

Understanding the Human Oral Microbiome with Next-Generation Sequencing Technologies

Next Generation Sequencing and Oralgen2.0

www.lanl.gov/bioscience
www.oralgen.lanl.gov

Patrick Chain




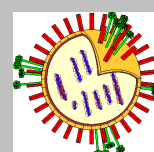



pchain@lanl.gov

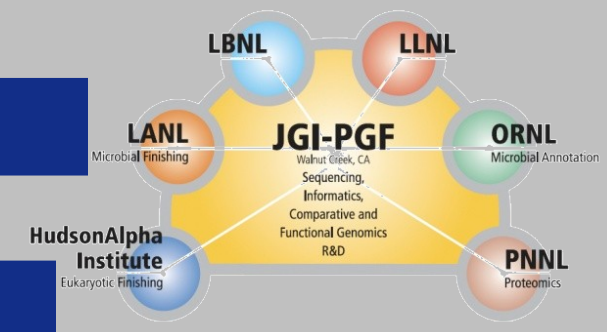
Metagenomics Applications Team, LANL
Metagenomics Program, JGI

IADR, Barcelona, Spain
July 2010

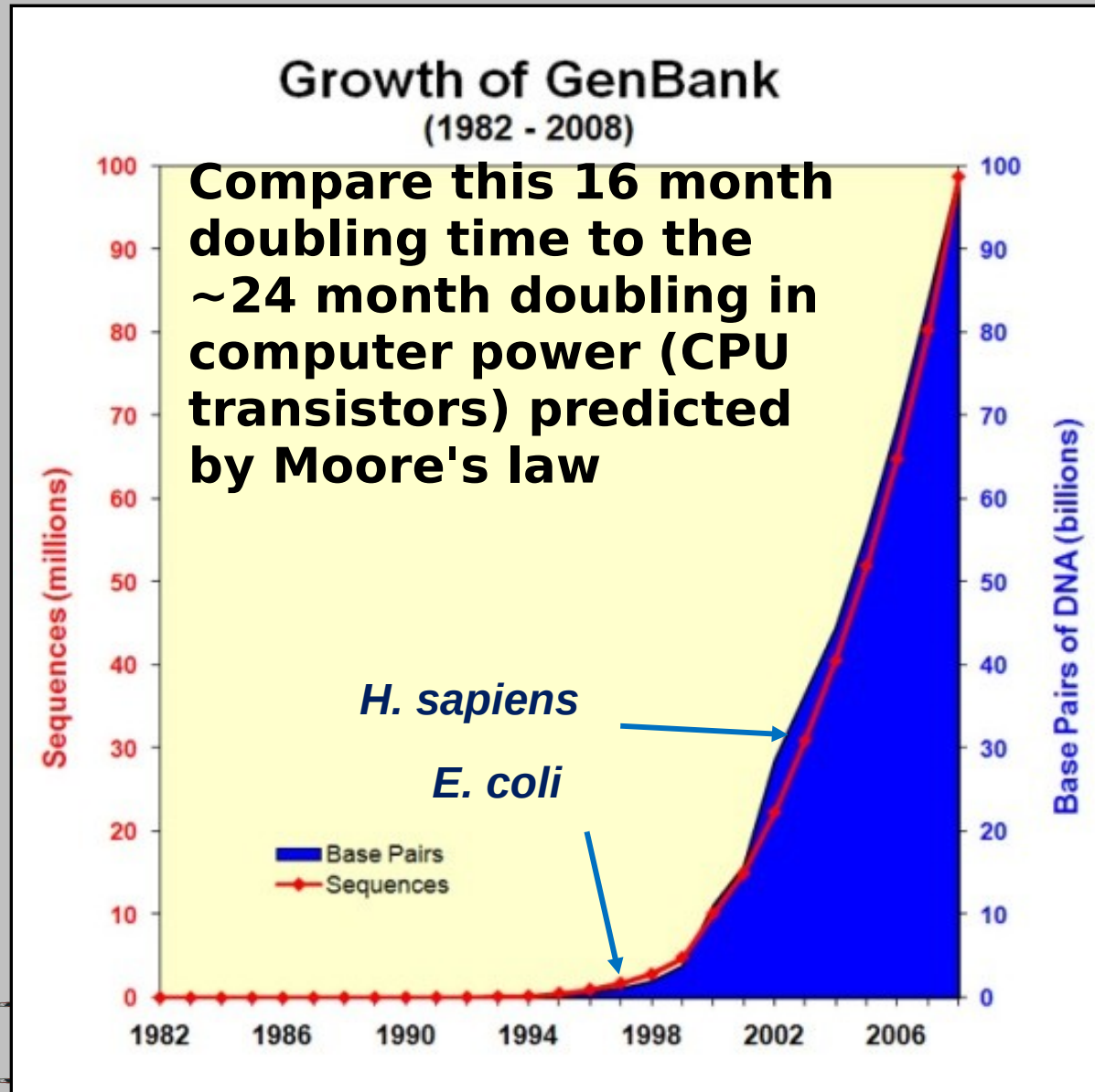


Genomics and Bioinformatics at LANL

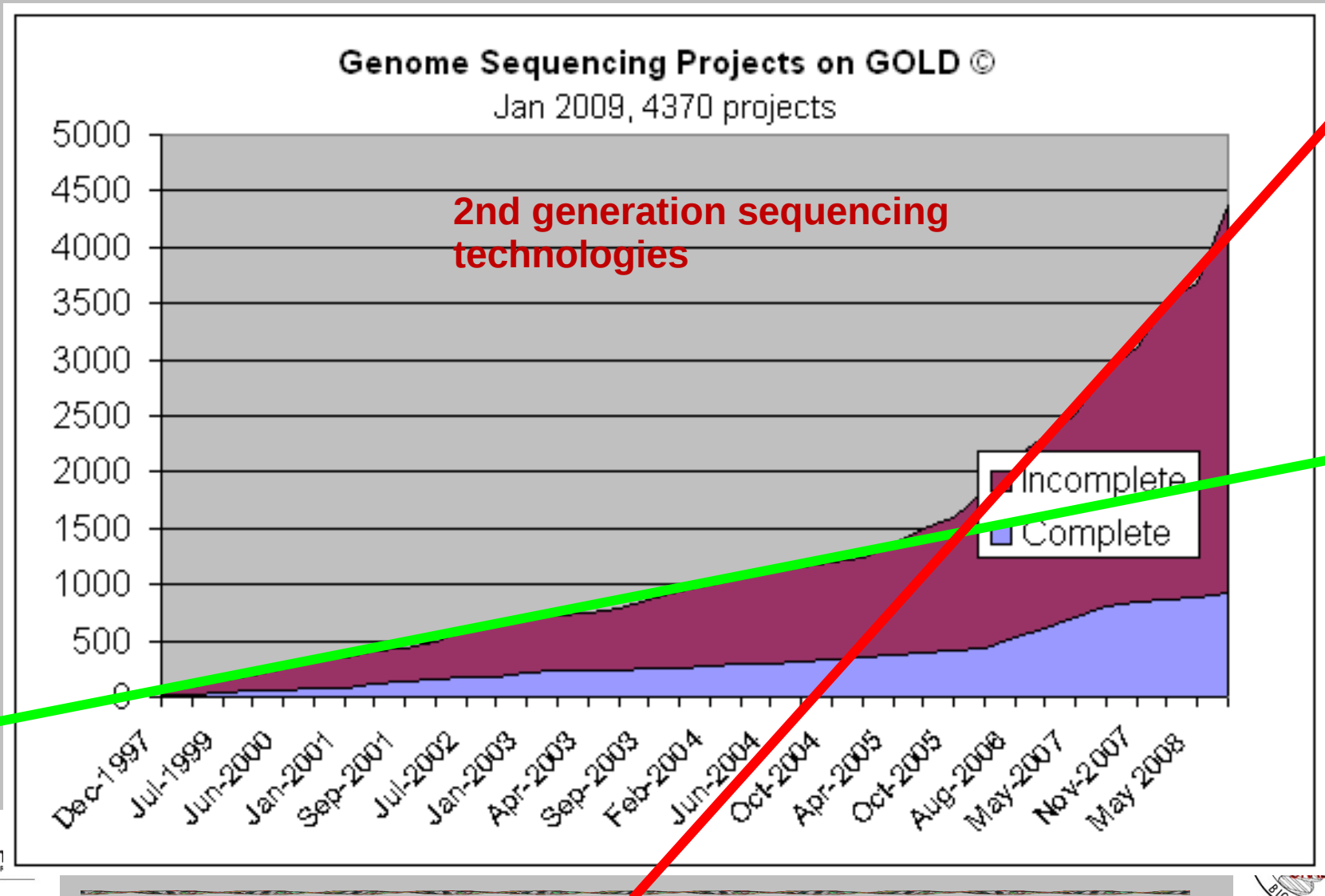
1982		GenBank
1988		The Los Alamos Center for Human Genome Studies
1990		HIV Sequence Database
1995		Influenza Sequence Database
1997		Joint Genome Institute
1998		STD Sequence Database
1999		Oral Pathogen Sequence Database



Growth of Sequence Databases

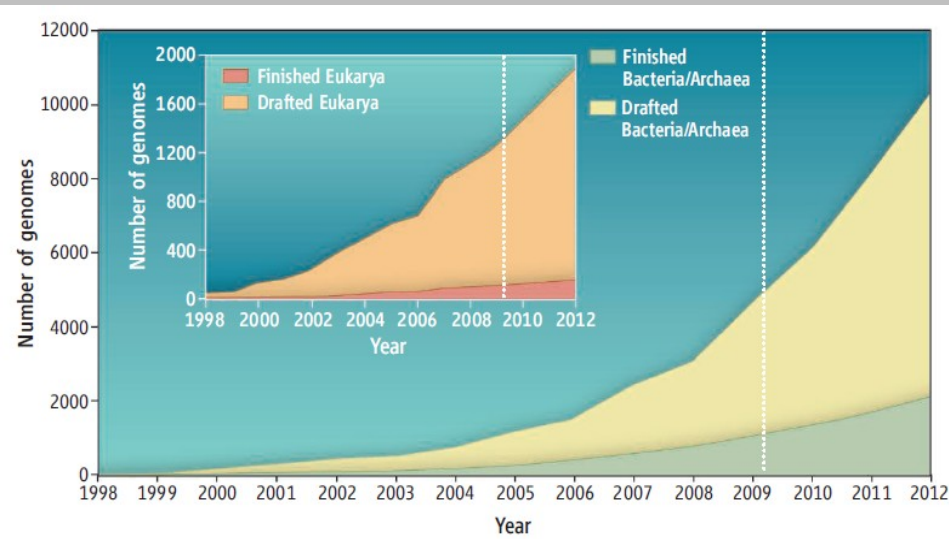


Genome projects: A growing problem



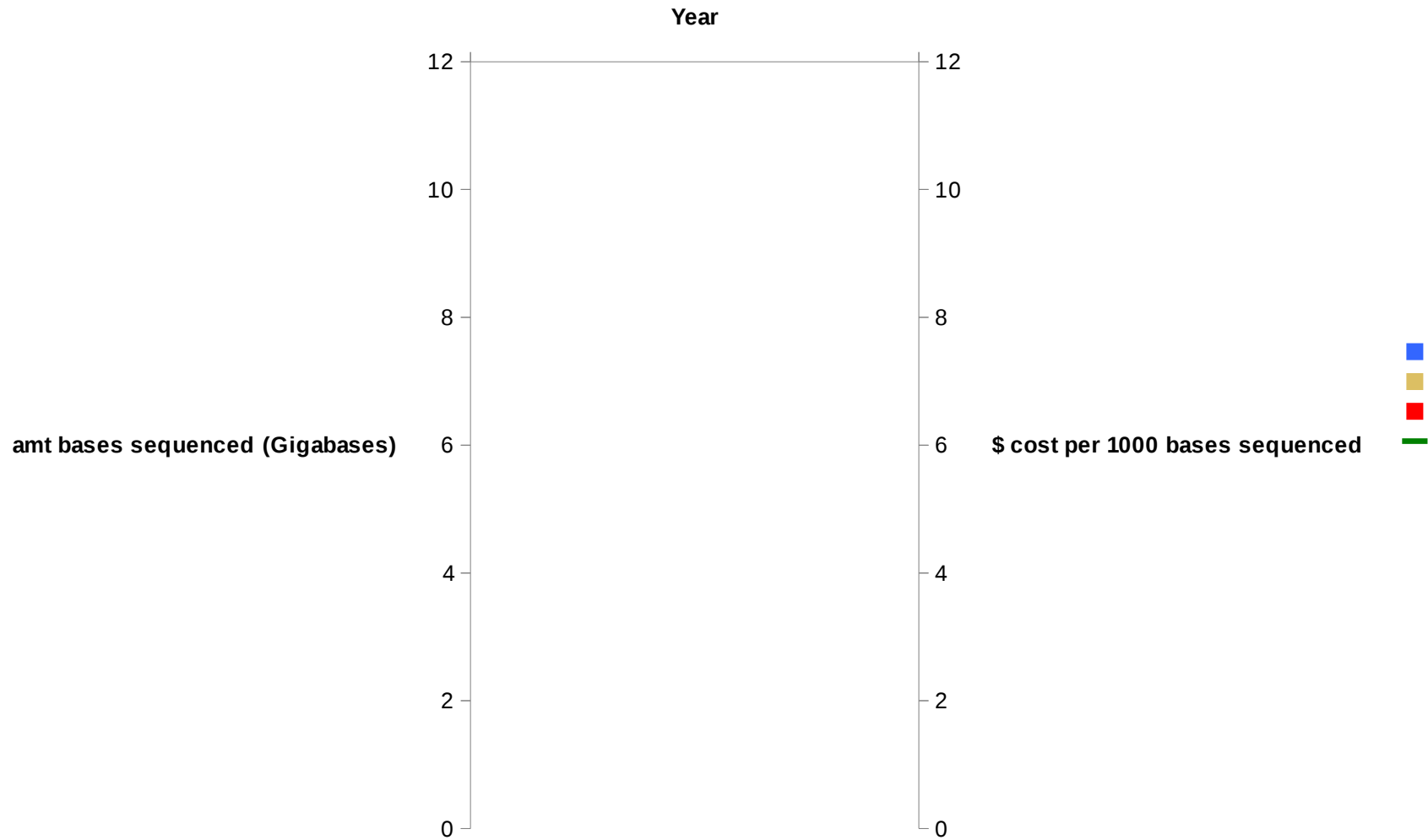
Behind the Surge: NGS Technologies

§ The advent of 2nd generation sequencing technologies has revolutionized the genomics field!



Driving the cost down: a new way to do business...

DOE JGI sequence output (\$/Kb) 2004-2009



When “more” isn’t just more!

- Sanger - 1975
- ABI gel “automated” - 1986
- ABI Capillary 1999 – current



Capillary Based Sequencer, 70 kb / run

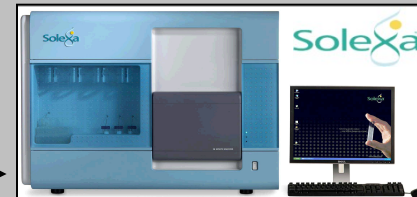
2005



1/2 day

30 > 100 > 400 mb / run
100bp > 250bp > 400bp
Pyrosequencing

2007



3 days

1.0 > 3 > 10 gb / run
25bp > 50bp > 75bp
Seq. by Synthesis

2008



7 days

1.0 > 10 > 20 gb / run
25bp > 35bp > 50bp
Seq. by Ligation

2010

What's

More on this at NGS workshop!! 11:45am Fri. Room 133



Oral pathogen sequence repository

www.oralgen.lanl.gov

Oralgen Databases - Windows Internet Explorer

http://oralgen/

File Edit View Favorites Tools Help

Oral Pathogen Sequence Databases

Home

Search

Search Oralgen Google

Oral Bacteria

- [Actinobacillus actinomycetemcomitans](#)
- [Actinomyces naeslundii](#)
- [Fusobacterium nucleatum](#)
- [Fusobacterium nucleatum polymorphum](#)
- [Porphyromonas gingivalis W83](#)
- [Porphyromonas gingivalis ATCC 33277](#)

Oralgen NEWS Help Desk

Oral Pathogen Sequence Databases

Oralgen NEWS Help Desk

These specialized databases are funded by the [National Institute of Dental and Craniofacial Research \(NIDCR\)](#) within the [National Institutes of Health](#), Bethesda Maryland.

The scope of the project includes molecular information pertaining to oral pathogens, bacterial and viral. The Oralgen project was renewed in May 2006, thus we anticipate continuing expansion over the coming five years. A full statement of our goals is forthcoming. In general we strive to provide one-stop shopping in which analysis and compilation go hand in hand. Currently, the included databases comprise:

Oral Bacteria

- [Actinobacillus actinomycetemcomitans](#)
- [Actinomyces naeslundii](#)
- [Fusobacterium nucleatum](#)
- [Fusobacterium nucleatum polymorphum](#)
- [Porphyromonas gingivalis W83](#)
- [Porphyromonas gingivalis ATCC 33277](#)

Los Alamos National Laboratory

Bioscience Division

ORALGEN

Goal to Maintain/Improve an Analytical Resource

- **Conduct Specialized Annotation and Other Analyses:**
 - ***Hypothetical protein ranking for *P. gingivalis* and *S. mutans****
 - ***Protein cellular localization prediction***
 - ***Metabolic pathways and transport capability analysis***
 - ***Recent gene duplication, insertion sequence, conjugative transposons, and “Genomic Island” predictions***
 - ***Small, non-coding RNAs (sRNAs) prediction***
 - ***Phylogenetic fingerprints***
 - ***etc.***
- **Advanced searching capabilities**
- **Provide comparative genomic tools (like StepToto DB)**
- **Links to other oral microbiome/pathogen resources**

New site up – slow and steady...

The screenshot shows a web browser window displaying the Oral Genomic and Metagenomic Database website. The browser's address bar shows the URL <http://hemisphere.lanl.gov/oralgen-tng/>. The website header features the Los Alamos National Laboratory logo on the left, the title "Oral Genomic and Metagenomic Database" in the center, and a search bar on the right with options to "Search Oralgen" or "Google". Below the header is a navigation menu with links for "My Oralgen", "Search", "Tools", "Community", and "Help".

The main content area is divided into several sections:

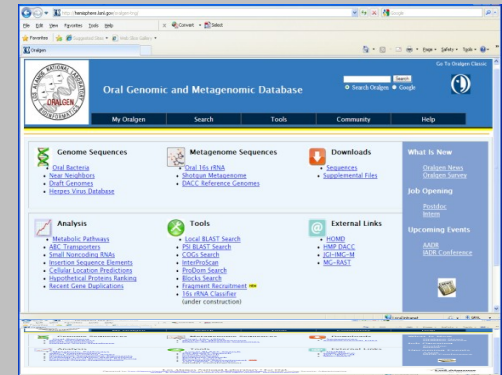
- Genome Sequences** (DNA double helix icon):
 - [Oral Bacteria](#)
 - [Near Neighbors](#)
 - [Draft Genomes](#)
 - [Herpes Virus Database](#)
- Metagenome Sequences** (Microarray icon):
 - [Oral 16s rRNA](#)
 - [Shotgun Metagenome](#)
 - [DACC Reference Genomes](#)
- Downloads** (Download arrow icon):
 - [Sequences](#)
 - [Supplemental Files](#)
- Analysis** (Line graph icon):
 - [Metabolic Pathways](#)
 - [ABC Transporters](#)
 - [Small Noncoding RNAs](#)
 - [Insertion Sequence Elements](#)
 - [Cellular Location Predictions](#)
 - [Hypothetical Proteins Ranking](#)
 - [Recent Gene Duplications](#)
- Tools** (Wrench and screwdriver icon):
 - [Local BLAST Search](#)
 - [PSI BLAST Search](#)
 - [COGs Search](#)
 - [InterProScan](#)
 - [ProDom Search](#)
 - [Blocks Search](#)
 - [Fragment Recruitment](#) NEW
 - [16s rRNA Classifier](#) (under construction)
- External Links** (@ icon):
 - [HOMD](#)
 - [HMP DACC](#)
 - [JGI-IMG-M](#)
 - [MG-RAST](#)
- What Is New** (Blue sidebar):
 - [Oralgen News](#)
 - [Oralgen Survey](#)
 - Job Opening**
 - [Postdoc](#)
 - [Intern](#)
 - Upcoming Events**
 - [AADR](#)
 - [IADR Conference](#)

The footer of the page includes the text "Operated by Los Alamos National Laboratory, Los Alamos, New Mexico" and the Los Alamos logo.

New technologies, new challenges!!

§ Keeping up! **Oralgen v2.0**

- How to feed the sequencing monsters
- How to handle the data (store, analyze)



§ What to do with a few hundred **tens of Million** reads?

- *de novo* sequencing (and finishing!)

How quickly things change...

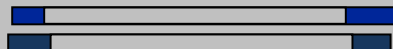
Sanger libraries

Paired ends

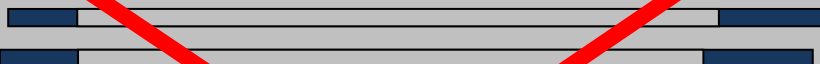
3 kb



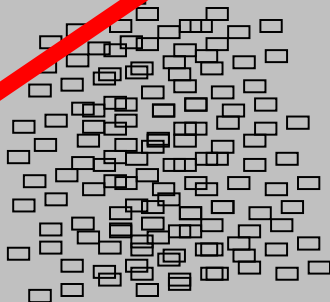
6-8 kb



and 40kb

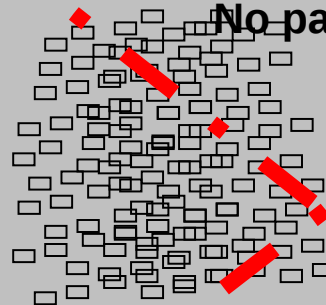


~~454 libraries (150-450 bp
No paired ends)~~

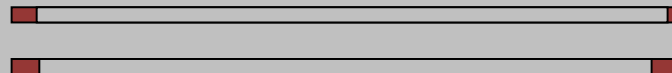


New Tech

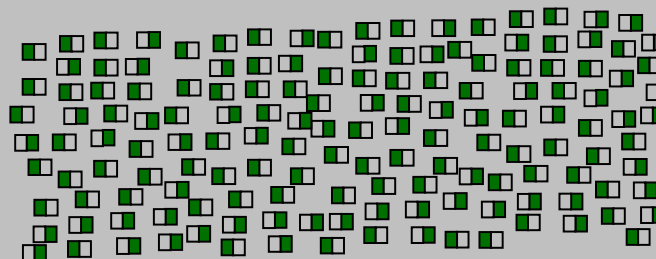
~~454 standard libraries (400 bases long
No paired ends)~~



~~454 Paired ends (Paired ends with reads of
150 bp - Average insert 15-20 kb)~~

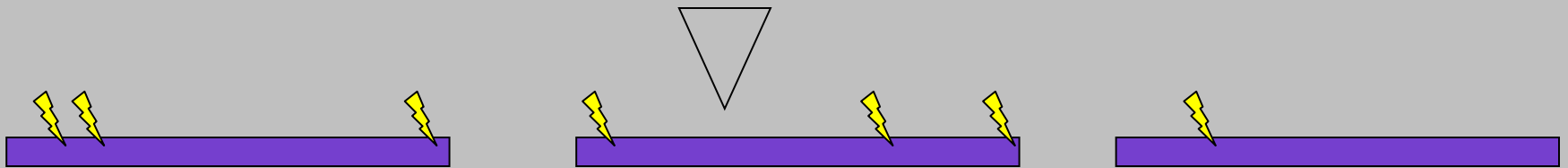
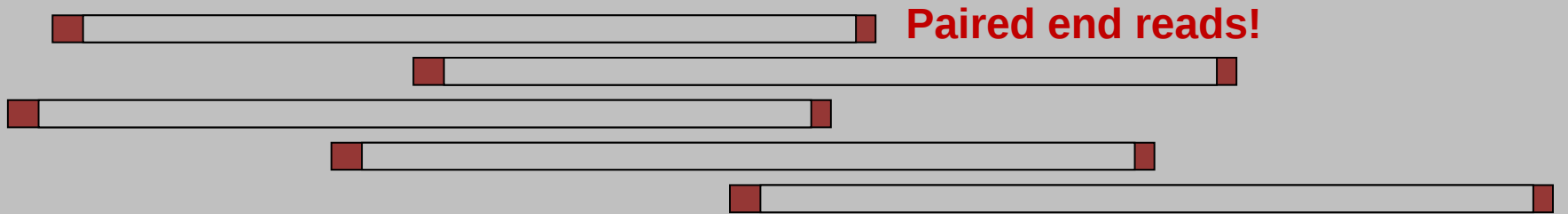
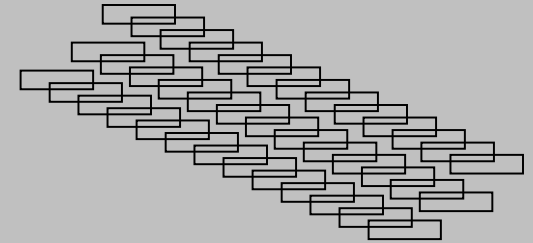
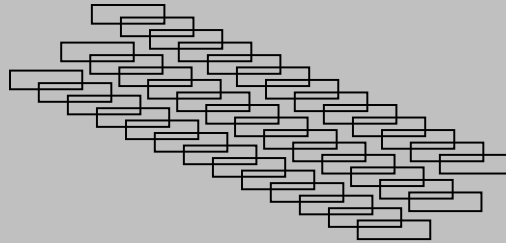
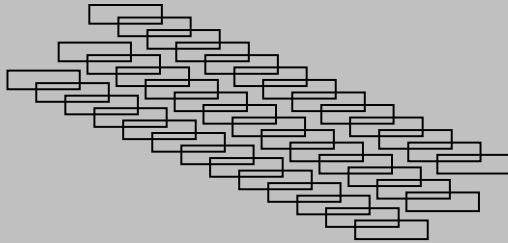


~~Solexa (Illumina) 150 bp inserts with reads
of 36 bp (one end)~~

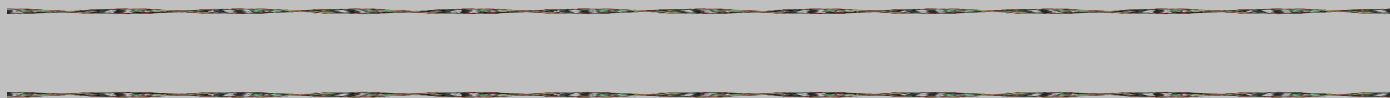


Genome Assembly with NGS technology

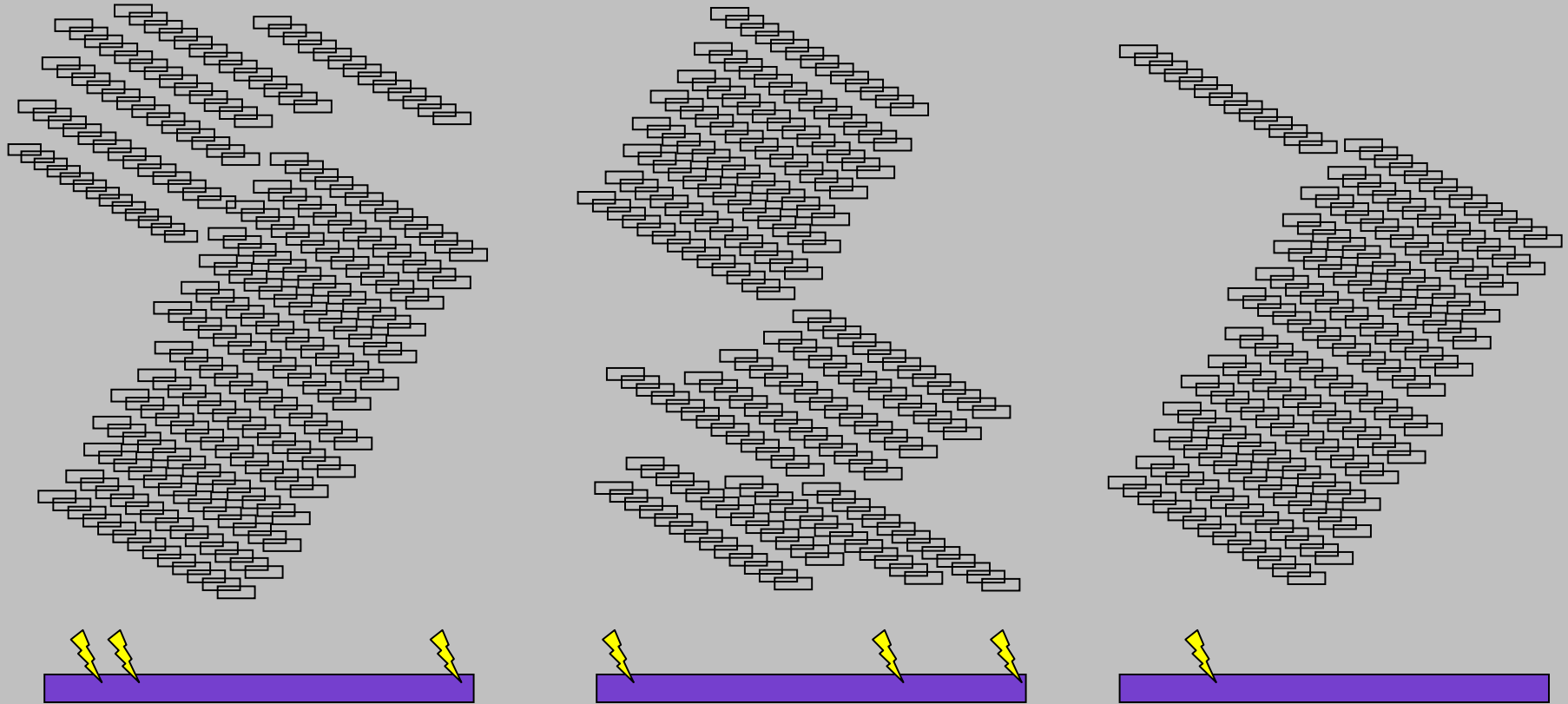
Short reads – but many of them! ■



Contigs (consensus sequences)



Genome Assembly with NGS technology

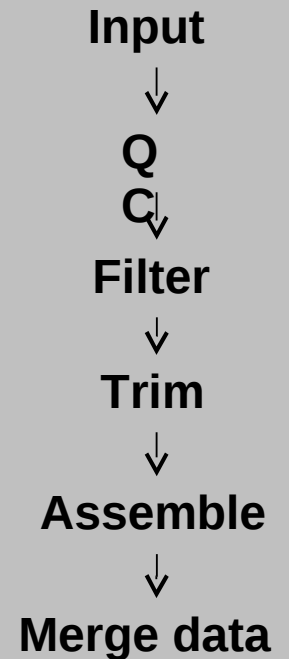


Contigs (consensus sequences)

****Corrected sequences, but still many gaps generally**

Improved NGS assembly

<i>Kingella kingae</i>	Original assembly	LANL improved assembly
Assembled_reads:		98.03%
Singleton:		2,388
Contigs_number:	180	101
N50:	21,244	54,948
N90:	6,210	11,656
Max:	76,431	188,302
Min:	103	110
Total_bases:	1,942,587	2,012,773

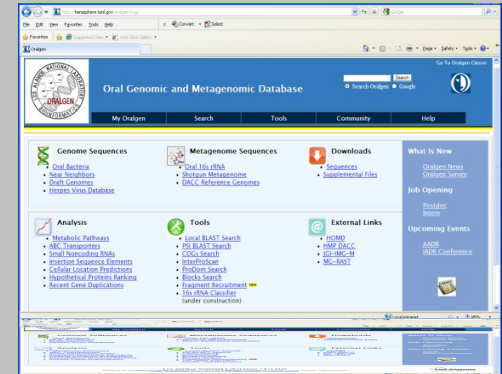


Can also suggest primers for genome closure
(ie. Finishing)

New technologies, new challenges!!

§ Keeping up! **Oralgen v2.0**

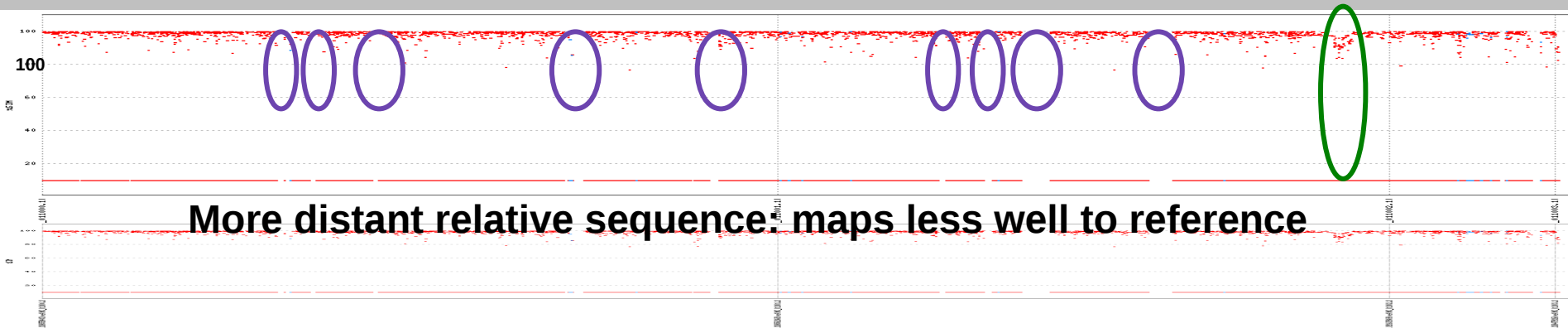
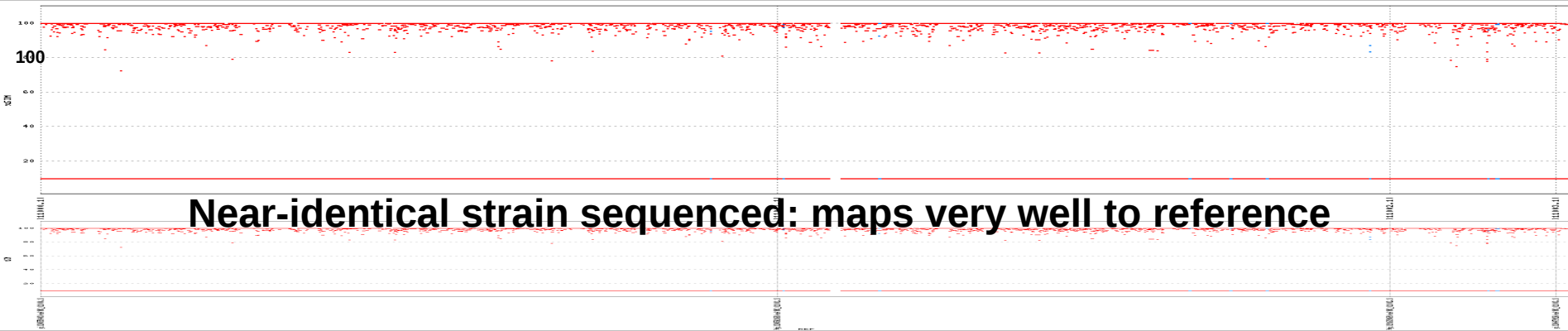
- How to feed the sequencing monsters
- How to handle the data (store, analyze)



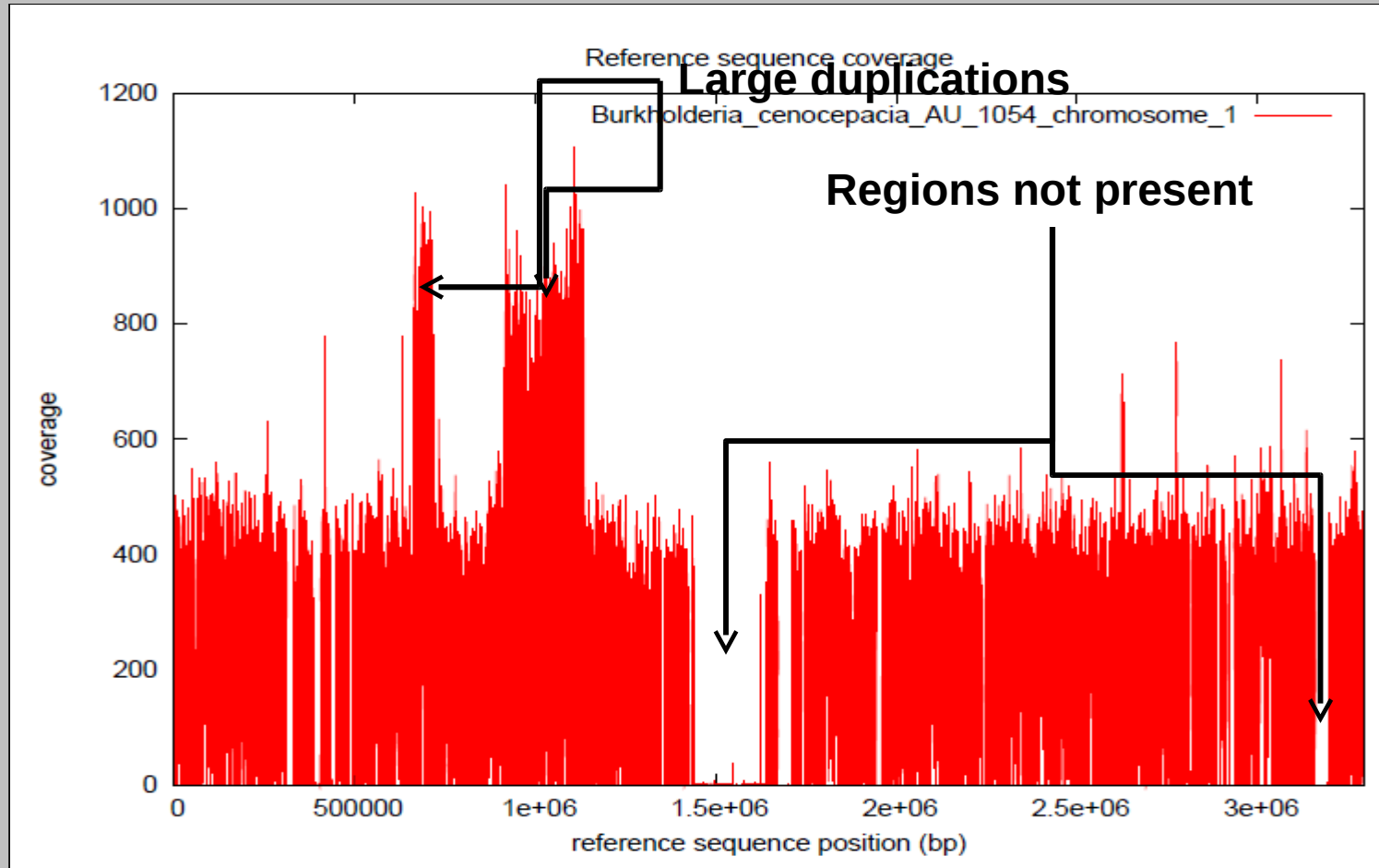
§ What to do with a few hundred **tens of Million** reads?

- *de novo* sequencing (and finishing!)
- Re-sequencing and Transcriptomics (RNAseq)

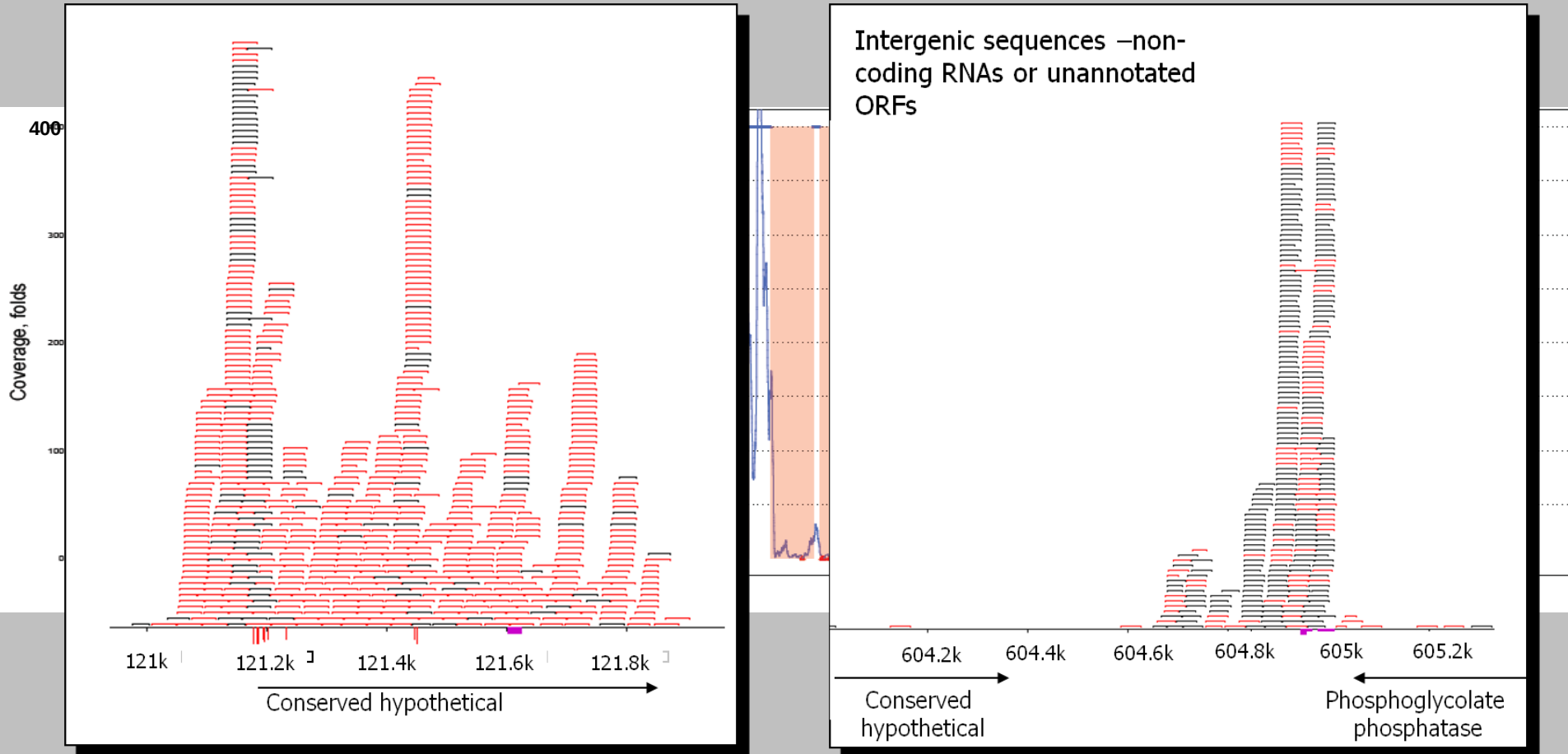
Mapping reads to a reference sequence



Uneven coverage revelations



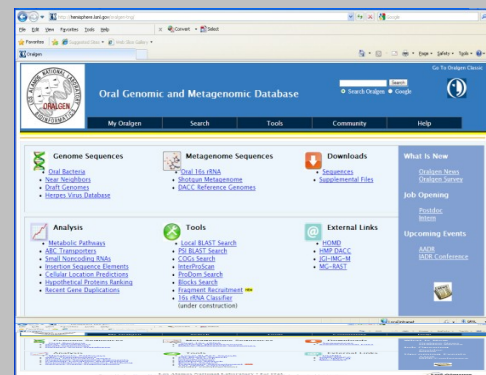
Transcriptome (RNA-sequencing)



New technologies, new challenges!!

§ Keeping up! **Oralgen v2.0**

- How to feed the sequencing monsters
- How to handle the data (store, analyze)



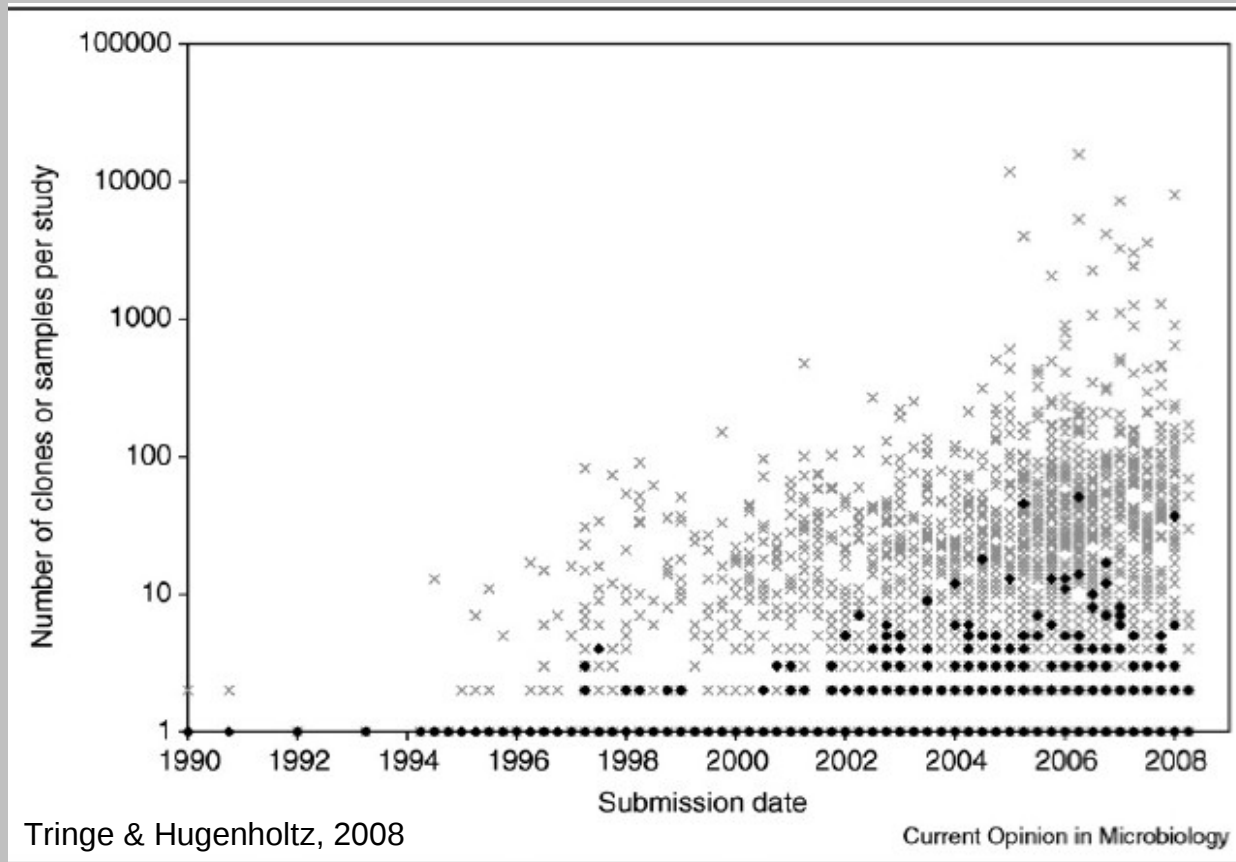
§ What to do with a few hundred ^{tens of} Million thousand reads?

- *de novo* sequencing (and finishing!)
- Re-sequencing and Transcriptomics
- Metagenomics (both rRNA surveys and shotgun “environmental” sequencing)

Metagenomic (16S) population surveys

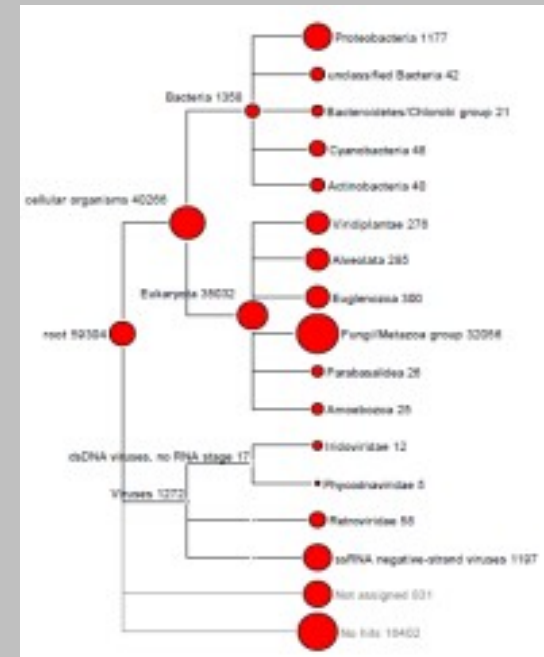
**“Pyrotags” gaining the upper hand:
200,000 to 1,000,000 reads per run**

- Much higher resolution than Sanger 16S surveys
- Has highlighted “rare biosphere”
- May even allow saturation of diversity in a given habitat



Improving Classification Tools

- § **RDP Naïve Bayesian Classifier:** a very fast and accurate taxonomy assignment algorithm
- § LANL has applied a Markov model to enhance accuracy while maintaining speed
- § Also, the LANL Markov NBC will be trained on oral microbes so may be more accurate
 - Will provide on Oralgen2.0!



However: only tells us who is there...

not what they are doing...

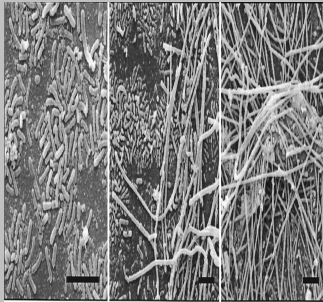
“Next Gen” metagenomics: Sequencing technology changing the landscape



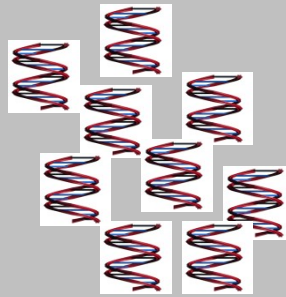
new way



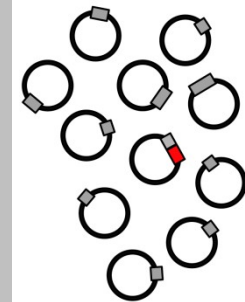
Sample



Extract DNA

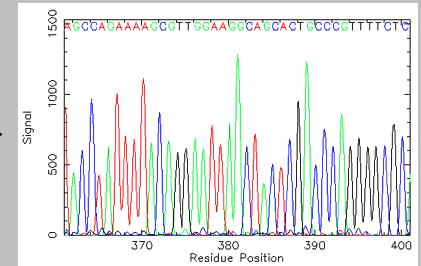


Shotgun clones



3, 8, 40 kb

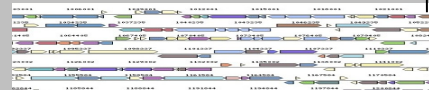
High throughput sequences



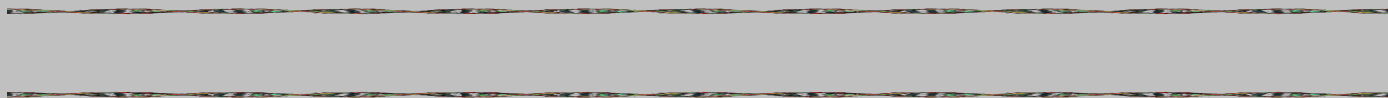
Assemble and Map reads



Gene calls/annotation

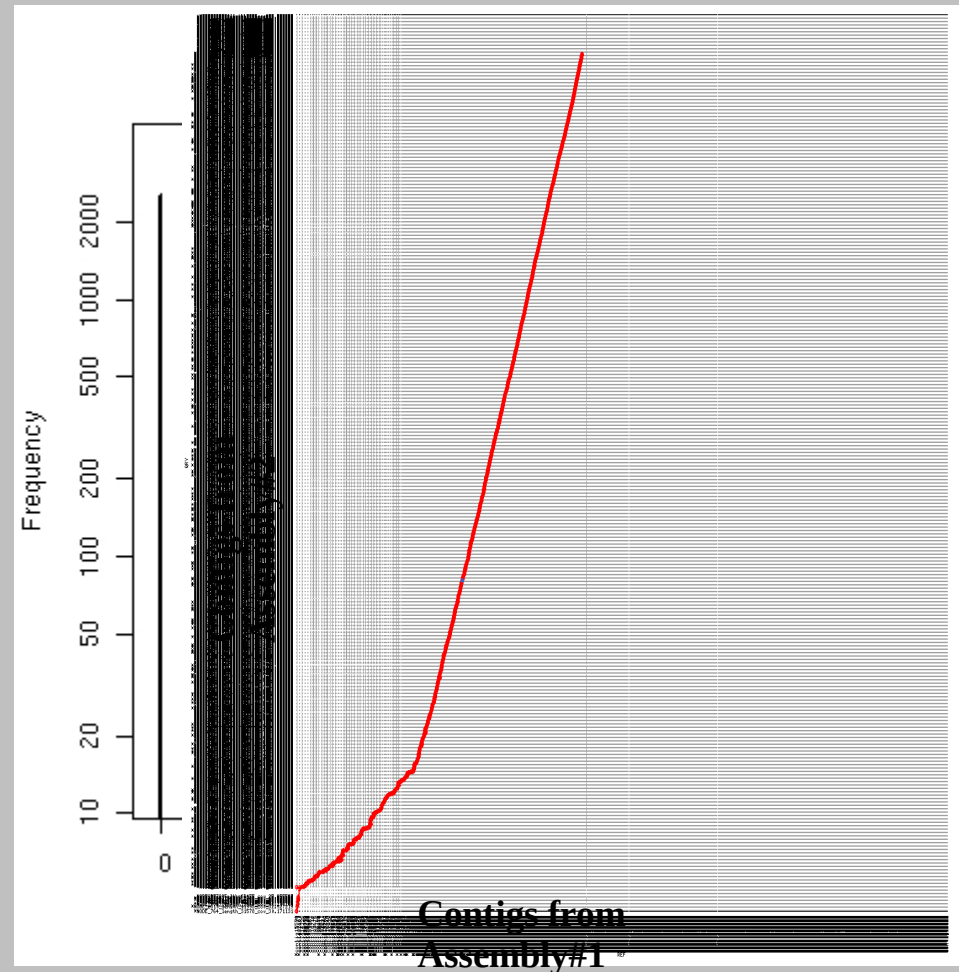


Bin fragments



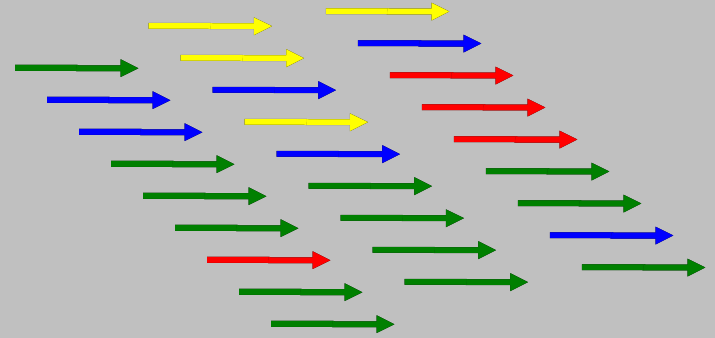
Current Impediments....

- § Different methods (assembly algorithms) and different technologies result in different contigs:
- 12,500 total (5kb max)
 - 287,000 total (30kb max)
 - For either: only ~50% of reads get assembled still!!
- § Very different answers when varying parameters...



Binning and other issues....

§ Binning methods do not yet work on short reads



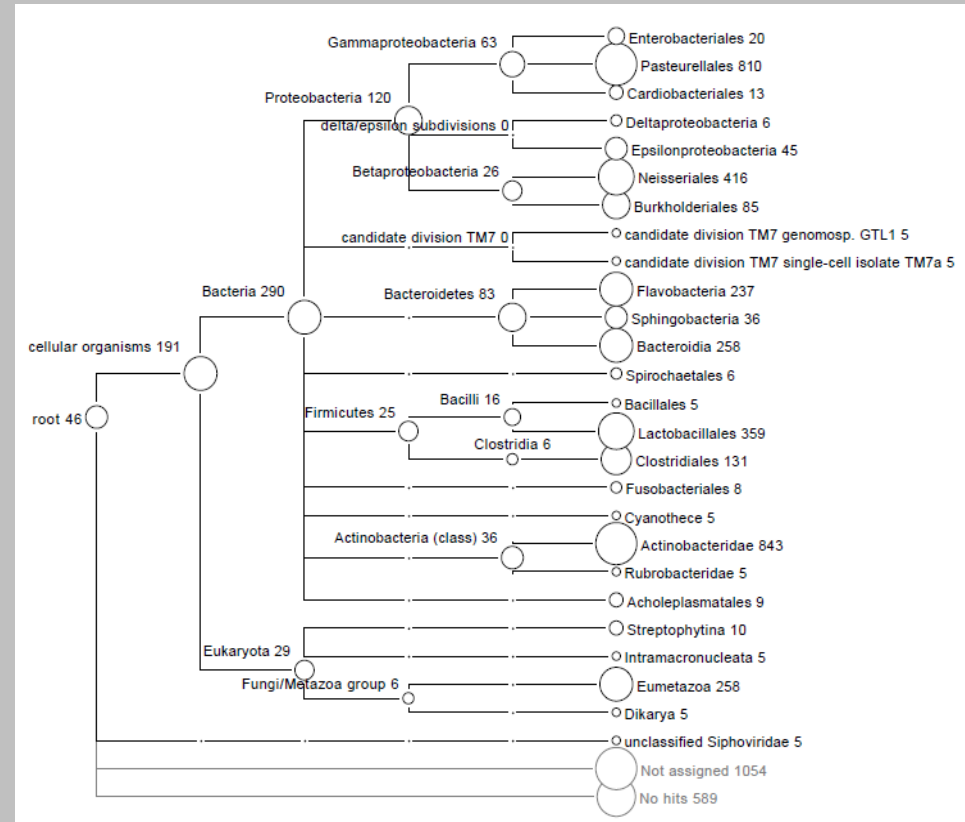
Binning and other issues....

§ Binning methods do not yet work on short reads

§ Annotation of metagenomic contigs takes time!

§ Read-based analyses take many CPU hours; so often look only at assembled contigs (ie. ~50% of the data)

Can one reduce complexity of the data/sample?



Binning via sequence recruitment...

Fragment Recruitment Service

[about fragment recruitment](#)

Select a Reference Genome: **Porphyromonas gingivalis ATCC 33277**
Selenomonas flueggei ATCC 43531
Prevotella melaninogenica ATCC 25845

Select a Query Metagenome: **HumanGutCommunitySubject7**
TM7b
TM7a

[Submit via file](#)

Los Alamos National Laboratory • Est 1943
Operated by [Los Alamos National Security, LLC](#) for the [U.S. Department of Energy Administration](#)
[Inside](#) | [© Copyright 2010 LANS LLC All rights reserved](#) | [Disclaimer](#)

Oral Genomic and Metagenomic Database

My Oralgen Search Tools Community Help

Genome Sequences

- Oral Bacteria
- Near Neighbors
- Draft Genomes
- Herpes Virus Database

Metagenome Sequences

- Oral 16s rRNA
- Shotgun Metagenome
- DACC Reference Genomes

Downloads

- Sequences
- Supplemental Files

Analysis

- Metabolic Pathways
- ABC Transporters
- Small Noncoding RNAs
- Insertion Sequence Elements
- Cellular Location Predictions
- Hypothetical Proteins Ranking
- Recent Gene Duplications

Tools

- Local BLAST Search
- PSI BLAST Search
- COGs Search
- InterProScan
- ProDom Search
- Blocks Search
- Fragment Recruitment
- 16s rRNA Classifier (under construction)

External Links

- HOMD
- HMP DACC
- JGI-IMG-M
- MG-RAST

What is New

- Oralgen News
- Oralgen Survey
- Job Opening
- Postdoc Intern
- Upcoming Events
- AADR IADR Conference

Los Alamos National Laboratory • Est 1943
Operated by [Los Alamos National Security, LLC](#) for the [U.S. Department of Energy's National Nuclear Security Administration](#)
[Inside](#) | [© Copyright 2008 LANS LLC All rights reserved](#) | [Disclaimer](#) | [Privacy](#)

What, of who, is present...

Fragment Recruitment test from GET no args - Mozilla Firefox

http://oralgen/FragmentRecruitment/cgi-bin/query_req.pl?reference=2&query=1&query=&new_file=0

Fragment Recruitment Service

[about fragment recruitment](#)

Reference Genome: Actinomyces sp. F0332
Query Metagenome: HumanGutCommunitySubject7

[tabular data](#)

Fragment Recruitment Service

[about fragment recruitment](#)

Reference Genome: Actinomyces sp. F0332
Query Metagenome: HumanGutCommunitySubject7

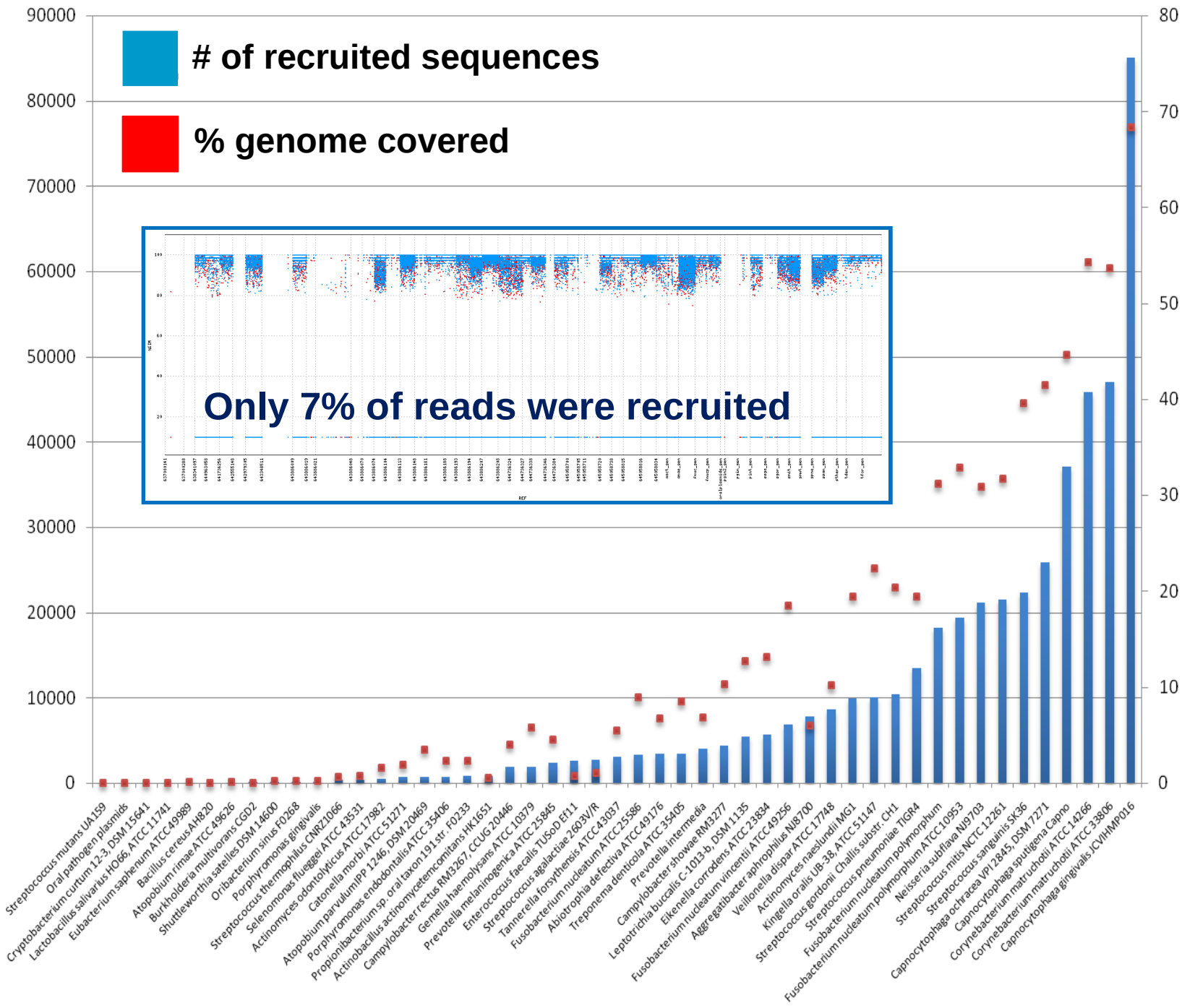
[tabular data](#)

38311.GG703879-GG703882.nuc

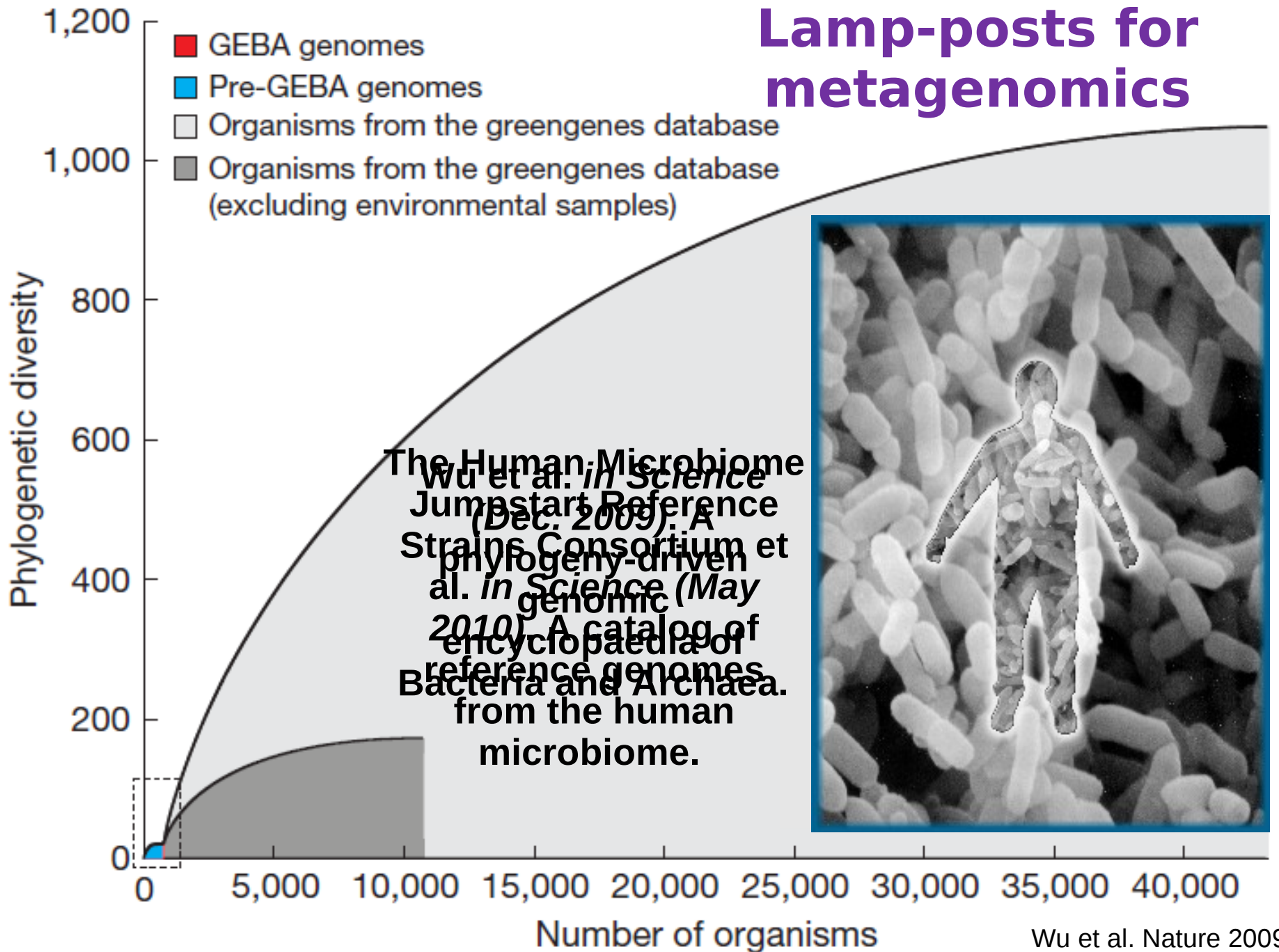
Fragment Recruitment Viewer displays the results from a NUCMER sequence comparison of an available microbial genome against selected metagenome sequence datasets. The reference genomes/contigs are arranged along the x-axis, as indicated by the tick marks. The percentage identities are arranged along the y-axis.

Can “recruit” reads and contigs to many reference genomes





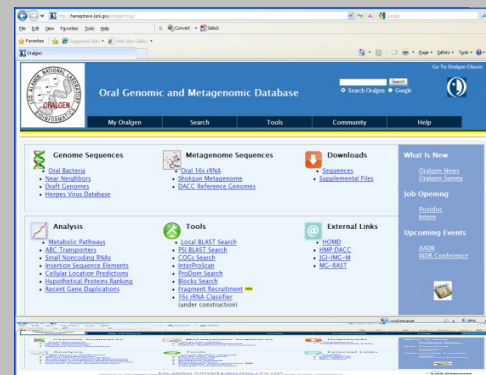
Lamp-posts for metagenomics



New technologies, new challenges!!

§ Keeping up! **Oralgen v2.0**

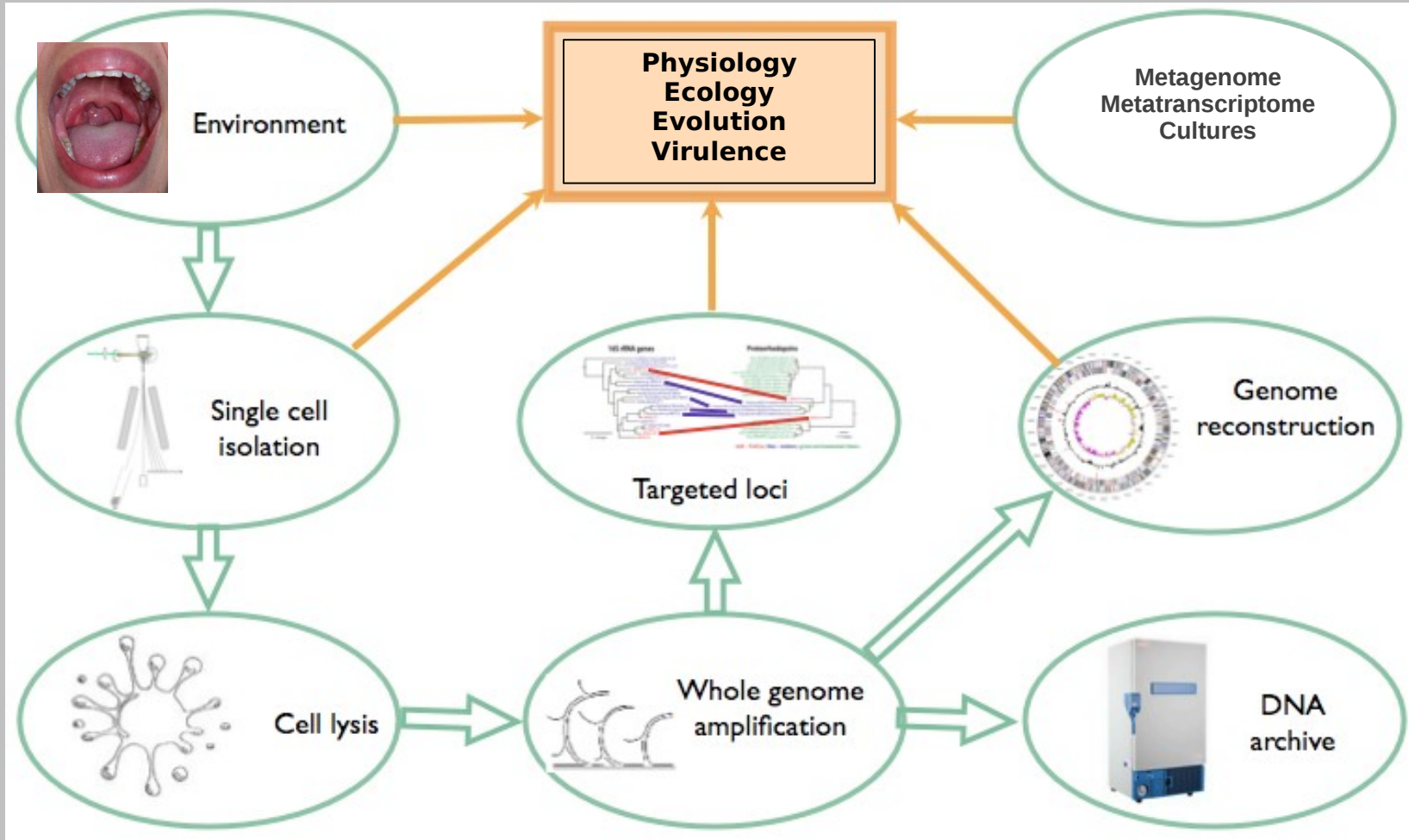
- How to feed the sequencing monsters
- How to handle the data (store, analyze)



§ What to do with a few hundred **tens of Million** reads?

- *de novo* sequencing (and finishing!)
- Re-sequencing and Transcriptomics
- Metagenomics (both rRNA surveys and shotgun “environmental” sequencing)
- **Single cell genomics**

Single cell genomics



Modified from <http://www.bigelow.org/files/8412/5850/8665/Pipeline.png>

Dealing with NGS Technology

1. Make *de novo* and re-sequencing of oral microbiome isolates and single cells a trivial task for investigators
Assembly suite
1. Provide a tailored metagenomic population survey resource and analysis suite
Markov NB Classifier
1. Develop and provide automated and semi-automated (meta)genomics analysis tools (coupled with manual-curation/enhancement)
*Refined annotations
Fragment recruitment*
1. Integrate the above, and work with experts to provide the community with a comprehensive genomics work environment
*Comprehensive database
Tailored analyses*

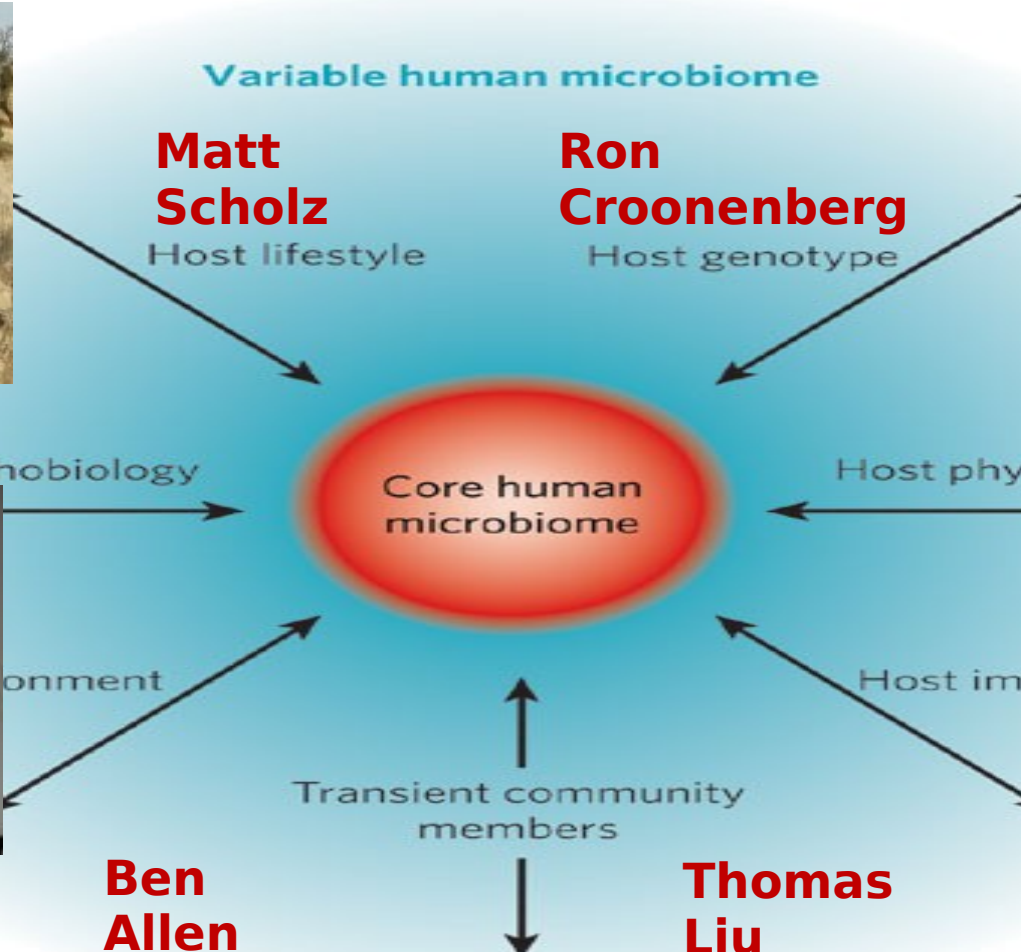
Acknowledgements



Chien-Chi Lo
Host pathobiology



Andrew Liu
Environment



Pavel Senin
Host physiology



Gary Xie
Host immunity

NIDCR National Institute of Dental and Craniofacial Research
National Institutes of Health