



# Computing for Genomics

Dan Rokhsar,  
DOE JGI/LBNL & UC Berkeley  
**NERSC User Meeting**  
4 February 2014

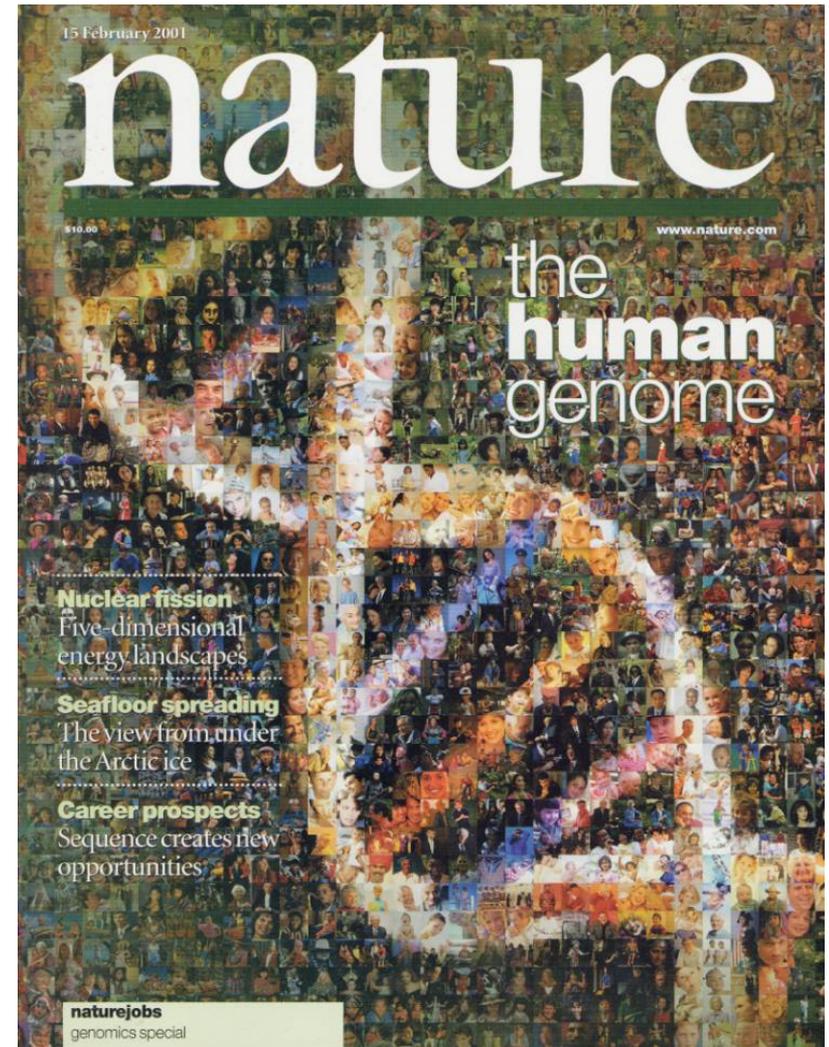
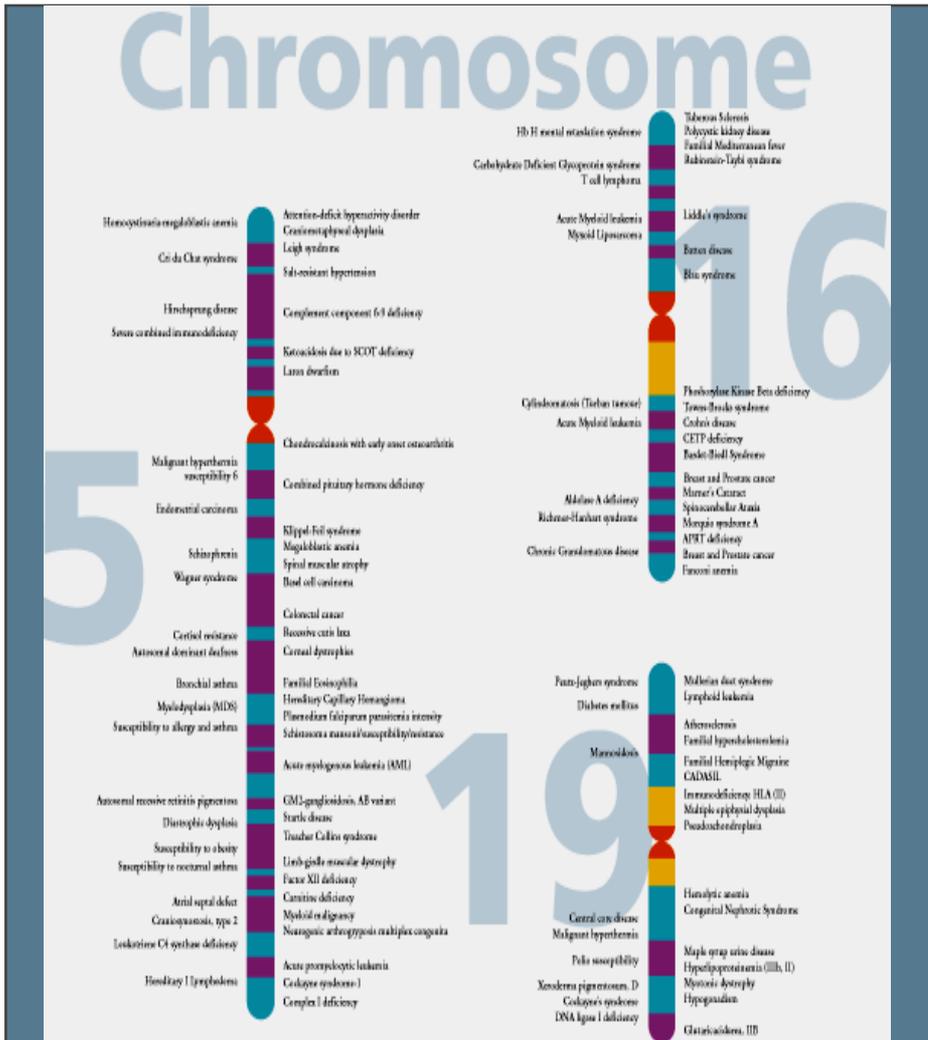
- **The DOE Joint Genome Institute**
  - Who we are, what we do, and why
- **Two computational problems in genomics**
  - “Assembling” genomes from “shotgun” sequence
  - “Annotating” “metagenomic” data

# DOE JGI- A Genomics User Facility



- In 1997, as part of the ramp-up for the human genome project, the DOE Office of Biological and Environmental Research created a virtual “Joint” Genome Institute
- DOE relevance: baseline for study of DNA damage from radiation.
- Goal: sequence 3 of the 23 human chromosomes (~11% of genome).
- “Joint” = LBNL, LLNL, LANL; in 1999 efforts were centralized at one location in Walnut Creek, CA

# By 2003, the “mission” was completed!



- **In the meantime: sequencing had become more efficient and cost-effective.**
- **And genomic information was becoming more and more central to studying all biological systems.**
- **As a continuation of the human genome project over a dozen animal genomes were sequenced by JGI for comparative purposes.**



- **In 2004, JGI became a DOE User Facility for Large Scale Genomics to Enable Bioenergy and Environmental Research**
- **Not “just” sequencing:**
  - New technologies & strategies
  - Novel computational methods and pipelines
  - Access to JGI scientific staff
  - DNA synthesis
  - Building user communities
- **Over 1,200 users ...**
- **In 2010, established partnership with NERSC to manage JGI computing infrastructure and collaborate on common scientific interests**



# SECOND "META- GENOME"

US DEPARTMENT OF ENERGY  
JOINT GENOME INSTITUTE (JGI)

**THIRD ANNUAL USER MEETING**

WALNUT CREEK, CA  
MARCH 26-28, 2008

**FOCUS:**  
Genomics of renewable energy strategies, biomass conversion to biofuels, environmental gene discovery, engineering of fuel-producing organisms. Featuring informatics workshops & tutorials.

**KEYNOTE:**  
**Steven Chu**  
Nobel Laureate, Director, Lawrence Berkeley National Laboratory

TO REGISTER:  
[WWW.JGI.DOE.GOV](http://WWW.JGI.DOE.GOV)

**JGI**  
Joint Genome Institute

**GENOMICS OF ENERGY & ENVIRONMENT**



# Science

# FIRST TREE GENOME

nature  
biotechnology

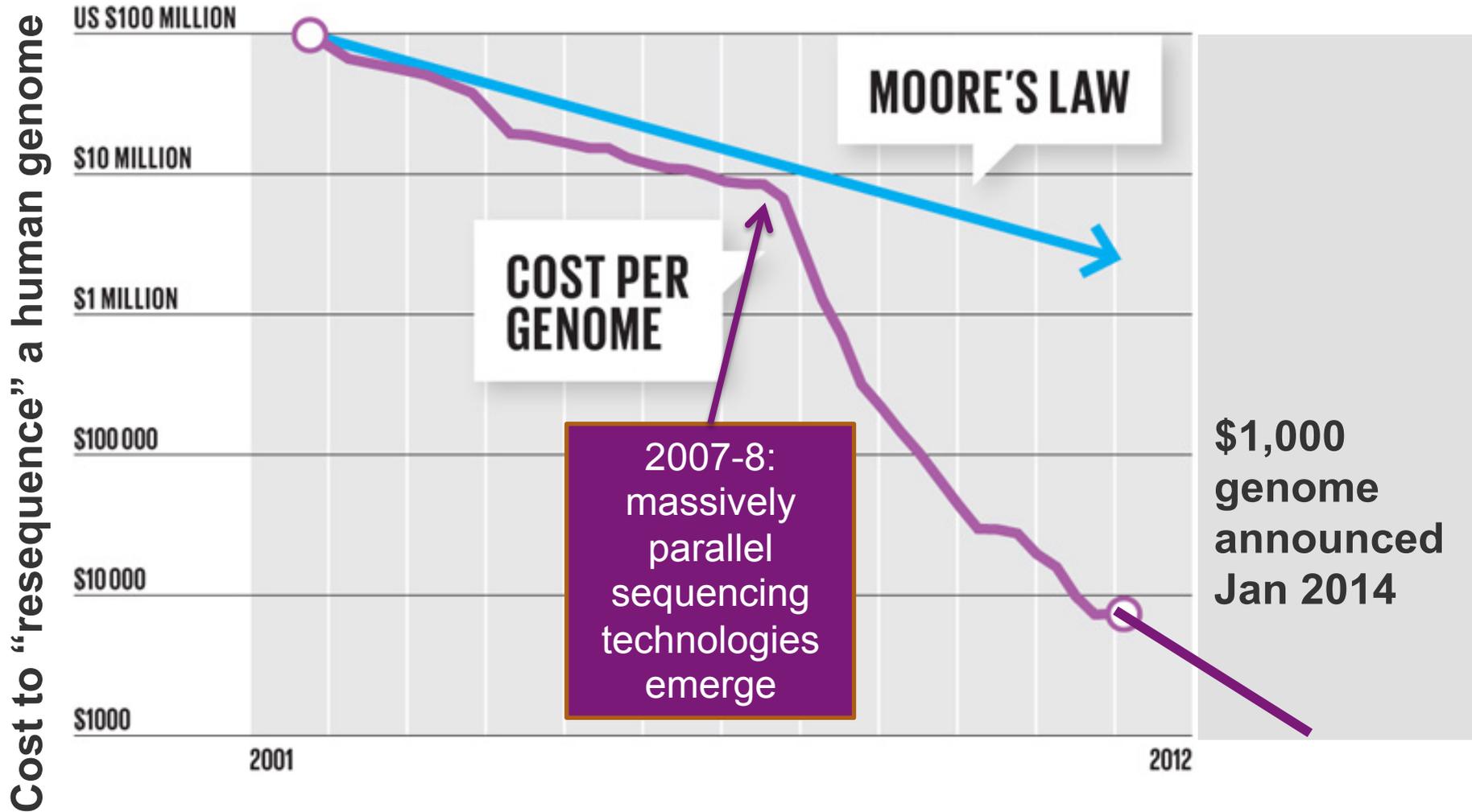
VOLUME 22 NUMBER 6 JUNE 2004  
[www.nature.com/naturebiotechnology](http://www.nature.com/naturebiotechnology)

**FIRST  
FILAMENTOUS  
FUNGUS  
GENOME**

Genome of a white-rot fungus  
Boosting essential oils in plants  
Imaging cancer drug effects on HER2

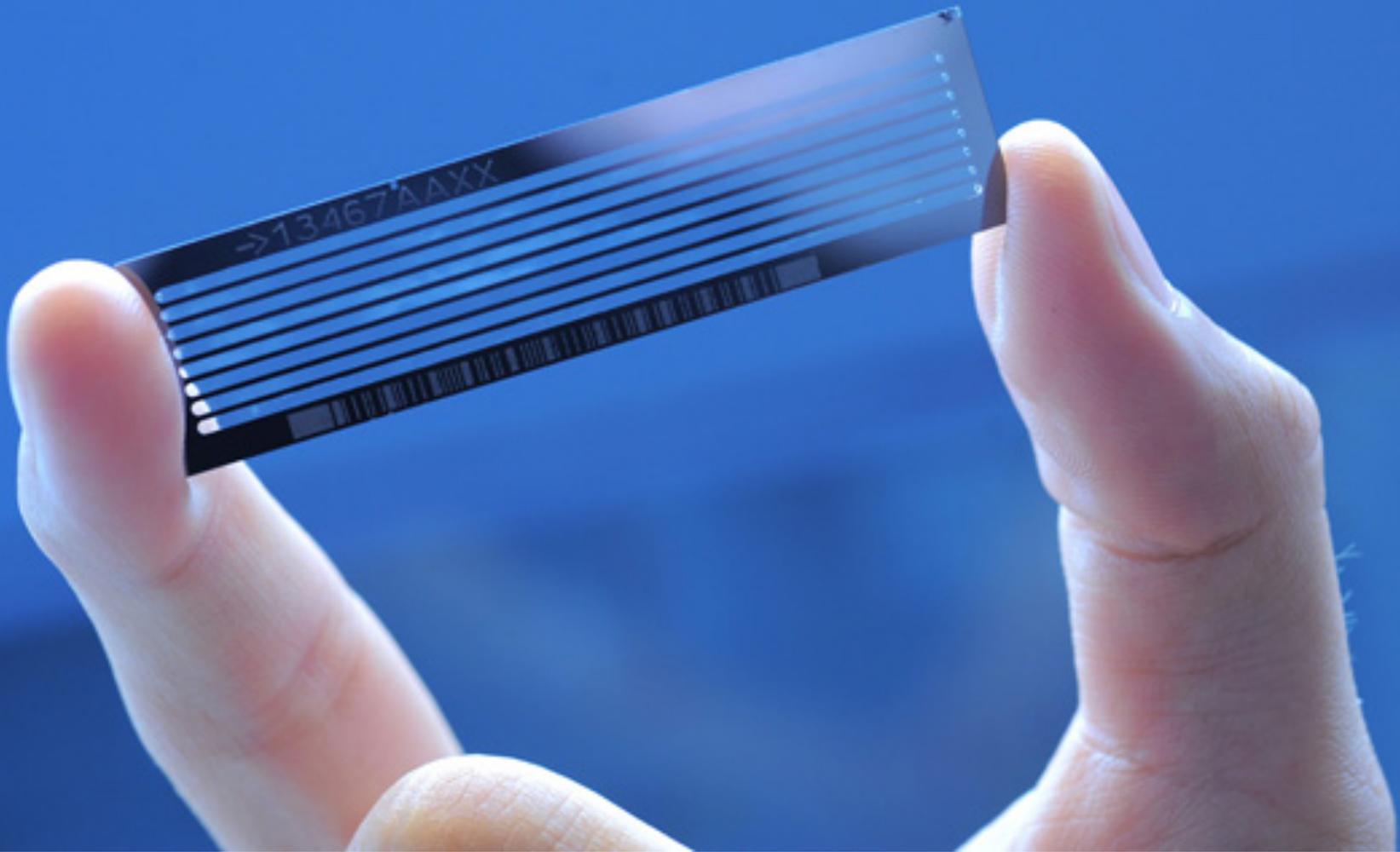


# Sequencing capacity is outstripping Moore's Law

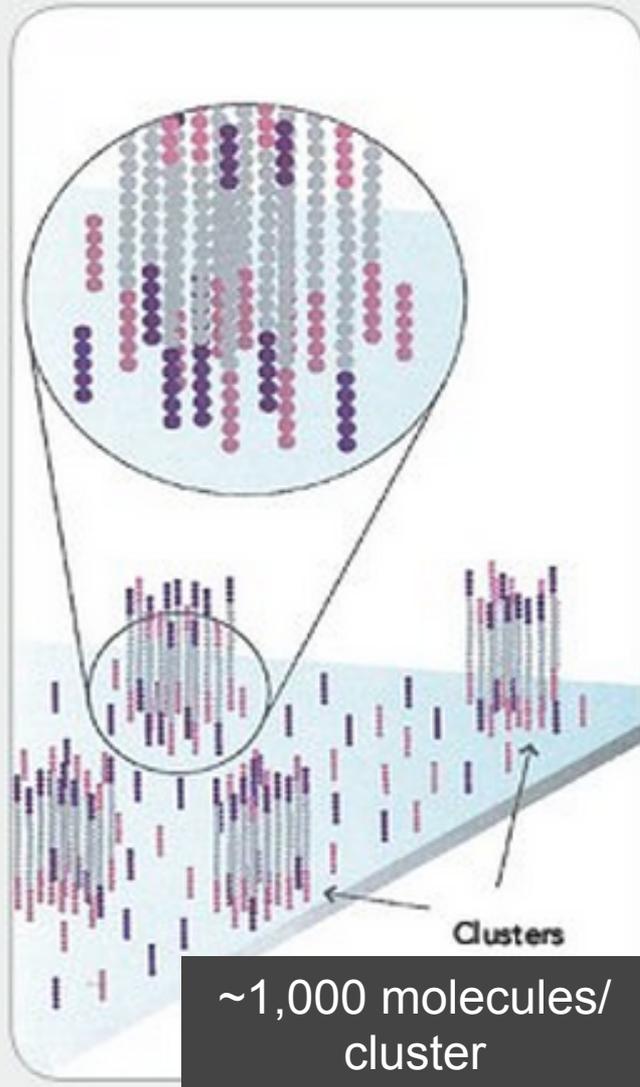


Source: National Human Genome Research Institute

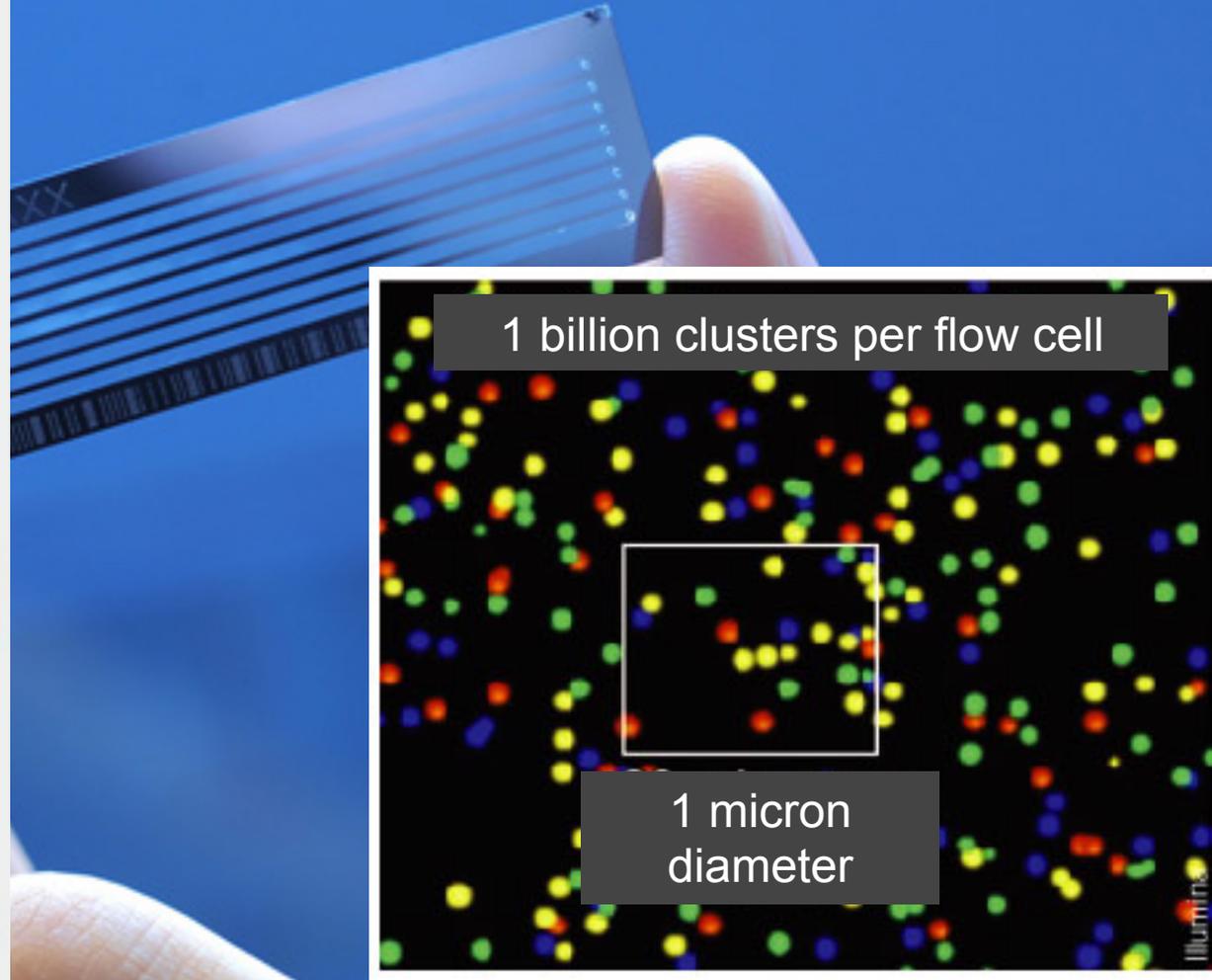
# Massively parallel sequencing



# Massively parallel sequencing

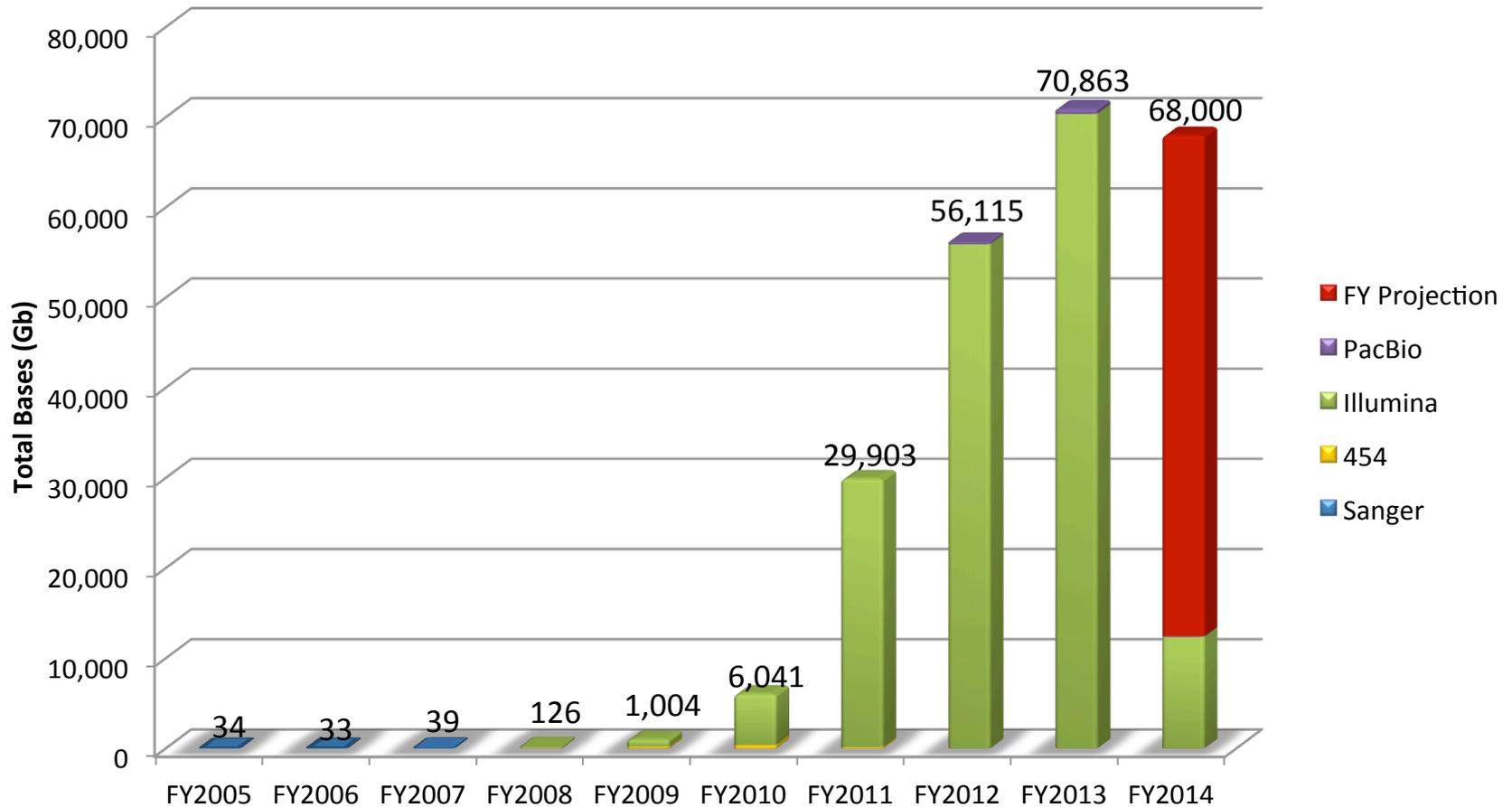


These methods produce billions of short sequences, each ~100-200 base pairs (bp) long



# Yearly JGI Sequencing Output

## FY Total Bases (Gb) Sequenced



**JGI Sequencing Output is Flat  
Reagents Costs are Increasing**

# JGI's Science Programs



## DOE Mission Areas



Bioenergy



Carbon Cycling

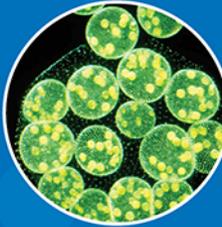


Biogeochemistry

## JGI Programs



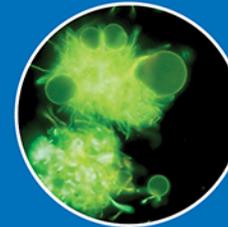
Metagenomes



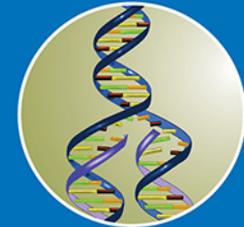
Plants



Fungi



Microbes



DNA Synthesis Science

## JGI Infrastructure



DNA Sequencing



Advanced Genomic Technologies

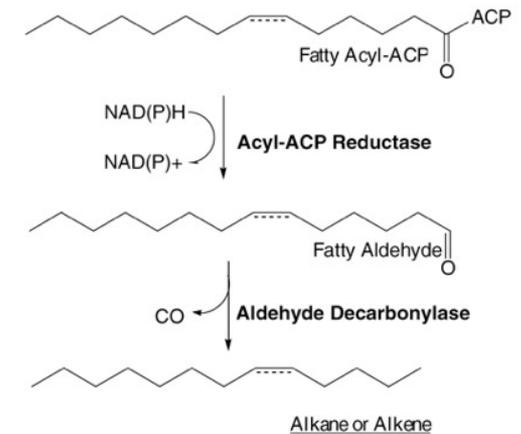
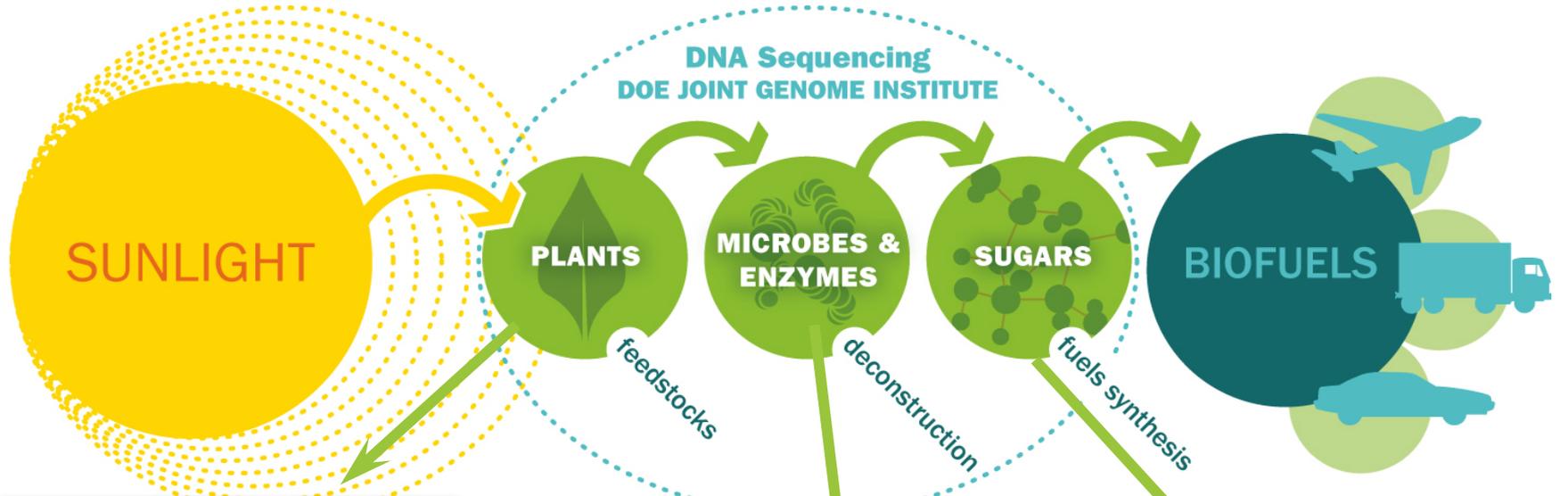


Computational Analysis



DNA Synthesis

# Genomics & Cellulosic Biofuels

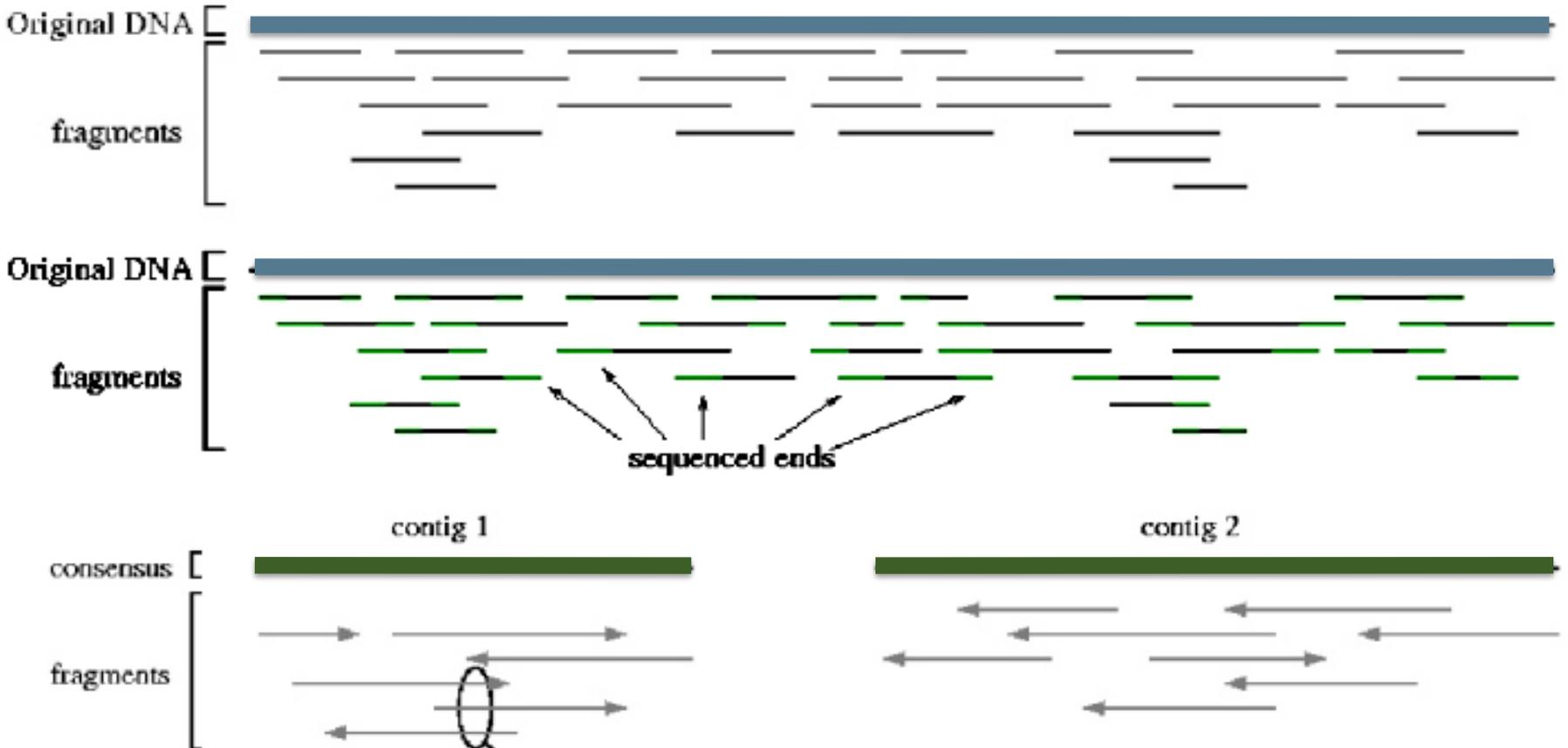




- **JGI:** Jarrod Chapman, Isaac Ho, Eugene Goltsman, Martin Mascher, Dan Rokhsar
- **NERSC/Berkeley/UCSB:** Evangelos Georganas, Veronika Strnadova, Aydin Buluc, Lenny Olikar, Joey Gonzales, John Gilbert, Stefanie Jegelka, Kathy Yelick
- **The “assembly” problem:**
  - We want the DNA sequence of chromosomes
  - Each chromosome is a single DNA molecule that can be represented as a string of the chemical “letters” A,G,C, and T
  - Chromosomes can be hundreds of thousands of “letters” (base pairs, nucleotides) long
- **But the data produced by modern sequencing instruments are billions of short (~hundred letters), redundant sequence fragments (“reads”)**

# “Shotgun” sequencing: break up genome into short pieces, determine sequence, and reassemble.

“DEPTH”



Aligned raw sequences produce “consensus”

```
AAAACTCGCCTGCTTATCAACCGATCCCCGCTACCTTCTACAGCCATCATT  
AAAACTCGCCTGCTTATCAACCGATCCCCGCTACCTTCTACAGCCATCATT  
AAAACTCGCCTGCTTATCAACCGATCCCCGCTACCTTCTACAGCCATCATT
```

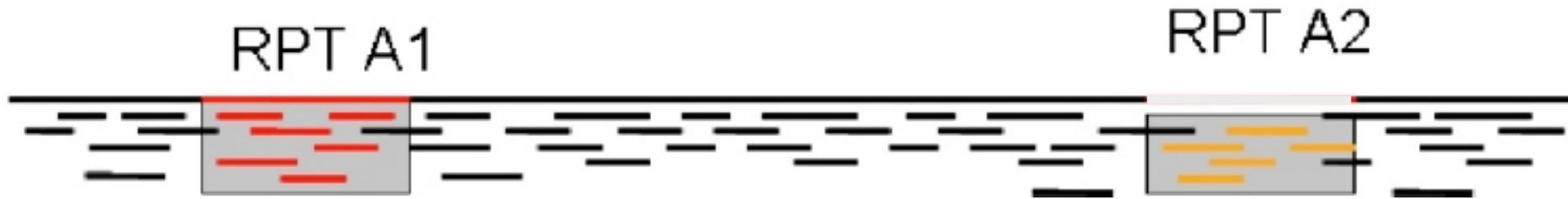


Figure 5. Two copies of a repeat along a genome. The reads colored in red and those colored in yellow appear identical to the assembly program.

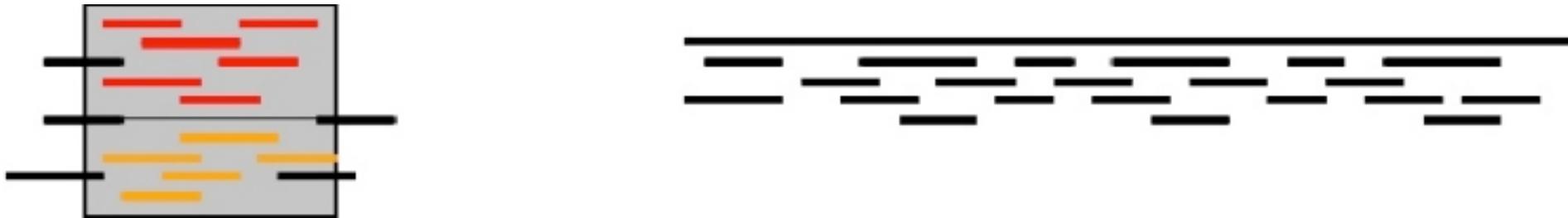
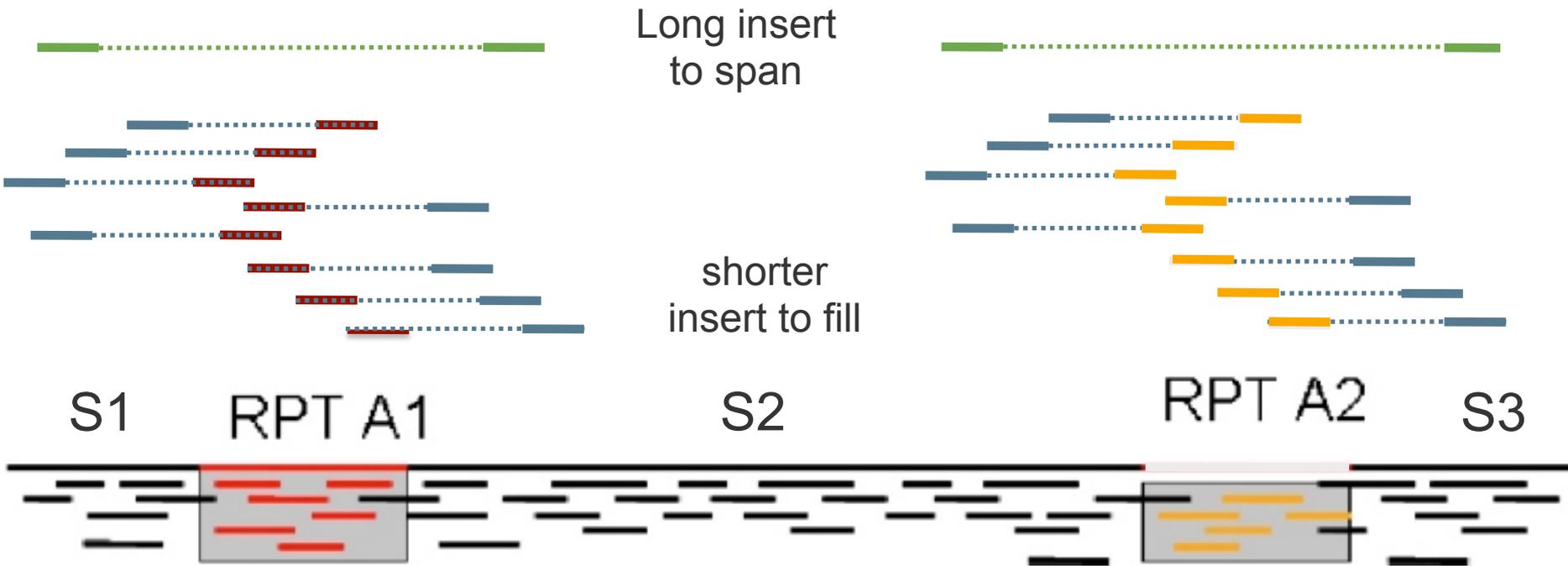


Figure 6. Genome mis-assembled due to a repeat. The assembly program incorrectly combined the reads from the two copies of the repeat leading to the creation of two separate contigs

# Paired-ends allow some shorter repetitive sequences to be skipped over and then back-filled

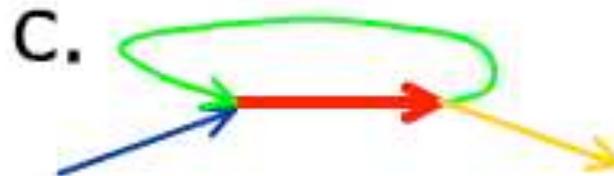
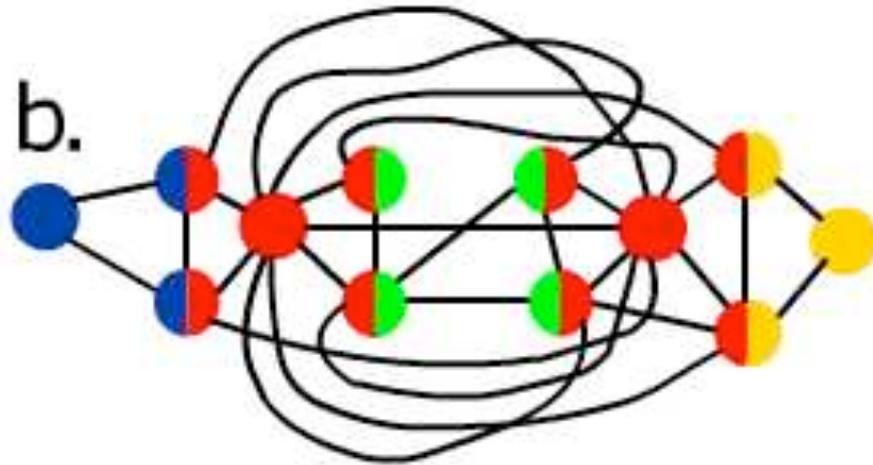
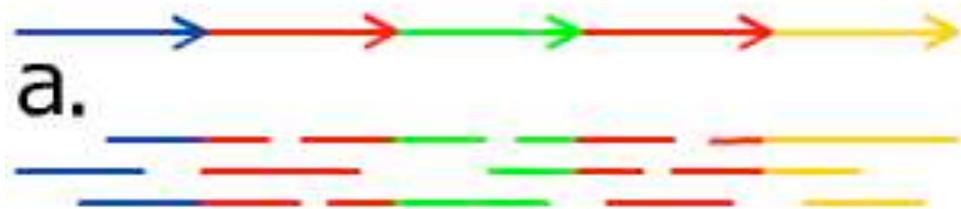


Also can use repeat-boundary spanning reads to define edges

Can represent unresolved assembly as a graph  
 $S1 > A > S2 > A > S3$

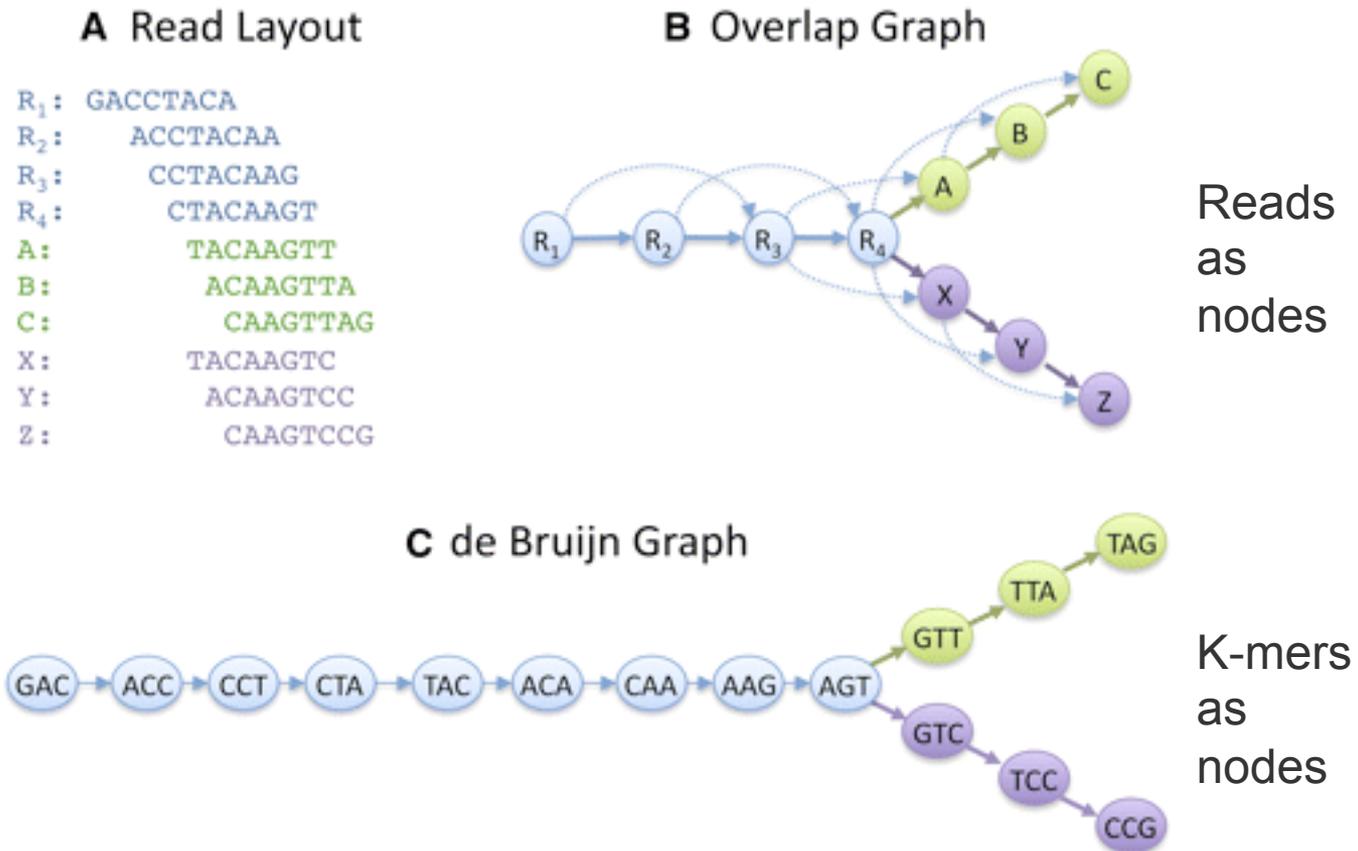
# Complex repetitive sequences

Can sometimes resolve path using flow constraints (repeats are higher apparent coverage) and paired-ends



# DeBruijn graph approach (Waterman, Idury, Pevzner, Myers, et al. )

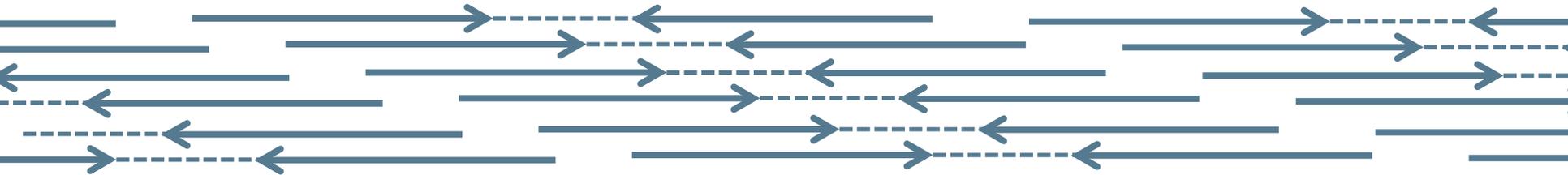
- Represent reads by overlapping k-mers
- Express assembly in terms of connectivity of k-mer graph



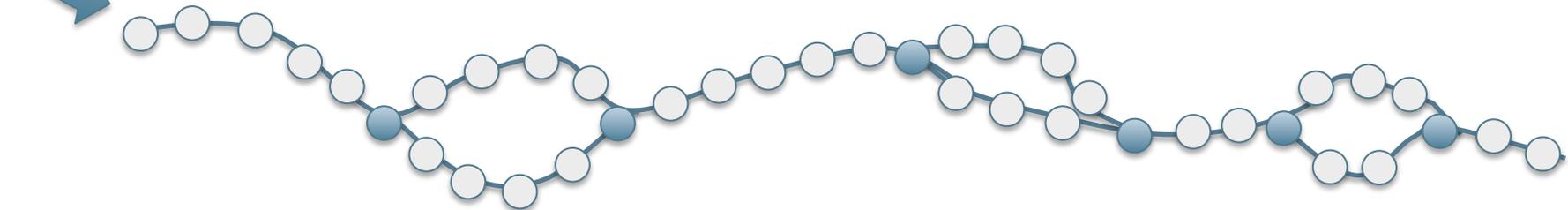
Schatz et al. (2010) Perspective: Assembly of Large Genomes w/2nd-Gen Seq. Genome Res.

# “mer-aculous” genome assembly

“shotgun sequencing” of a diploid genome



Represent each read by a set of overlapping words (**k-mers**) that are long enough to be mostly unique in genome, but short enough to be unlikely to contain a sequencing error



A “stringy” graph with bubbles representing polymorphisms

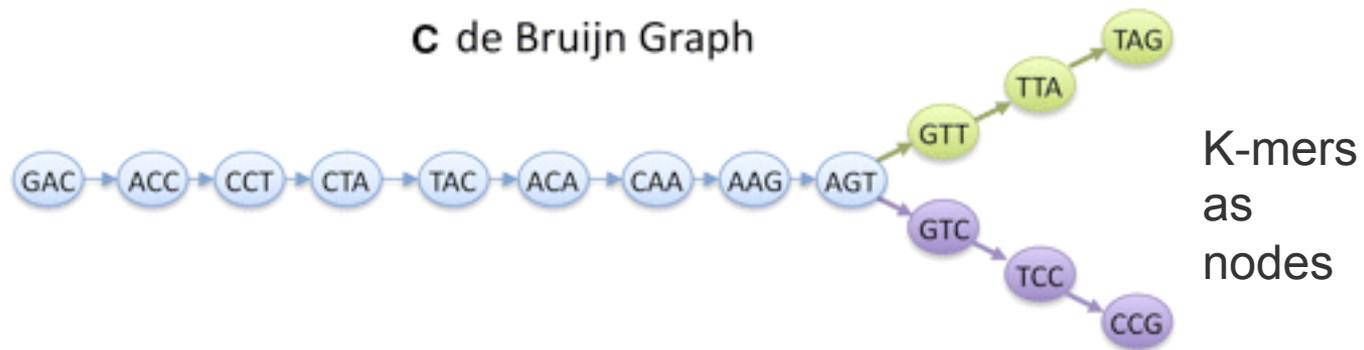
# Counting k-mers is computationally simpler than comparing reads to each other to find overlap

Counting k-mers is dependent on total data size, not depth, so it avoids all-vs-all alignments of reads.

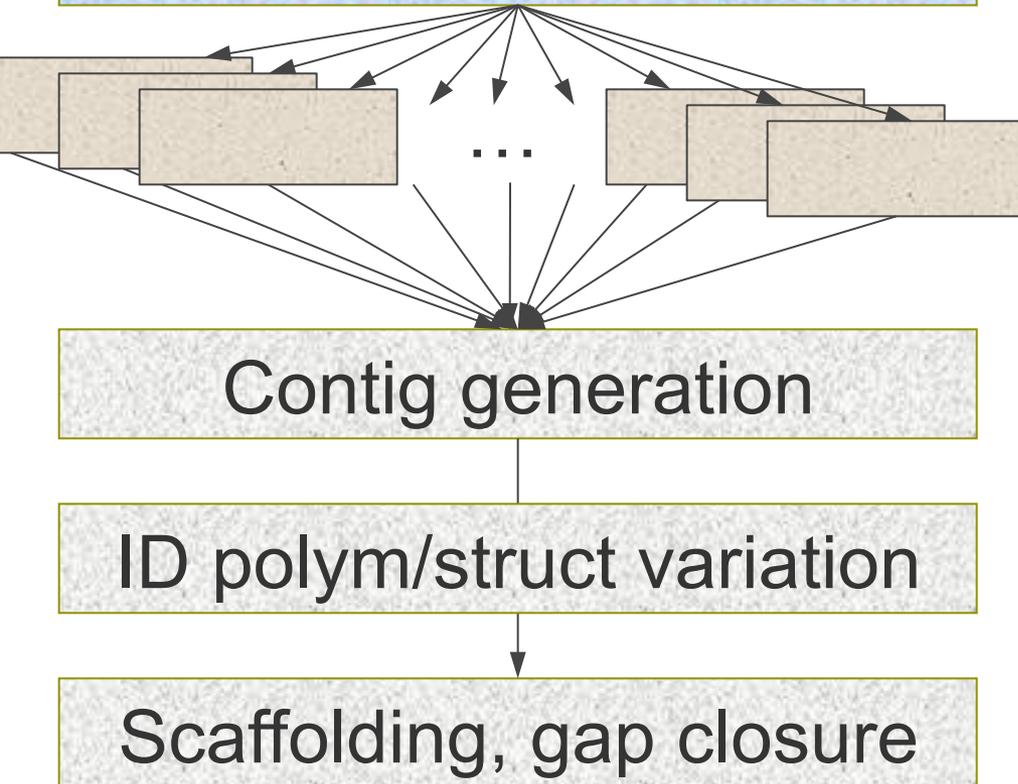
In unique genomic sequence, each (error-free) read contains the same k-mers. So data compression is achieved.

But info is lost in converting reads into strings of adjacent k-mers.

This info must be recovered elsewhere in the algorithm. These are longer-range (beyond nearest neighbor) connections in the graph.



Input short-read dataset  
**>10<sup>9</sup> Illumina reads**



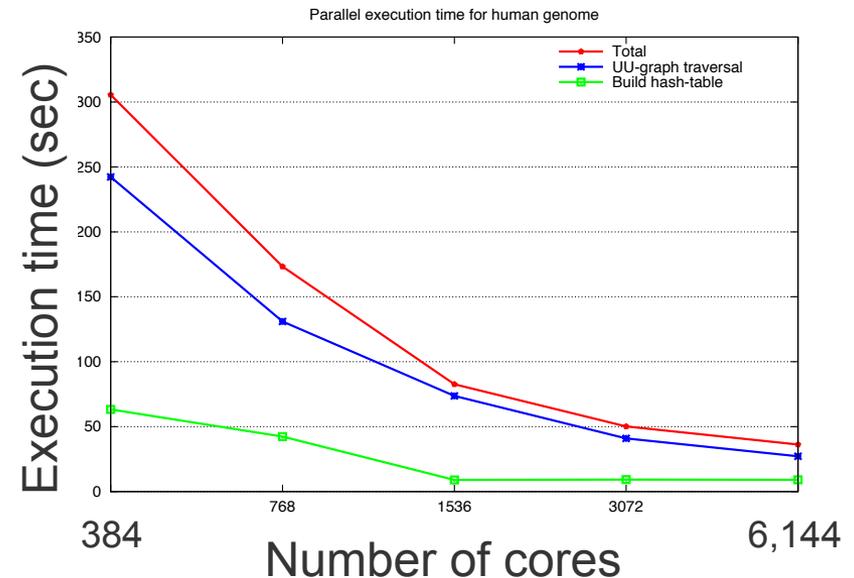
- Parallelized calculation of mer-graph, traversal, etc.
- Load balancing depend on genome, data quality, and uniformity.
- led Assemblathon I and II in several key categories.
- Ongoing collaboration with NERSC to produce distributed parallel version using UPC/PGAS.

Special features of metagenomes: sample-specific depth and variation profiles. **Must be learned from the data.**

# Towards a highly parallelized mer-acious assembler

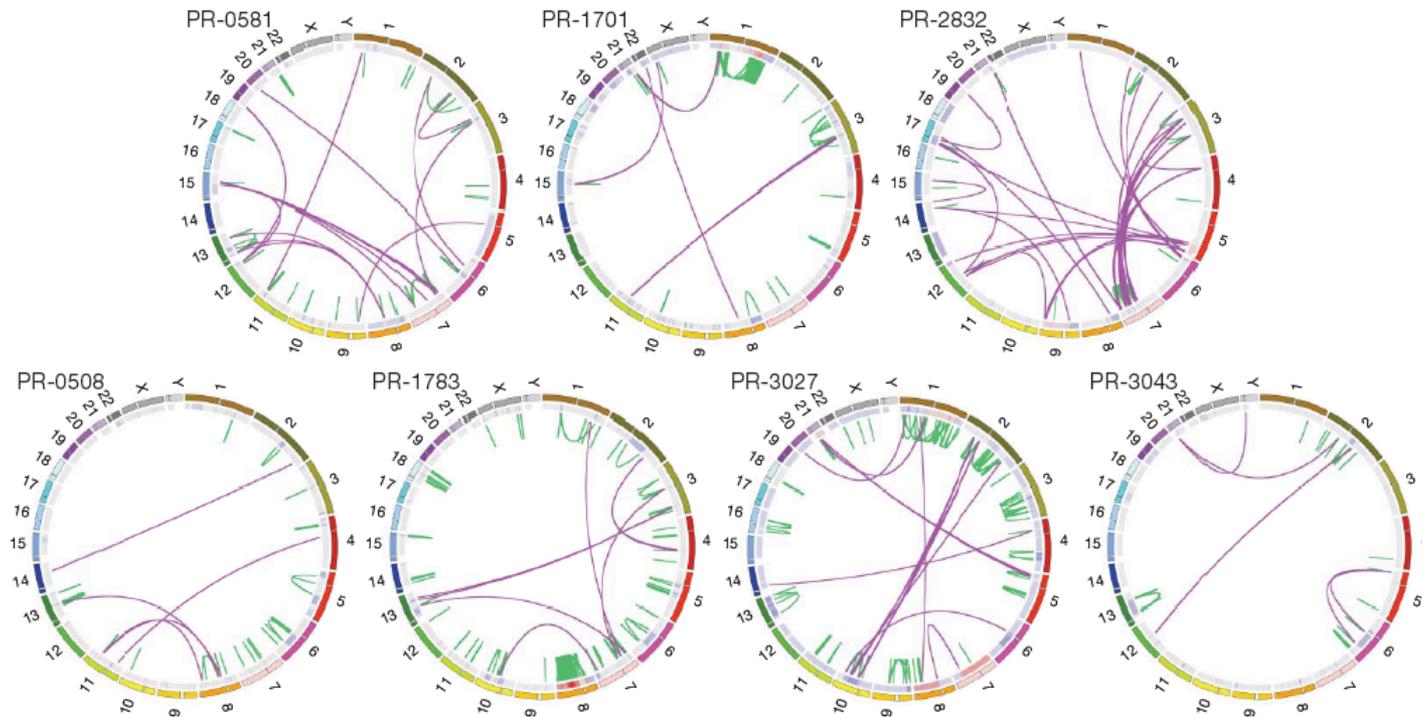
— Jarrod Chapman, Evangelos Georganas, Aydin Buluc, Kathy Yelick, Dan Rokhsar

- **Two critical steps in Meracious have been translated to UPC (Unified Parallel C)**
- **Realized thousand-fold speedup!**
- **Allows access to arbitrary memory footprint for large assembly problems**
- Remaining steps are also being translated to UPC.



# Genomic rearrangements in prostate cancers

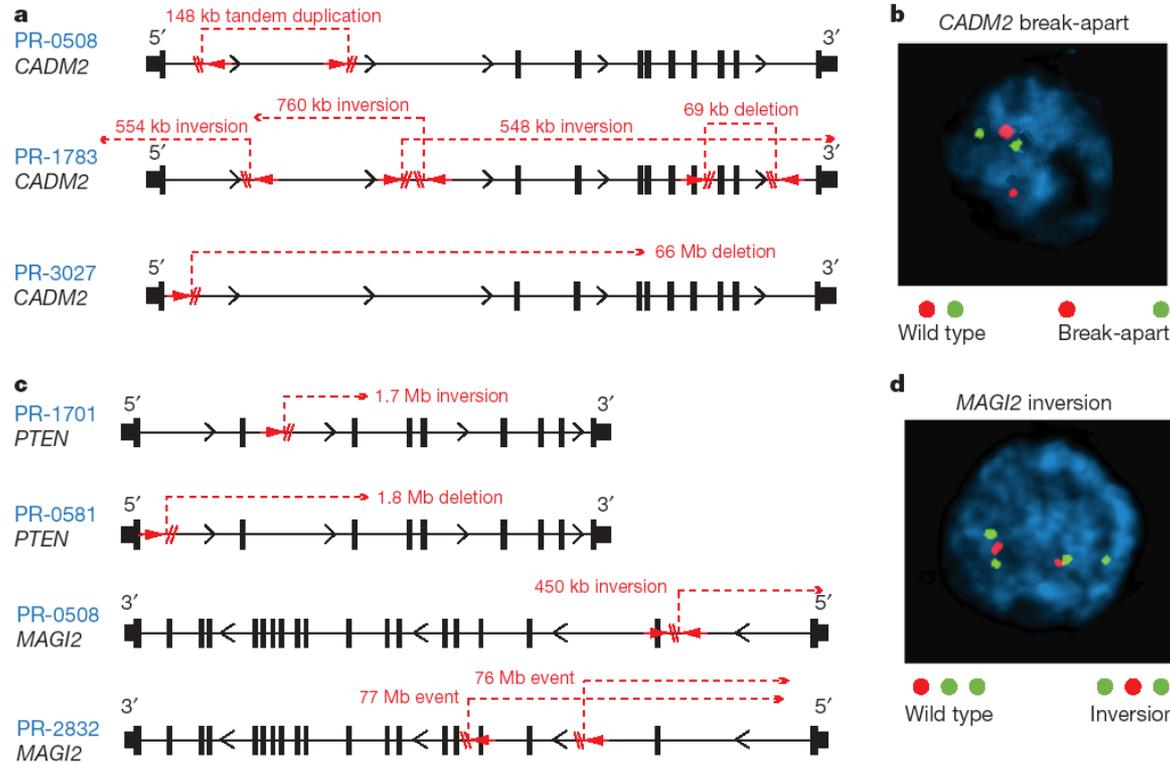
With souped-up UPC meraculous, we'll be able to assemble complex genomes like human genomes in minutes ....



**Figure 1 | Graphical representation of seven prostate cancer genomes.** Each Circos plot<sup>12</sup> depicts the genomic location in the outer ring and chromosomal copy number in the inner ring (red, copy gain; blue, copy loss). Interchromosomal translocations and intrachromosomal rearrangements are

shown in purple and green, respectively. Genomes are organized according to the presence (top row) or absence (bottom row) of the *Tmprss2-ERG* gene fusion.

# Some rearrangements break up genes. Some of these recur independently across prostate cancers.



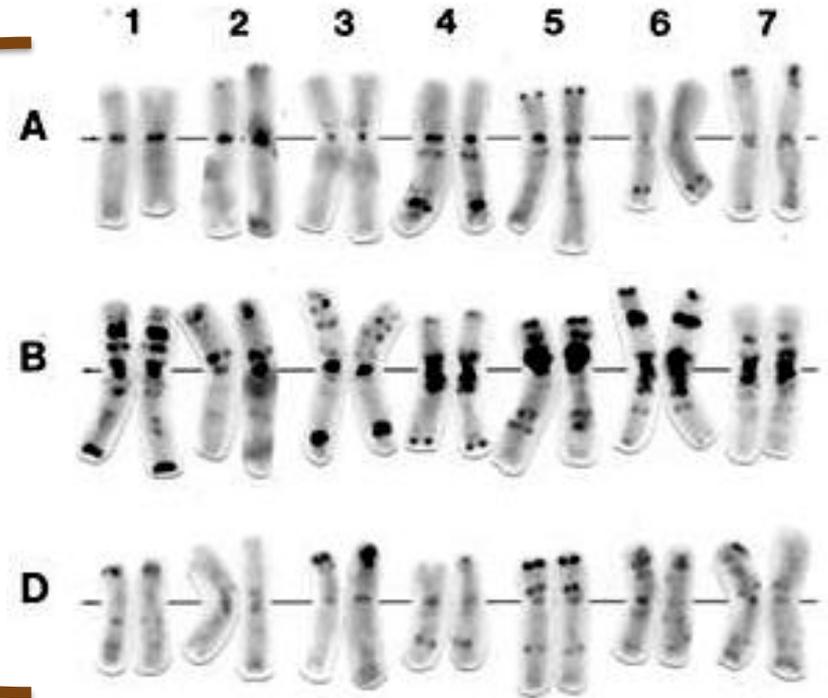
**Figure 4 | Disruption of *CADM2* and the *PTEN* pathway by rearrangements.** a, Location of intragenic breakpoints in *CADM2*. b, *CADM2* break-apart demonstrated by FISH in an independent prostate tumour. c, Location of intragenic breakpoints in *PTEN* (top) and *MAGI2* (bottom).

d, *MAGI2* inversion demonstrated by FISH in an independent prostate tumour, using probes flanking *MAGI2* (red and green) and an external reference probe also on chromosome 7q (green). The probes and strategy for detecting novel rearrangements by FISH are shown in diagram form in Supplementary Fig. 8.

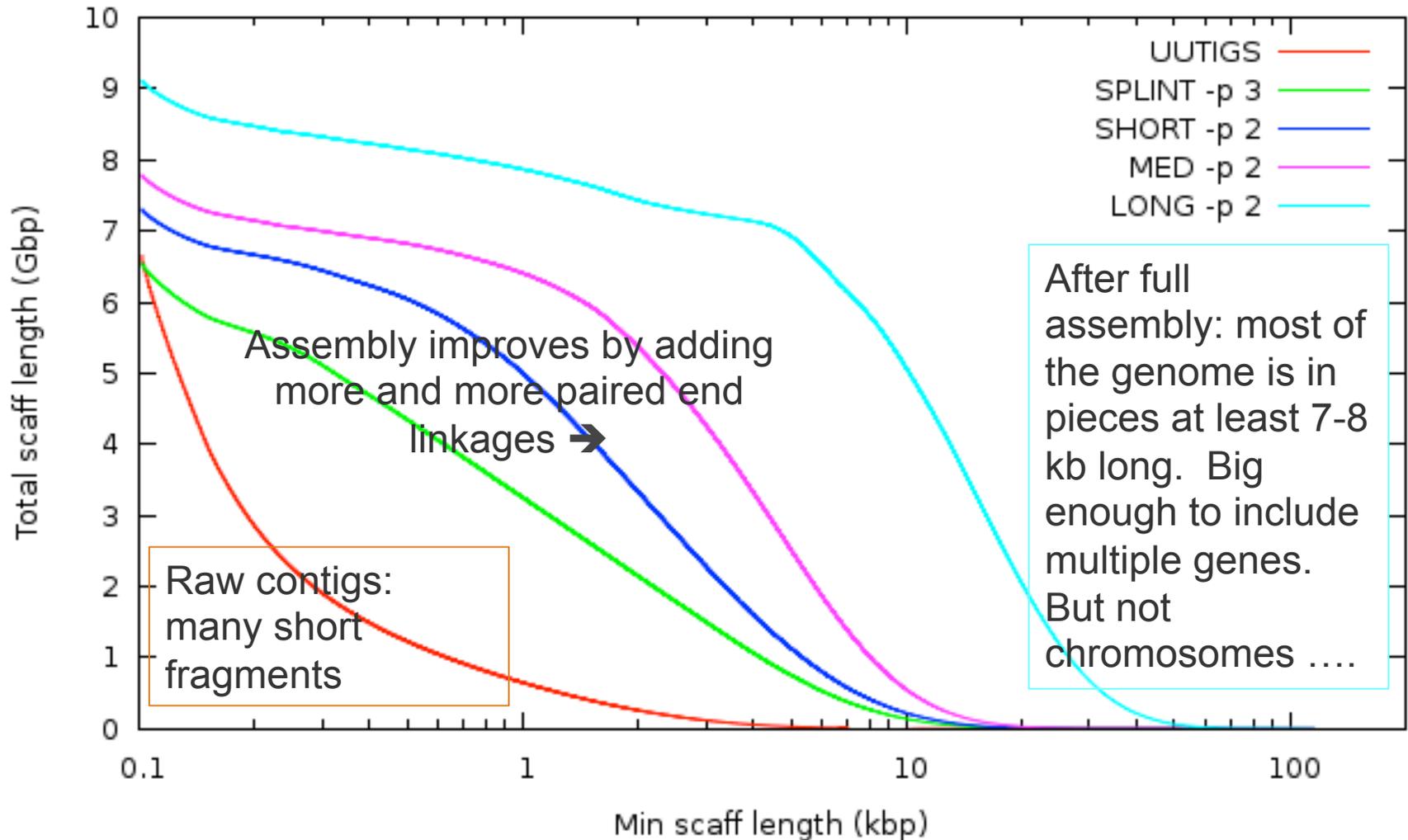
**Can we access these kinds of structural changes through a "diff" that acts on a pair of mer-graphs?**

# Assembly humongous genomes

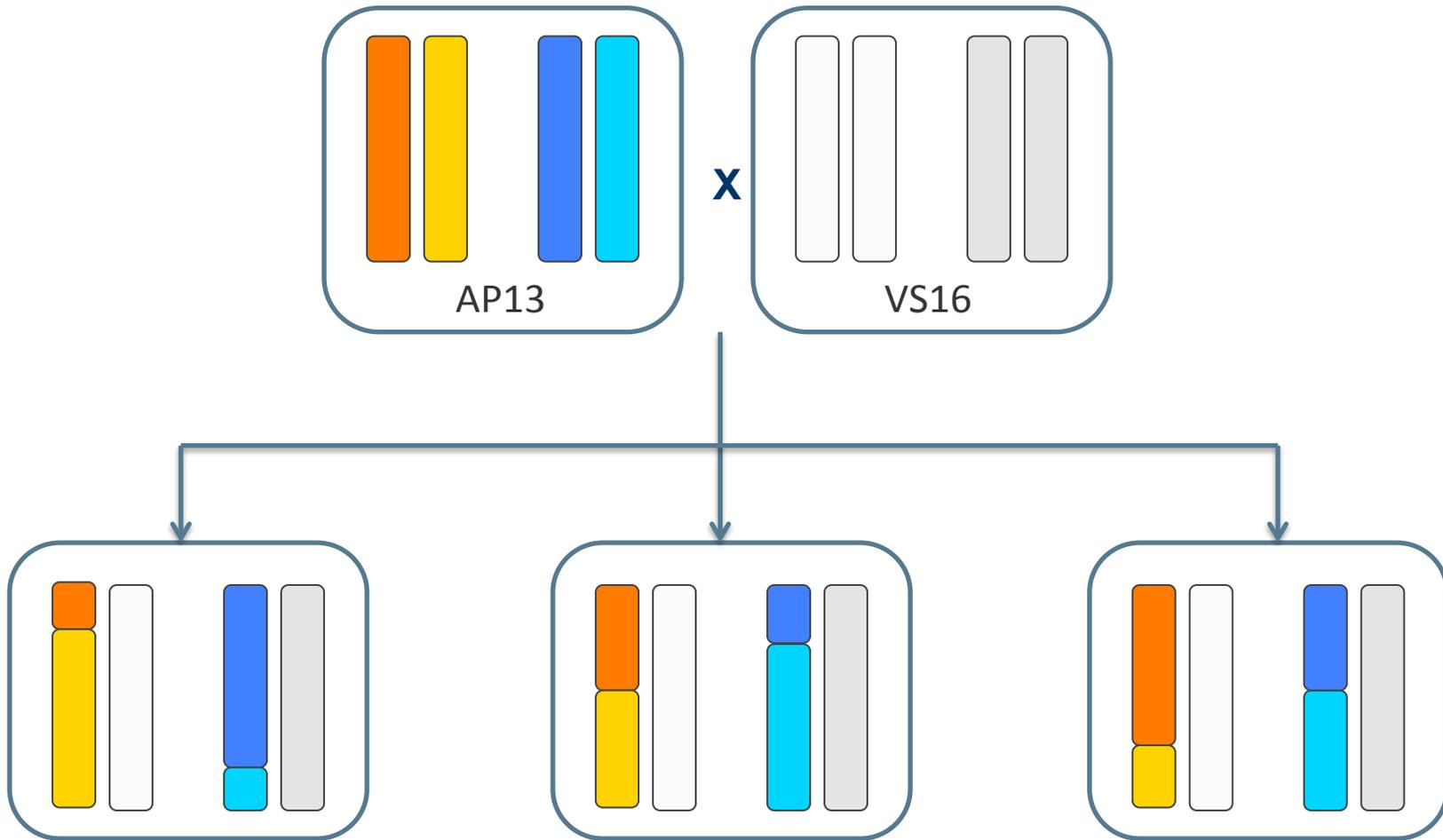
- Human genome: 3 giga-base pairs
- Maize: 2.4 Gbp
- Switchgrass: 1.4 Gb
- Miscanthus: 2.5 Gb
- ...
- Barley genome: 7 Gbp
- **Wheat genome: 17 Gbp**
- Pine genome: 20 Gbp
- Salamander: 20-30 Gbp



# Progress towards assembling wheat genes (hybrid C/UPC version)



# Genetic mapping with millions of markers

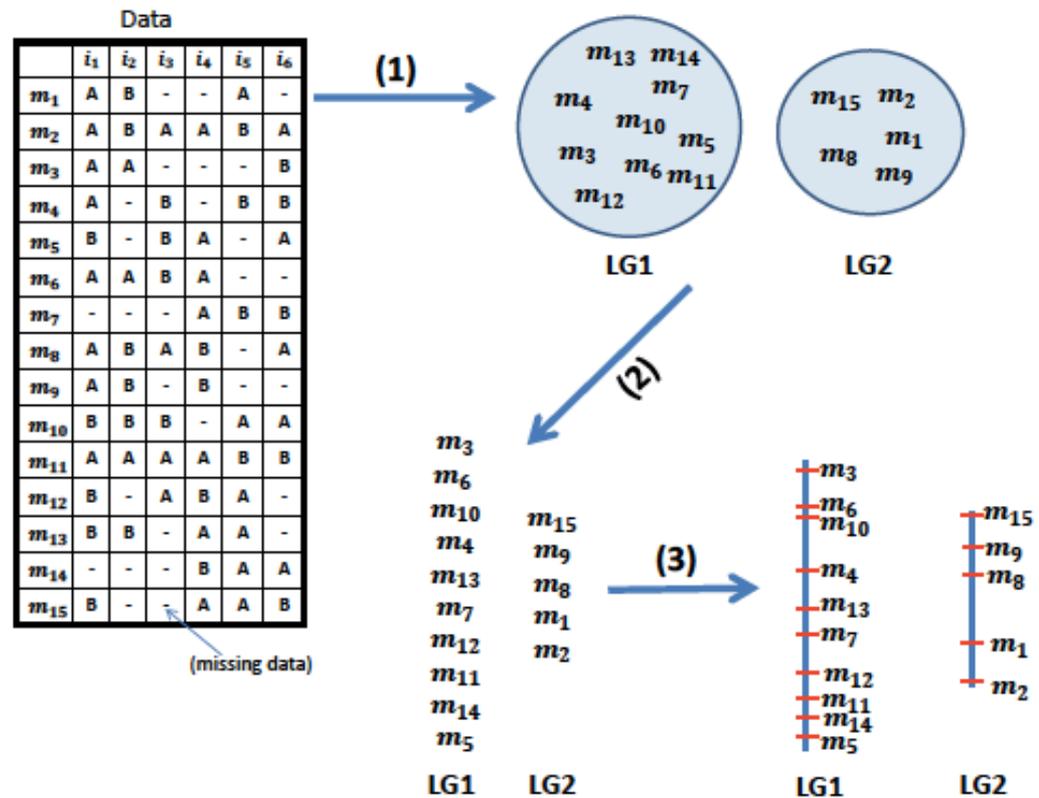


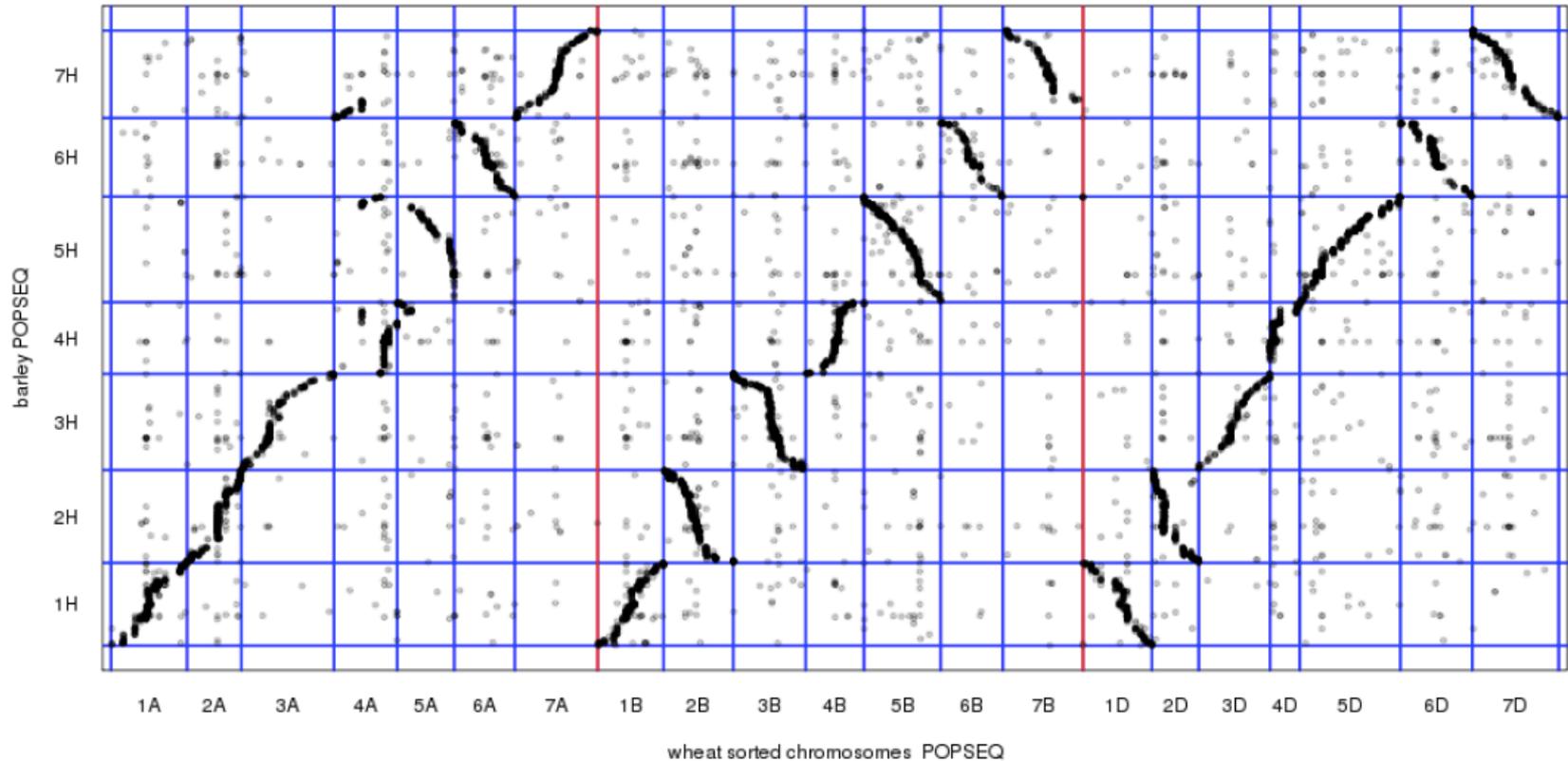
F1 recombinants track “orange” vs “yellow” in offspring

# Efficient and Accurate Clustering for Large-Scale Genetic Mapping \*

Veronika Strnadova<sup>†</sup>    Aydın Buluç<sup>‡</sup>    Jarrod Chapman<sup>§</sup>    John R. Gilbert<sup>¶</sup>  
 Joseph Gonzalez<sup>||</sup>    Stefanie Jegelka<sup>\*\*</sup>    Daniel Rokhsar<sup>††</sup>    Leonid Olikier<sup>‡‡</sup>

- Given a collection of “markers” (short sequences that can be either “yellow” or “orange”)
- Compare markers with every other marker, across all family members.
- Look for pairs of markers such that when one is yellow, the other is almost always yellow; when one is orange, the other is almost always orange.
  - “Color” is unknown at the start. Like a (trivial) Ising gauge factor.
  - Amount of mis-correlation is related to distance along the chromosome.
- Need clever ways to organize the calculation, taking advantage of the inherent linearity of chromosomes.
- Can handle 10’s of millions of markers, far more than other genetic mapping codes.





Within a week or two we will have a de novo assembly and map of the wheat genome, after the algorithmic dust settles!

# “Annotating” DNA sequence



- **JGI:** Amrita Pati, Marcel Huntemann, Nikos Kyrpides
- **NERSC:** Seung-Jin Sul, Kjersten Fagnan, Shane Canon
  
- **The “annotation” problem:**
  - When sequencing a “metagenome” from a microbial community, we may sample a billion gene fragments, derived from the constituent microbes.
  - Q. How do we computationally infer the function of these fragments?
  - A. By grouping related gene sequences together, and detecting similarity with genes of known function.

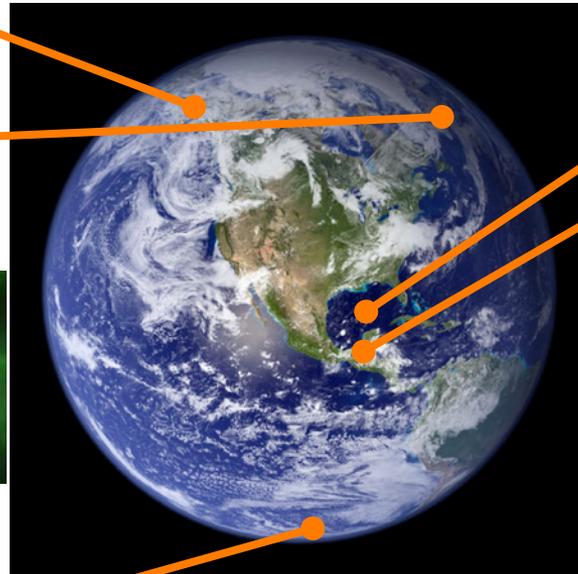
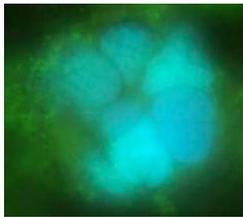
# Microbial communities (“meta” genomes)

Metagenomic analysis of a permafrost microbial community reveals a rapid response to thaw. *Nature*



Metagenome, metatranscriptome and single-cell sequencing reveal microbial response to Deepwater Horizon oil spill. *ISME J*

The metagenome of a marine anammox bacterium illustrates role in the global nitrogen cycle. *Environmental Microbiology*



Global transcriptome response to ionic liquid by a tropical rain forest soil bacterium, *Enterobacter lignolyticus*. *PNAS*



The genome of the polar eukaryotic microalga *coccomyxa subellipsoidea* reveals traits of cold adaptation. *Genome Biology*

# Growth of metagenomic data



IMG/M Home

Find Genomes

Find Genes

Find Functions

**DECEMBER 2012**



The data deluge  
...is here



Comparative analysis of  
**1,000,000,000+**

Genes in IMG/M

Powered by **NERSC**



INTEGRATED MICROBIAL GENOMES  
with MICROBIOME SAMPLES

**February 2014**

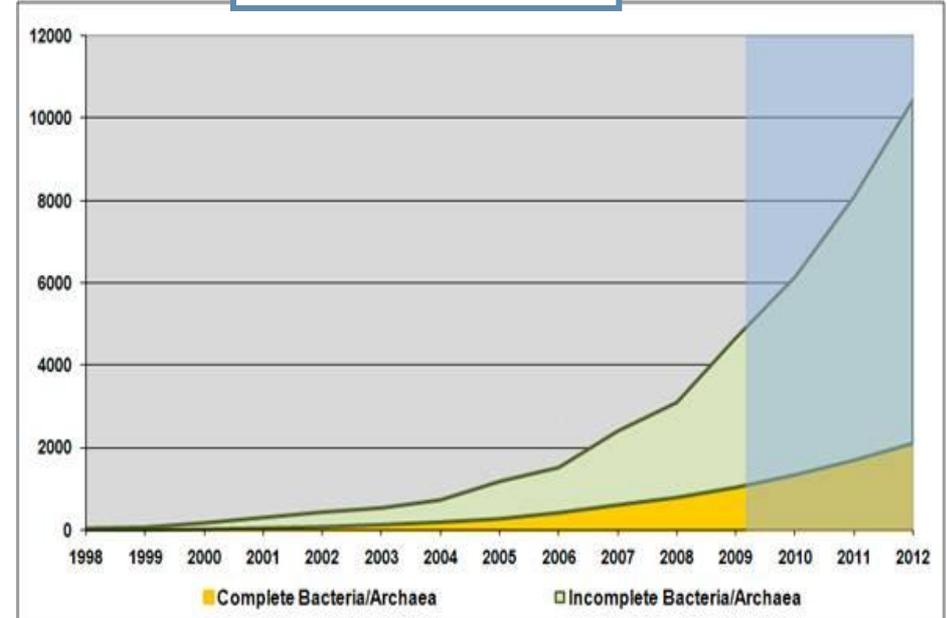
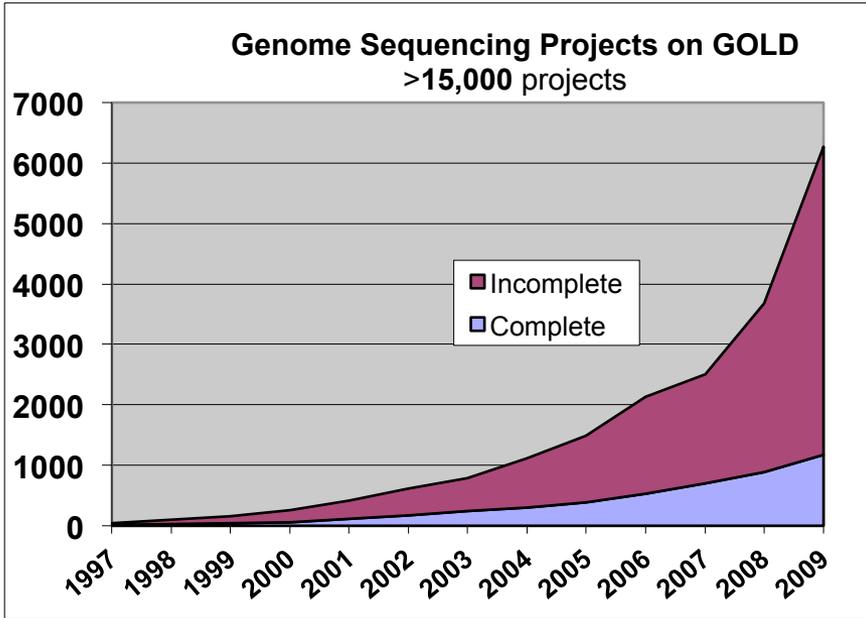
Samples	<b>3,854</b>
DNA (bps)	<b>3,9 Tb</b>
Genes	<b>23.4 Billion</b>

**CLUSTER CAPACITY**

minimum	<b>250 M genes /wk</b>
maximum	<b>1.3 Billion genes /wk</b>

# Microbial Genome Projects

P. Chain *et al.* Science, 2009



	1995-2009	2010-2014
<b>Finished</b>	1000	3,000
<b>Draft</b>	1000	10,000
<b>Genes</b>	6 Million	39 Million

	now
<b>Finished</b>	2,975
<b>Draft</b>	12,379
<b>Genes</b>	52 Million

# Metagenomics analysis challenges



POWERED BY

**NERSC**

- **How do you compare 15 Billion genes (all vs all)?**

- 1. Cluster all genes from “isolate” genomes**

- 22 Million Genes

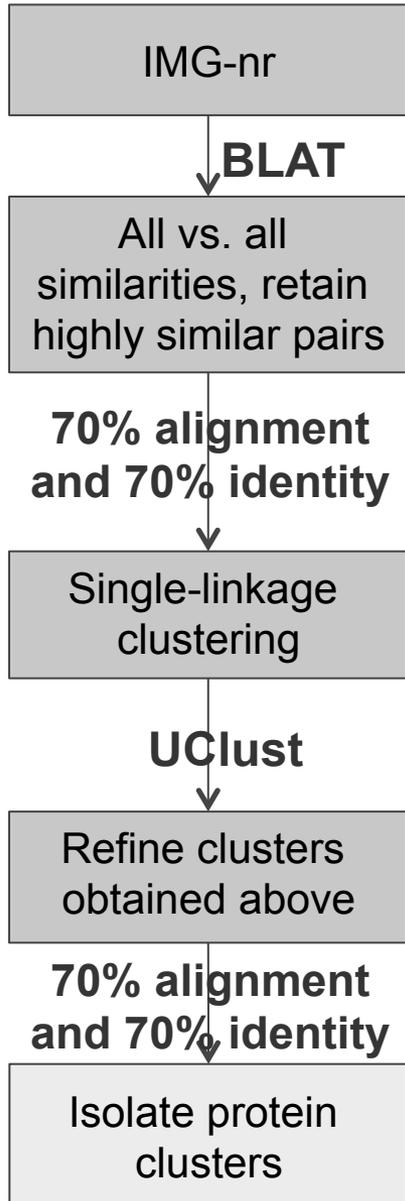
- 2. Pledge all genes of the metagenomes to the clusters of the isolate genomes**

- 15 Billion genes

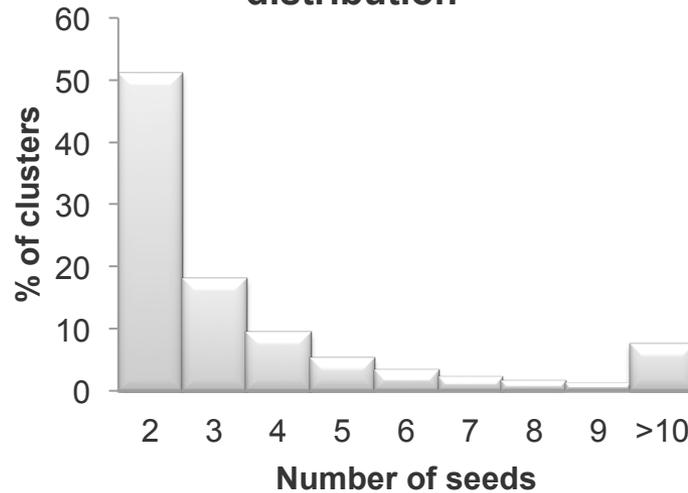
- 3. Cluster all the metagenomes genes that did not pledge and create additional (purely metagenomics) clusters.**

- These may have interesting unknown functions

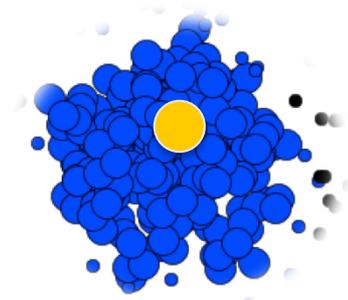
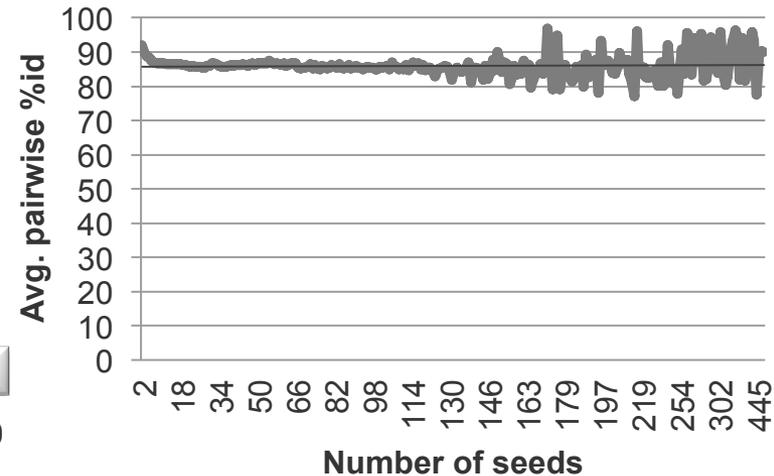
# 1. Generation of isolate protein clusters



Cluster seed set size distribution



Avg pairwise %id vs. #seeds



5448778.1  
 MGRISSGGMMFKAITTVAALVIATSAMAQDDLTISSLAKGETTKAAFNMVQGHKLP  
 KGGTYTPAQTVTLGDETYQVMSACKPHDCGSORIAVMWSEKSNQMTGLESTIDEKTS

## Isolate Genomes

Total Genes	<b>22 Million</b>
NR Genes	<b>16 Million</b>
Clusters	<b>1.8 Million</b>
Genes in clusters	<b>8 Million</b>
Singletons	<b>8 Million</b>

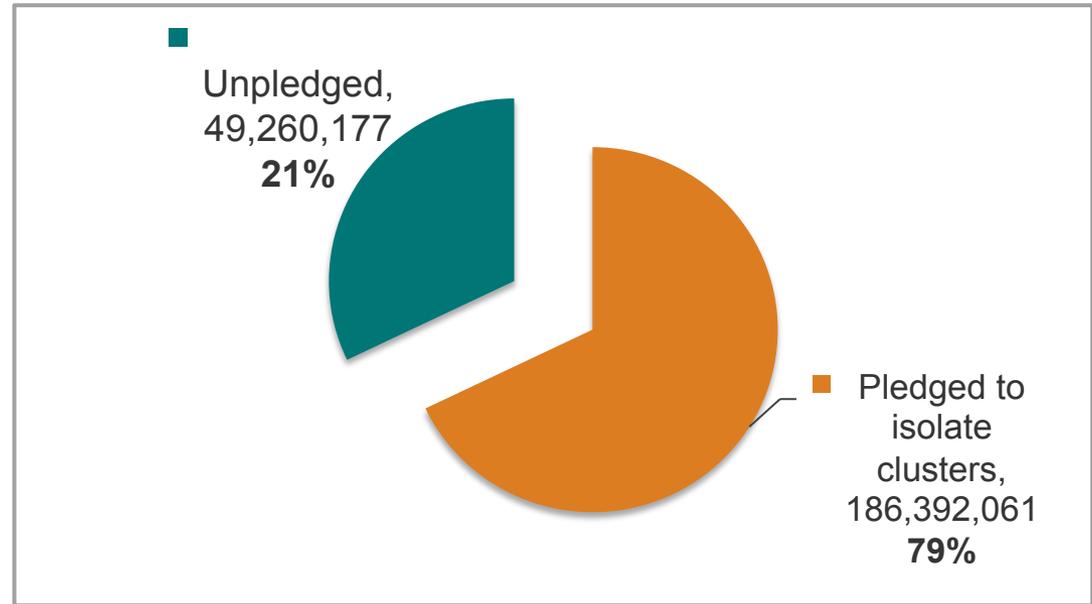
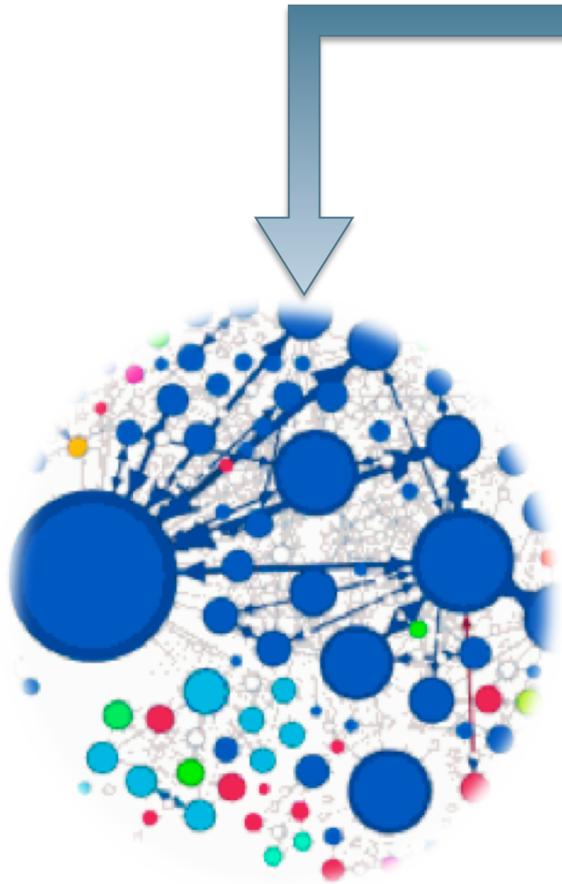
## Create new clusters from metagenomic genes that don't have any hits to gene clusters from isolates

Pledge all metagenomic genes to the clusters of isolates



Create new clusters from the unpledged metagenomic genes

# Metagenomic Clustering from assembled metagenomes

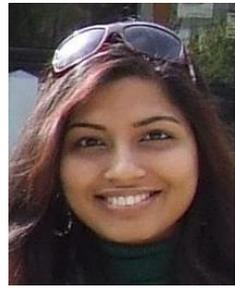


Metagenome Genes	
Total Genes	<b>49 Million</b>
Clusters	<b>1.3 Million</b>
Genes in clusters	<b>25 Million</b>
Singletons	<b>24 Million</b>

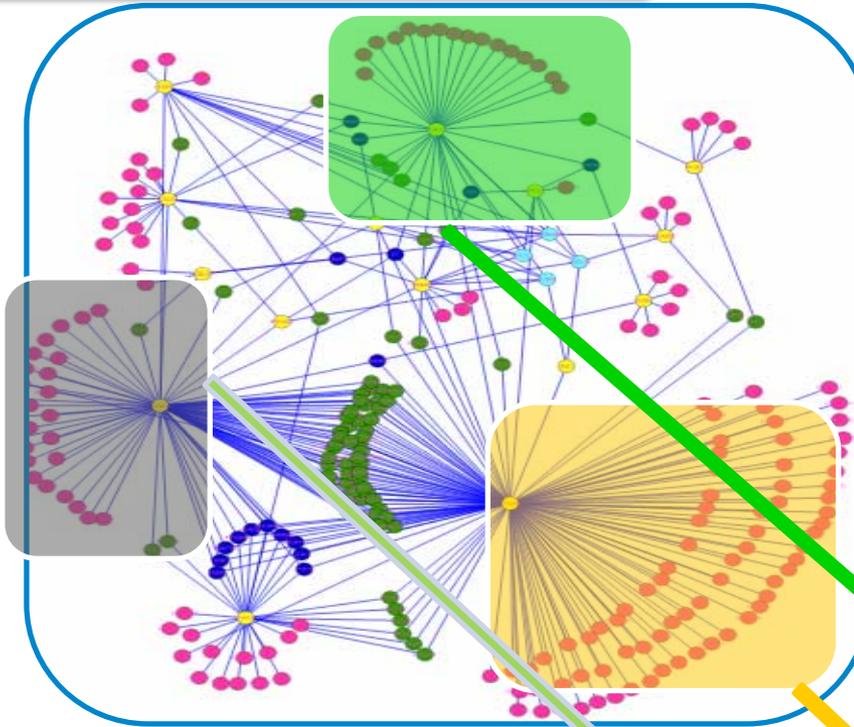
## Metagenome clusters 1.27M clusters

Within a cluster linkages have 70%id and 50% alignment length across both members

# IMG/M: 23 Billions Genes



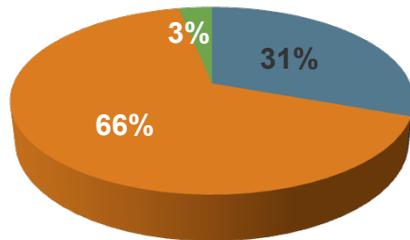
## Gene Clustering



## Metagenome Classification

- Metagenomic Samples (1930)
  - Environmental (711)
    - Aquatic (580)
      - + Non-marine Saline and Alkaline (70)
      - + Marine (264)
      - + Thermal springs (0)
      - + Freshwater (152)
    - Terrestrial (119)
      - + Rock-dwelling (subaerial biofilms) (5)
      - + Soil (113)
      - + Deep subsurface (1)
    - Air (7)
      - + Indoor Air (2)
      - + Outdoor Air (4)
    - + Host-associated (1122)
    - + Engineered (97)

Habitat distribution of metagenomes



■ Environmental ■ Host-associated ■ Engineered



Marcel Huntemann  
Amrita Pati (PI)

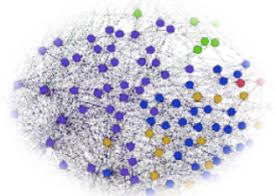
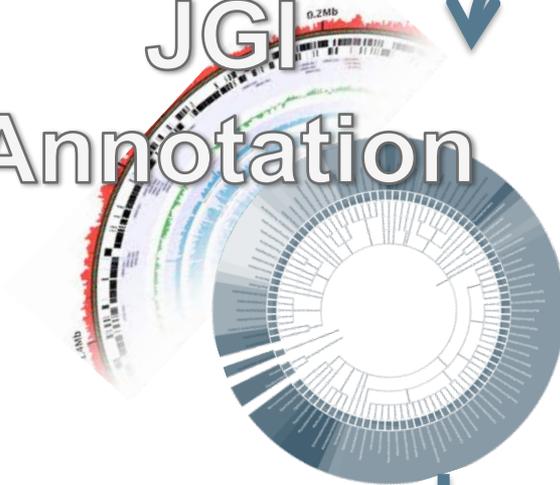
Seung-Jin Sul  
Shane Canon (PI)

- Terabytes of sequence data
- JGI engineers
- Distributed computational paradigms
- Tools for bioinformatics sequence analysis

- Massive computational systems
- Consultants to help build interfaces with custom hardware

HOPPER  
GENEPOOL  
JESUP TESTBED

**JGI**  
**Annotation**



**R&D towards generating scalable computational frameworks for big sequence data**

**Functional and phylogenetic annotation for 12 Billion metagenome genes in 12 months**

**Massive backlogs cleared**

- **The DOE Joint Genome Institute**
  - Who we are, what we do, and why
- **Two computational problems in genomics**
  - “Assembling” genomes from “shotgun” sequence
  - “Annotating” “metagenomic” data
- **We have a User Group meeting too!**



# GENOMICS OF ENERGY AND ENVIRONMENT

Meeting

March 18 - 20, 2014  
Walnut Creek, CA

Invited presentations, workshops and tutorials on sequence-based bioinformatics, and data management systems.

DNA Synthesis & Synthetic Biology

Single-Cell Genomics for Bioprospecting

Biofuel Traits in Biomass Feedstocks

HPC for Next-Gen Sequencing Applications

Functional Metagenomics



## Assembly

- **JGI:** Jarrod Chapman, Isaac Ho, Eugene Goltsman, Martin Mascher, Dan Rokhsar
- **NERSC/Berkeley/UCSB:** Evangelos Georganas, Veronika Strnadova, Aydin Buluc, Lenny Oliker, Joey Gonzales, John Gilbert, Stefanie Jegelka, Kathy Yelick

## Metagenome annotation

- **JGI:** Amrita Pati, Marcel Huntemann, Nikos Kyrpides
- **NERSC:** Seung-Jin Sul, Kjersten Fagnan, Shane Canon

**Thanks to DOE OBER, ASCR, and LBNL for support**

**Also: Thanks to NERSC for JGI computing partnership!**