



Efficient and Scalable Workflows for Genomic Analyses

Subho S. Banerjee, Arjun P. Athreya, Liudmila S. Mainzer, C. Victor Jongeneel, Wen-Mei Hwu, Zbigniew T. Kalbarczyk, Ravishankar K. Iyer

A decorative background for the bottom section of the slide, featuring a complex network diagram. The diagram consists of numerous light blue circular nodes of varying sizes, interconnected by thin, light blue lines. The nodes are scattered across the width of the slide, with some clusters and some isolated nodes. The overall effect is a dense, interconnected web of data points.

CSL.ILLINOIS.EDU



Summary

Contributions:

- **Common mathematical kernels:** Static analysis of genomic analyses algorithms
- **Performance Pathologies:** Measurement driven diagnosis of performance bottlenecks
- **IGen:** A scalable genomic data analytics framework which overcomes observed inefficiencies

“Variant Calling and Genotyping” Workflow as the driver

	Baseline Runtime	IGen Accelerated Runtime **	Speedup
Blue Waters – Single Node (CPU)	59 hr	28 hr	2.1x
IBM Power 8 – Single Node (CPU + GPU + FPGA)	36 hr	11 hr	5.3x, 3.2x
Blue Waters – 10 Nodes (CPU)	-	2 hr	22x



Outline

- Genomics Primer: Variant Detection
- Kernels for Genomics
- Performance Pathologies in State of the Art Genomics Pipelines
- IGen: The **I**llinois **G**enomics Execution **E**nvironment



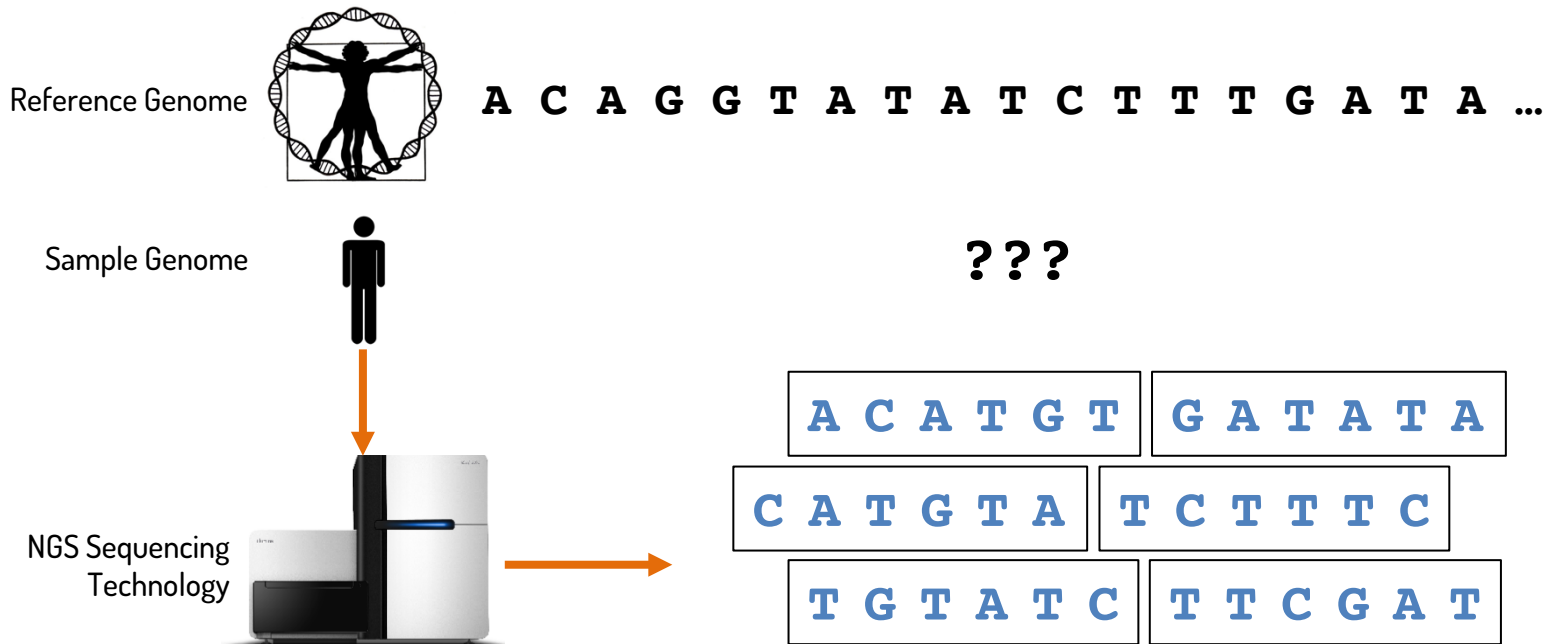
Overview: Variant Calling and Genotyping

Detecting and characterizing mutations in a sample genome



Overview: Variant Calling and Genotyping

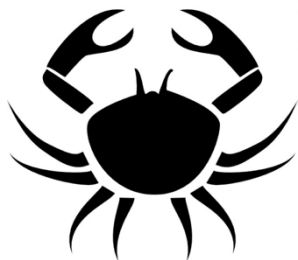
Detecting and characterizing mutations in a sample genome



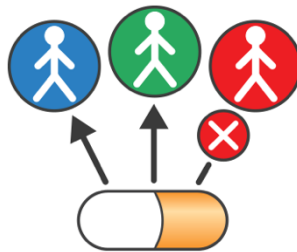


Overview: Variant Calling and Genotyping

Detecting and characterizing mutations in a sample genome



Diagnosis
e.g., Cancer



Personalized Medicine

Reference G

Sample G

~3GB

T A ...

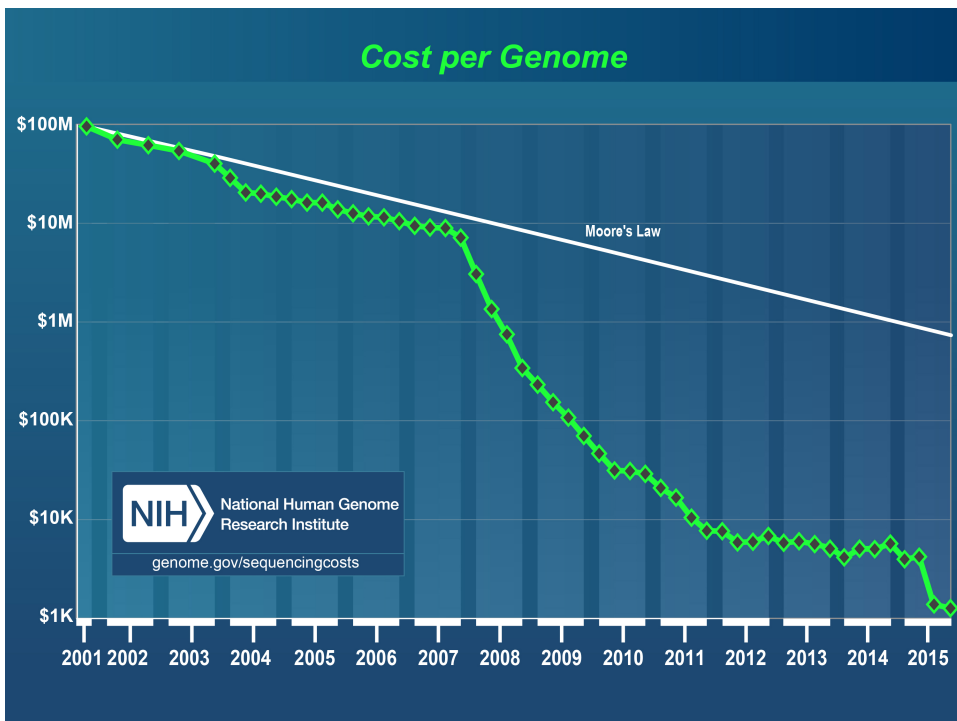
A T A

~200GB (compressed)

Noisy Data Polyploid Samples



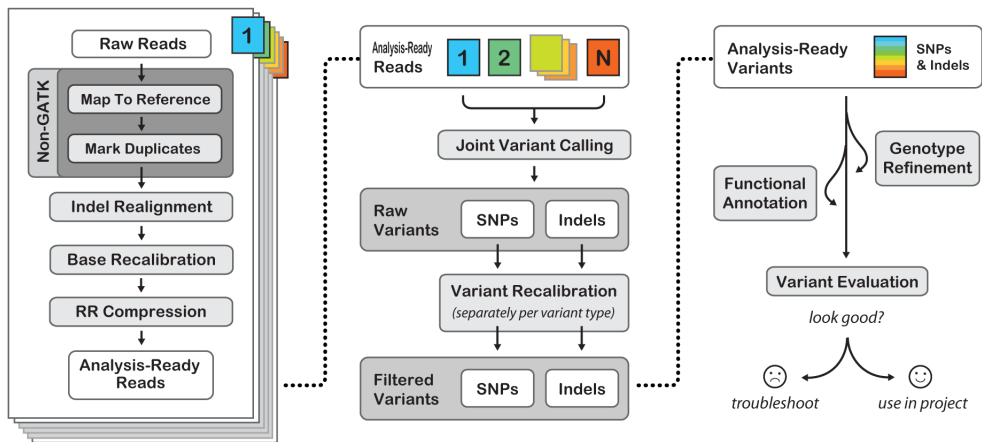
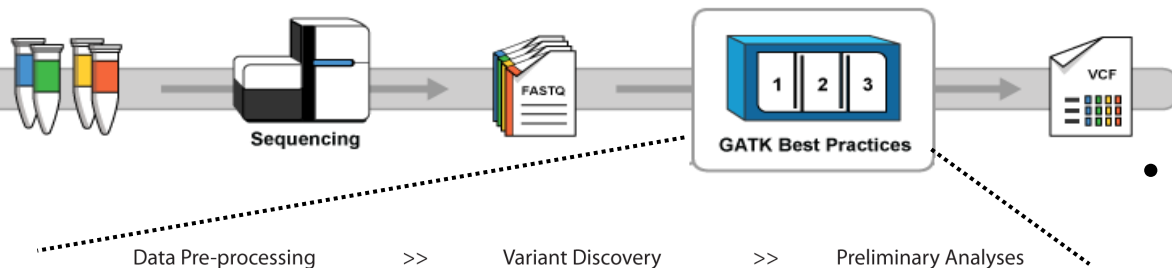
The Genomic Data Deluge



- Falling costs
 - **Capex**: Cost of buying sequencing machines
 - **Opex**: Cost of sequencing genomes
- Potential for large amounts of sequence data to be generated over a short span of time
- Societally important problem
 - Scope for personalized medicine changing healthcare delivery

<http://www.genome.gov/sequencingcosts>

Variant Discovery as a Workflow



- The Broad Institute Best Practices Guidelines

- Tools come from disparate sources
 - Designed for workstations
 - **Few** are performance tuned
 - Do not fit well in traditional HPC



Outline

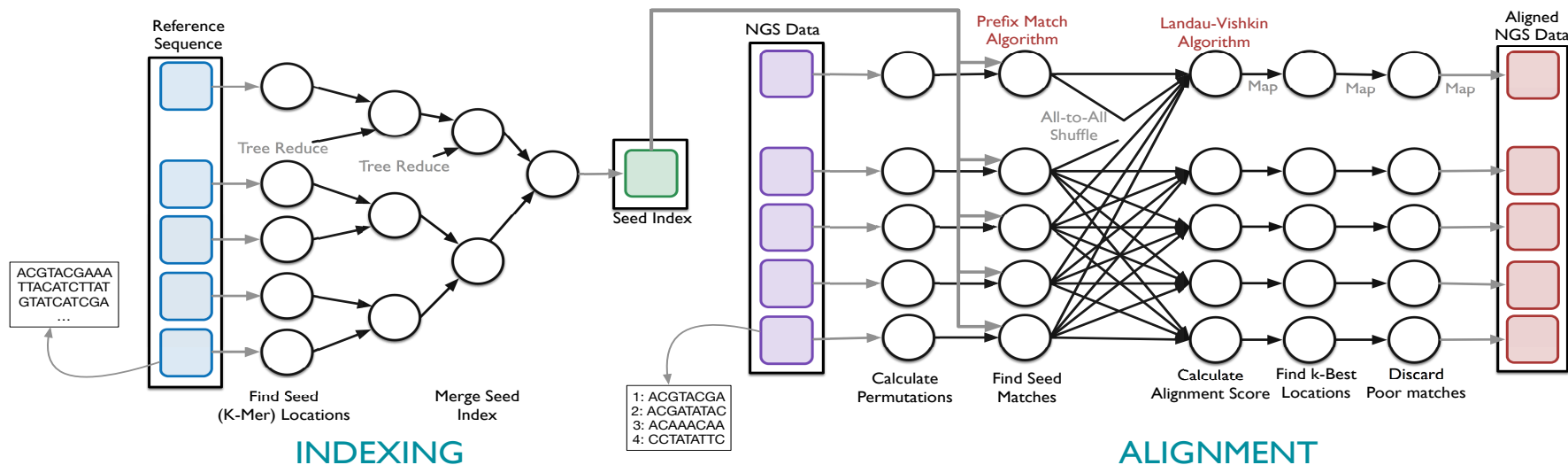
- Genomics Primer: Variant Detection
- Kernels for Genomics
- Performance Pathologies in State of the Art Genomics Pipelines
- IGen: The **I**llinois **G**enomics Execution **E**nvironment



Computational Kernels

- **Computational Kernels:** Basic Mathematical Operations common to large number of bioinformatics analysis
- Kernels enable system level optimizations effecting a large number of tools
- Clearly show commonalities between different tools performing the same analysis
- Provide an interface between algorithm designers and system designers
 - Future benchmarks for data-intensive HPC machines
- Defines a simple data-flow abstraction for non-expert programmers (biologists)

Kernels



Single Ended NGS Read Alignment as a DFG



Kernels

Repeated kernel usage across tools/stages

Workflow stage	Kernels
Error Correction	K-mer computation
Alignment	K-mer computation, Prefix Tree, Edit-distance computation
Indel Re-Alignment	Edit-distance computation
Re-Calibration	Yates correction
Variant Calling	Entropy, Convolution, Assembly, Edit-distance, Pair-HMM, Bayesian inference

See paper for common kernels across multi-sequence alignment, metagenomics and phylogeny

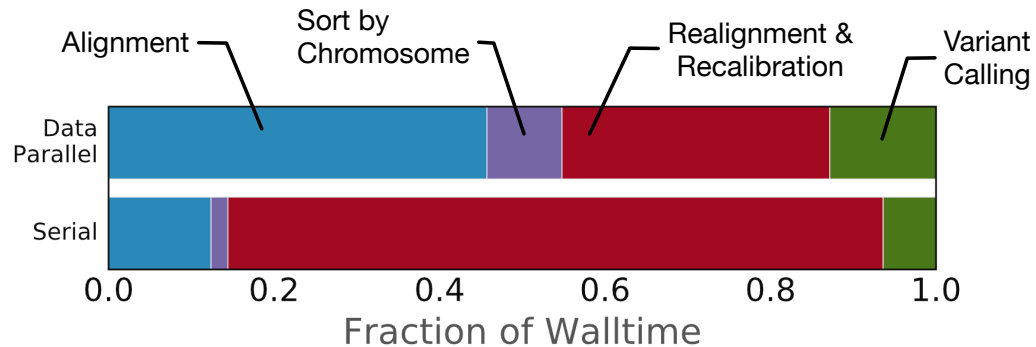
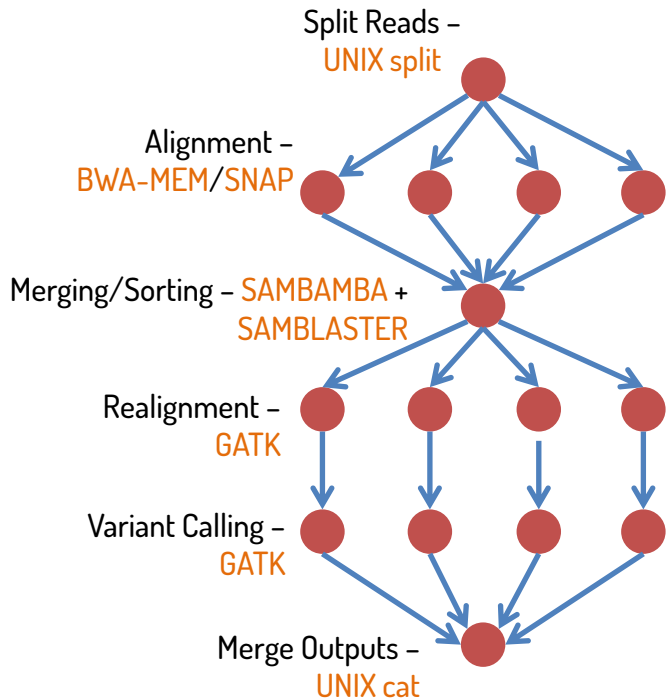


Outline

- Genomics Primer: Variant Detection
- Kernels for Genomics
- Performance Pathologies in State of the Art Genomics Pipelines
- IGen: The **I**llinois **G**enomics Execution **E**nvironment

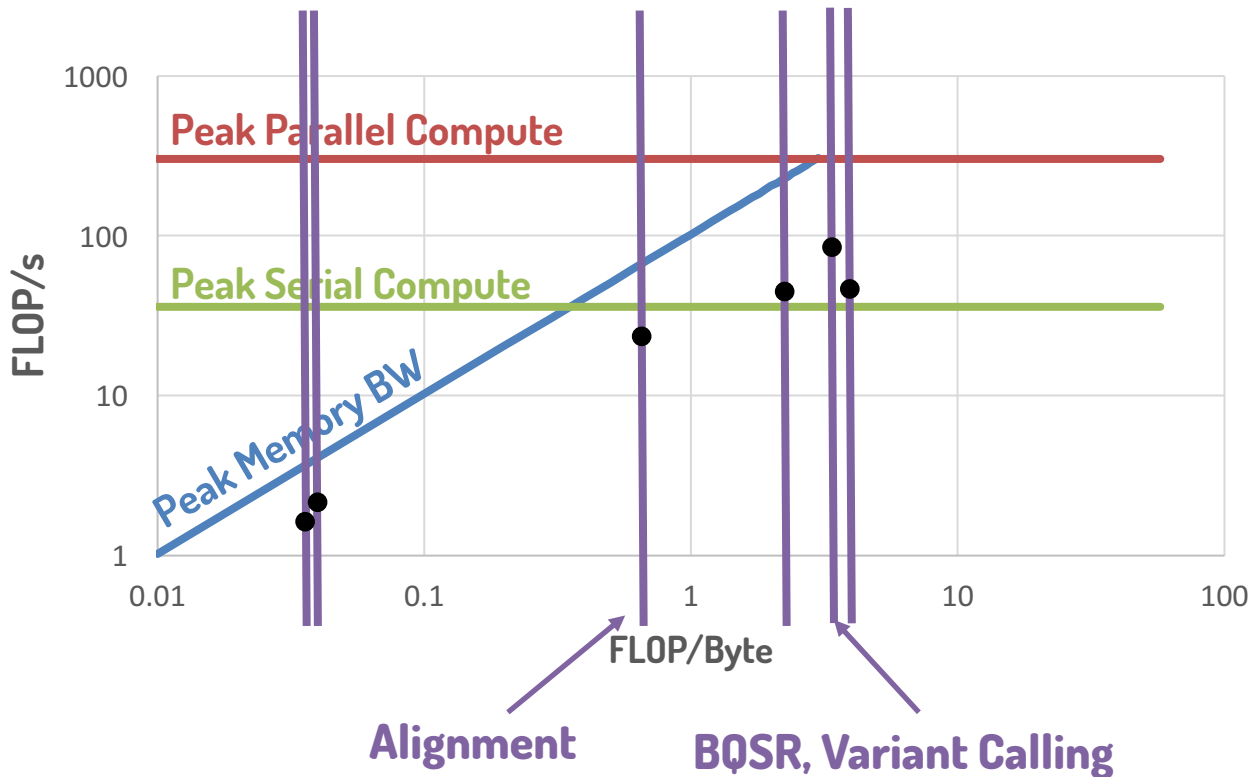


Deploying genomics workflows on parallel systems



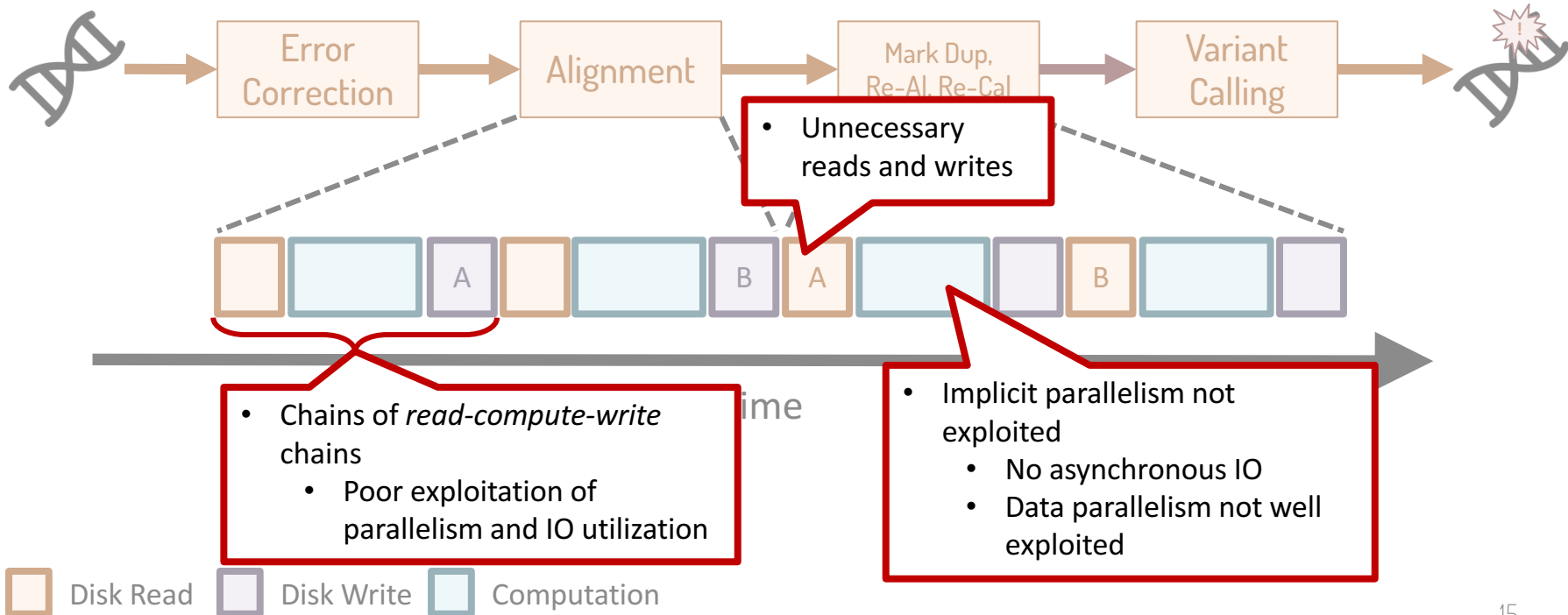


Tools are not well suited for HPC machines





Understanding Performance Issues

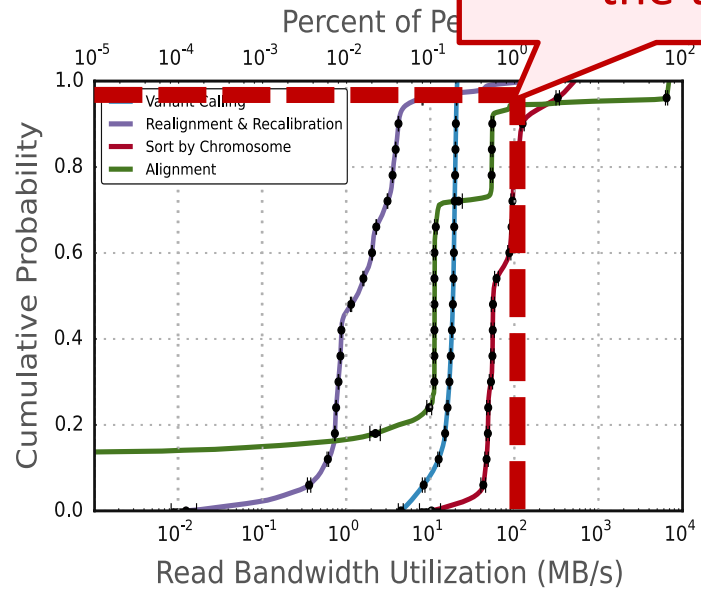
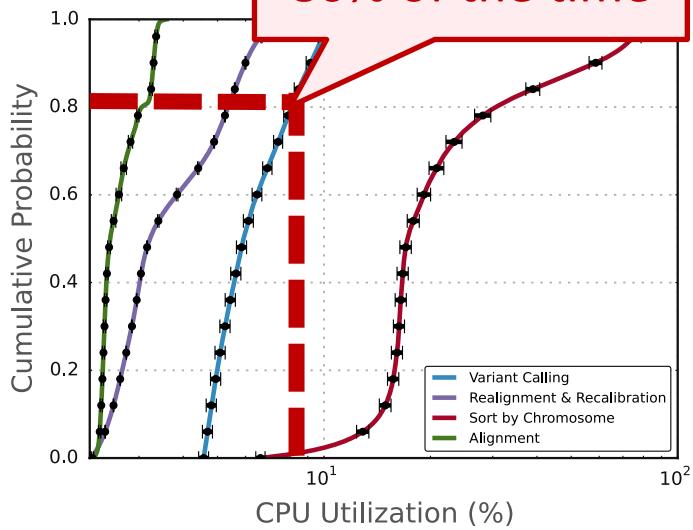


Understanding Performance

Measurement on Blue-Waters

Less than 10% CPU Utilization 80% of the time

Less than 1% IO utilization 95% of the time



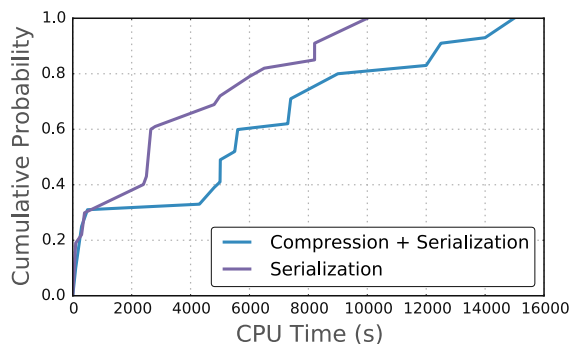
System resource utilization for phases of the Broad Institute Best Practices Guide Workflow



Variant Calling on Spark

ADAM and Avocado

- Best performance – Best place to start!
- ADAM: Extremely efficient data formats for parallel compute



Time Spent in Serialization for ADAM based file formats

- Several Problems
 - Serialization takes a lot of time
 - Easy to program \neq Good performance
 - Single Node performance quite poor, Great Scalability
 - Non-trivial (12.3 %) amount of time spent in fault-tolerance related computation/messaging
 - JVM - Garbage collection



Outline

- Genomics Primer: Variant Detection
- Kernels for Genomics
- Performance Pathologies in State of the Art Genomics Pipelines
- IGen: The **I**llinois **G**enomics Execution **E**nvironment

Sequences to Systems

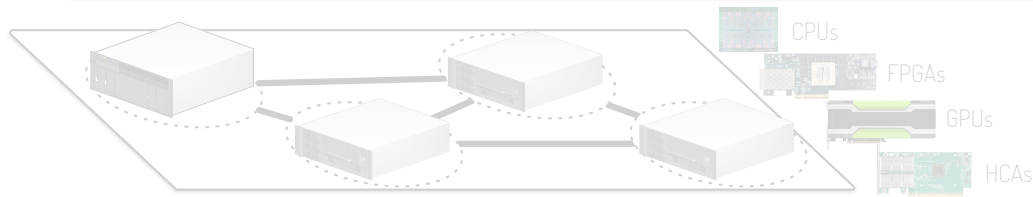
Key Idea: Decouple **algorithms**, **schedule** and **accelerators**

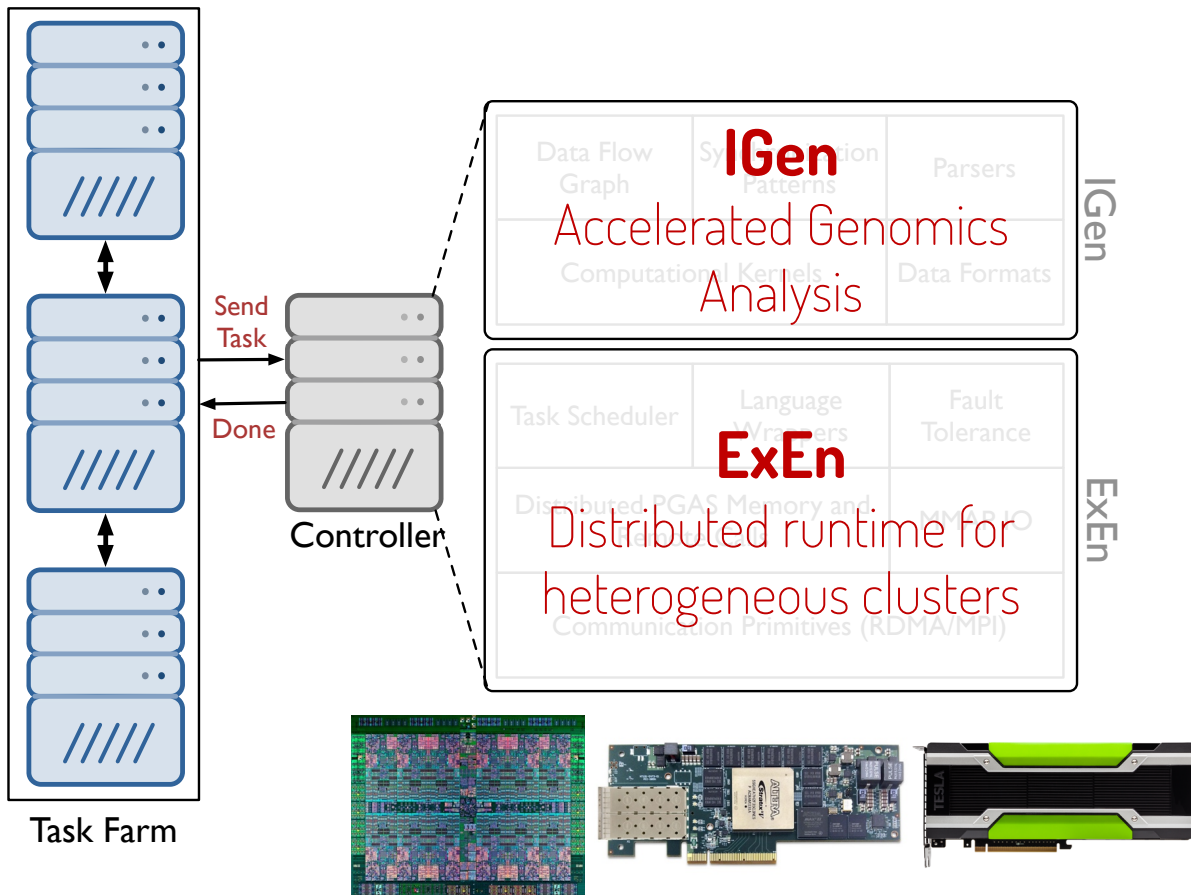
- Algorithm: What is computed
- Schedule: Where and when it is computed
- Accelerators: How it is computed

• Distributed Scheduling

Hardware Layer

- Massively parallel processors
- Specialized hardware

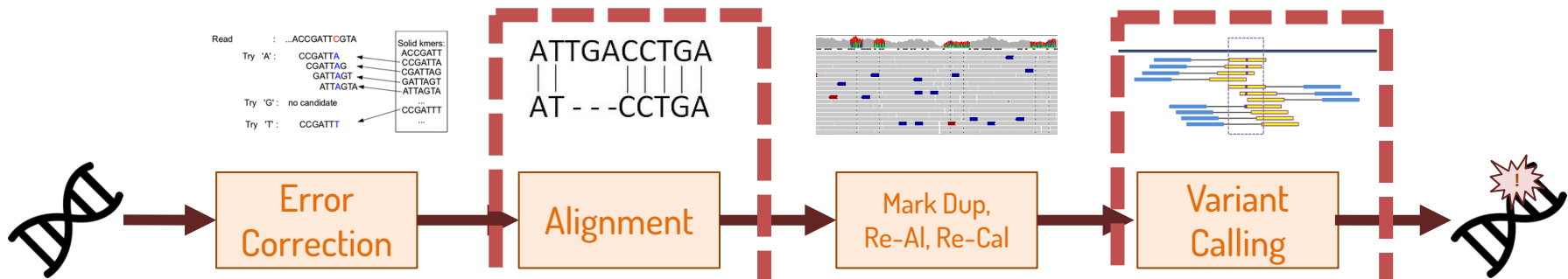






Variant Calling and Genotyping in IGen

Detecting and characterizing mutations in a sample genome



Blue Waters - 22 Nodes (CPU)	-	14.3 hr	3.8	2.3	0.2	4.8	3.1 hr
Blue Waters - 10 Nodes (CPU)	-	48 min	63 min*			12 min	

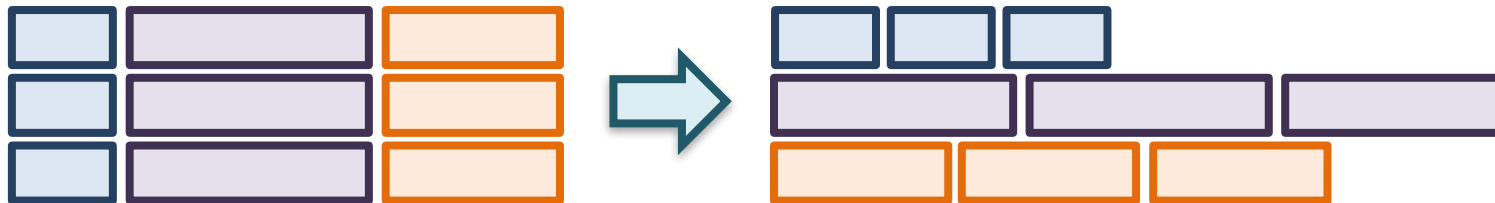
Broad Inst. Best Practices
IGen Accelerated

* Whole Human Genome @ 60x coverage
 ** Default tool parameters



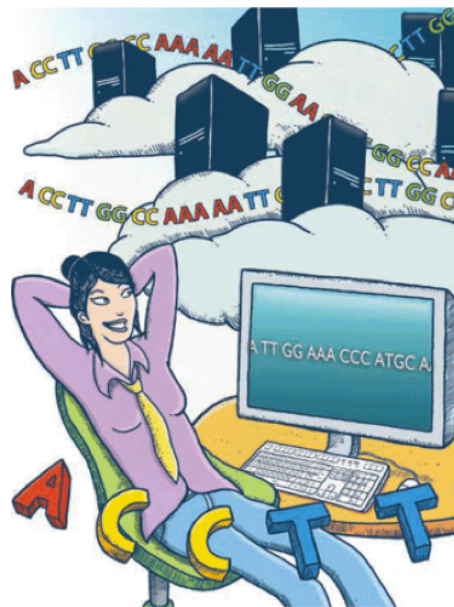
Enhancements in IGen

- Succinct data representations
 - All tools use ASCII based in-memory representations
 - Use 2 and 4 bit representation for Nucleotides/CIGAR
- Asynchronous File IO
- Column based data-structures to improve locality and aid vectorization



- Compiler assisted and SIMD intrinsic based implementations of kernels

Conclusions



- Bringing computer systems and analytics to precision medicine
 - ExEn and IGen for accelerated NGS analysis
 - NEAT and AssembleSV for quality control of NGS pipelines
 - Statistical analysis for deriving actionable intelligence



CompGen Machine

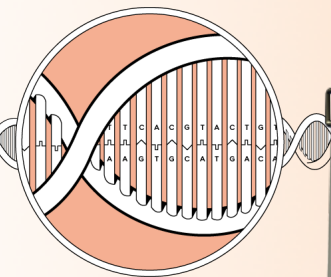
Medical Knowledge



Continuous Monitoring



*omics Data



Timely Diagnosis

Personalized Drugs

Model Drug Response

New Biological Insight



Medical Devices



Patient Records