# Use of Globus in GenePattern at UCSD

## Ted Liefeld

Mesirov Lab, School of Medicine
University of California, San Diego

# > 10,000 active tools for bioinformatics

# GenePattern wraps software tools in an accessible visible format

```
> java -Djava.awt.headless=true
-Dwin=cluster.exe -Dmac=clusterMac
-Dlinux=clusterLinux
-Dlinux64=clusterLinux64 -cp
hcl.jar/legacy-gp-modules.jar/ant.jar
org.genepattern.modules.hcl.RunCluster -f
input.filename log.transform row.center
row.normalize column.center
column.normalize -u output.base.name -e
column.distance.measure -g
row.distance.measure -m clustering.method
```

Standard "command-line" method for running analysis



Corresponding GenePattern visual representation

# GenePattern

## Analysis interface



www.genepattern.org



Gene expression heatmap



3D principal components
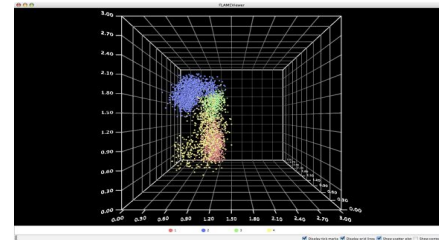


Cell populations

# Hundreds of genomics analysis tools

**Machine learning**
•Clustering, classification, dimension reduction

**Gene expression analysis**
• DESeq2, BWA, HISAT2, HT-Seq, Salmon, Kallisto, Cufflinks, etc.

**Single-cell RNA-seq analysis**
•Seurat, Scanpy, STREAM, CONOS

**Cancer genomics**
•GISTIC, MutSigCV, HAPSEG,

**Gene Set Enrichment Analysis**
•GSEA, ssGSEA, GSEAPreranked

**Collaborative projects**
• OpenCRAVAT *Karchin Lab*
• AMARETTO *Pochet Lab*
• CoGAPS *Fertig Lab*
• MutPanning *Van Allen Lab*
• NDEx *Ideker Lab*
• Next-Generation Clustered Heatmaps *Weinstein Lab (beta)*

**Other**
• Proteomics, Flow Cytometry, Network Analysis, Data import and formatting utilities, etc

# The GenePattern Notebook Environment



- Integrates GenePattern with Jupyter Notebook

- Access hundreds of GenePattern genomic analyses from within a notebook without the need for code

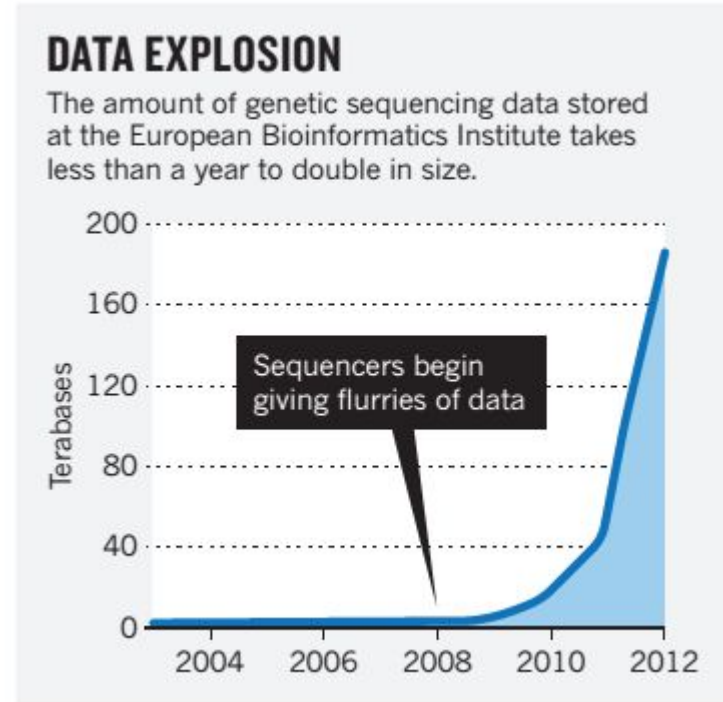notebook.genepattern.org



Reich et al.,
*Cell systems,* 2017

# Genomic Dataset Growth

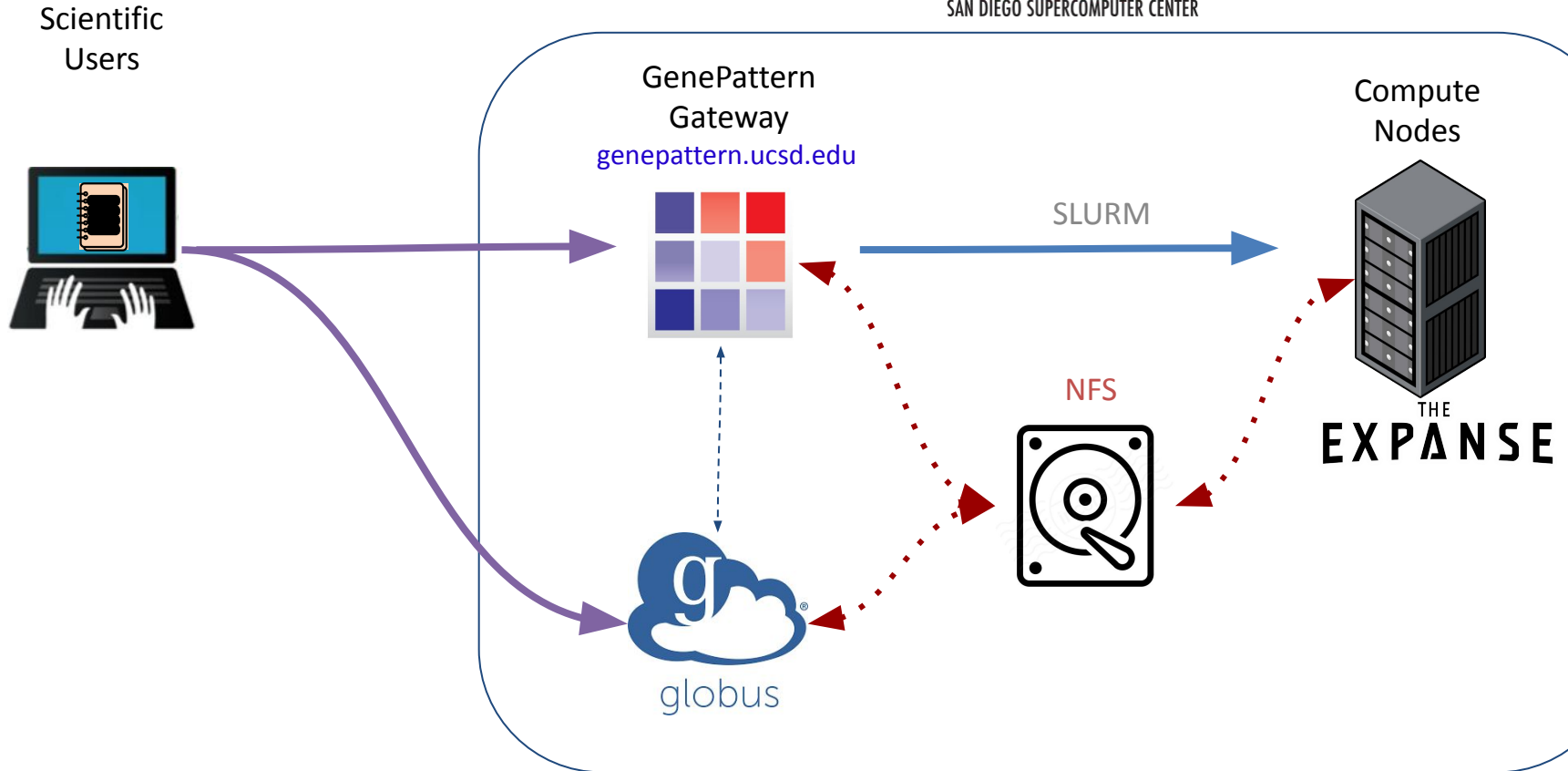~2000,    100 samples, 10000 gene transcripts
     <10 MB file sizes

~2015,   3000 samples,  60000 transcripts,
     10 - 100MB

~2022, 320000 samples, ~60000 transcripts,
     40 - 100 **GB**



**DATA EXPLOSION**
The amount of genetic sequencing data stored
at the European Bioinformatics Institute takes
less than a year to double in size.

Sequencers begin
giving flurries of data

Kashish et al, 2018

# High-level GenePattern Architecture

# GenePattern and Globus



**Sign in to GenePattern**          Click to Register

Username: ted

Password: ••••••••••

Sign in

Forgot your password?

Sign on using your Globus account
You may also use this link to sign in with Google or institutional (for many universities) credentials via Globus.

# GenePattern and Globus

# GenePattern and Globus

# GenePattern and Globus

# GenePattern and Globus

# Acknowledgements

**Mesirov Lab UCSD**
Anthony Castanza
David Eby
Edwin Huang
Forrest Kim
Michael Reich
Thorin Tabor
Helga Thorvaldsdottir
Alexander Wenzel

**SDSC**
Mahidhar Tatineni
Richard Wagner

**Globus**
Brigitte Raumann

www.genepattern.org