# Enabling Economical Genome Analyses through Optimization and Scalable Workflows

ParslFest 2024
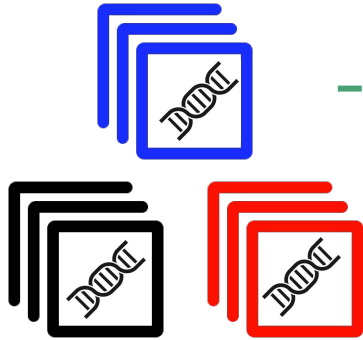
Akila Perera, Jason J. Pitt

Cancer Science Institute of Singapore

# Non-linear Multithreading of Bioinformatics Tools

# SWAG: Scalable Workflows for Analyzing Genomes

- Transparent parallelism

- Portability

- Robust to failures

- Scalability

# SWAG - GDC Workflow



- Genomic raw fastq data split into multiple chunks for parallelization

- BAM files split by region for parallel variant calling

- Compute considerations due to chromosome size variability

  - Dynamic resource allocation

  - Multiple executors with different walltimes supported by underlying queues

# SWAG - Optimizing Workflow Performance



**General Executor**

- bamtofastq
- co_clean_indel_realign
- co_clean_apply_bqsr
- split_fastq

- Strelka2
- Mutect2
- make_readgroups
- split_contigs
- validate_contigs
- merge_vcf

normal queue on ASPIRE2A
(ncpus=128, mem=440GB)

1-4 nodes, 96 workers, ~1cpu each

**Compute-optimized Executor**

- bwa_sort
- bwa_align
- mark_duplicates
- merge_bam_files
- co_clean_realign_target

normal queue on ASPIRE2A
(ncpus=128, mem=440GB)

1-2 nodes, 32 workers, 4cpus each

**Long Tasks Executor 1**
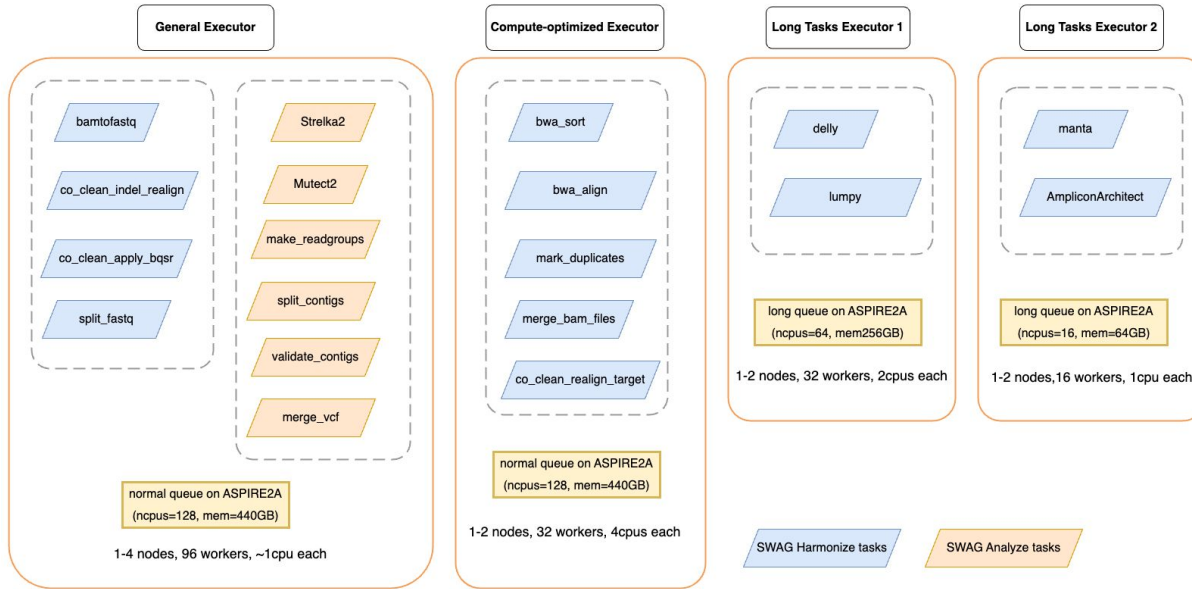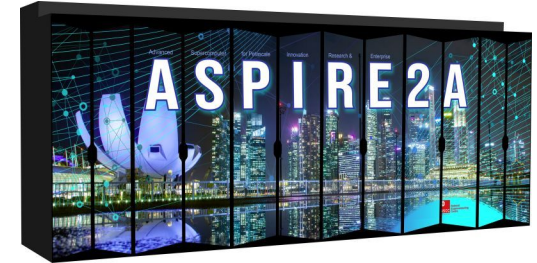
- delly
- lumpy

long queue on ASPIRE2A
(ncpus=64, mem256GB)

1-2 nodes, 32 workers, 2cpus each

**Long Tasks Executor 2**

- manta
- AmpliconArchitect

long queue on ASPIRE2A
(ncpus=16, mem=64GB)

1-2 nodes,16 workers, 1cpu each

SWAG Harmonize tasks     SWAG Analyze tasks

ASIPRE2A – Singapore's National
Petascale Supercomputer
(~800 nodes & ~100,000 cores)

- Canned configurations for Parsl executors for each workflow

- Easily to switch between various HPC/compute infrastructure

- Task profiling to understand resource utilization patterns

- Dynamic assignment of tasks to executors based on scalability profile

# Key Results to Date

Related manuscripts

- Processed 250+ whole genome sequencing (WGS) of isogenic cell lines and mouse tissue samples
  - ~164,000 cpu core hours on ASPIRE1 (NSCC)
  - ~50 tb of raw data

- Processed 12,000+ whole exome sequencing (WXS) samples from TCGA
  - ~2M cpu core hours on ASPIRE1 (NSCC)
  - ~800 tb of raw data

- Processed 250+ deeply sequenced WGS samples from TCGA
  - ~200,000 cpu core hours on ASPIRE2A (NSCC)
  - ~250 tb of raw data
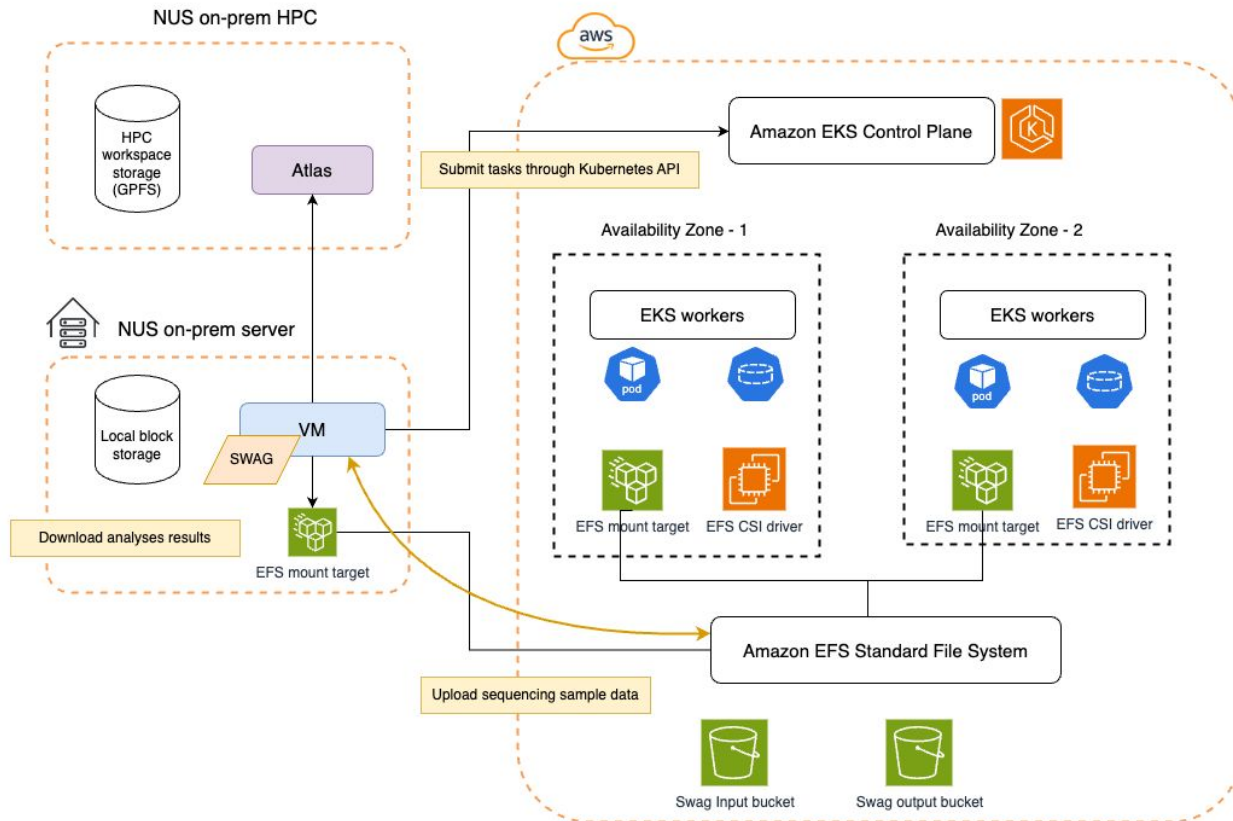
Kong LR… Pitt JJ, Venkitaraman AR. *Cell* 2024

Perera A... Pitt JJ. *Communications Biology* 2024

Wu A… Pitt JJ. *Briefings in Bioinformatics* 2023

Wong H… Pitt JJ. *In preparation*

Wu A… Pitt JJ. *In preparation*

# SWAG Roadmap - Multi-site Genomic Workflows

# Improving Pipeline Efficiency by Advancing Parsl Capabilities

- Task failures due to biological variability
  - Chromosomes vary widely in length
  - WGS sample data size varies from 100Gb - 300Gb
  - Variability in runtime causes task timeouts
- Inefficient distribution of tasks
  - Whole node allocation requirement for consistent performance in HPC
  - Resource waste due to idling cpu cores
  - Guaranteed failure when allocated to an old node

- Dynamic assignment of tasks
  - Duration-sensitive task allocation
  - Feedback loop through a user provided function on retries
- Flexible task allocation to reduce costs
  - Cost-effective task packing strategy
  - Node-aware scheduling