WASHINGTON, DC

Advancing data intensive science at GWU

Clark Gaylord Director, Research Technology Services George Washington University cgaylord@gwu.edu

Thanks to: Glen MacLachlan, Adam Wong, Jarda Flidr, Terrence Lewis

WASHINGTON, DC

Research Technology Services Portfolio

- TRADITIONAL HPC CLUSTER: PEGASUS 210 NODE, HYBRID CPU/GPU, SLURM, IB, ETC
- TEACHING HPC CLUSTER: CERBERUS 12 NODE, FOCUS ON USE IN CLASSROOM
- HTCONDOR: "PILOT" 8 NODE DEPLOYMENT
 - Entre to, working with Open Science Grid
- CAPITAL AREA RESEARCH AND EDUCATION NETWORK (CAAREN) REGIONAL INTERNET2
- CLOUD SERVICES: MOSTLY ENGINEERED, PURPOSE-BUILT SOLUTIONS FOR SPECIAL USE CASES
- MANY BOUTIQUE SOLUTIONS FOR RESEARCH GROUP CLUSTERS, DATABASES, VISUALIZATION
- DATA MANAGEMENT, PROTECTED DATA ENVIRONMENTS
 - REDCAP
 - GLOBUS WITH GOOGLE DRIVE, BOX, S3 CONNECTORS
- Solution consultancy

A good fit for HTC

WASHINGTON, DC

- Pegasus is heavily utilized often over 85% allocated.
- LIKE MANY CLUSTERS, MOST JOBS HAVE RELATIVELY MODEST RESOURCES
 - More than 91% of our jobs finish in under 1 hour.
 - More than 97% of our CPU-based jobs run on a single node.
 - Approx 1/2 of a recent sample of our jobs require less than 10 GB of RAM.
 - [EVERYTHING IS PARETO DISTRIBUTED]
- This profile of JOBS FIT NICELY INTO THE HTC MOLD.

Current configuration

WASHINGTON, DC

GWU HTCONDOR IS AN ON-PREM 9-NODE DELL CLUSTER (POWER EDGE R730)

- 1 Physical manager node
- 8 Physical worker nodes
- Each node features:
 - 28 INTEL XEON E5-2680 CORES
 - 128 GB MEMORY
 - 150 GB SCRATCH SSD STORAGE

Current growth

WASHINGTON, DC

Addition of P100/V100 GPUs repurposed from Pegasus login nodes, etc. Augmenting the original HTCondor deployment for higher storage deployment.

- Additional VMs highly scalable deployment (begin with 8)
 - OPENSTACK HOSTS GIVES FLEXIBILITY FOR OS "PSEUDO-BARE METAL"
- Storage
 - 1.2PB GLUSTER OVER ZFS
 - 3.2GB/s aggregate I/O
- IPv6 connectivity through CAAREN (Equinix facility)
- 4 additional GPUs

Pegasus HPC



WASHINGTON, DC

[insert Pegasus photo here]



WASHINGTON, DC

A Science Driver from Genomics

- Prof. Anelia Horvath from Milken Institute School of Public Health at GW Functional Analysis of SNVs (Single Nucleotide Variants) from single-cell RNA sequencing data.
- Single cell RNA sequencing (cell-level transcriptome studies) in cancer genomics.
 - Aim: develop methods to assess functionality of SNVs through their relationships with dynamic transcriptome traits from cancer scRNA-seq data.
 - Data source: generated by the 10x Genomics platform a high-throughput single-cell sequencing platform.
 - Datasets: consisted of $\sim 9K$ cells per individual, and $\sim 150K$ sequencing reads per cell, and ~ 150 nt/read; which produces a single file of $\sim 1-1.5$ billion sequencing reads.
 - Software: R/Python/cellRanger(from 10x Genomics).

WASHINGTON, DC

A Science Driver from Genomics

- CELLRANGER IS A SET OF ANALYSIS PIPELINES THAT PROCESS CHROMIUM SINGLE-CELL RNA-SEQ OUTPUT TO ALIGN READS, GENERATE FEATURE-BARCODE MATRICES AND PERFORM CLUSTERING AND GENE EXPRESSION ANALYSIS.
- The pipelines of tools are:
 - CELLRANGER MKFASTQ CONVERT RAW DATA FASTQ FORMAT DATA
 - CellRanger Count Perform Tasks such as Alignments, filtering and data clustering
 - \circ cellRanger aggr aggregate results from multiple runs of cellRanger count
 - CellRanger Reanalyze Perform Secondary Analysis by Fine-Tuning Parameters

WASHINGTON, DC

A Science Driver from Genomics

- The problem size (~1.5 billion sequencing reads) of this genomic data analysis demands for a great amount of computational resources.
- The data processing nature of such computational task makes it easily be broken into numerous independent prices (i.e. HTC-able) and thus it is a good candidate for a HTCondor cluster!
- A CELLRANGER WORKFLOW WOULD FIT WELL TO HTCONDOR'S WORKFLOW MODEL WHERE ITS DAGMAN (DIRECTED ACYCLIC GRAPH MANAGER) MANAGES DEPENDENCIES BETWEEN JOBS AUTOMATICALLY.
- These CellRanager workflows of computation allow us to experience and learn about the advanced feature of HTCondor in workflow construction.

WASHINGTON, DC

A Science Driver from Physics

- PROF. WILLIAM J. BRISCOE FROM THE PHYSICS DEPARTMENT AT GW IS WORKING ON AN INTERMEDIATE ENERGY NUCLEAR PHYSICS EXPERIMENT AT THE THOMAS JEFFERSON NATIONAL ACCELERATOR FACILITY (JLAB) USING THE CEBAF LARGE ACCEPTANCE SPECTROMETER (CLAS).
- JLAB CAN PROVIDE EITHER CONTINUOUS ELECTRON AND PHOTON BEAMS WITH ENERGIES UP TO 12 GIGA ELECTRON VOLTS (12 GeV) TO STUDY PROPERTIES OF QUARKS IN HADRONS AND MESONS.
- EXPERIMENTS ARE EXTREMELY EXPENSIVE AND TIME-CONSUMING FOR HUNDREDS OF PHYSICISTS. CAREFUL SIMULATION ARE PARAMOUNT.

WASHINGTON, DC

- RESEARCHER SIMULATES PART OF AN EXPERIMENT USING GEANT4 MONTE CARLO (GEMC)
- GEMC IS A C++ APPLICATION USED BY PHYSICISTS TO MODEL THE INTERACTION OF CHARGED PARTICLES PASSING THROUGH MATTER SUCH AS INCIDENT, SCATTERED, AND RECOIL PARTICLES PASSING THROUGH TARGETS AND DETECTORS.

A Science Driver from Physics



WASHINGTON, DC

A Science Driver from Physics

- CLAS delivers GEMC and related data as both Docker and Singularity images via CVMFS.
- The simulation needed to be completed in 2-3 months and required approx 4.3 million core hours to complete which is a little more than GW could provide locally or what OSG could provide opportunistically.
- OSG ALLOCATED ADDITIONAL RESOURCES AND SIMULATION COMPLETED ON SCHEDULE.

WASHINGTON, DC

Omics data, Medicine and Public Health

PROF. ERIC VILAIN FROM THE CHILDREN'S NATIONAL HOSPITAL (DC) WITH THE INSTITUT NATIONAL POUR LA RECHERCHE BIOMÉDICAL OF THE DEMOCRATIC REPUBLIC OF THE CONGO.

The objective is to ascertain microbial role in the Konzo disease and then to understand microbiota in settings such us different towns, cities and regions in Congo. Consequently, actionable genomics experiments could be performed to advance precision medicine in DRC

WASHINGTON, DC

- GWU/ CHILDREN'S NATIONAL HOSPITAL
 - JONATHAN LOTEMPIO, BS
 - D'ANDRE SPENCER, BA, MPH
 - NEERJA VAHIST, BS
 - MATT BRAMBLE, PHD
 - ERIC VILAIN, MD, PHD
- INSTITUT NATIONAL POUR LA RECHERCHE BIOMÉDICAL
 - KIZITO MOSEMA, MD, MPH
 - KEVIN KARUME, BS
 - Désiré Tshala Katumbay, MD, PhD
 - JEAN-JACQUES MUYEMBE, PHD
 - DIEUDONNE MUMBA, MD, PHD

United States – DR Congo Collaboration





WASHINGTON, DC

Vast amount of omics data

To increase representation of Congolese individuals, about 100-1000 high quality references genomes will be generated via the advanced sequencing technologies including:

- HIGH COVERAGE ILLUMINA SEQUENCING
- HIGH COVERAGE PACBIO HIFI SEQUENCING
- BIONANO OPTICAL MAPPING

WASHINGTON, DC

Huge computation demand

- Mapping of sequences need to be done on 25 TB of data by end of 2021 and another 100 TB of data by the following 1-2 years.
- Common alignment tools such as Samtools and Minimap2 will be used.
- These can be carried out efficiently in a HTCONDOR cluster.

WASHINGTON, DC

Serving the research community

- Working with HTCondor community (and OSG) leverages collaboration to bring resources to the table that are outside scope for HPC
- Continuing with HTCondor deployment to contribute to and partner with OSG
- FLAGSHIP HPC CLUSTER CONTINUES TO BE STRATEGIC FOR MANY RESEARCH APPLICATIONS
 NOT ONLY MPI...
- Much of GWU research is naturally high throughput which often *can* adapt to HPC but this isn't necessarily ideal

WASHINGTON, DC



We really appreciate the effort from the HTCONDOR Team in supporting our HTCONDOR project at GW. A special thanks goes to John M Knoeller (TJ) for meeting regularly with us and in helping us to solve any Condor related problems.