# High Performance Computing in Life Sciences

## Part I
### HPC Introduction

## PartII
### BioComputing Sofware Introduction

Oleksandr Moskalenko
om@ufl.edu

Matt Gitzendanner
magitz@ufl.edu
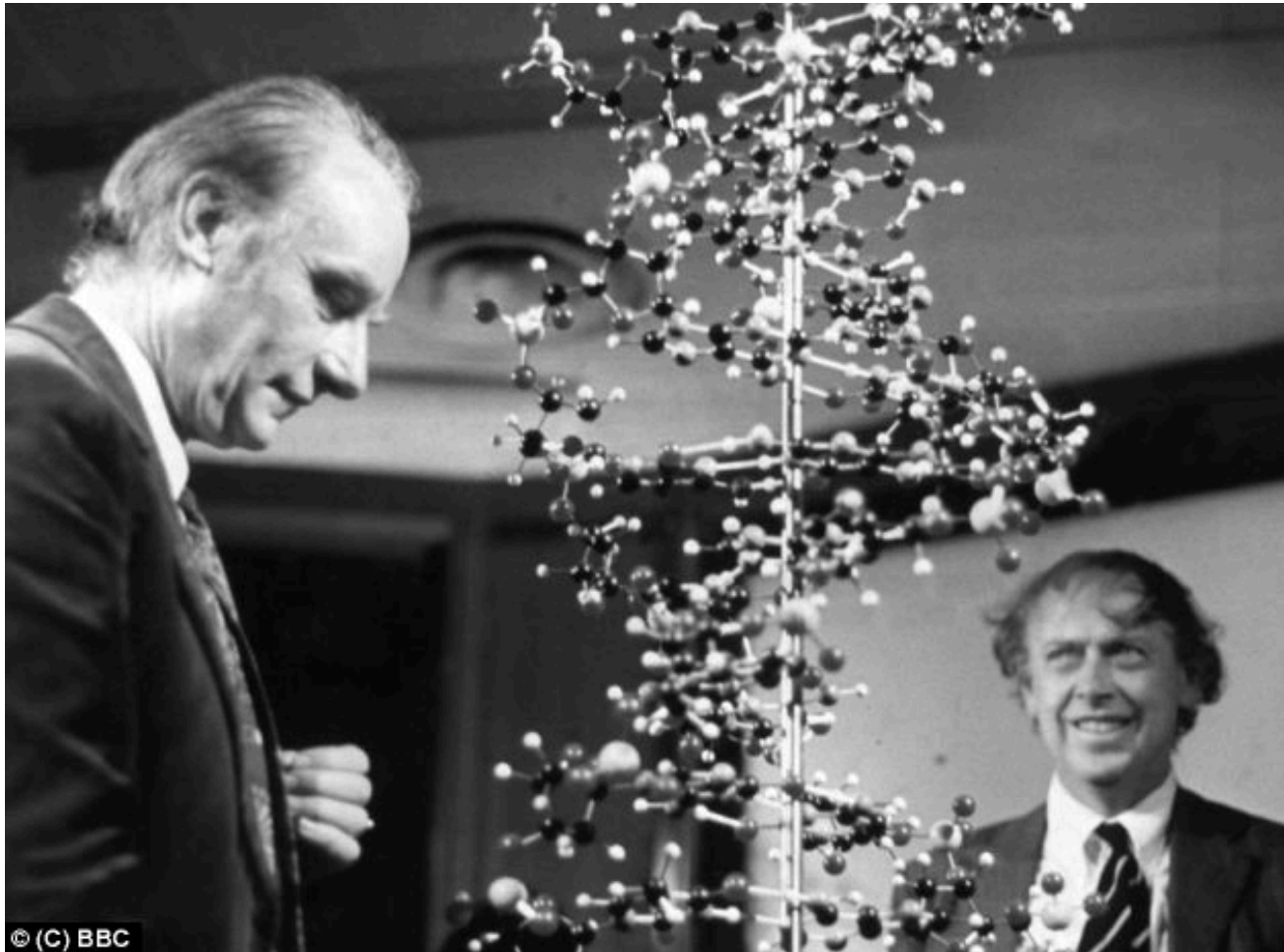
**UF** | Research Computing
Information Technology

# Summary

- ◦ The scale of biocomputing challenges
- ◦ The evolution of High-Performance Computing
- ◦ Current state of the traditional computing
- ◦ Parallelizing analyses
  - • Traditional multiprocessing
  - • Hadoop
  - • Specialized approaches
- ◦ The interfaces
  - • GUI vs. Web vs. Batch (comman-line)
- ◦ Biocomputing Software (Part II)

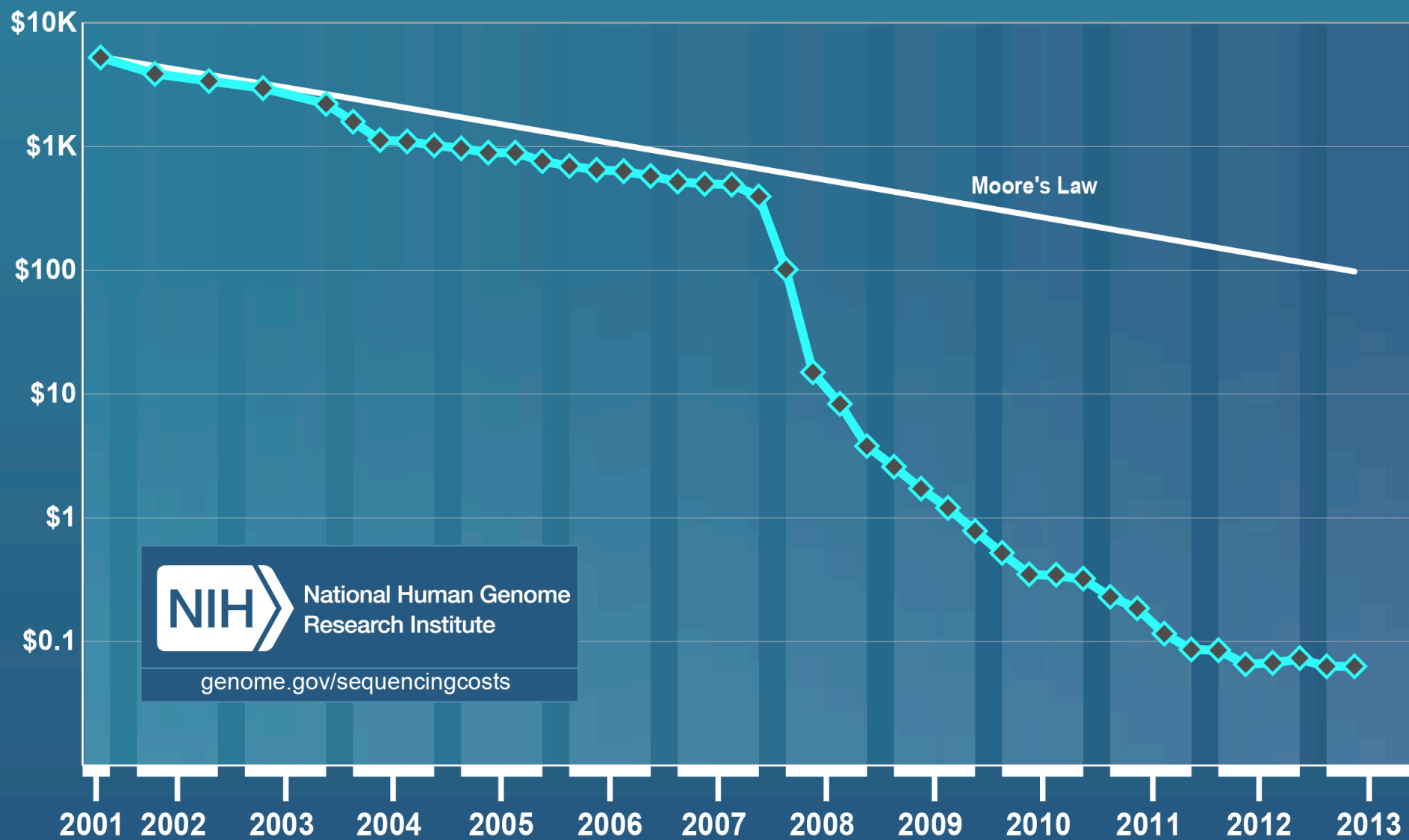# Historical Perspective

**From a molecule to millions of genomes**
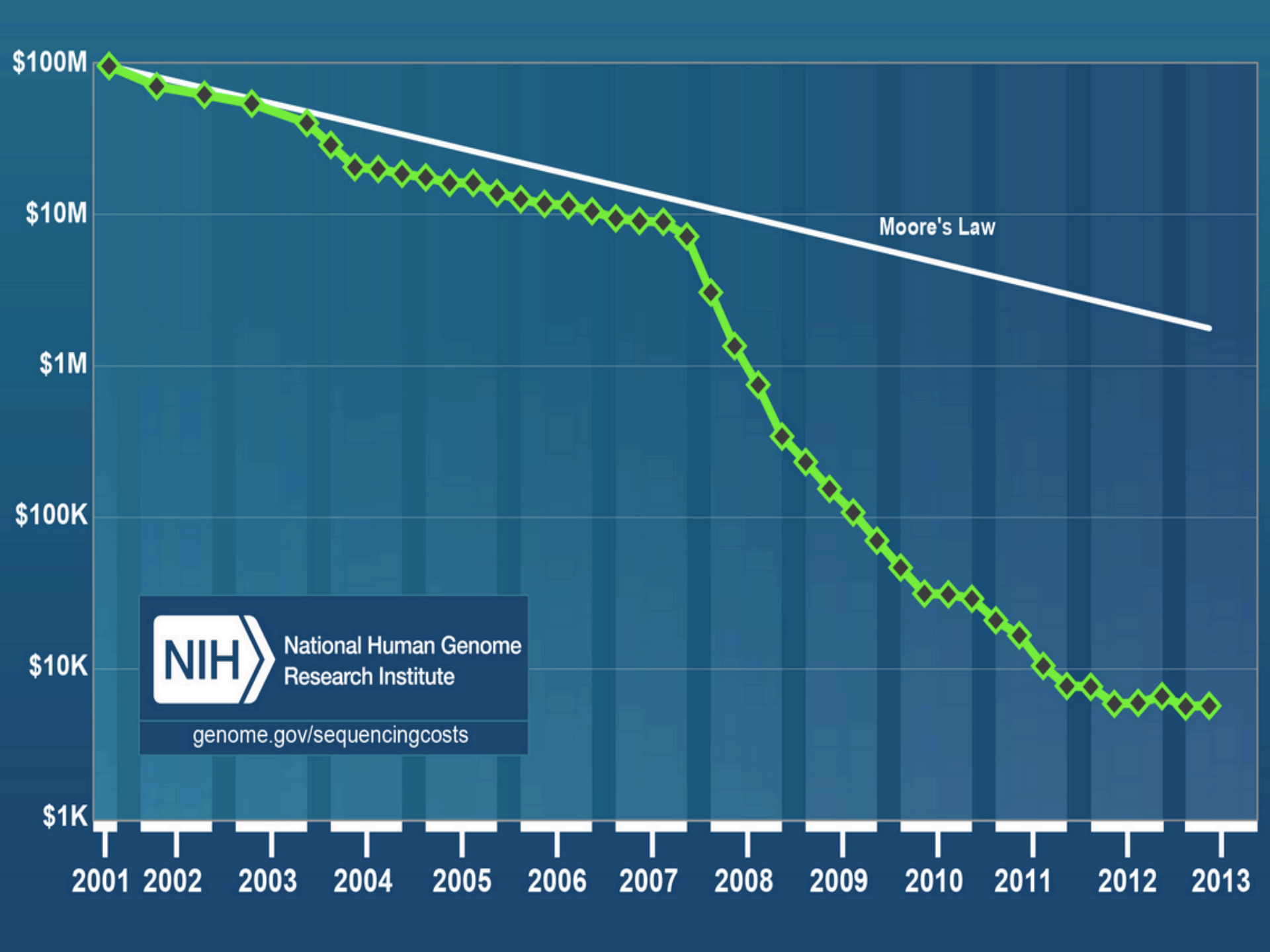
# The Beginning



© (C) BBC
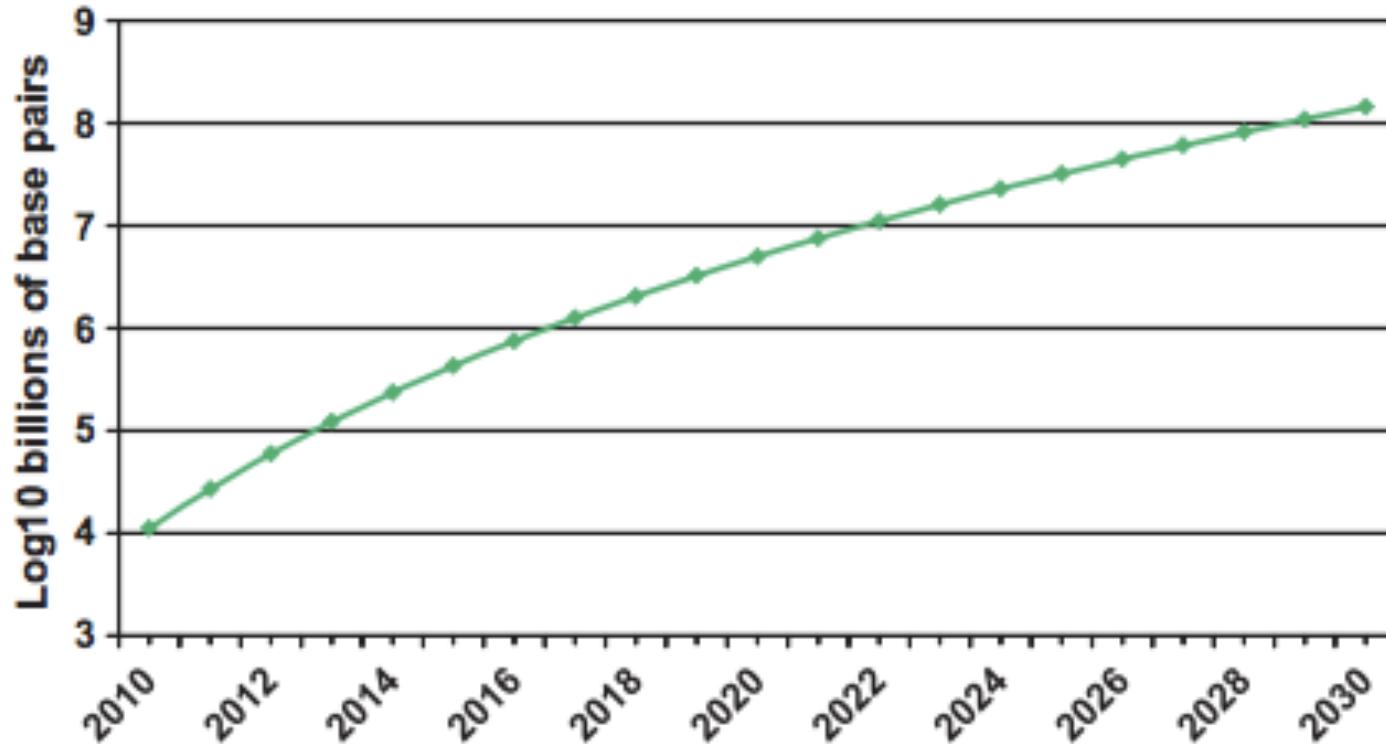
# Sequencing Data Scaling

- ◦ Genome Size * Coverage
  - Viral – 1-100kbp
  - Bacteria, Archaea – 1-10Mbp
  - Simple Eukaryotes – 10-100 Mbp
  - Animals, Plants – 100Mbp - > 100Gbp
- ◦ Sequencing Coverage
  - ~10x in the Sanger Shotgun WGS times
  - ~30x for an average analysis
  - ~100x for metagenomic studies
  - Up to ~1000x for low-frequency SNP analysis in mixed samples

# Growth of Sequencing Data



$10^6$ (Mb) -> $10^9$ (Gb) -> $10^{12}$ (Tb) -> $10^{15}$ (Pb) -> $10^{18}$ (Eb) -> $10^{21}$ (Zb)
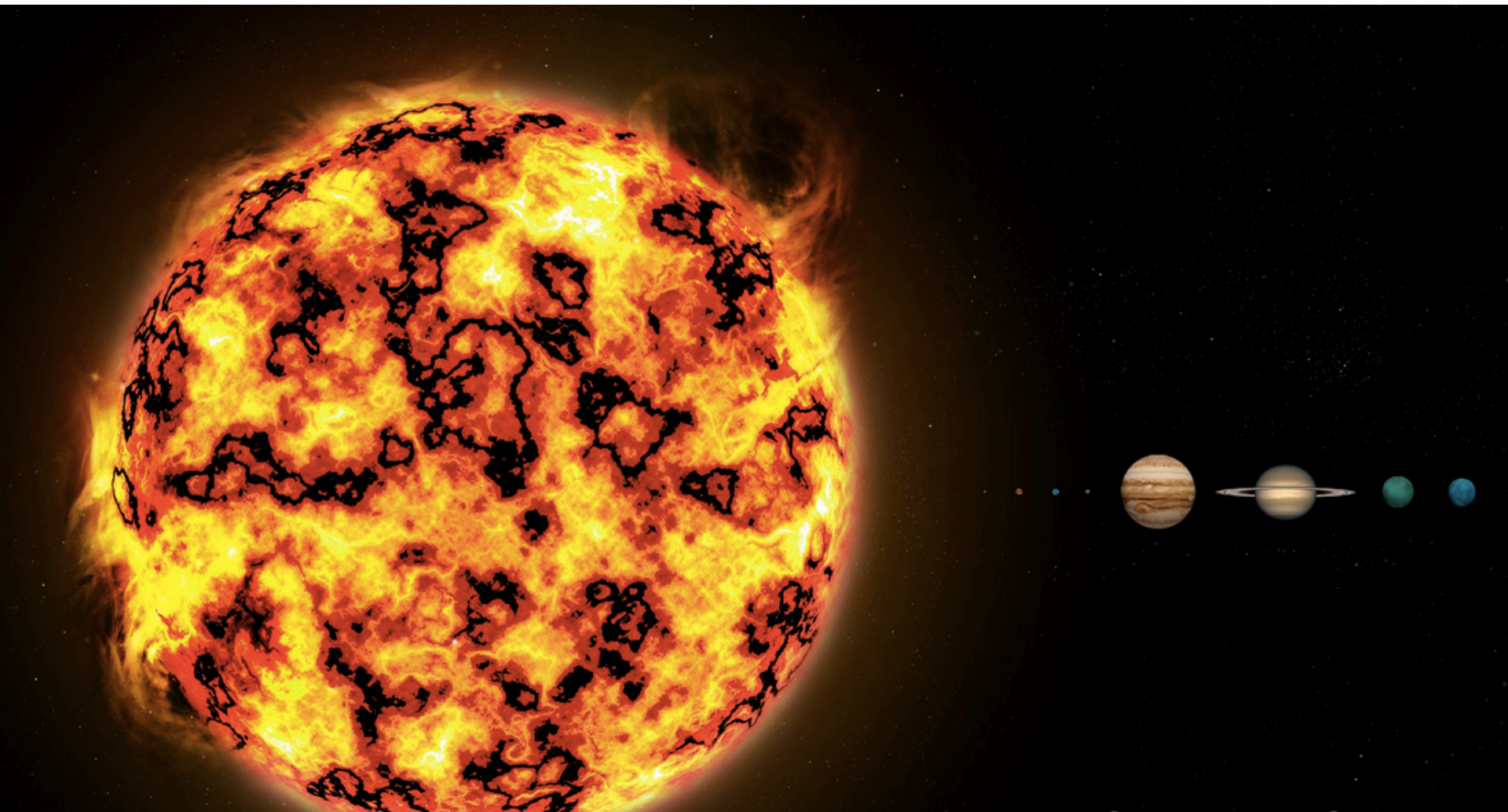
Grossman et al. (2011)

# Growth of Sequencing Data

- 1 Gigabyte: A pickup truck filled with paper OR A symphony in high-fidelity sound OR A movie at TV quality
- 10 Terabytes: The printed collection of the US Library of Congress
- 2 Petabytes: All US academic research libraries
- 5 Exabytes: All words ever spoken by human beings.
- 2.7 Zettabytes: the total amount of global data in 2012 (IDC).

$10^6$ (Mb) -> $10^9$ (Gb) -> $10^{12}$ (Tb) -> $10^{15}$ (Pb) -> $10^{18}$ (Eb) -> $10^{21}$ (Zb)

Grossman et al. (2011)

# BioComputing Growth - NGS

# Evolution of HPC

## From Local to Global

# "Local" BioComputing

# Early Grid BioComputing

HiPerGator

The University of Florida Supercomputer for Research
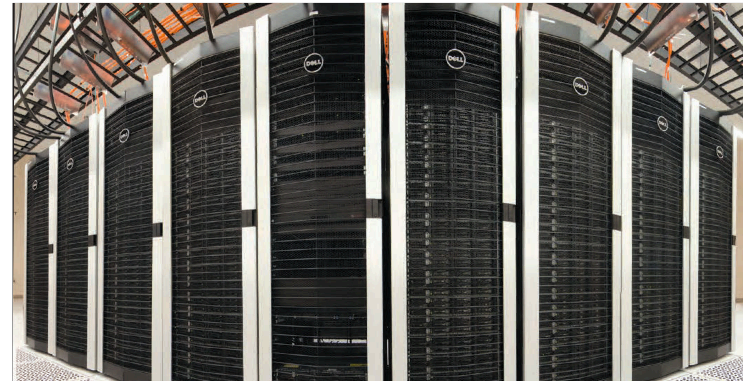
# Contemporary Cluster Specs

- Storage and Networking:
  - 2Pb – Lustre parallel file system
  - 100Gbit networking, Infiniband Fabric
- Computing nodes:
  - 64 x 2.4GHz AMD Abu Dhabi cores
  - 254gb of usable memory
  - 1TB of local storage
- Big memory nodes:
  - 512Gb and 1TB memory with 48-80 cores
- GPU nodes:
  - Tesla, Fermi, Kepler GPU classes

# HPC Considerations

▸ Scale

# HPC Considerations

▸ Computational capacity vs. power and cooling

# UF Data Center

▸ UF Data Center on Eastside Campus
  ◦ 10,000 sq.ft and 1.75 MW total
  ◦ 5,000 sq. ft. space for Research Computing
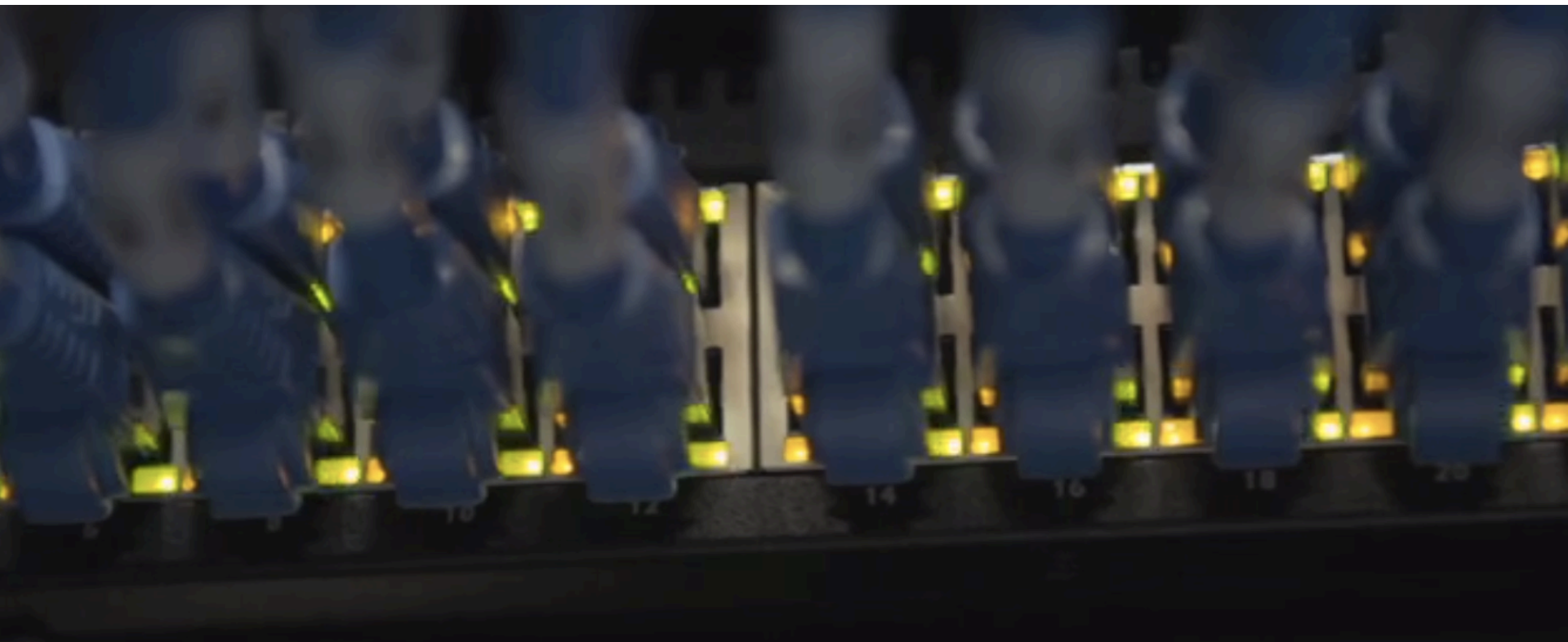
# HPC Considerations

▸ Interconnects
▸ Networking

**Internet2 Network**

▸ Internet2 Innovation Platform
  ◦ 100 Gpbs connectivity
  ◦ Campus Research Network now 200 Gbps

# HPC Considerations

▸ Storage
▸ Parallel file systems
▸ High I/O storage
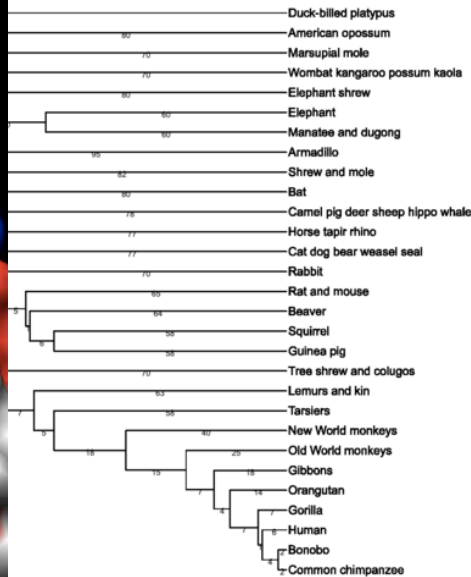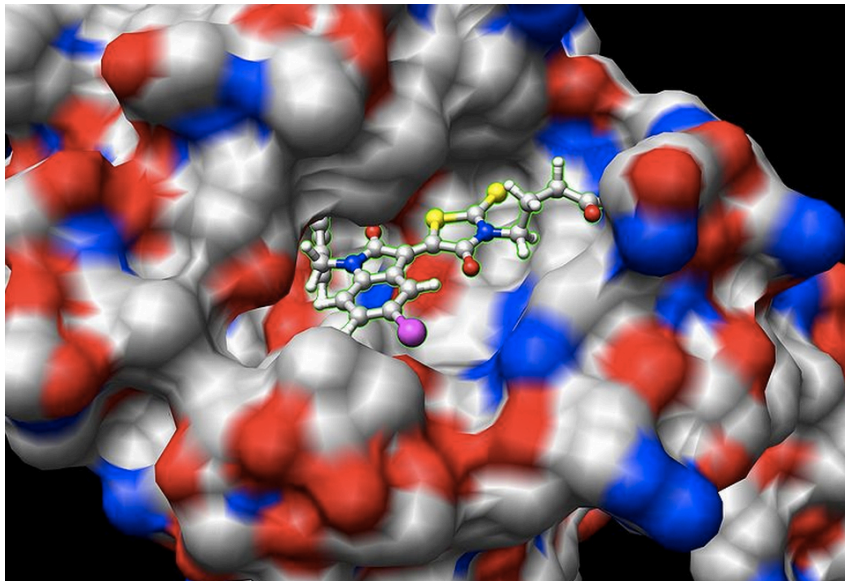▸ Distributed storage

# Scaling the HPC

## The power of many

# Computational Power

◦ Modeling, phylogenetics, simulations

# Traditional Computation

- ◦ *De*-novo genome assembly
- ◦ Short-read mapping
- ◦ RNA-Seq
- ◦ BS-Seq
- ◦ CHIP-Seq
- ◦ SNP calling
- ◦ Pathway analysis
- ◦ …

- ◦ Why? Poor parallelization
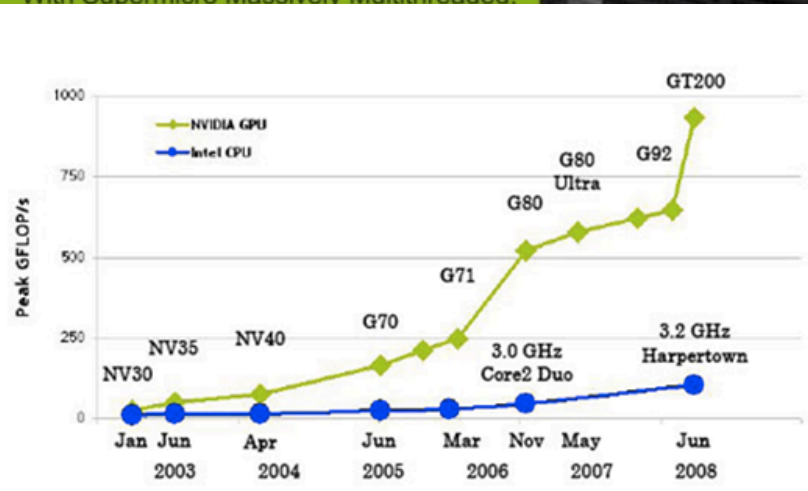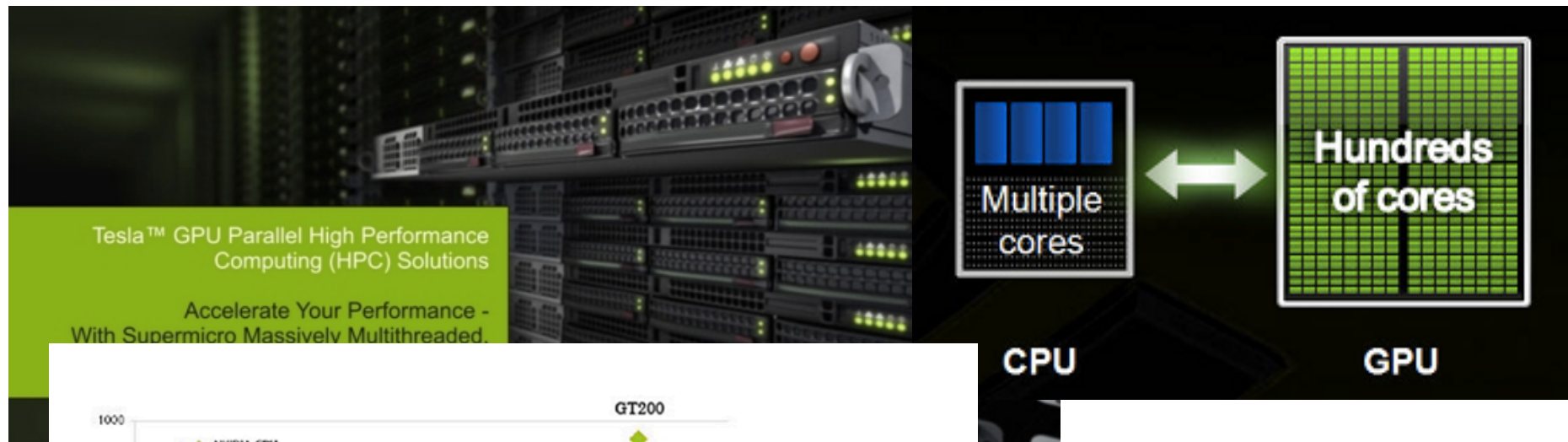
# Circumventing the Moore's Law

## Divide and conquer

# Traditional Parallel Computing

- Split analyses manually, run separately
- Multi-core (SMP) analyses with enabled software
- Multi-node (MPI) analyses with specially constructed software

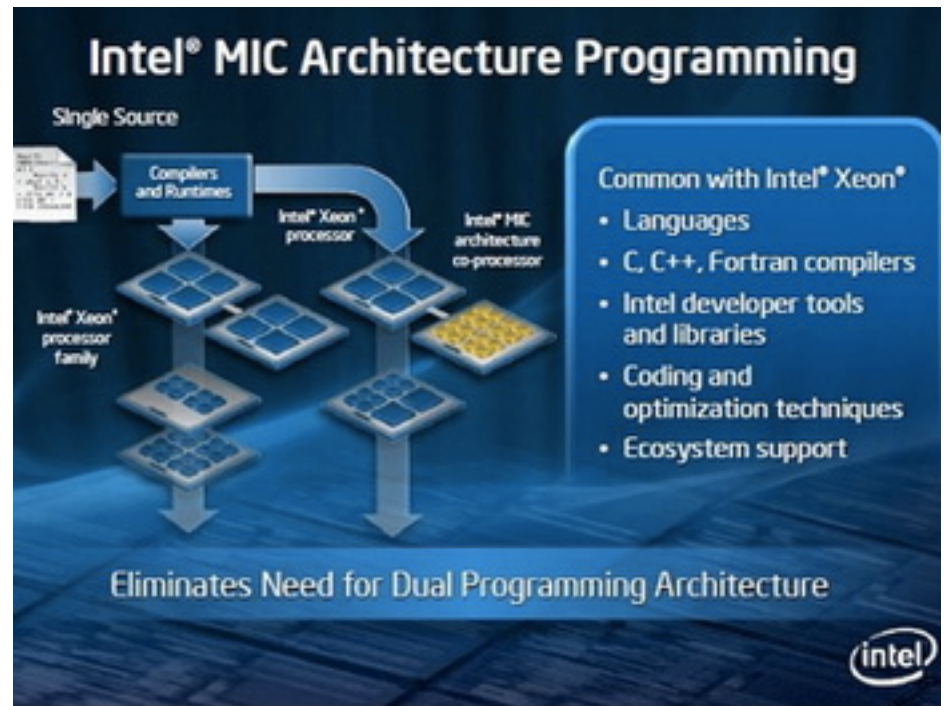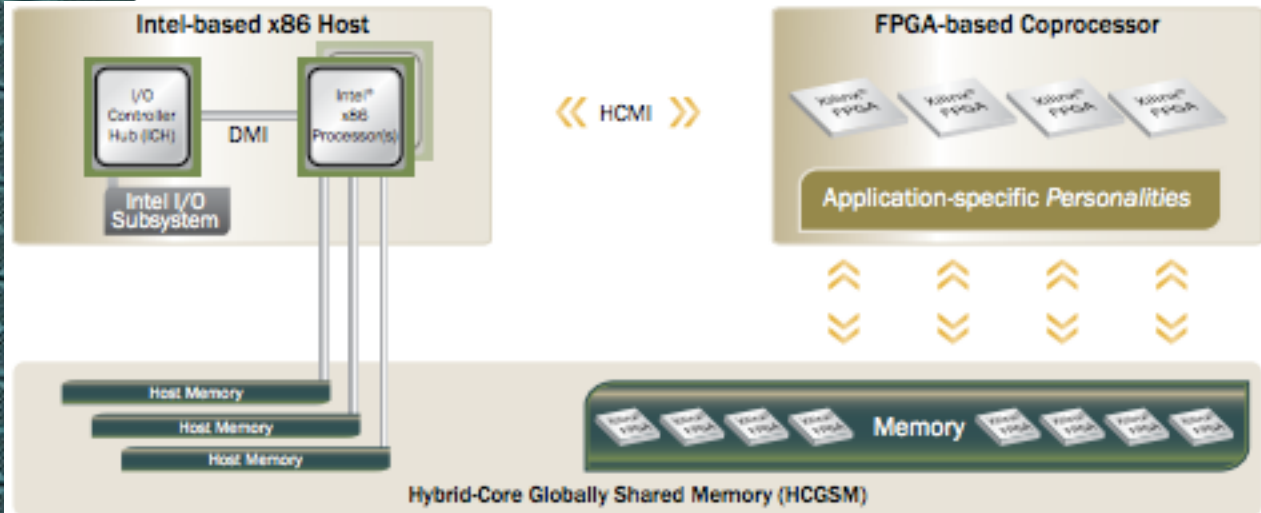# GPU Computing

◦ Highly Parallelizable



Need the code!

CUDA

# MIC Computing

- ◦ Highly Parallelizable
- ◦ Standard x86 cores
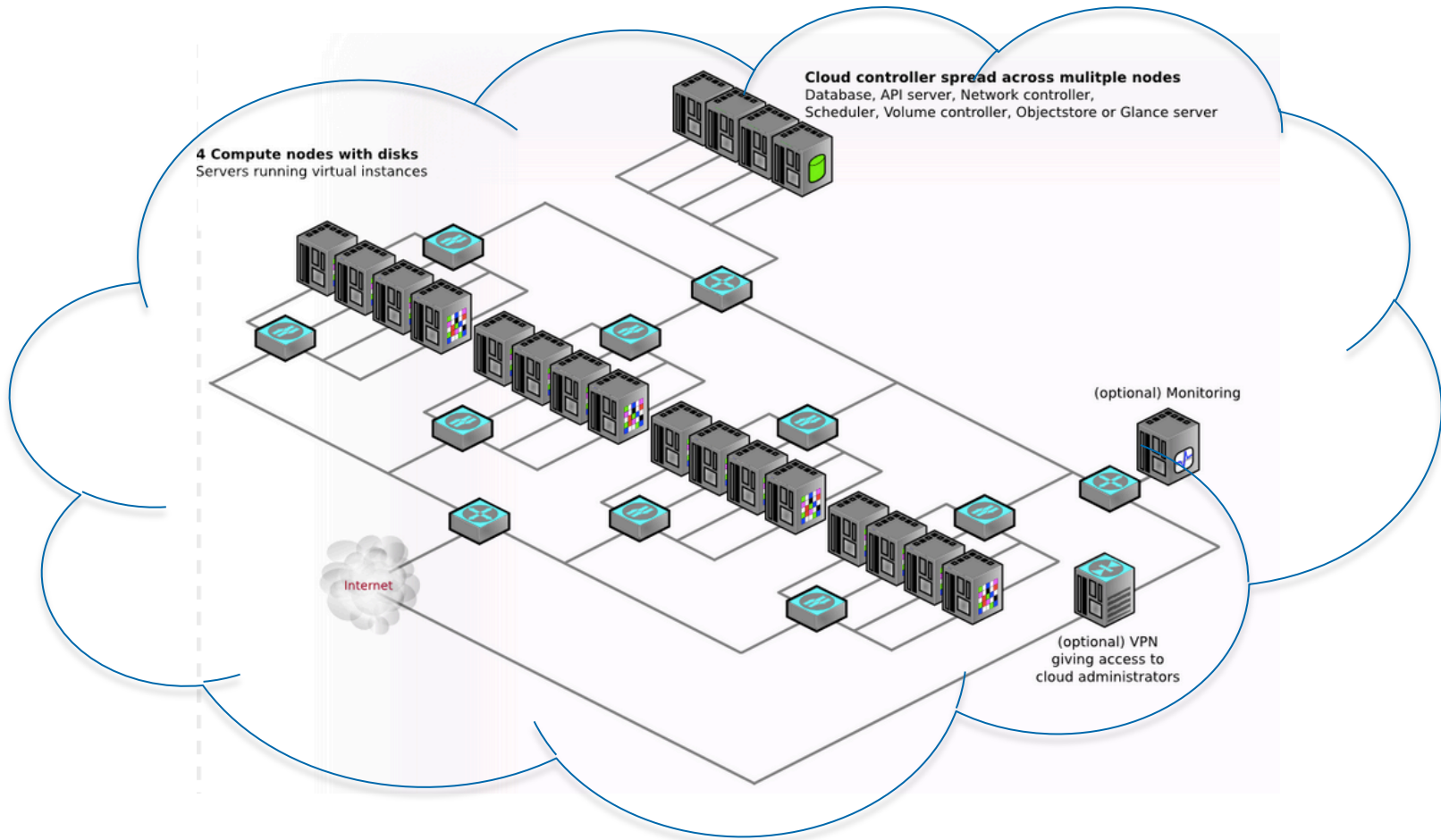- ◦ No need for learning a different programming paradigm ???





Intel® MIC Architecture Programming

# Specialized Processing

# Distributed Computation (Hadoop)



BIG DATA

Results Hypotheses Patterns

Map-Reduce Approach

# Biocomputing Cloud 9 ???

# Interfaces

**Interfaces, Interfaces, Interfaces!!!**

# What the Future May Bring

# Graphical User Interfaces

# Graphical User Interfaces

▸ Proprietary applications
  ◦ Graphical User Interface
  ◦ Integrate multiple tools, pipelines
  ◦ User friendly-wizards for analyses
  ◦ Many can tie into servers or clusters
  ◦ Often highly optimized
  ◦ Expensive
  ◦ Limited flexibility
  ◦ Limited scalability
  ◦ Proprietary algorithms

# Web Interfaces

# Web Interfaces

▸ Galaxy
  ◦ Free, Open Source
  ◦ Public or private instance, physical or cloud-based
  ◦ Web interface
  ◦ Most applications can be integrated
  ◦ User made pipelines
  ◦ Moderately scalable
  ◦ Integrating applications time consuming
  ◦ User made pipelines—where to start? reliability?

# Batch Processing

# Batch Processing

**User interaction**



Login node
(Head node)

**Scheduler**



Tell the scheduler what you want to do

**Compute resources**



Your job runs on the cluster

# Batch Processing

▸ The Linux Command Line
  ◦ Maximum flexibility
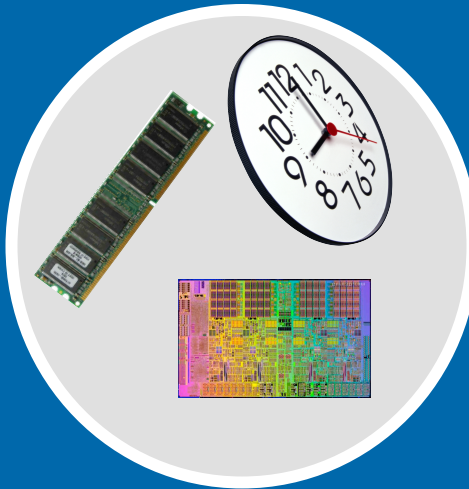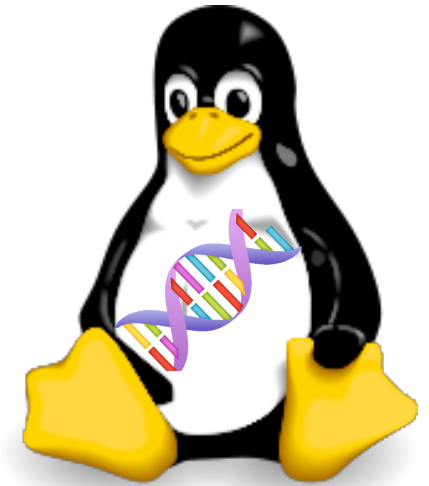  ◦ Most informatics tools run under Linux
  ◦ Write your own tool, or script
  ◦ Maximum scalability
  ◦ Learning barrier of entry

```
Last login: Thu Jul 25 12:03:00 on ttys001
You have mail.
FLMNH-SOL-MAC1:~ gitz$
```

gitz — bash — 100×50 — ⌘2

# Batch processing

▸ Submission Script

```
#!/bin/bash
#PBS -N My_Job_Name
#PBS -M Joe_Shmoe@ufl.edu
#PBS -m abe
#PBS -o My_Job.log
#PBS —e My_Job.err
#PBS -l nodes=1:ppn=1
#PBS -l walltime=00:05:00
#PBS —l pmem=900mb


cd $PBS_O_WORKDIR
date
module load test_app
test_app —i file.txt
```
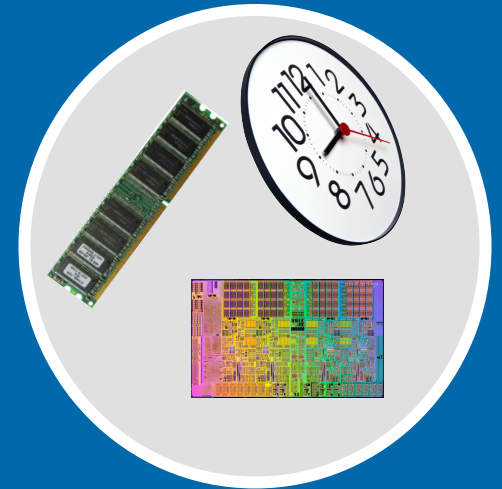
**Scheduler**



Tell the scheduler what you want to do

Compute resources

Your job runs on the cluster

# Accessing software via environment modules

▸ `module load trinity`

▸ Automatically:
  ◦ Sets, `$HPC_TRINITY_DIR`
    • To run Inchworm, simply type

      `inchworm --reads reads.fa --run_inchworm [opts]`

  ◦ Loads Bowtie and Allpaths, two Trinity dependencies
    • You don't need to hunt those down, or worry if they are in your path or not

# It's all in the software!

Matt Gitzendanner

UF Research Computing

# Questions?

Thank you!