913 |
| Google AutoML | 0.615982 |
| Erkut & Mark,Google AutoML,Sweet Deal | 0.615858 |
| Sweet Deal | 0.615766 |
| Arno Candel @ H2O.ai | 0.615492 |
| ALDAPOP | 0.615040 |
| 9hr Overfitness | 0.614371 |
| Shlandryn | 0.614132 |
| Erin (H2O AutoML 100 mins) | 0.612657 |

# To do ML in production, in addition to the actual ML...
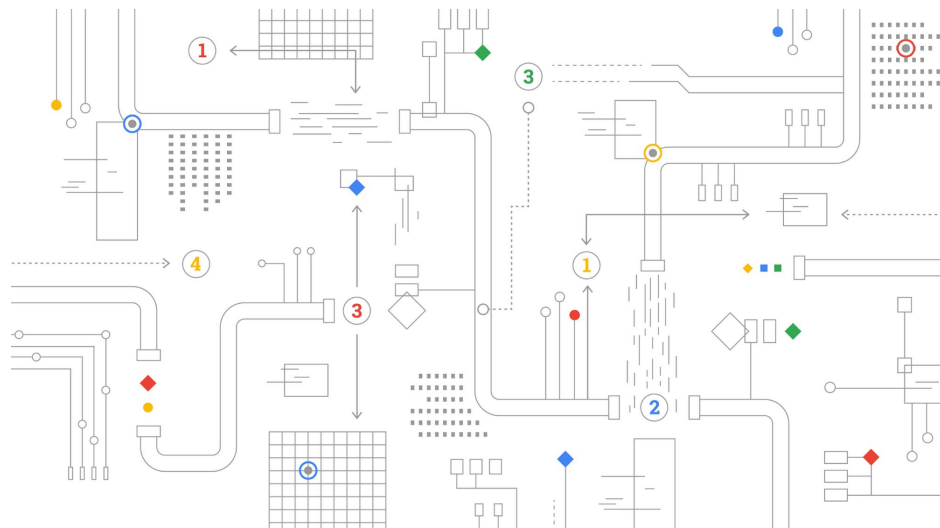
# ...you have to worry about so much more.

Google Cloud

# Operatize ML pipelines, not ML models

- Input validation
- Model retraining
- Reusable and shareable components
- ML microservices
- Serverless

# AI Hub

## Public Content

## + Private Content

### By Google

Unique AI
assets by Google

AutoML, TPUs, kaggle
Cloud AI Platform, etc.

Research at Google

DeepMind

### By Partners

Created, shared
& monetized
by anyone.

### By Customers

Content shared Securely within
and with other organizations.

Google Cloud

# AI Platform Overview

| | Data Readiness | Feature Engineering | Training/ HP-Tuning | Model Management | Prediction | Understanding / Tuning | Productionize |
|---|---|---|---|---|---|---|---|
| **Fully Managed** | **AI Platform** Labeling | | Training | Hub Metadata | Prediction | Explainability | |
| **Semi-Managed** | | DL Environment (DL VM + DL Container) | | | | | |
| | | Notebooks | | | | | |
| **On Prem & DIY** | | | **Kubeflow** Training | Hub / Metadata | Prediction | Catib | Serving |
| | | | Pipelines | | | | |

# Choose operational complexity based on level of control needed



**Cloud AI Platform**

Turnkey serverless training and batch/online predictions

**Kubeflow pipelines on Google Kubernetes Engine**

Fully customizable

**Deep Learning VM Cloud Functions Cloud Dataflow**

Customized training, online/batch predictions

Google Cloud

# Takeaways

### GCP for HTC and AI

GCP well suited for high throughput computing with many partners, schedulers and cost effective solutions.

### AI capabilities both for quick prototyping as well as scaled training

Cloud is ideal both for quick prototyping, sharing and reusing code and ML models, as well as for reproducible workloads (both large-scale training and high-throughput inference)

### Scale and hybrid approach

You can use GCP capabilities to achieve enormous scale to add them to your existing on-premise resources

Google

Thank you

Google

Karan Bhatia
karanbhatia@google.com
Leonid Kuligin
kuligin@google.com