# High Performance Computing in Life Sciences

## Part I
**HPC Introduction**
**Introduction**

## PartII
**BioComputing Sofware**

Oleksandr Moskalenko
om@ufl.edu

Matt Gitzendanner
magitz@ufl.edu

UF | Research Computing
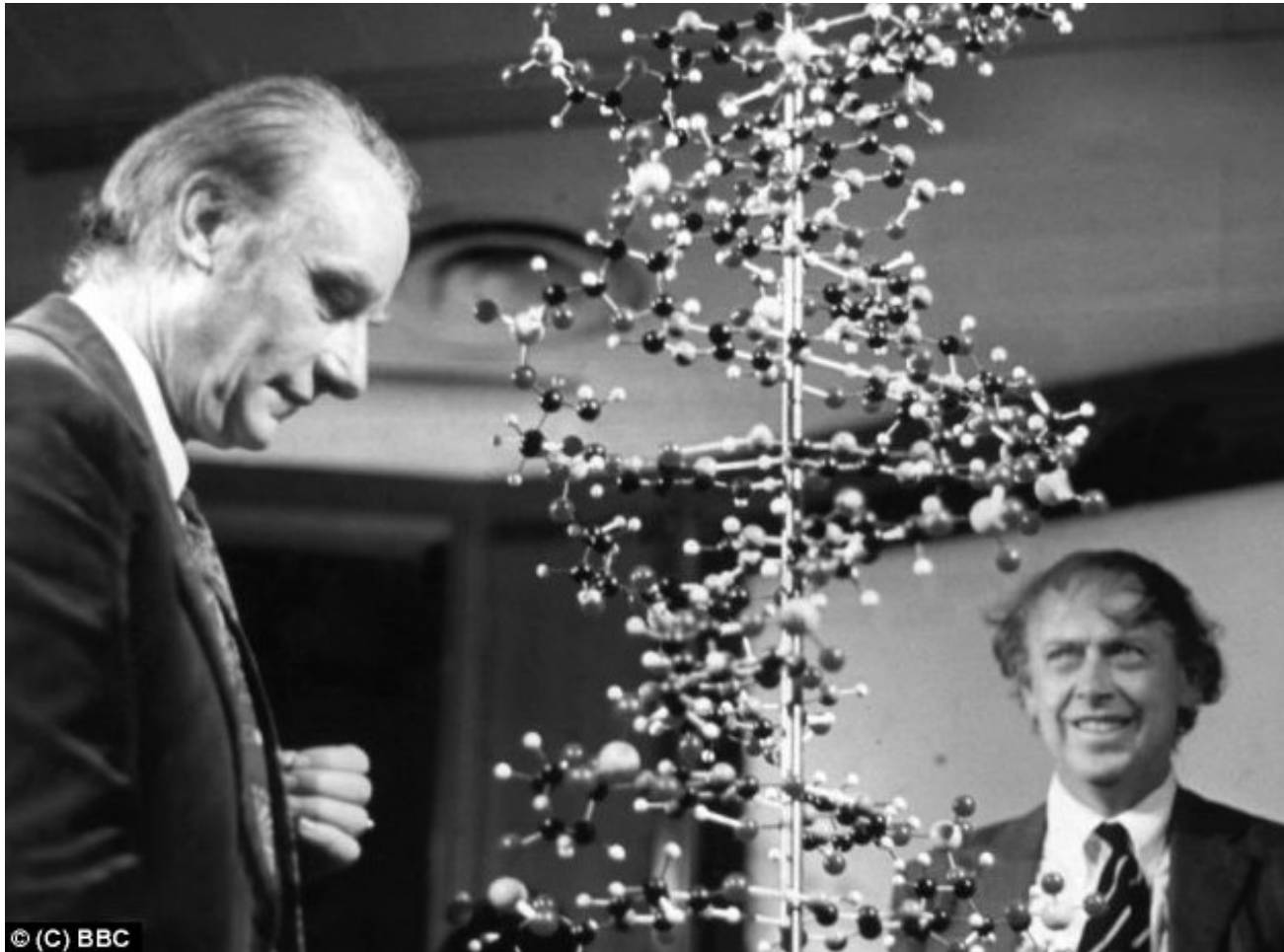Information Technology

# 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)

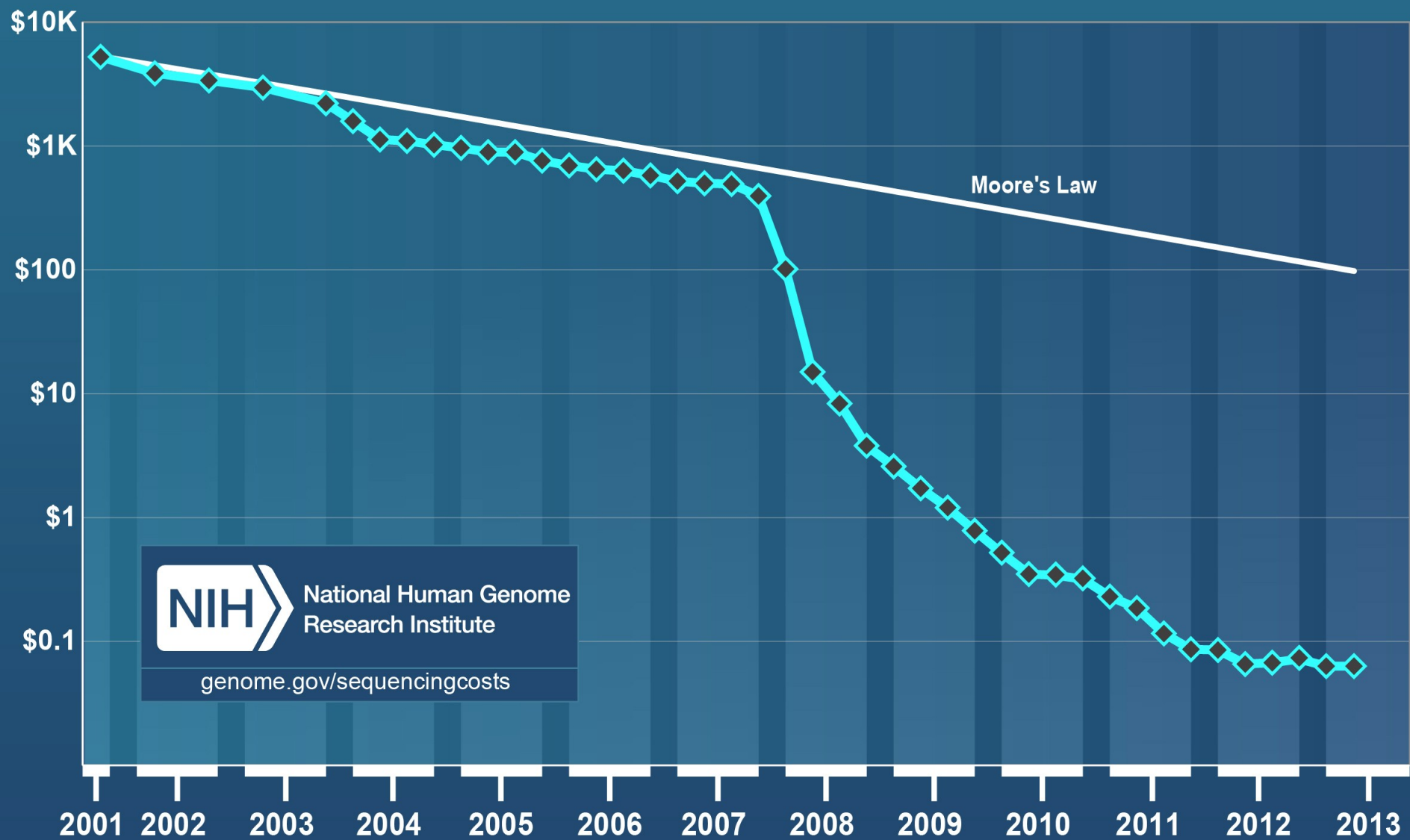# Historical Perspective

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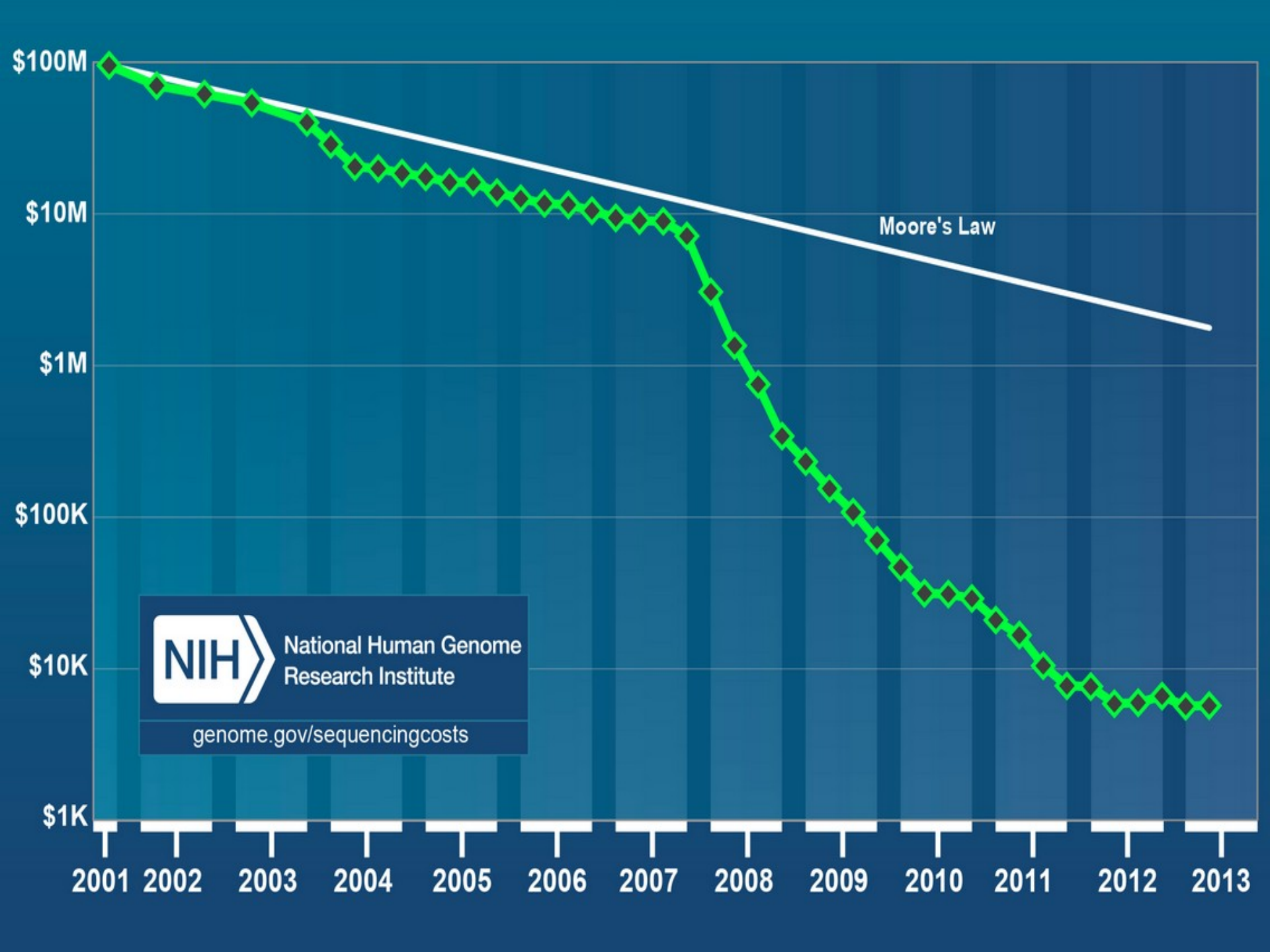
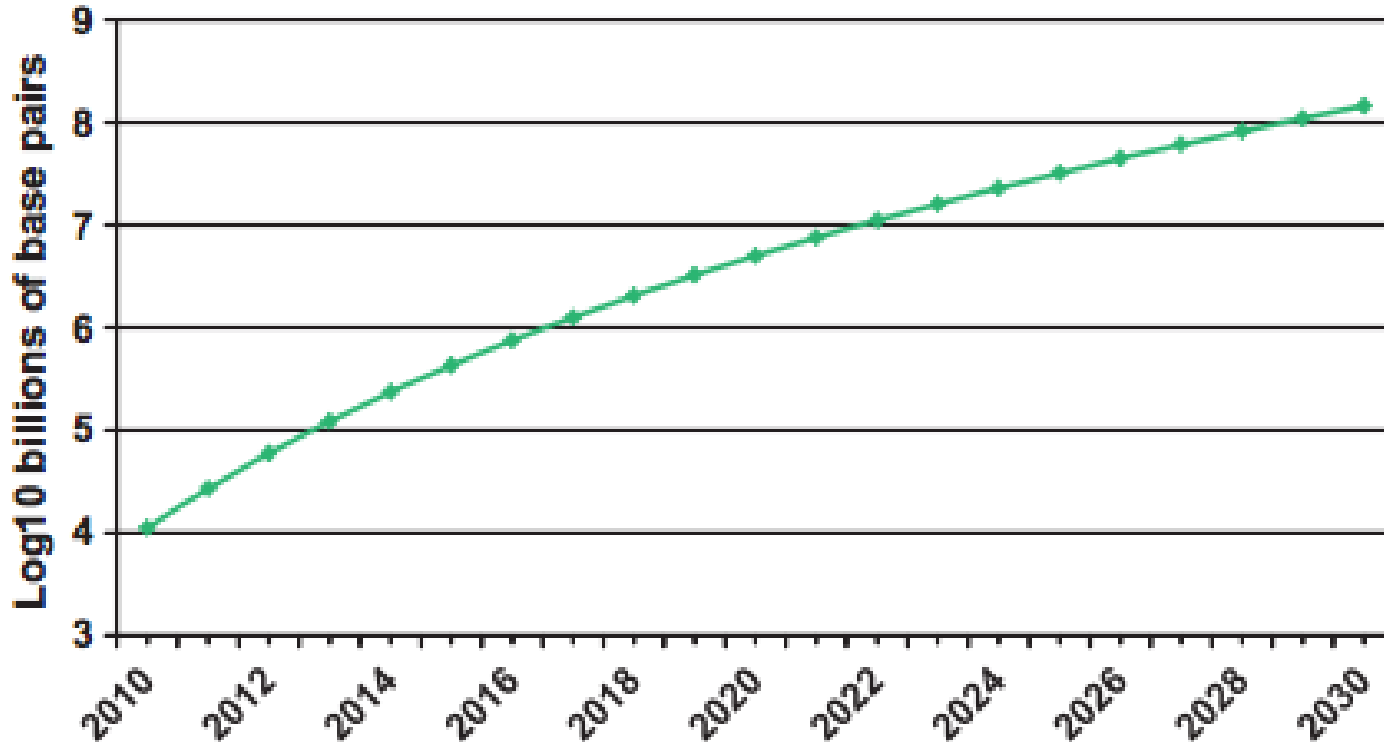# The Beginning


© (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

# Cost per Raw Megabase of DNA Sequence



Moore's Law

**NIH** National Human Genome Research Institute

genome.gov/sequencingcosts

# Growth of Sequencing Data



106 (Mb) -> 109 (Gb) -> 1012 (Tb) -> 1015 (Pb) -> 1018 (Eb) -> 1021 (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).

106 (Mb) -> 109 (Gb) -> 1012 (Tb) -> 1015 (Pb) -> 1018 (Eb) -> 1021 (Zb)

Grossman et al. (2011)

# BioComputing Growth - NGS

# Evolution of HPC

## From Local to Global

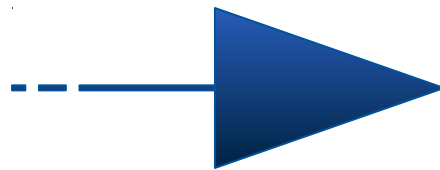# "Local" BioComputing

# 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 logo]

▸ 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
- …

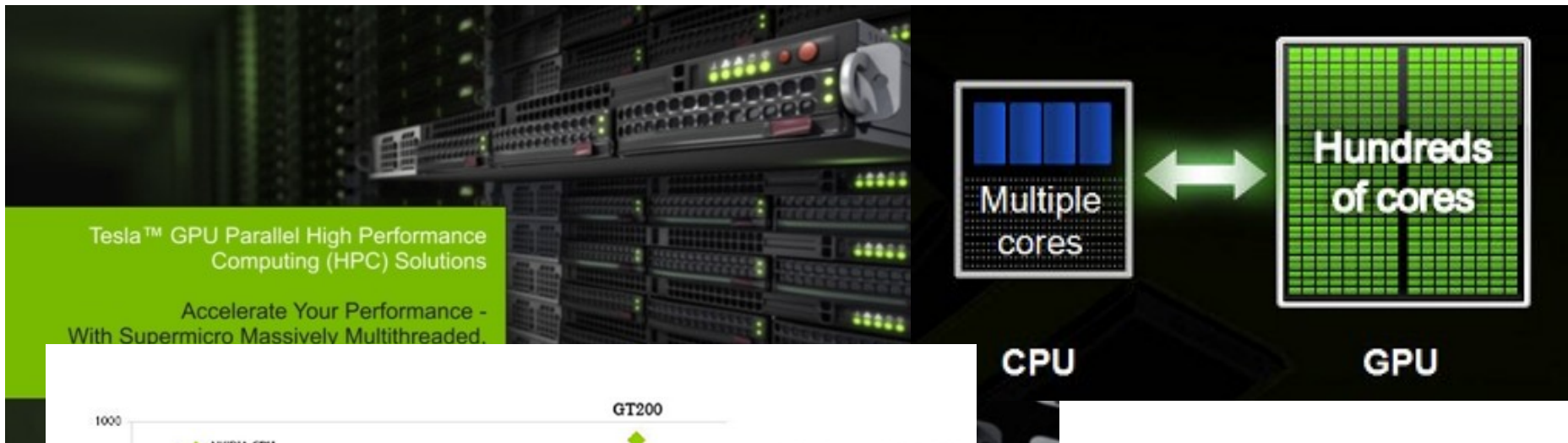- Why? Poor parallelization

# 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
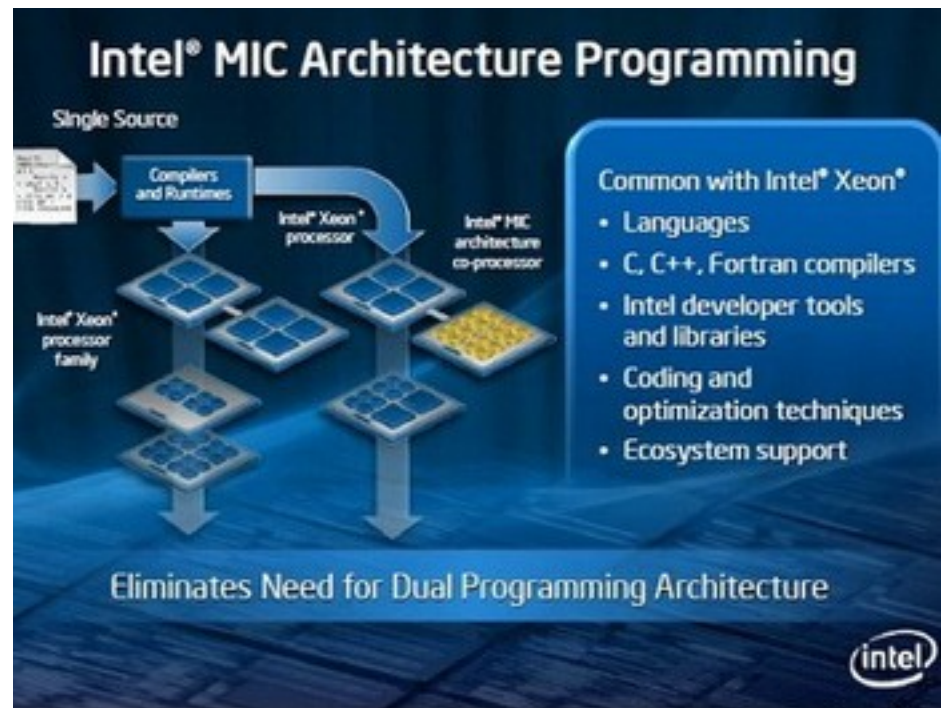
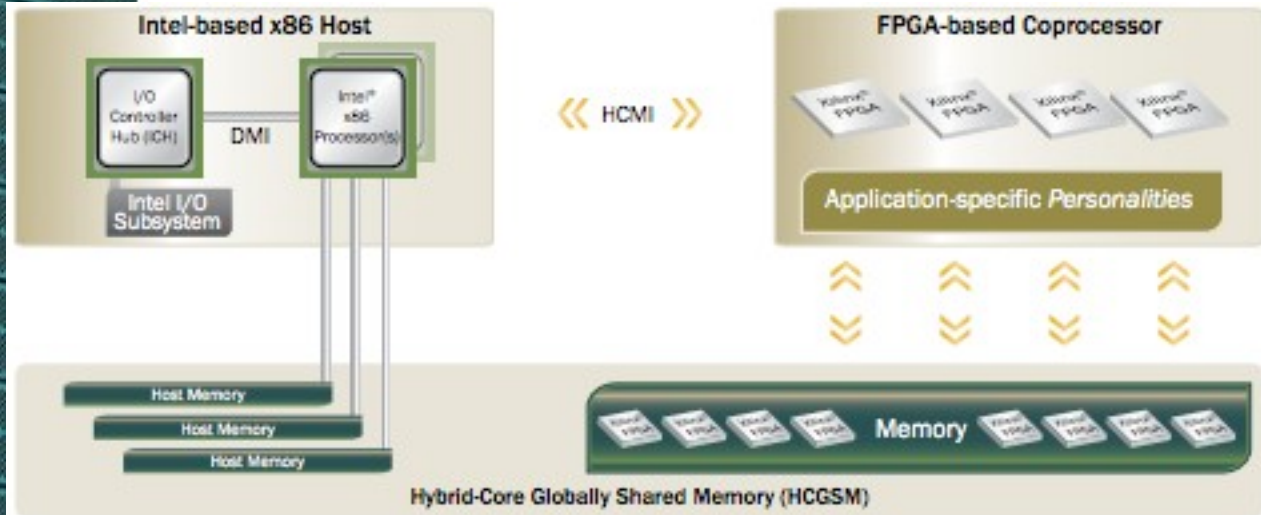# GPU Computing

◦ Highly Parallelizable



Need the code!

CUDA

# MIC Computing

- Highly Parallelizable
- Standard x86 cores
- No need for learning a different programming paradigm ???
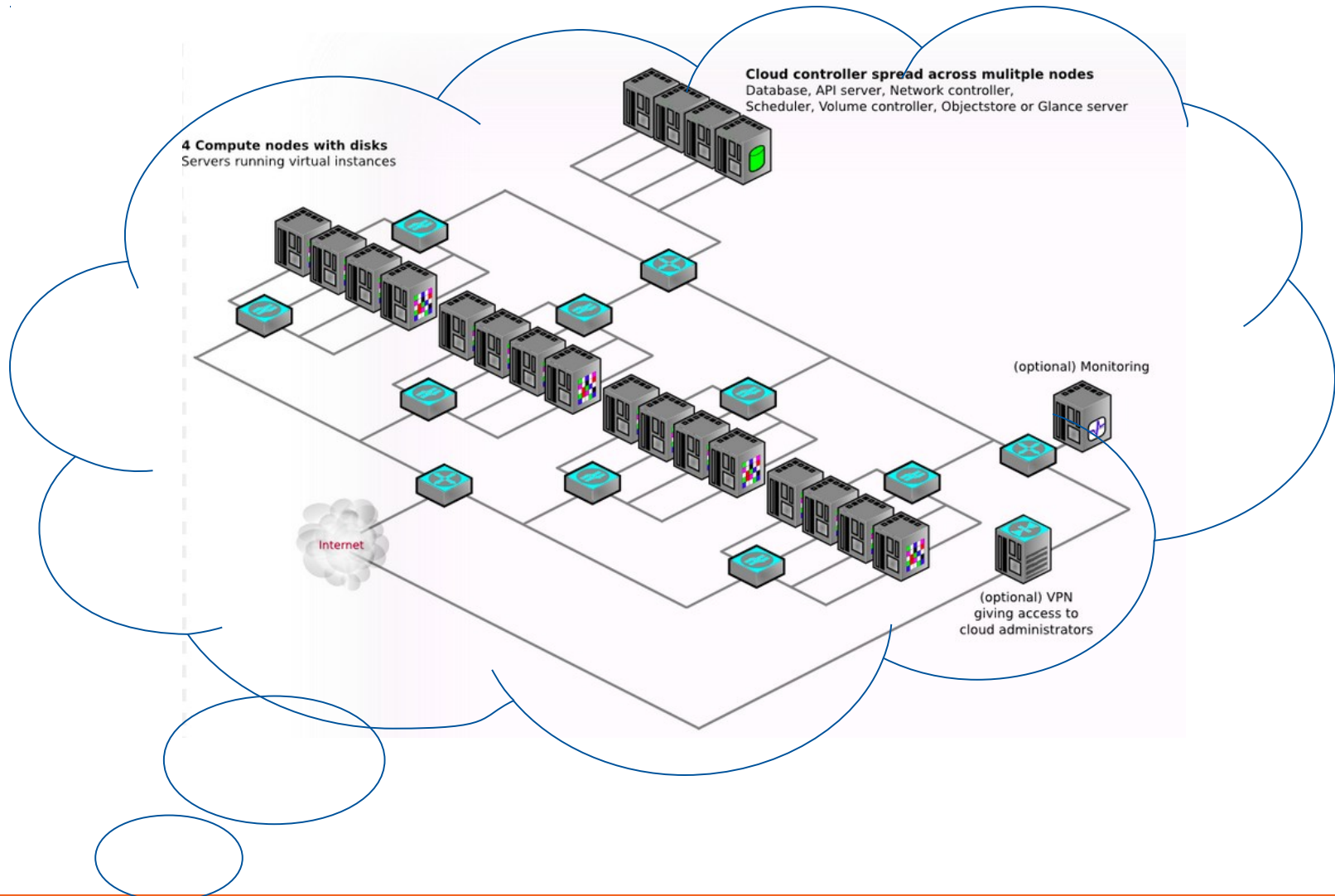
# Specialized Processing

# Distributed Computation (Hadoop)



BIG DATA
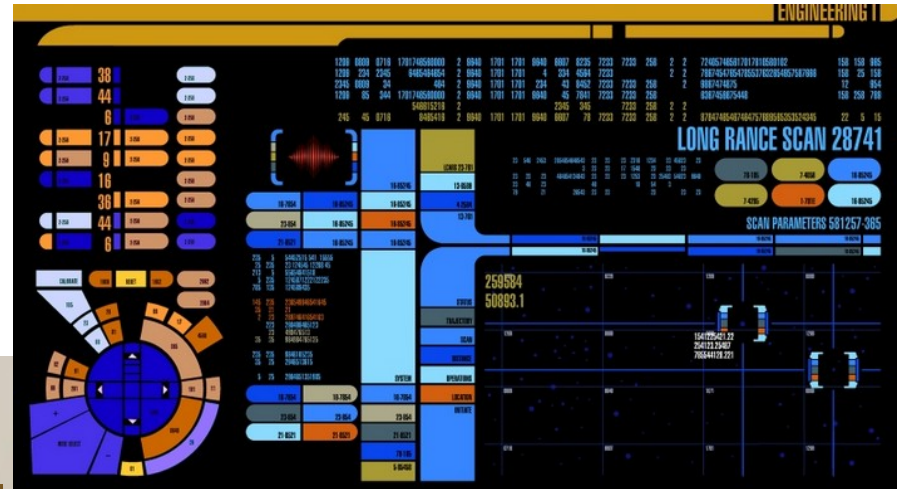
Results Hypotheses Patterns

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

CLC bio
Accelerating Scientific Research

gx

NextGENe®
2nd Generation Sequence Analysis Software

Partek®
Genomics Suite™

# Web Interfaces



- Click to edit Master text styles
  - Second level
    - Third level
      - Fourth level
        - Fifth level

# 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
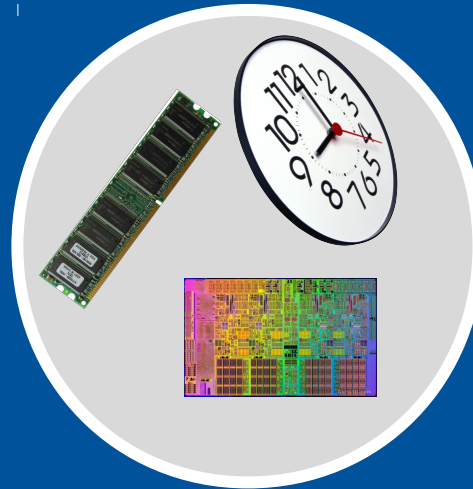- ◦ User made pipelines—where to start? reliability?

# Batch Processing

# Batch Processing

## User interaction
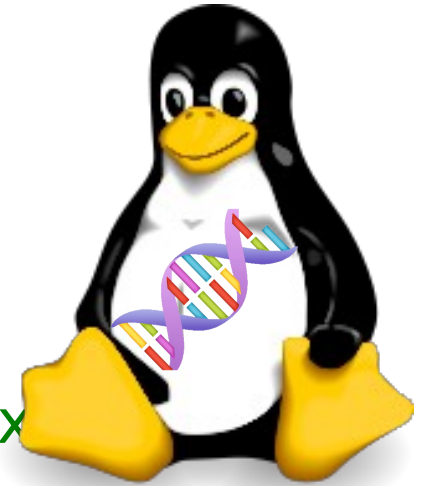
**Galaxy**

Login node
(Head node)

## Scheduler

Tell the scheduler what you want to do

# 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$
```

# 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
```
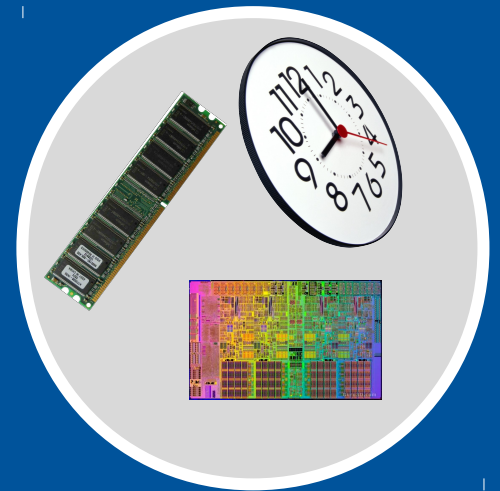
Scheduler



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

► Click to edit Master text styles
  ○ Second level
    ▪ Third level
      ▪ Fourth level
        ▪ Fifth level
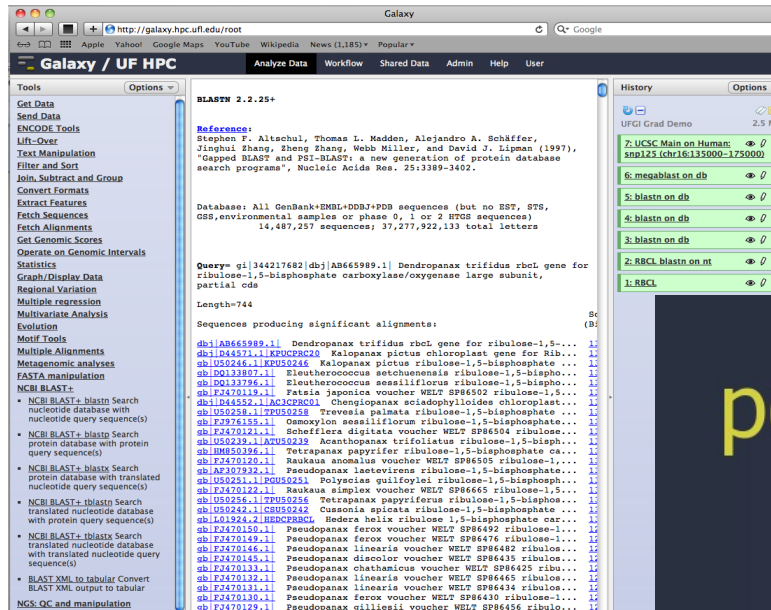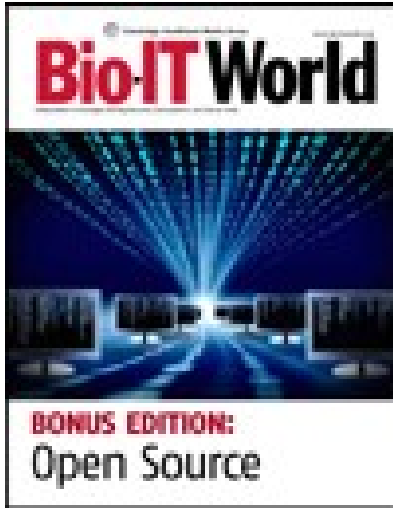
# Galaxy: Data intensive biology *for everyone*

▸ Accessible, reproducible, transparent computational biology

▸ galaxy.hpc.ufl.edu
  ◦ Local instance of Galaxy
    ▫ Faster access to storage, easier upload
    ▫ Local compute resources
    ▫ Local control

# What is Galaxy?

# Cluster basics

## User interaction

**Galaxy**

Login node
(Head node)

## Scheduler

Tell the scheduler what you want to do

# Cluster basics

## User interaction

Galaxy

Login node
(Head node)

## Scheduler

Tell the scheduler what you want to do

# Pond *et al*. 2009, *Genome Research*

**Resource**

# Windshield splatter analysis with the Galaxy metagenomic pipeline

Sergei Kosakovsky Pond,[1,2,6,9] Samir Wadhawan,[3,6,7] Francesca Chiaromonte,[4] Guruprasad Ananda,[1,3] Wen-Yu Chung,[1,3,8] James Taylor,[1,5,9] Anton Nekrutenko,[1,3,9] and The Galaxy Team[1]

**Figure 3.** (*A*) Galaxy history pane showing all steps of a metagenomic analysis described in the study. (*B*) Workflow representation of analysis. Using workflow functionality, the user can rerun analyses in their entirety.

# Galaxy demo

## http://galaxy.hpc.ufl.edu

**Tools**

## Filter (version 1.1.0)

**Filter:**

7: Join two Datasets.. and data 6          ▴▾

Dataset missing? See TIP below.

**With following condition:**

c6/c2>0.5

Double equal signs, ==, must be used as shown above. To filter for an arbitrary string, use the Select tool.

**Number of header lines to skip:**

0

**Execute**

⚠ Double equal signs, ==, must be used as *"equal to"* (e.g., **c1 == 'chr22'**)

ⓘ **TIP:** Attempting to apply a filtering condition may throw exceptions if the data type (e.g., string, integer) in every line of the columns being filtered is not appropriate for the condition (e.g., attempting certain numerical calculations on strings). If an exception is thrown when applying the condition to a line, that line is skipped as invalid for the filter condition. The number of invalid skipped lines is documented in the resulting history item as a "Condition/data issue".

ⓘ **TIP:** If your data is not TAB delimited, use *Text Manipulation->Convert*

### Syntax

The filter tool allows you to restrict the dataset using simple conditional statements.

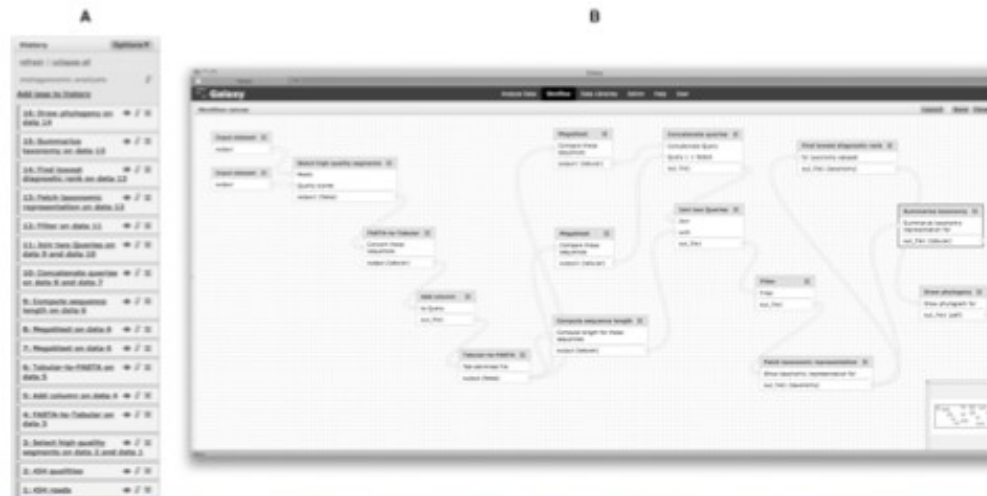Columns are referenced with **c** and a **number**. For example, **c1** refers to the first column of a tab-delimited file
Make sure that multi-character operators contain no white space ( e.g., <= is valid while < = is not valid )
When using 'equal-to' operator **double equal sign '==' must be used** ( e.g., **c1=='chr1'** )
Non-numerical values must be included in single or double quotes ( e.g., **c6=='+'** )
Filtering condition can include logical operators, but **make sure operators are all lower case** ( e.g., **(c1!='chrX' and c1!='chrY') or not c6=='+'** )

### Example

**c1=='chr1'** selects lines in which the first column is chr1
**c3-c2<100*c4** selects lines where subtracting column 3 from column 2 is less than the value of column 4 times 100
**len(c2.split(','))** < 4 will select lines where the second column has less than four comma separated elements
**c2>=1** selects lines in which the value of column 2 is greater than or equal to 1
Numbers should not contain commas – **c2<=44,554,350** will not work, but **c2<=44554350** will
Some words in the data can be used, but must be single or double quoted ( e.g., **c3=='exon'** )

**History** 🔄 ⚙

**Metagenomics**

7.9 MB                                    ✂ 📄

**7: Join two Datasets on data 5 and data 6** 👁 ✏ ✖
18,638 lines
format: tabular, database: ?
💾 ⓘ ⬆ 📊                              ✂ 📄

| **1** | **2** | **3** | **4** | |
|---|---|---|---|---|
| 2 | 78 | 2 | gi\|288887617\|gb\|CP001891.1\| |
| 2 | 78 | 2 | gi\|206564770\|gb\|CP000964.1\| |
| 10 | 167 | 10 | gi\|220939440\|emb\|FP017181.6\| |
| 10 | 167 | 10 | gi\|74136715\|gb\|AC132854.8\| |
| 10 | 167 | 10 | gi\|66841675\|gb\|AC140489.2\| |
| 10 | 167 | 10 | gi\|48374116\|emb\|BX629345.5\| |

**6: Compute sequence length on data 4** 👁 ✏ ✖

**5: megablast on db** 👁 ✏ ✖

**4: Rename sequences on data 3** 👁 ✏ ✖

**3: Select high quality segments on data 1 and data 2** 👁 ✏ ✖

**2: Trip_B.fasta** 👁 ✏ ✖

**1: Trip_B.qual** 👁 ✏ ✖

Analyze Data · Workflow · Shared Data · Visualization · Admin · Help · User

Using 98.1 GB

**Tools**

Evolution
Phylogenetics
Motif Tools
Multiple Alignments
Metagenomic analyses

- dnaclust Cluster sequences into OTUs using DNAclust
- fastaselectclust Get Fasta file of cluster centres from DNAclust output
- dnaclust2tab Convert dnaclust to tabular
- cutClust Remove clusters below a certain depth
- count_clustersize Get cluster size DNAclust output
- Fetch taxonomic representation
- riboPicker Easy identification and removal of rRNA-like sequences.
- Summarize taxonomy
- Draw phylogeny
- Find diagnostic hits
- Find lowest diagnostic rank
- Poisson two-sample test

NCBI BLAST+
FASTA manipulation
NGS: QC and manipulation
NGS: Assembly
NGS: Picard (beta)
NGS: Mapping
NGS: Indel Analysis

**Find lowest diagnostic rank (version 1.0.1)**

**for taxonomy dataset:**

10: Fetch taxonomic r..n on data 9

**require the lowest rank to be at least:**

Family

Subphylum
Superclass
Class
Subclass
Superorder
Order
Suborder
Superfamily
Family

T... e lowest taxonomic rank for which a mategenomic sequencing read is diagnostic. It takes
... *Fetch Taxonomic Ranks* tool (aka Taxonomy format) as the input.

S... reads, **read_1** and **read_2**, with the following taxonomic profiles (scroll sideways to see the
e...

```
read_1 1 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum1 subphylum1 superclass1 cla
read_1 2 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum1 subphylum1 superclass1 cla
read_2 3 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum3 subphylum3 superclass3 cla
read_2 4 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum4 subphylum4 superclass4 cla
```

For **read_1** taxonomic labels are consistent until the genus level, where the taxonomy splits into two branches, one
ending with *subspecies1* and the other with *subspecies2*. This implies **that the lowest taxomomic rank read_1
can identify is SUBTRIBE**. Similarly, read_2 is diagnostic up until the **superphylum** level. As a results the output of
this tool will be:

```
read_1 2 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum1 subphylum1 superclass1 cla
read_2 3 root superkingdom1 kingdom1 subkingdom1 superphylum1 n        n          n          n
```

where, **n** means *EMPTY*.

**What's up with the drop down?**

Why do we need the *require the lowest rank to be at least* dropdown? Let's look at the above example again.
Suppose you need to find only those reads that are diagnostic on at least phylum level. To do this you need to set
the *require the lowest rank to be at least* to **phylum**. As a result your output will look like this:

**History** ⟳ ✿

**Metagenomics**

10.7 MB

✿ **10: Fetch taxonomic representation on data 9**

**9: Convert on data 8**
16,877 lines
format: tabular, database: ?

| 1 | 2 | 3 | 4 | 5 | | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 2 | 78 | 2 | gi | 288887617 | gb | CP001891.1 |
| 2 | 78 | 2 | gi | 206564770 | gb | CP000964.1 |
| 14 | 68 | 14 | gi | 386794017 | gb | CP001925.1 |
| 14 | 68 | 14 | gi | 383101383 | gb | CP002291.1 |
| 14 | 68 | 14 | gi | 374356928 | gb | CP003109.1 |
| 14 | 68 | 14 | gi | 349736152 | gb | CP003034.1 |

**8: Filter on data 7**

**7: Join two Datasets on data 5 and data 6**
18,638 lines
format: tabular, database: ?

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 2 | 78 | 2 | gi\|288887617\|gb\|CP001891.1 |
| 2 | 78 | 2 | gi\|206564770\|gb\|CP000964.1 |
| 10 | 167 | 10 | gi\|220939440\|emb\|FP017181. |
| 10 | 167 | 10 | gi\|74136715\|gb\|AC132854.8 |
| 10 | 167 | 10 | gi\|66841675\|gb\|AC140489.2 |
| 10 | 167 | 10 | gi\|48374116\|emb\|BX629345.5 |

Analyze Data | Workflow | Shared Data ▾ | Visualization ▾ | Admin | Help ▾ | User ▾

Using 98.1 GB

## Tools

**Evolution**
**Phylogenetics**
**Motif Tools**
**Multiple Alignments**
**Metagenomic analyses**

- **dnaclust** Cluster sequences into OTUs using DNAclust
- **fastaselectclust** Get Fasta file of cluster centres from DNAclust output
- **dnaclust2tab** Convert dnaclust to tabular
- **cutClust** Remove clusters below a certain depth
- **count_clustersize** Get cluster size DNAclust output
- **Fetch taxonomic representation**
- **riboPicker** Easy identification and removal of rRNA-like sequences.
- **Summarize taxonomy**
- **Draw phylogeny**
- **Find diagnostic hits**
- **Find lowest diagnostic rank**
- **Poisson two-sample test**

**NCBI BLAST+**
**FASTA manipulation**
**NGS: QC and manipulation**
**NGS: Assembly**
**NGS: Picard (beta)**
**NGS: Mapping**
**NGS: Indel Analysis**
**NGS: RNA Analysis**
**NGS: SAM Tools**
**NGS: GATK Tools (beta)**
**NGS: Peak Calling**

## Draw phylogeny (version 1.0.0)

**Draw phylogram for:**

11: Find lowest diagn.. on data 10

**show ranks from root to:**

Family

Choosing to show entire tree may produce very large PDF file disabling your viewer

**select font size:**

Normal

**maximum number of leaves:**

0

set to 0 to show all

**Execute**

### What it does

Given taxonomy representation (produced by *Taxonomy manipulation–>Fetch Taxonomic Ranks* tool) this utility produces a graphical representations of phylogenetic tree in PDF format.

### Example 1: Fake data

Suppose you have the following dataset:

```
Species_1 1 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum1 subphylum1 superclass
Species_2 2 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum1 subphylum1 superclass
Species_3 3 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum3 subphylum3 superclass
Species_4 4 root superkingdom1 kingdom1 subkingdom1 superphylum1 phylum4 subphylum4 superclass
```

Drawing the tree with default parameters (without changing anything in the interface) will produce this tree:

(for explanation of colors and numbers on the tree scroll to the bottom of this help section)

Here *Class* rank represent terminal nodes (leaves) of the tree because it is the default setting of the "*show ranks from root to*" drop-down. Changing the drop-down to "*Subspecies*" will produce this:

## History

**Metagenomics**
10.7 MB

- 🕐 **11: Find lowest diagnostic rank on data 10**
- ✳ **10: Fetch taxonomic representation on data 9**
- **9: Convert on data 8**
- **8: Filter on data 7**
- **7: Join two Datasets on data 5 and data 6**
- **6: Compute sequence length on data 4**
- **5: megablast on db**
- **4: Rename sequences on data 3**
- **3: Select high quality segments on data 1 and data 2**
- **2: Trip_B.fasta**
- **1: Trip_B.qual**

100%    50%    0%

Viruses:4 ──────────────────────────────────────────── Flaviviridae:4
Archaea:1 ── Thaumarchaeota:1 ──────── Nitrosopumilales:1 ── Nitrosopumilaceae:1

Eukaryota:54
  Fungi:1 ── Dikarya:1 ── Ascomycota:1 ── Pezizomycotina:1 ── Eurotiomycetes:1 ── Eurotiomycetidae:1 ── Eurotiales:1 ── Trichocomaceae:1
  Metazoa:44
    Chordata:7 ── Craniata:7 ── Gnathostomata:7
      Actinopterygii:1 ── Cypriniformes:1 ── Cyprinoidea:1 ── Cyprinidae:1
      Mammalia:6
        Euarchontoglires:1 ── Primates:1 ── Haplorrhini:1 ── Cebidae:1
        Laurasiatheria:5 ── Ruminantia:5 ── Bovidae:2
                                             Cervidae:3
    Arthropoda:37 ── Hexapoda:37 ── Insecta:37 ── Neoptera:37
      Neuroptera:1 ── Mantispidae:1
      Diptera:6
        Nematocera:1 ── Chironomoidea:1 ── Simuliidae:1
        Brachycera:5
          Empidoidea:1 ── Dolichopodidae:1
          Syrphoidea:2 ── Syrphidae:2
          Ephydroidea:1 ── Drosophilidae:1
          Tephritoidea:1 ── Tephritidae:1
      Hymenoptera:1 ── Chalcidoidea:1 ── Agaonidae:1
      Coleoptera:1 ── Polyphaga:1 ── Curculionoidea:1 ── Curculionidae:1
      Hemiptera:26
        Membracoidea:1 ── Membracidae:1
        Aphidoidea:25 ── Aphididae:25
      Amphiesmenoptera:2 ── Lepidoptera:2 ── Glossata:2
        Noctuoidea:1 ── Noctuidae:1
        Tortricoidea:1 ── Tortricidae:1
  Viridiplantae:8 ── Streptophyta:8
    rosids:5 ── Malpighiales:4 ── Salicaceae:4
               Malvales:1 ── Malvaceae:1
    Liliopsida:1 ── commelinids:1 ── Poales:1 ── Poaceae:1
    Coniferopsida:2 ── Coniferales:2 ── Pinaceae:2

root:603

Bacteria:544
  Fusobacteria:1 ── Fusobacteriia:1 ── Fusobacteriales:1 ── Leptotrichiaceae:1
  Firmicutes:98 ── Bacilli:98
    Lactobacillales:76
      Lactobacillaceae:2
      Streptococcaceae:23
      Enterococcaceae:24
      Leuconostocaceae:27
    Bacillales:22 ── Paenibacillaceae:22
  Actinobacteria:2 ── Actinobacteria:2 ── Actinobacteridae:2 ── Actinomycetales:2
    Streptomycineae:1 ── Streptomycetaceae:1
    Frankineae:1 ── Geodermatophilaceae:1
  Proteobacteria:442
    Alphaproteobacteria:1 ── Rhizobiales:1 ── Rhizobiaceae:1
    Gammaproteobacteria:441
      Pseudomonadales:12
        Moraxellaceae:2
        Pseudomonadaceae:10
      Chromatiales:1 ── Chromatiaceae:1
      Enterobacteriales:428 ── Enterobacteriaceae:428

# Reference-based mapping

▸ Map NGS reads onto a reference genome
- ◦ Identify SNPs
- ◦ RNA-seq
- ◦ ChIP-seq
- ◦ Etc.

# Bowtie (Langmead *et al*. 2009)

▶ Pre-built reference genome index
- ○ Burrows-Wheeler transform
- ○ Index needs to be computed prior to mapping
  - ▯ Either build your own: bowtie-build
  - ▯ Or ask for index to be installed for you

▶ Important parameters
- ○ -v vs. –n
  - ▯ Two mapping modes

# Bowtie (Langmead *et al*. 2009)

- ▸ Mapping mode
  - ○ -v: map reads that have less than *v* mismatches
    - Ignores quality scores
    - -v can be 0-3

Number of mismatches for SOAP–like alignment policy (–v):

`-1`

–1 for default MAQ–like alignment policy

Reference    ATGCGTAGTACGTCAACGTGTCACGTGACAGACAGT
Read         CGAAGTACGACAACGGGTCAC

If number of mismatches
<= v, read maps

# Bowtie (Langmead *et al*. 2009)

▶ Mapping mode

○ -n: map using quality scores

- -n: Mismatches in seed (0-3), ignores quality

- -l: seed length (default 28bp)

- -e: max quality score of mismatches across read (default 70)

- Quality scores range from 0-40

Reference ATGCGTAGTACGTCAACGTGTCACGTGACAGACAGT
Read      CGAAGTACGACAACGGGTCAC

Seed: -l 7
-n 1

If sum of quality scores on the mismatches is <=e, read maps here, otherwise not

# Bowtie (Langmead *et al.* 2009)

► Mapping mode
  ◦ -n: map using quality scores
    □ -n: Mismatches in seed (0-3), ignores quality
    □ -l: seed length (default 28bp)
    □ -e: max quality score of mismatches across read (default 70)

**Maximum number of mismatches permitted in the seed (–n):**

2

May be 0, 1, 2, or 3

**Maximum permitted total of quality values at mismatched read positions (–e):**

70

**Seed length (–l):**

28

Minimum value is 5

# Bowtie (Langmead *et al*.  2009)

- Dealing with multiple mappings
  - -k: report up to *k* good alignments per read (1)
  - -a: report all alignments for a read (slow!)
  - -m: don't report if more than m alignments exist
  - -M: like –m, but report 1 random alignment
  - --best: guarantees alignment is in best stratum
  - --strata: don't report suboptimal strata

# Bowtie (Langmead *et al*.  2009)

- Keeping unmapped/mapped reads
  - --un <filename> unmapped reads
  - --al <filename> mapped reads
  - Can be helpful for downstream analyses

- Use –S for SAM output
  - Most likely will process output using SAM anyway

- -p: Bowtie is threaded, can run using multiple cores on *one* node
  - E.g.: nodes=1:ppn=8

# Bowtie2 (Langmead & Salzberg 2012)

- Adds gapped read alignment (indels)
- Faster than Bowtie for reads longer than 50bp
- Supports local alignment
  - Can trim ends that don't map
- Can map reads over Ns in reference
- No colorspace option

# Bowtie2 (Langmead & Salzberg 2012)

▸ Presets for both global and local
  ◦ --very-fast(-local)
  ◦ --fast(-local)
  ◦ **--sensitive(-local) Defaults**
  ◦ --very-sensitive(-local)

# Other mapping applications

- BWA
- Lastz
- Maq
  - Bowtie is generally faster
- Mosaik
  - Handles gapped alignments relative to reference
- PerM
- SRMA

# de Novo Assembly

- ▸ No reference genome
- ▸ Assemble contigs from reads
  - ◦ Assemble scaffolds using paired-end data
- ▸ Most short-read assemblers are de Buijn graph-based



Original sequence
GTAGTATAGTCAGTATCA

Sequence reads
GTAGTA  TAGTAT  AGTATA
   GTATAG  TATAGT
ATAGTC  TAGTCA  AGTCAG
   GTCAGT  TCAGTA
CAGTAT  AGTATC  GTATCA

k-mers (2-mers)
GT  TA  AG  AT  TC  CA

Consensus overlap assembly
GTAGTA
 TAGTAT
  AGTATA
   GTATAG
    TATAGT
     ATAGTC
      TAGTCA
       AGTCAG
        GTCAGT
         TCAGTA
          CAGTAT
           AGTATC
            GTATCA
GTAGTATAGTCAGTATCA

de Bruijn graph

Nature Reviews | Microbiology

# kmers

▶ A kmer is a sequence of length *k*
- Longer kmer
  - More unique
  - Fewer reads/kmer
- Shorter kmer
  - Less unique
  - More reads/kmer

▶ The kmer you use does matter!
- Try different kmers



Schatz *et al*, 2010

# Velvet (Zerbino & Birney 2008)

- Two stages
  - velveth
    - Creates the hash table of kmers
  - velvetg
    - Uses the de Bruijn graph to create contigs & scaffolds
- kmer is critical
  - Default maximum value is 31
  - If you need longer kmer, let us know

# Velvet (Zerbino & Birney 2008)

- Can use multiple types of sequencing inputs
  - Short, long
  - Paired, single
  - Different insert sizes
  - Reference
- A mix of library types is typically needed for de novo genome assembly
- Many helpful scripts distributed with Velvet
  - VelvetOptimiser—helps pick best kmer

# Other de novo assembly applications

- Abyss
- ALLPATHS-LG
  - Has very specific requirements for library types and coverage
- Metavelvet
  - Modified version of Velvet for metagenomics
- Newbler
  - Provided by Roche (454), but can use Illumina data
- SOAPdenovo
- For RNA-seq
  - Oases (builds on after Velvet)
  - SOAPdenovo-TRANS
  - Trinity

# Galaxy demo

## http://galaxy.hpc.ufl.edu

# Questions?

Thank you!