



Hybrid Worlds of Cloud and HPC

and the Challenges for System Software

Kevin KISSELL, Google Cloud Office of the CTO

October 13, 2020

An Agricultural Revolution on the HPC Server Farm

Stresses on the Ecosystem

- End of Moore's Law

- Limits of symmetric exascale HPC

Adaptive Mutations

- New computational models and paradigms

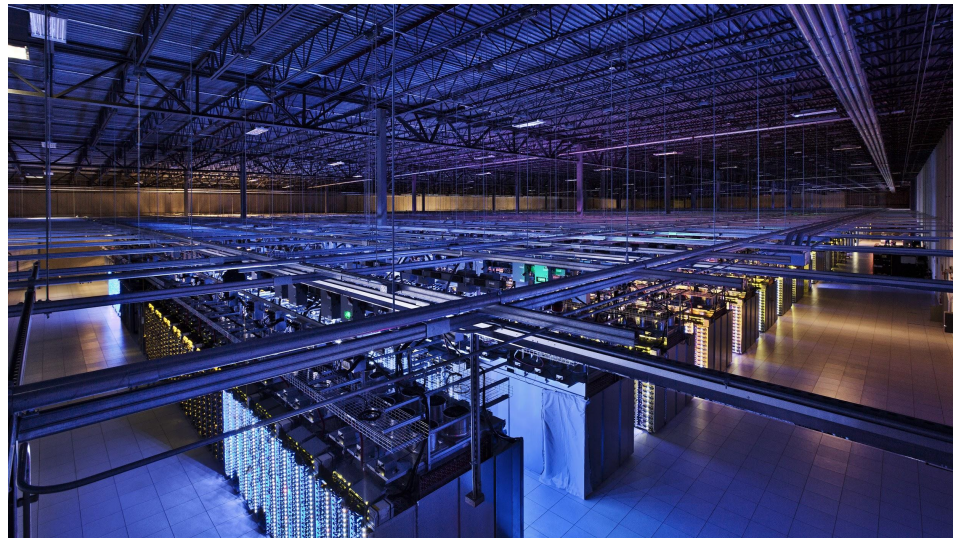
- New models for storage

Changes to the Climate

- Multi-level virtualization

- Robust planetary networks

**We are moving away from monoculture,
and breeding hybrids**



Hybrid Deployment On-Premises/Cloud

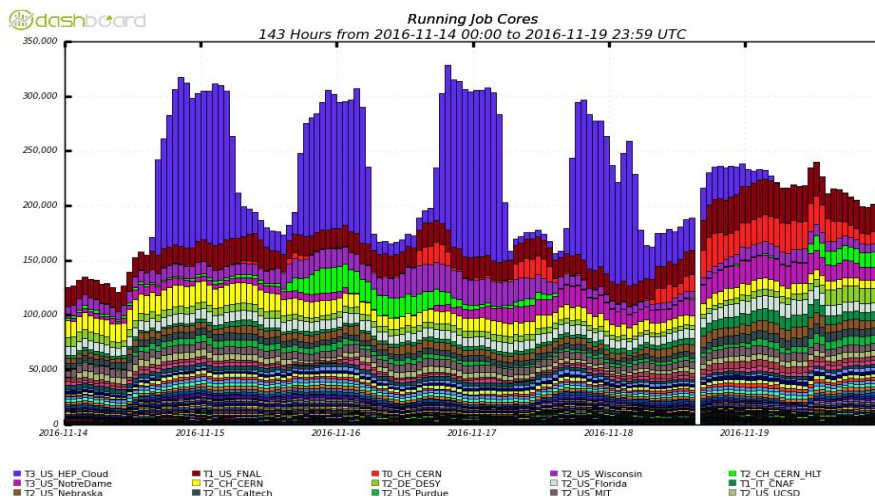
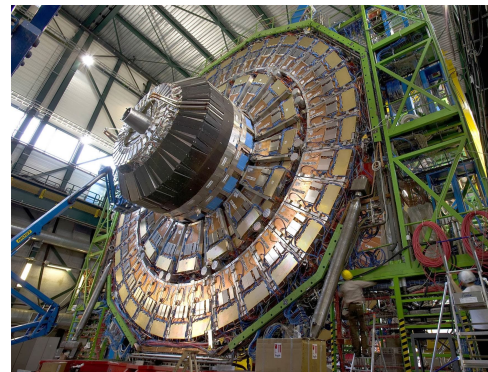
CERN LHC Experiments 2016

Fermilab/CMS

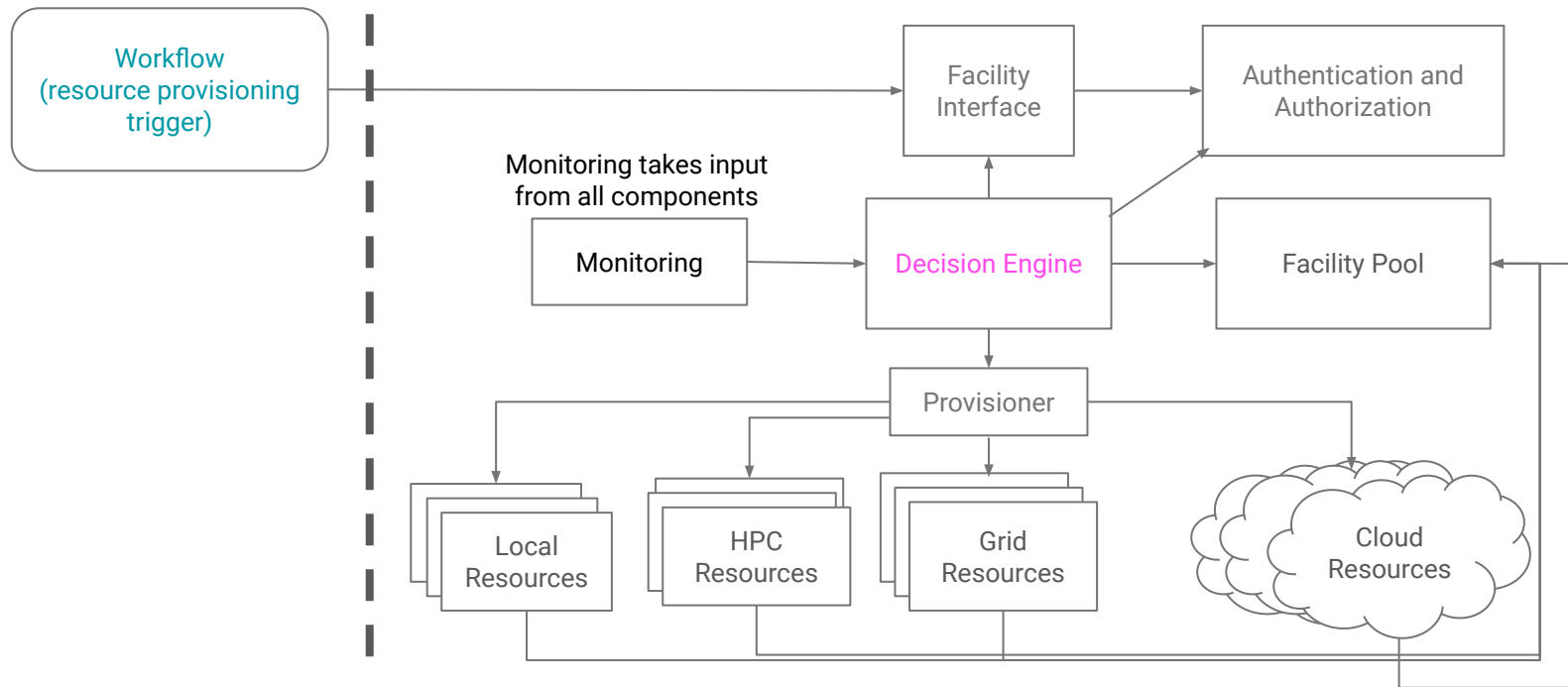
HTCondor made Google Cloud aware

Added 160 000 virtual cores to HEPCloud

Roughly doubled HEPCloud capacity during SC16



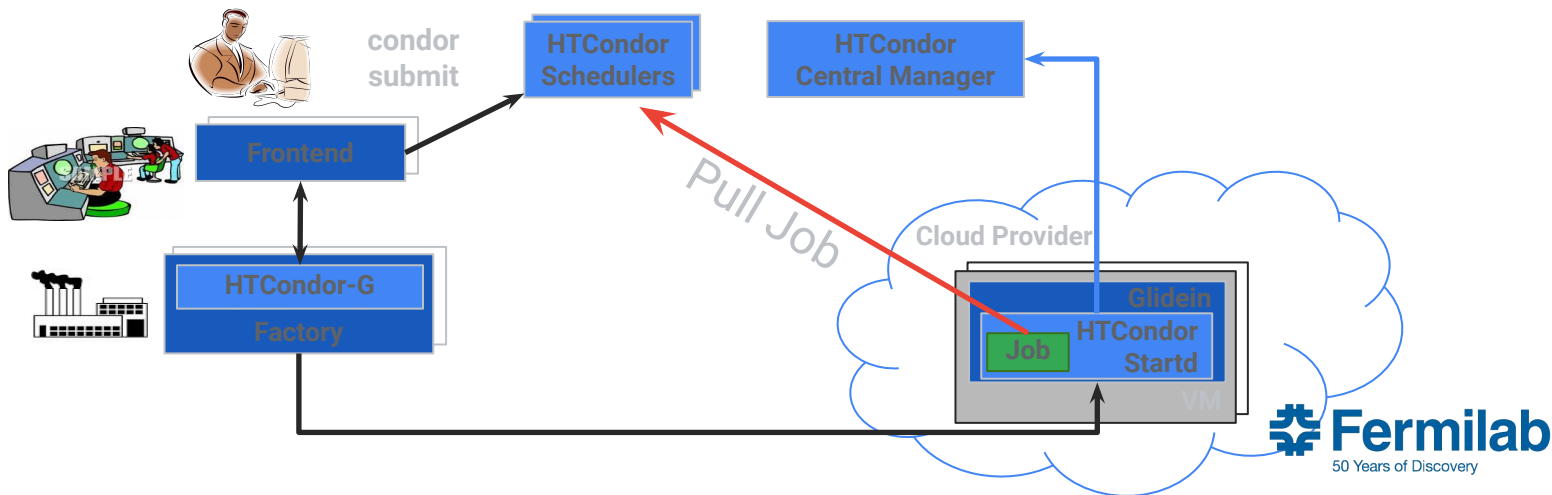
Hybrid Deployment HEPCloud Architecture



Hybrid Deployment

Provisioning remote resources via glideinWMS

- GlideinWMS submits “pilot jobs” to compute resources based on demand
- Pilot jobs execute on the resource and fetch user jobs from a queue
 - Pilot jobs hide heterogeneity of compute from the user and validate environment (will not start user jobs on bad resources)



Hybrid Deployment On-Premises/Cloud

CERN LHC Experiments 2019

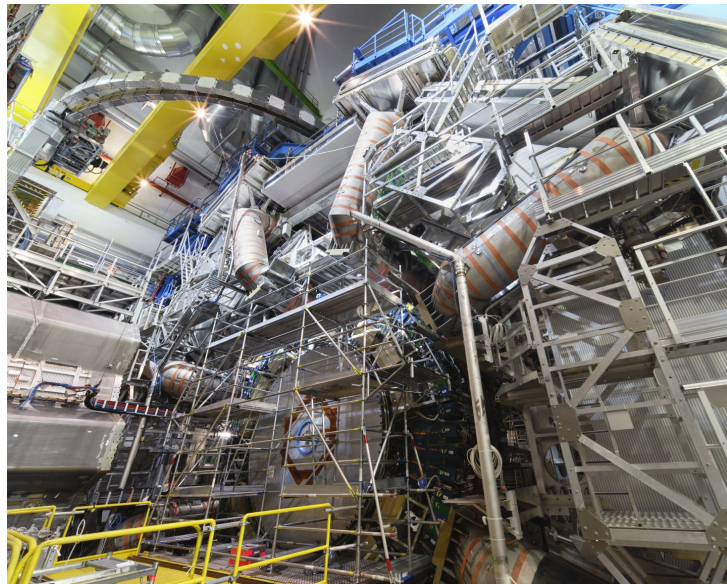
CERN KubeCon 2019

Reproduced discovery of Higgs boson using
Kubernetes, Google Container Engine (GKE)

Containers allowed use of 2010 binaries

70TB Data, 20 000 cores

Setup and run completed “live” during talk!



Hybrid Algorithms

ML for HPC

Google Brain Research for US NOAA

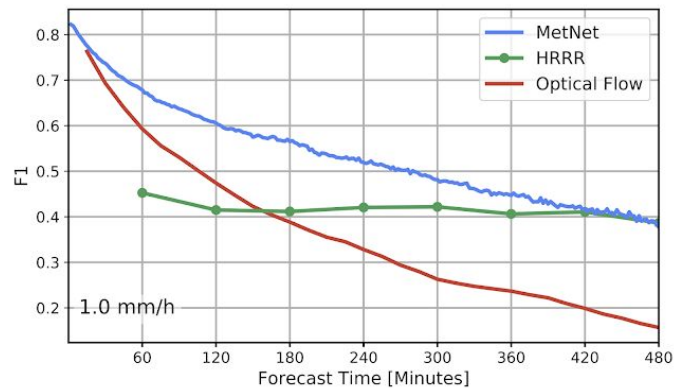
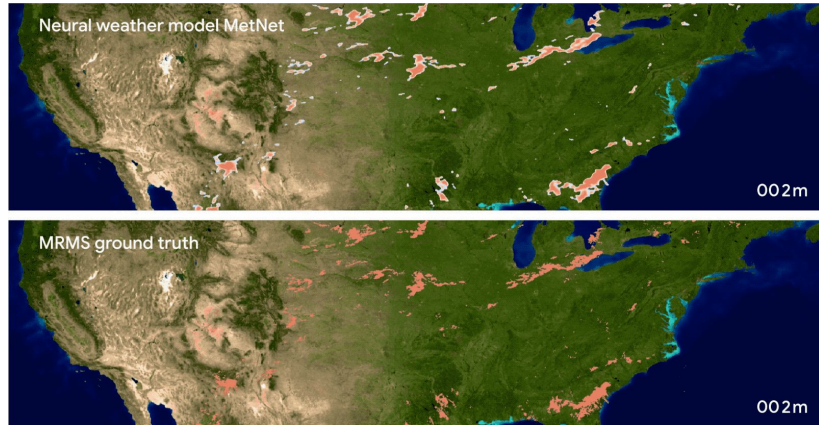
MetNet Neural Weather Model

Runs on 256 Google TPUs

Outperforms current physics based models for speed and accuracy out to 8 days

Parallel scaling allows prediction for entire US in seconds.

NOAA Environmental Data Set now part of freely accessible Google Public Data Sets



<https://arxiv.org/abs/2003.12140>

Hybrid Algorithms

TPUs

Matrix Multiples Dominate ML Computation

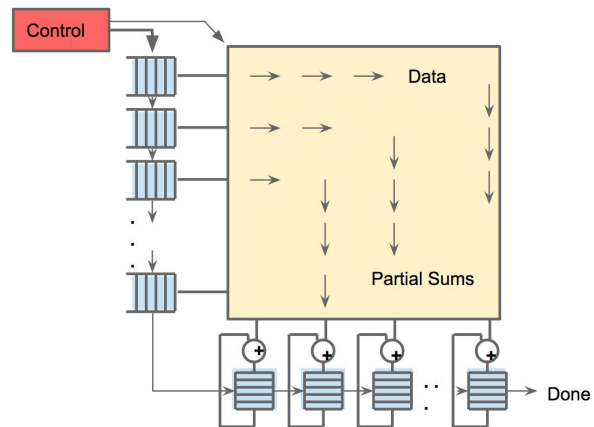
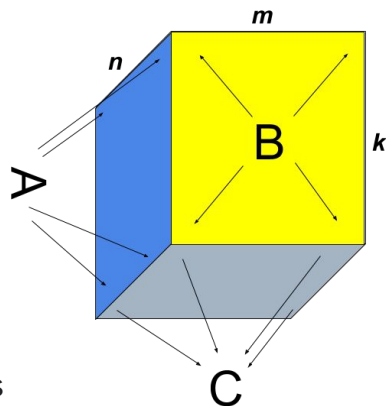
Systolic Array Multiplier Architecture Maximizes Parallelism with Minimal Architectural State

> 100 TOPs per chip

> 100 POPs per pod in a toroidal mesh

Support processors handle communication, perform JIT compilation on XLA from RPCs

Only programmable in TensorFlow and PyTorch

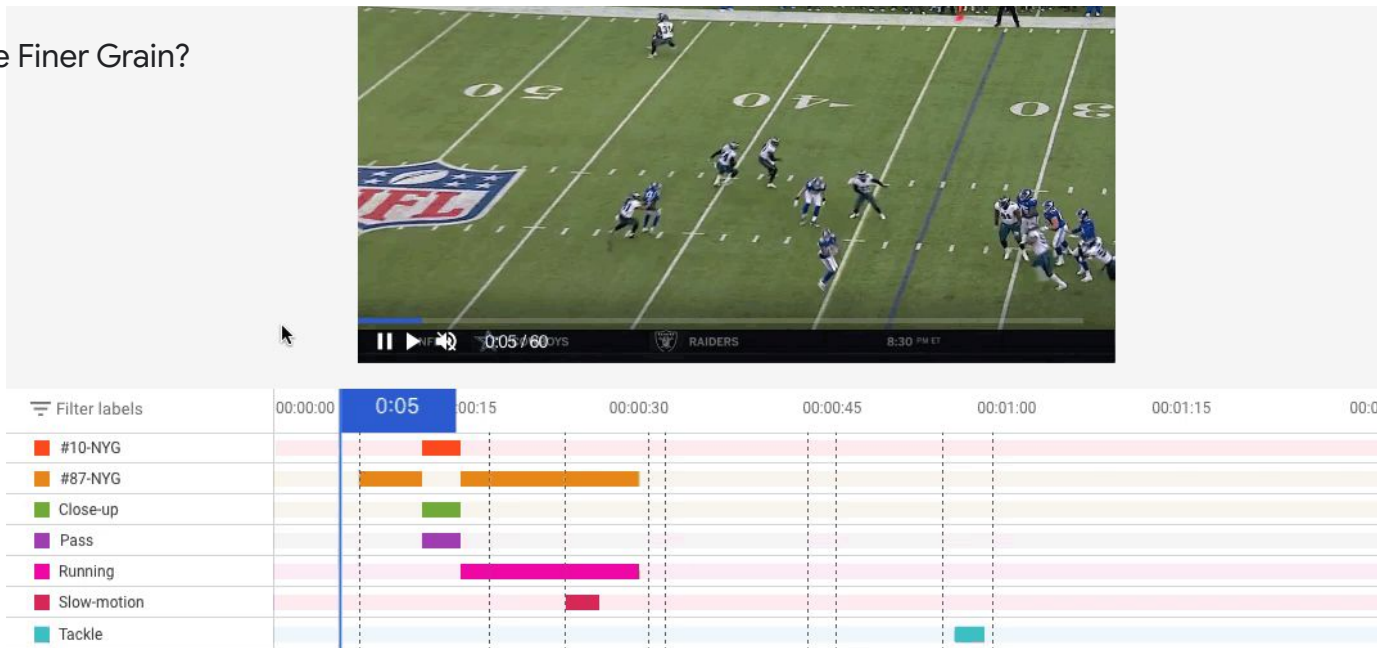


Hybrid Algorithms

Applications or APIs?

Very Loose Coupling Leads to APIs
on Deep Complex Stacks

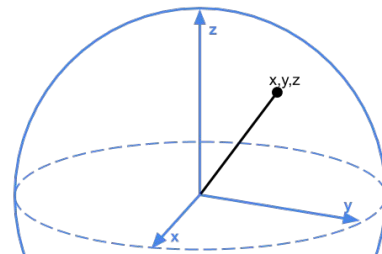
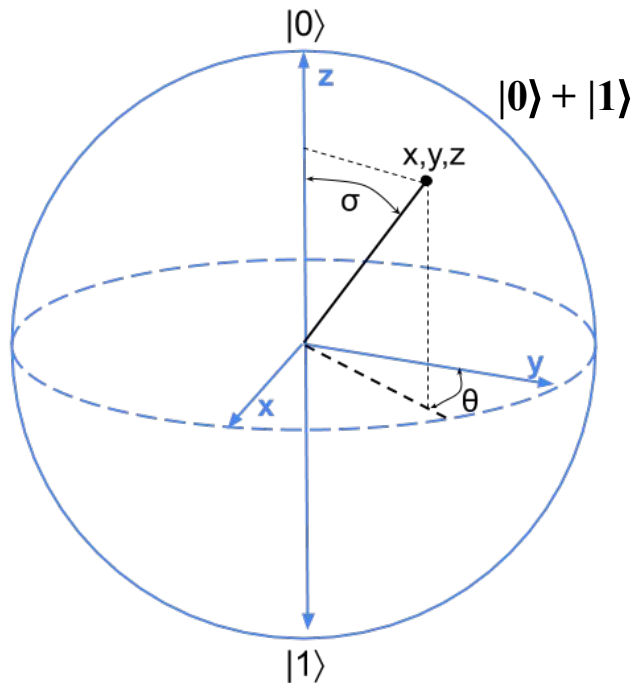
How to enable Finer Grain?



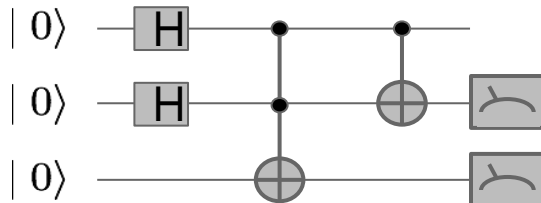
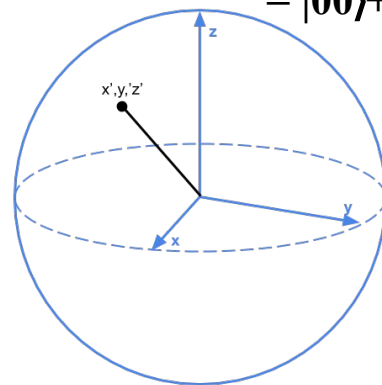
Hybrid Processing Quantum Computing

Not just binary,
Not just probabilistic,
But Turing Complete

“All” one needs to do
is master
Quantum Particles



$$(|0\rangle + |1\rangle)^2 = |00\rangle + |01\rangle + |10\rangle + |11\rangle$$

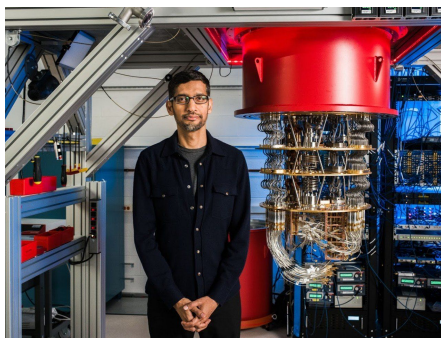
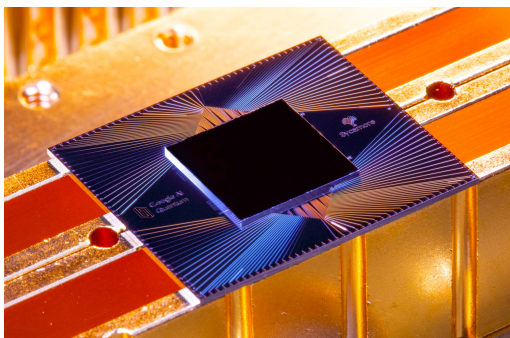


Hybrid Processing

Google Sycamore superconducting qubit platform

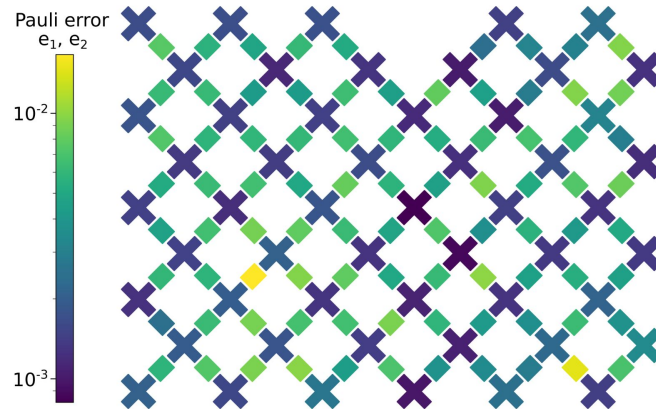
Sycamore platform has 54 planar transmon qubits

tunably coupled in square lattice array



Pauli and measurement errors

Average error	Isolated	Simultaneous
Single-qubit (e_1)	0.15%	0.16%
Two-qubit (e_2)	0.36%	0.62%
Two-qubit, cycle (e_{2c})	0.65%	0.93%
Readout (e_r)	3.1%	3.8%



Hybrid Processing

Variational Quantum Methods

Quantum computing is a disruptive technology.

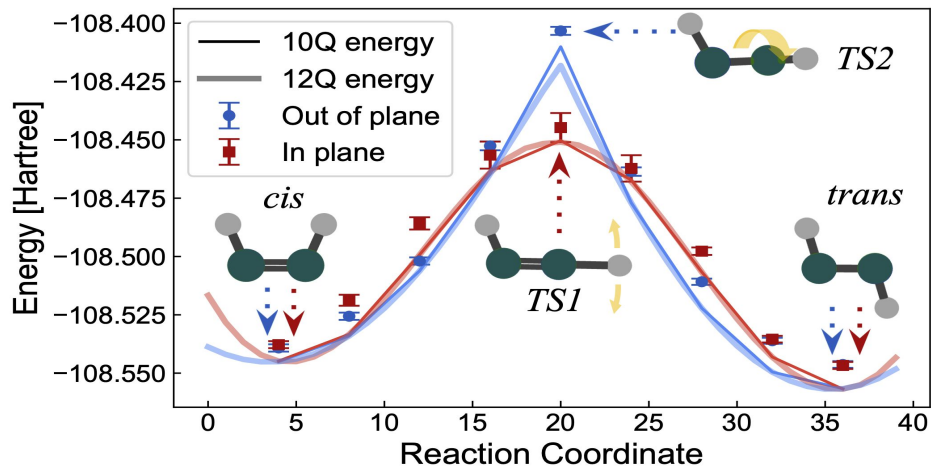
Many applications will require large scale error-corrected machines.

Variational quantum methods are a quantum/classical hybrid for near-term quantum devices.

A parameterised quantum model of a system is iteratively optimised using “classical” computing algorithms.

First demonstrated by Google in 2016 for H_2

With Sycamore quantum processor in 2020, modeled H_2N_2 (Diazene) with sufficient precision to simulate isomerisation.



<https://arxiv.org/abs/2004.04174>



Cirq



OpenFermion

Google Cloud

Copyright 2017-2020 Google LLC

Hybrid Processing

What is a Quantum Operating System?

Cirq = Python + Qubits

TensorFlow = Python + Tensors + ML

TensorFlow Quantum = Python + Qubits + Tensors + Quantum ML

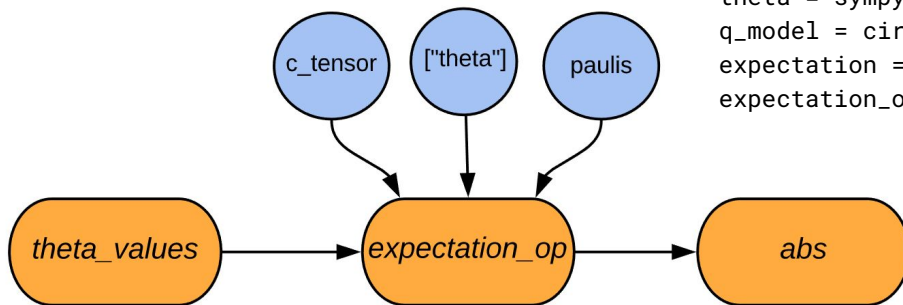
But what's under the Python?

Quantum State cannot be copied, loaded, or stored

Quantum computations are statistical, many runs needed

```
import cirq, random, sympy
import numpy as np
import tensorflow as tf
import tensorflow_quantum as tfq

qubit = cirq.GridQubit(0, 0)
# Quantum data labels
expected_labels = np.array([[1, 0], [0, 1]])
# Random rotation of X and Z axes
angle = np.random.uniform(0, 2 * np.pi)
# Build the quantum data
a = cirq.Circuit(cirq.Ry(angle)(qubit))
b = cirq.Circuit(cirq.Ry(angle + np.pi/2)(qubit))
quantum_data = tfq.convert_to_tensor([a, b])
# Build the quantum model
q_data_input = tf.keras.Input(shape=(),
dtype=tf.dtypes.string)
theta = sympy.Symbol('theta')
q_model = cirq.Circuit(cirq.Ry(theta)(qubit))
expectation = tfq.layers.PQC(q_model, cirq.Z(qubit))
expectation_output = expectation(q_data_input)
```



Hybrid Geography

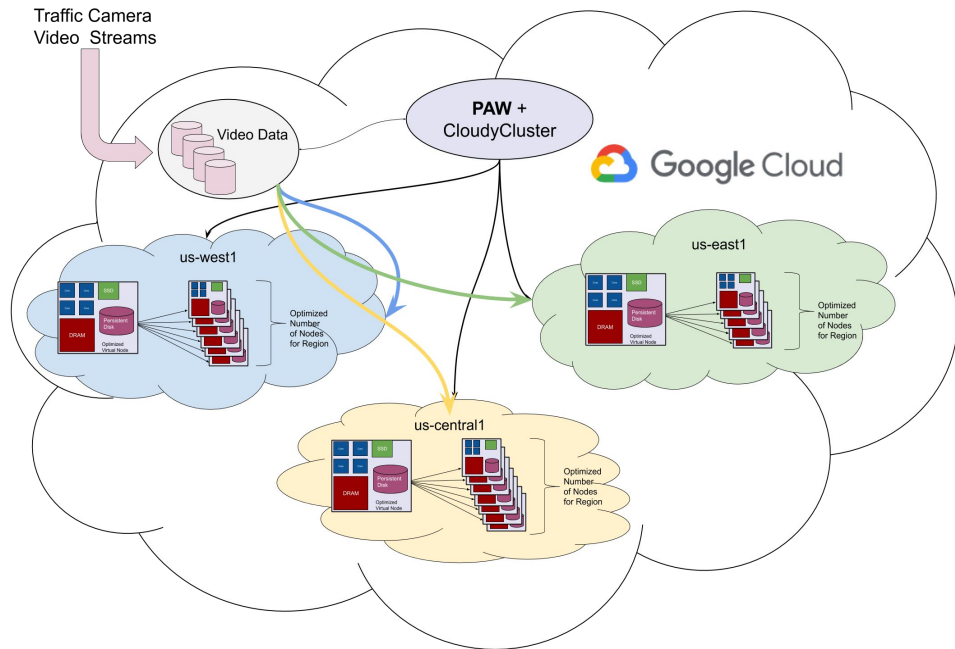
Clemson Experiment

Machine Vision at Scale for Emergency Evacuation Management

2.14 Million Virtual CPUs in 133,573 virtual machines

Orchestrated across six geographical regions, 19 GCP zones

Ephemerally, #5 on Top 500



<https://conferences.computer.org/sc19w/2019/#!/toc/21>

“On-Demand Urgent High Performance Computing Utilizing the Google Cloud Platform”

Hybrid Cloud Services

Caltech IceCube Experiment

IceCube Neutrino Observatory in Antarctica

Calibration requires accurate model of photon propagation

Numerically intensive, massively parallel computation problem

Distributed across 51,000 GPUs

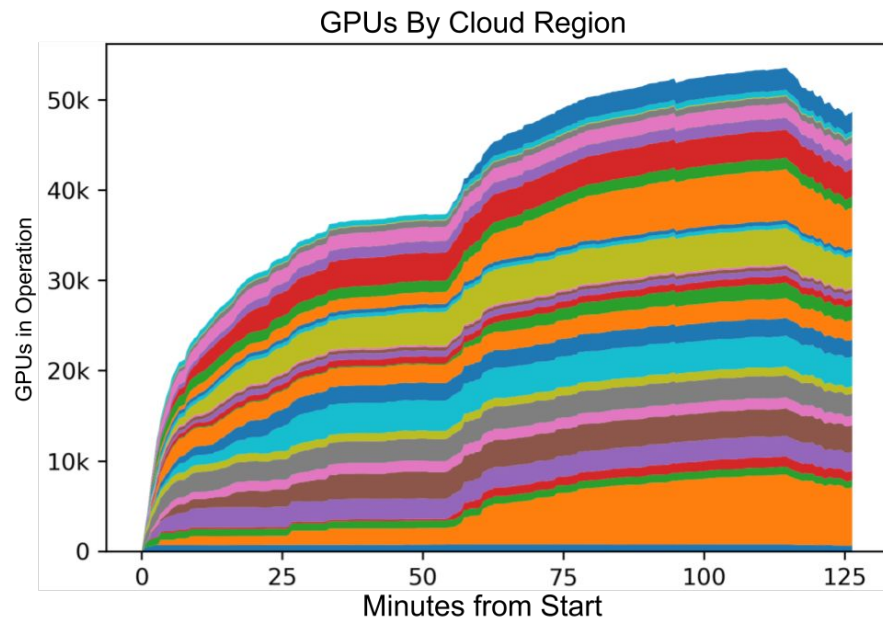
3 Public Cloud Vendors, Including Google

28 cloud regions/zones

380 PFLOP32 nominal peak

90% of Summit - for this sort of problem

HTCondor!



<https://arxiv.org/abs/2002.06667>

What Could Serverless Supercomputing Mean?

Container Orchestration of Microservices
Becoming Dominant Paradigm
For “Enterprise” Computing

Flexible, Open-sourced Abstractions:
Build & Deploy (Knative, Cloud Run) on top of
Service Mesh (Istio) on top of
Container Orchestration (Kubernetes) on top of
Infrastructure

What is the future model for HPC Applications?

