

Lessons learned from operating a big metagenomics resource

Folker Meyer, PhD

Argonne National Laboratory
and
University of Chicago



Introduction

Sequence quality

Assembly

Gene prediction

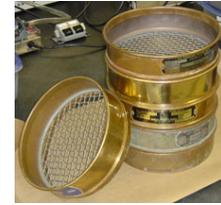
Functional annotation

Current MG-RAST work

Metagenomics

Environmental clone libraries (“functional metagenomics”)

- Sanger sequencing of BAC clones with env. DNA
- low throughput but supports in vitro screens
- **Discovery of novel functions**



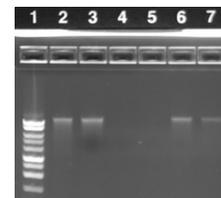
Amplicon studies (single gene studies, 16s rDNA)

- Sequencing of PCR amplified ribosomal genes
- sequence quality limited (→ rare biosphere debate)
- often can't distinguish between individual strains
- **Ecology**

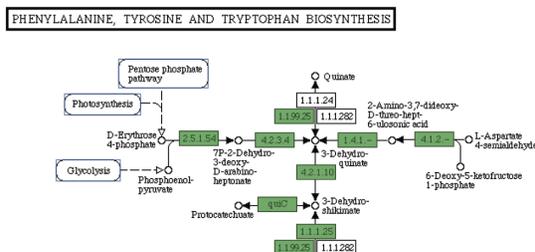


Shotgun metagenomics

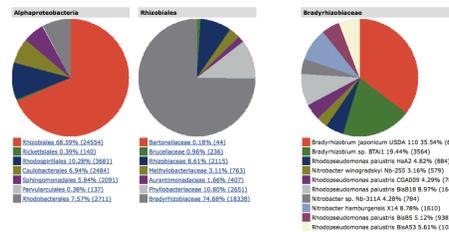
- “random shotgun DNA sequencing applied directly to environmental samples”
- **Consensus sequences & Ecology**



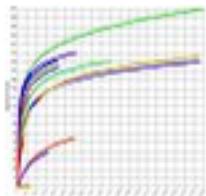
What are they doing?



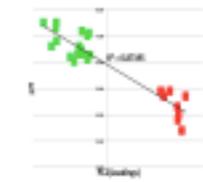
Who are they?



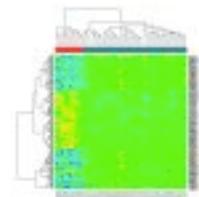
Why I am talking about this?



rarefaction

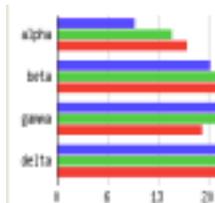


PCoA

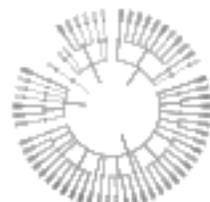


heatmap

# of metagenomes	111,797
# base pairs	44.18 Tbp
# of sequences	380.7 billion
# of public metagenomes	16,791



barchart

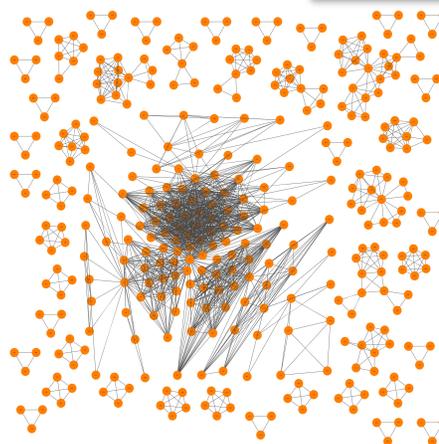
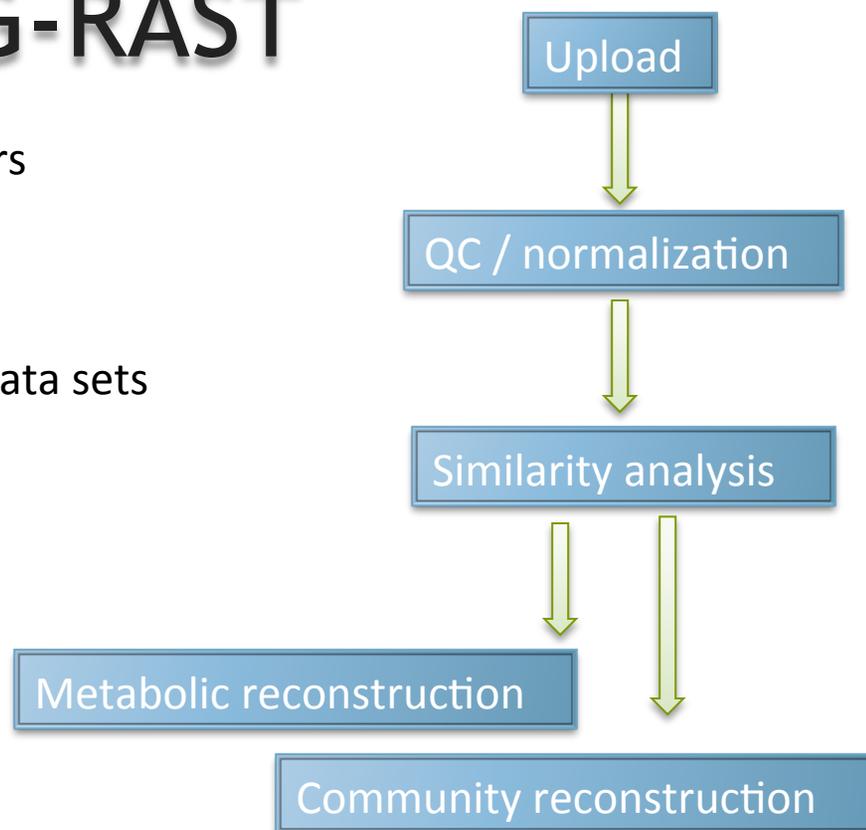


tree

13,000 users
>100,000 data set
>~1.5 Terabases per month

Brief history of MG-RAST

- December 2007 (v1)
 - 100+ groups and ~250 data submitters
 - 100+ data sets, ~**10+ GBp** total size
- October 2009 (v2)
 - Pre-publication sharing available
 - ~1500 data submitters, ~300 public data sets
 - 6000+ data sets
 - **200+ GBp**
- March 2011 (release v3)
 - 2500+ data submitters
 - ~2000 public data sets
 - 25,000 data sets total
- March 2012 (v 3.1.2)
 - **12 Terabasepairs** (10^{12} bp)
- May 2012 (3.20)
 - **13.8 TBp** (10^{12} bp)
 - 128 billion sequences
- December 2013 (3.3.7.3)
 - **39 TBp**
 - 340 billion sequences
 - 100,000 jobs / 14,000 public



Metabolic Reassembly 1.0.0	ID: 4442083.0.0	ID: 4441643.3.0	ID: 4441548.3.0	ID: 4441878.3.0
Amino Acids and Derivatives	10000	10000	10000	10000
Carbohydrates	10000	10000	10000	10000
Cell Division and Cell Cycle	10000	10000	10000	10000
Cell Wall and Capsule	10000	10000	10000	10000
Clustered Isotopologues	10000	10000	10000	10000
Co-factors, Vitamins, Prosthetic Groups, Signaling	10000	10000	10000	10000
DNA Metabolism	10000	10000	10000	10000
Fatty Acids and Lipids	10000	10000	10000	10000
Macromolecular Synthesis	10000	10000	10000	10000
Metabolic Transport	10000	10000	10000	10000
Metabolism of Aromatic Compounds	10000	10000	10000	10000
Metabolism	10000	10000	10000	10000
Motility and Chemotaxis	10000	10000	10000	10000
Nitrogen Metabolism	10000	10000	10000	10000
Nucleosides and Nucleotides	10000	10000	10000	10000
Protein Metabolism	10000	10000	10000	10000
Proteolysis	10000	10000	10000	10000
Respiration	10000	10000	10000	10000
Signal Transduction	10000	10000	10000	10000
Stress Response	10000	10000	10000	10000
Transcription	10000	10000	10000	10000
Translation	10000	10000	10000	10000
Unclassified	10000	10000	10000	10000
Vitamins	10000	10000	10000	10000

simplified



Introduction

Sequence quality

Assembly

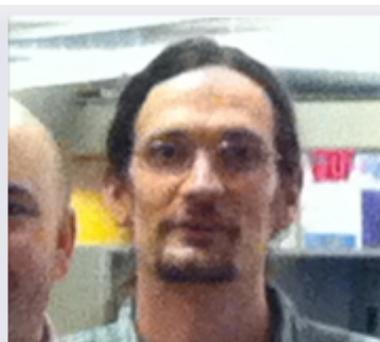
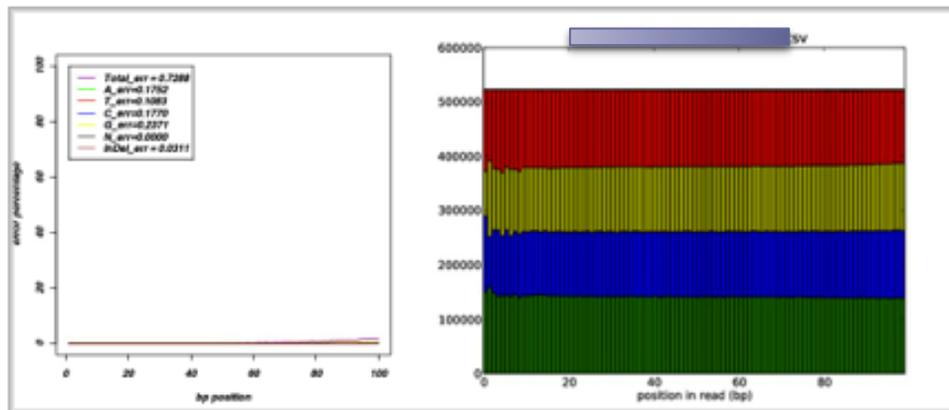
Gene prediction

Functional annotation

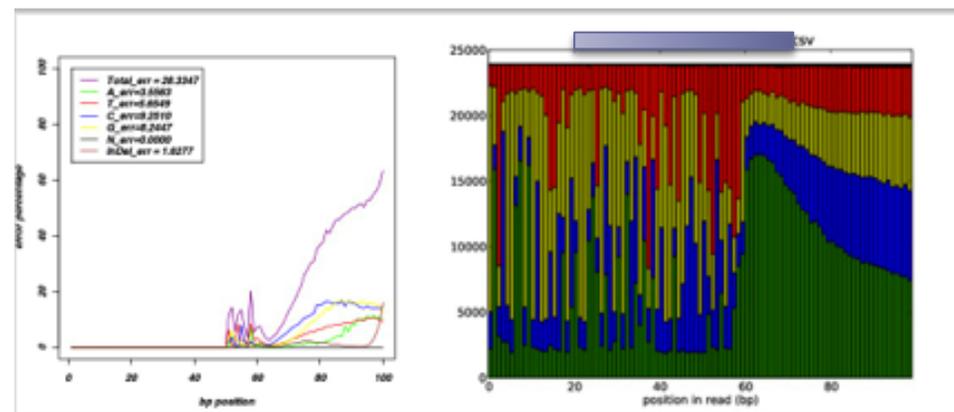
Functional annotation

Current MG-RAST work

Not all sequence data created equal



DRISEE: Keegan, et al, PLoS Comp. Bio, 2012



Every platform has unique issues

- Not all samples from one platform are the same
- Massive operator variation
- Most users don't understand sequence quality
- Majority of users lacks tools to gain insight into "sequence space"
- Sequence providers work hard to obscure issues
- FAQ: Shouldn't my student do the QC locally?
 - → sure he will do a much better job than us, we have only seen 10,000 similar data sets



Introduction

Sequence quality

Assembly

Gene prediction

Functional annotation

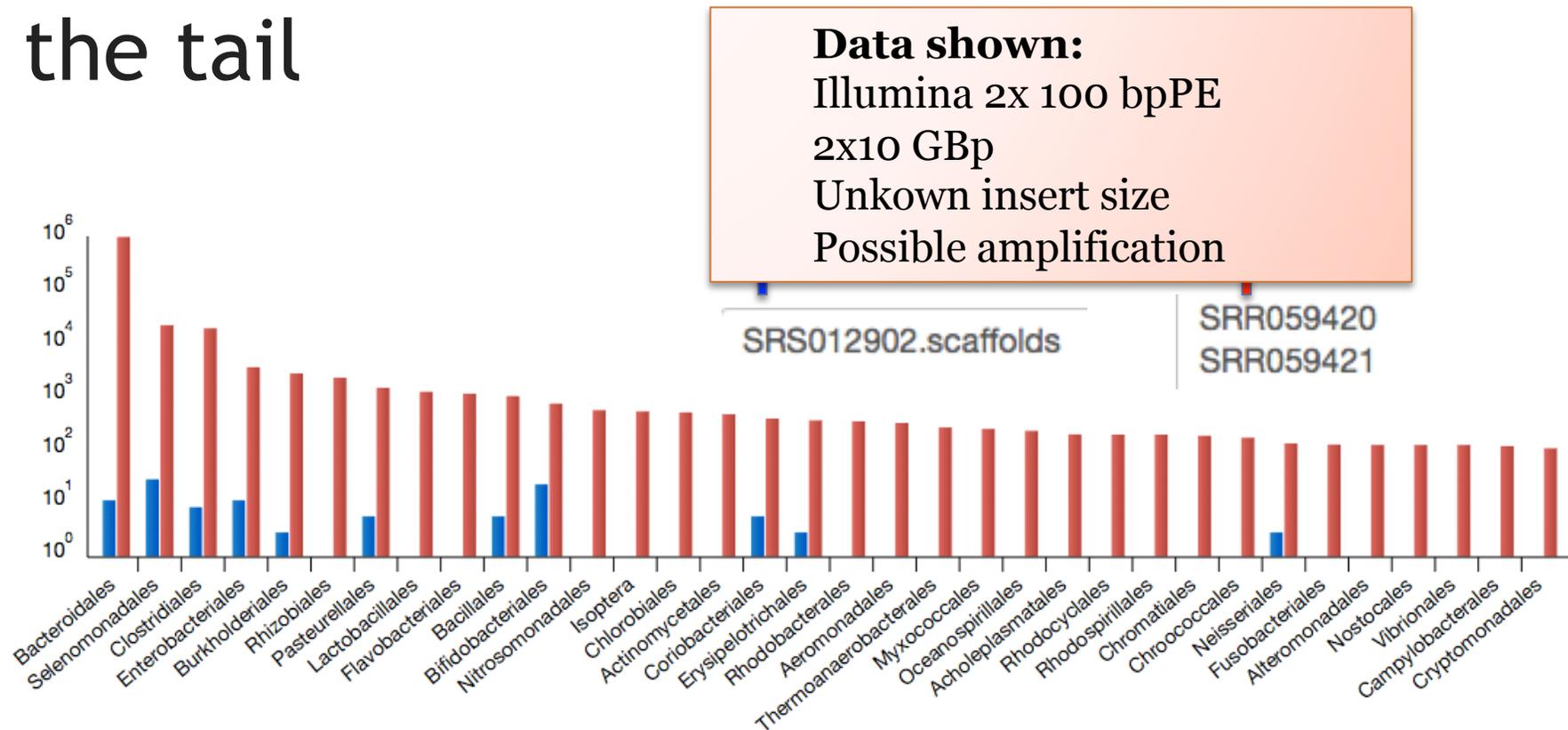
Current MG-RAST work

Two provoking examples:

- **Example 1:**
- **high diversity host microbiome sample**
 - One physical sample
 - two libraries, two sequencing runs
 - unexpected results
- **Example 2:**
- **Medium diversity environmental sample**
 - One physical sample
 - One library, one sequencing run
 - Very unexpected results

Example 1: Human microbiome

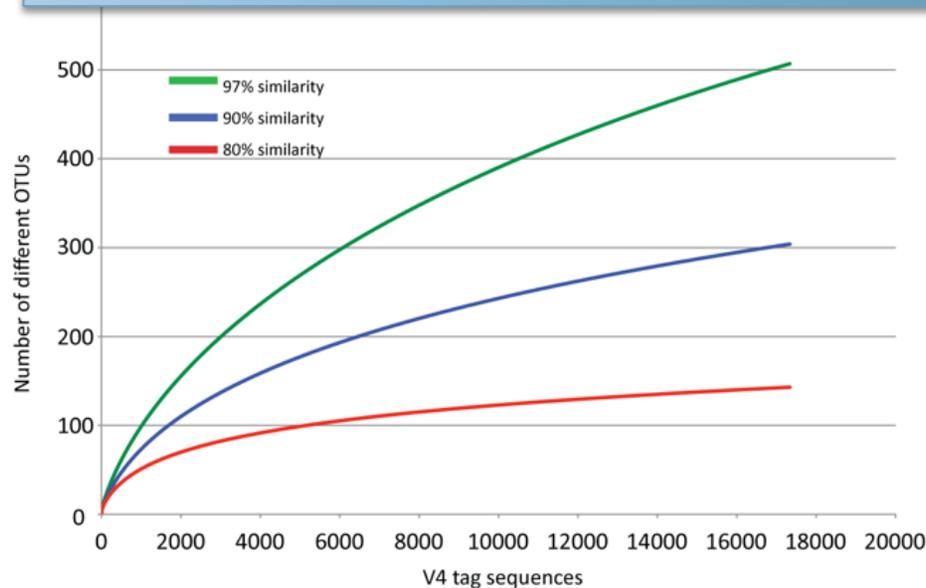
We are missing more than the rare taxa at the tail



I spent less time finding a sample with this problem than plotting the graph.

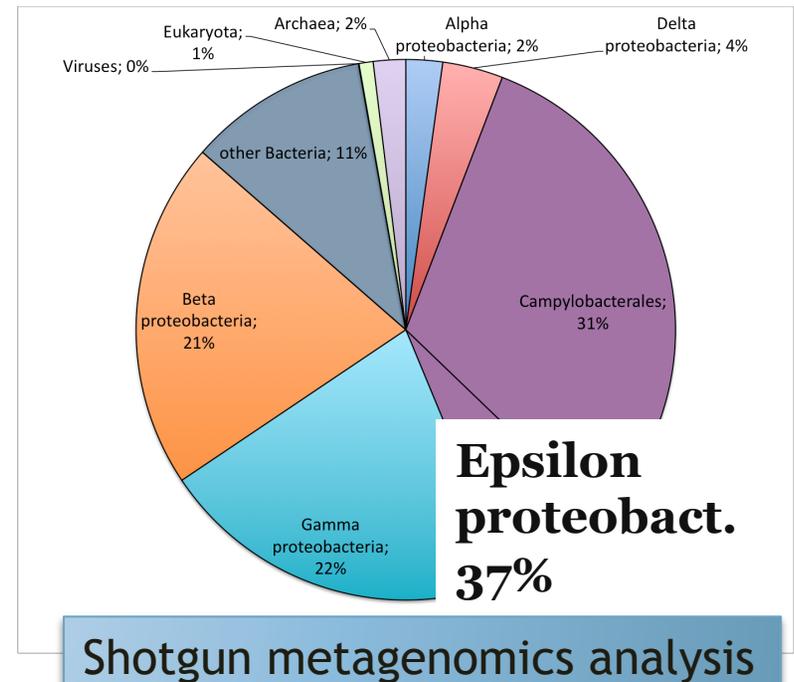
How complex is the community?

16s Rarefaction analysis using RDP classifier



- Based on 17,334 V4 tag sequences
 - 500+ "species" inferred at 97% sequence similarity to define the OTU
 - **Proteobacteria (86.2%)**
 - **genus Sulfuricurvum (46.9%)**
- Assumption: easy assembly

RDP: Cole, et al, NAR, 2007



- Shotgun metagenome data confirm
 - **~37% epsilon proteobacteria**
 - no Sulfuricurvum genome sequence
- Assumption: easy assembly

MG-RAST: Meyer, et al, BMC Bioinformatics, 2008

- *Sulfuricurvum* reads == ~38% of all reads
- Approx. 3.45Gbp are *Sulfuricurvum* reads → 45 million reads
- Assuming 3MBp genome size

Assembly should be easy

Expected:

Length of genome	G = 3000000 bp
Number of sequences	N = 45000000
Average length of sequences	L = 75 bp
Minimum overlap	T = 30 bp
Redundancy of coverage	= 1125.0
Expected number of contigs	= 1.0
The expected number of seqs in a contig	= 45000000.0
The expected length of each contig	= 3000000
Percentage of coverage of genome	= 100%

→ 1100X coverage of genome

Observed:

Resulting contigs:

244,782 Contigs

67Mbp

N50: 227bp

Long contigs >= 1000

5,967 Contigs

No long contigs from
Epsilon proteobacteria

Sequence quality the source of the problem?

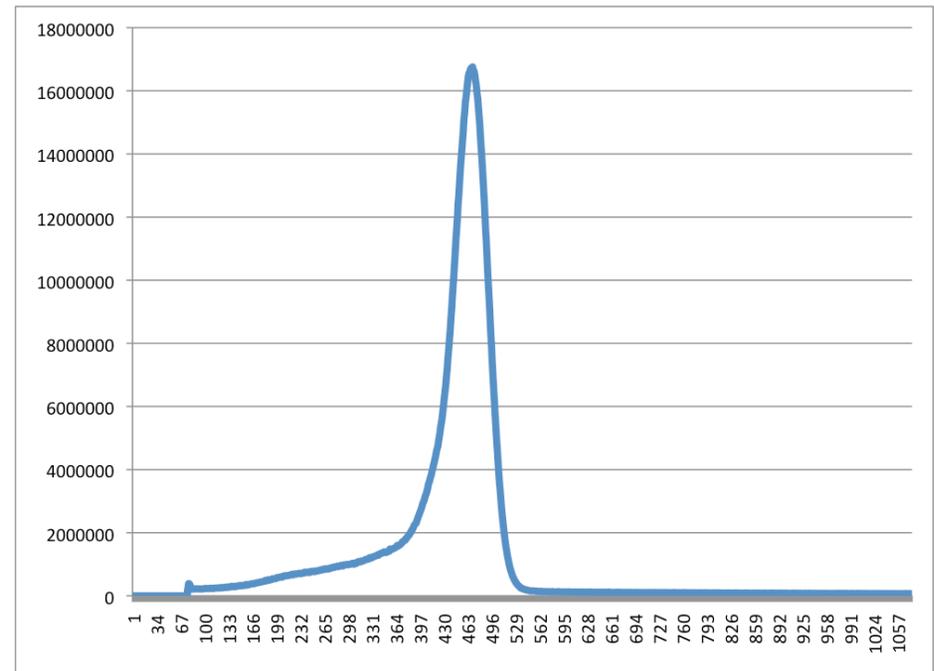
Two arguments against that:

- a) Very good DRISSE score
INSERT SCORE HERE
- b) Mate pairs map nicely to contigs
- mate percentage: 86.015%
 - average distance: 446
 - As designed
 - Using exact matches (300-1000 window)

➔ NO! Look for other reasons

Distribution of mate-pair “insert lengths”

mps



Insert size (bp)

Using exact bowtie matches

Assembly tool(s) or parameters?

Two options:

1) Velvet parameters are incorrect

→ we varied the k-mer length parameter

2) Velvet does not work properly?

→ Will other assemblers produce similar result?

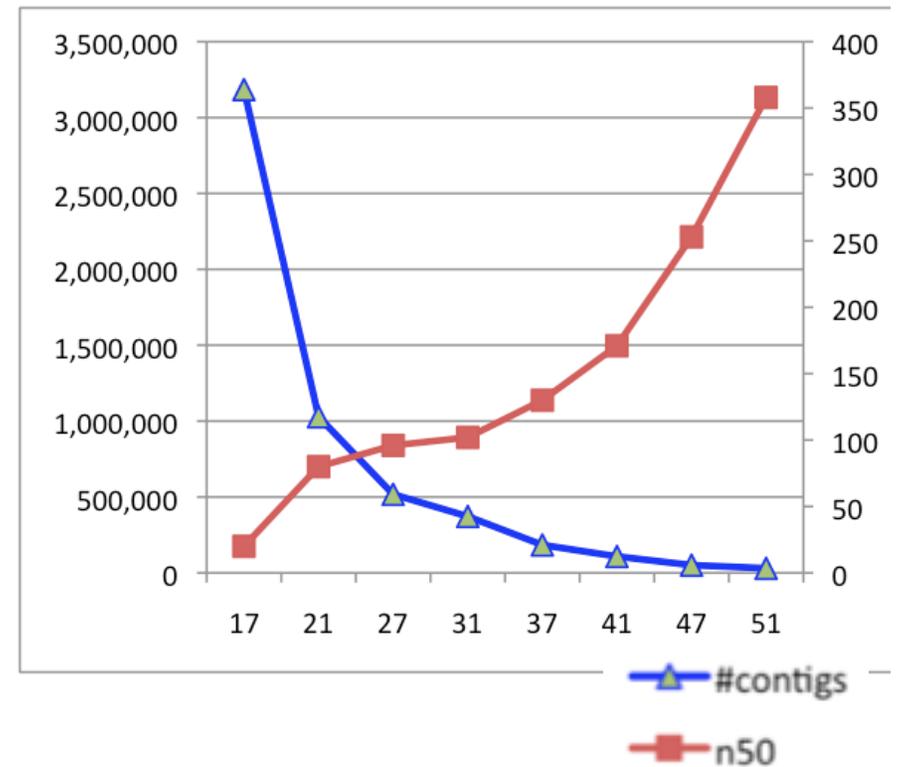
→ We tried: SOAP de-novo, Newbler, Abyss, AMOS, AllPaths

- We learned many valuable lessons
- Requires significant computational resources
- Most tools can't handle 45 million reads

→ similar results with other assemblers

Mihai Pop in 2010 @ HMP in St. Louis, MO
“metagenomes can't be assembled”
and
“all assemblers are equally bad”

Effects of changing k-mer lengths on contig number and sizes



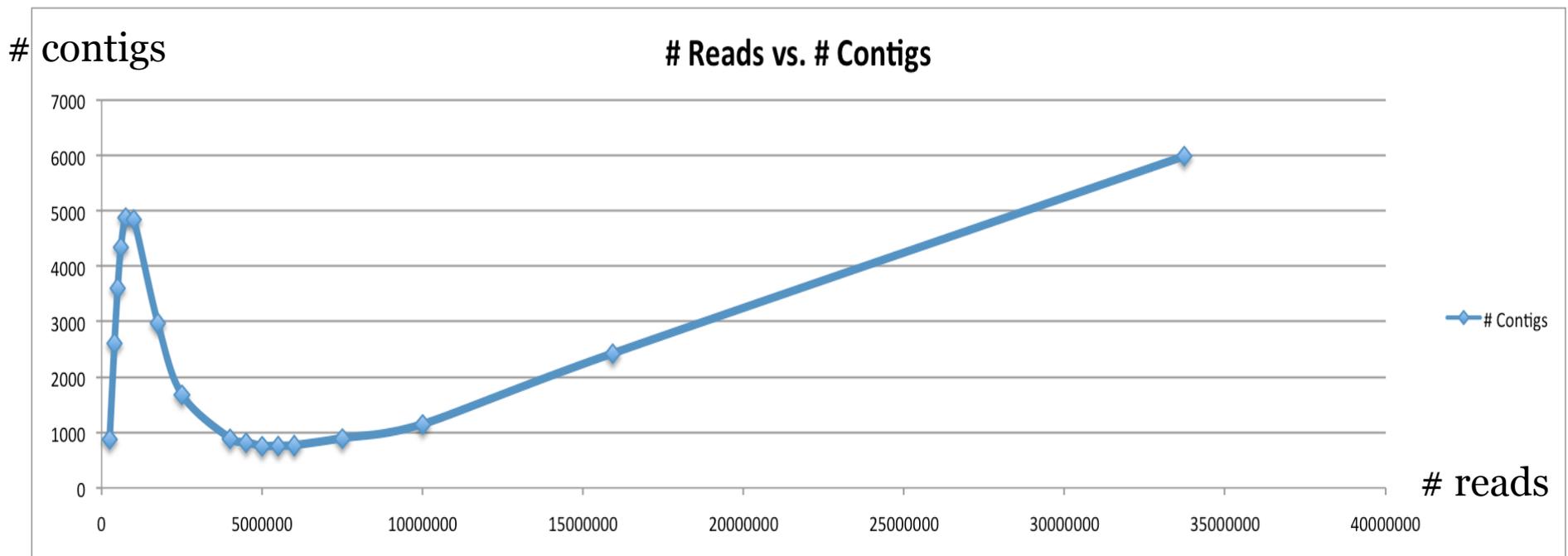
Short k-mers == many contigs but short
Long k-mers == few long contigs

Too much coverage?

- Observations:
 - Contigs yield is very low
 - varying parameters yields different contigs
 - Genome yield is low for all cases
- What if...
 - Coverage was not too low, but too high?
- Hypothesis:
 - Strain variation is confusing assembly tools
 - IFF high coverage confuses assemblers by exposing strain variation, lowering coverage might help.
 - Test strategy: Subsampling

In addition: Integrating different parameter sets will overcome parameter set artifacts

The effects of subsampling



Plotting number of contigs vs. number of reads

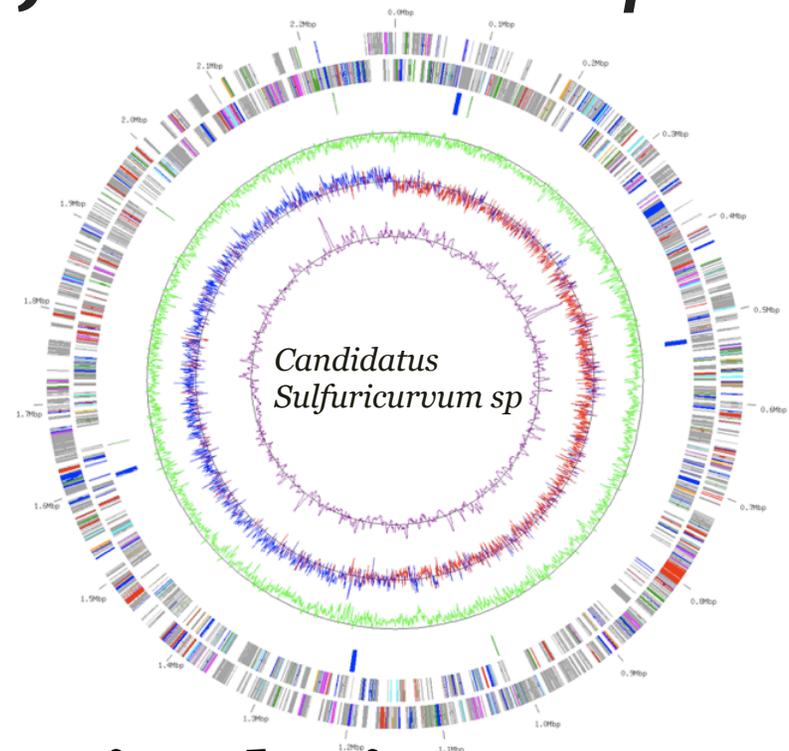
→ More reads with variation confuse velvet (and the other tools)

→ Existing assembly tools have been built for clonal situations

→ Even modern tools like metaVelvet, meta-IDBA, metaSOAP will succumb to too much strain variation

Result: *Candidatus Sulfuricurvum* sp.

- new genus
 - Epsilon proteo bacteria
 - *Campylobacteriales*
 - 2.3 MBp genome
 - 1 contig
 - 43% GC
 - High quality
 - **Verified by many overlapping by intact mate-pairs**
- likely to not be chimeric**



Genome annotation with RAST: Aziz et al, BMC Genomics, 2008
Genome published: Handley et al, AEM, 2014 (in press)



State of the art

(yes I am being provocative)

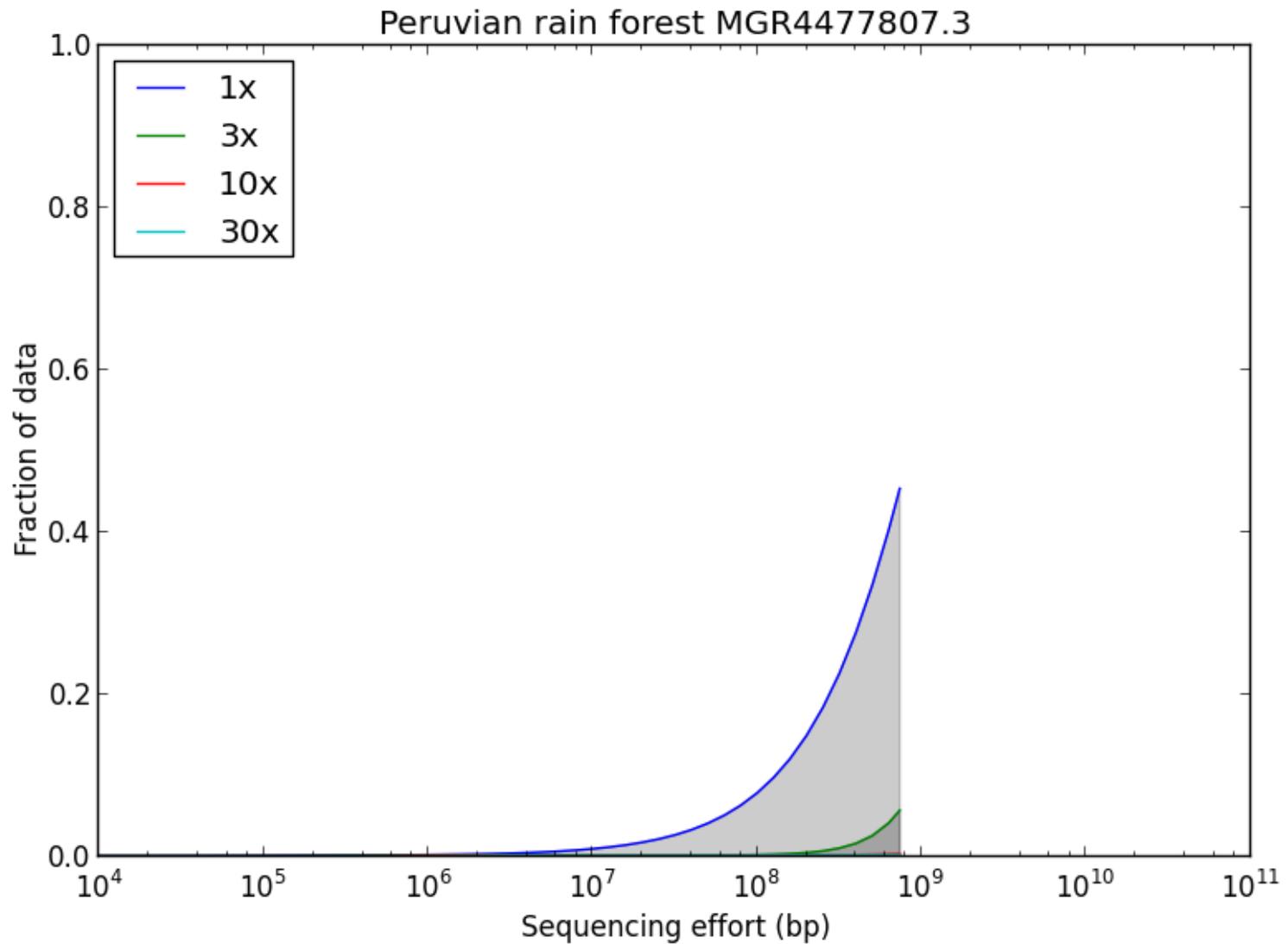
- Apply a single algorithm
 - Assembly, clustering, etc.
- Assume optimal results
- Proceed
- Data sets are sizable, trying to reduce them to manageable size quickly

Alternative:

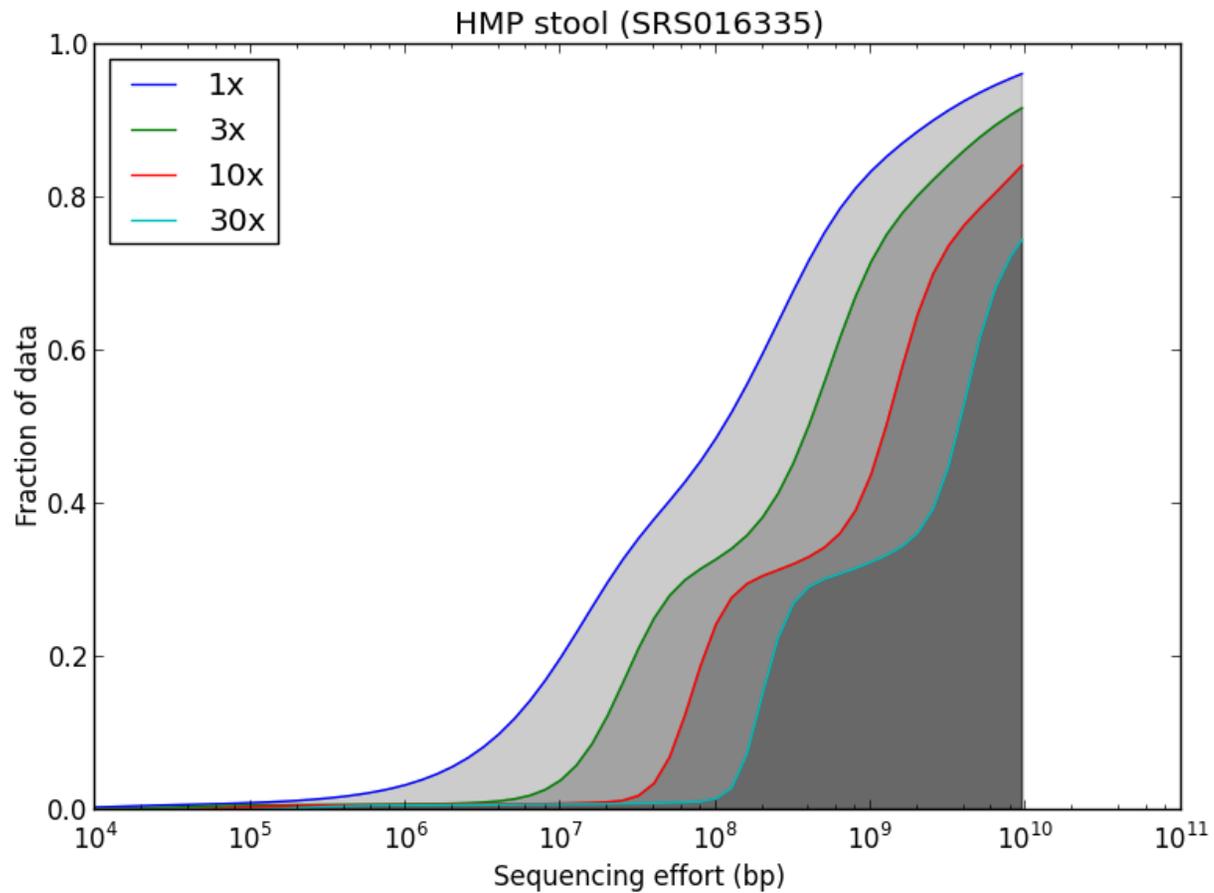
Using nonpareil-K (npK) to understand coverage strata first

- npK evolved from
 - Nonpareil (Rodriguez & Konstantinidis, Bioinformatics, 2013)
 - and Kmer-spectrum-analyzer (Williams et al, BMC Genomics 2014)
- npK: does not use sub-sampling, instead uses Kmer index
- Allows rapid classification “strata”
 - Next slides explain the term 😊

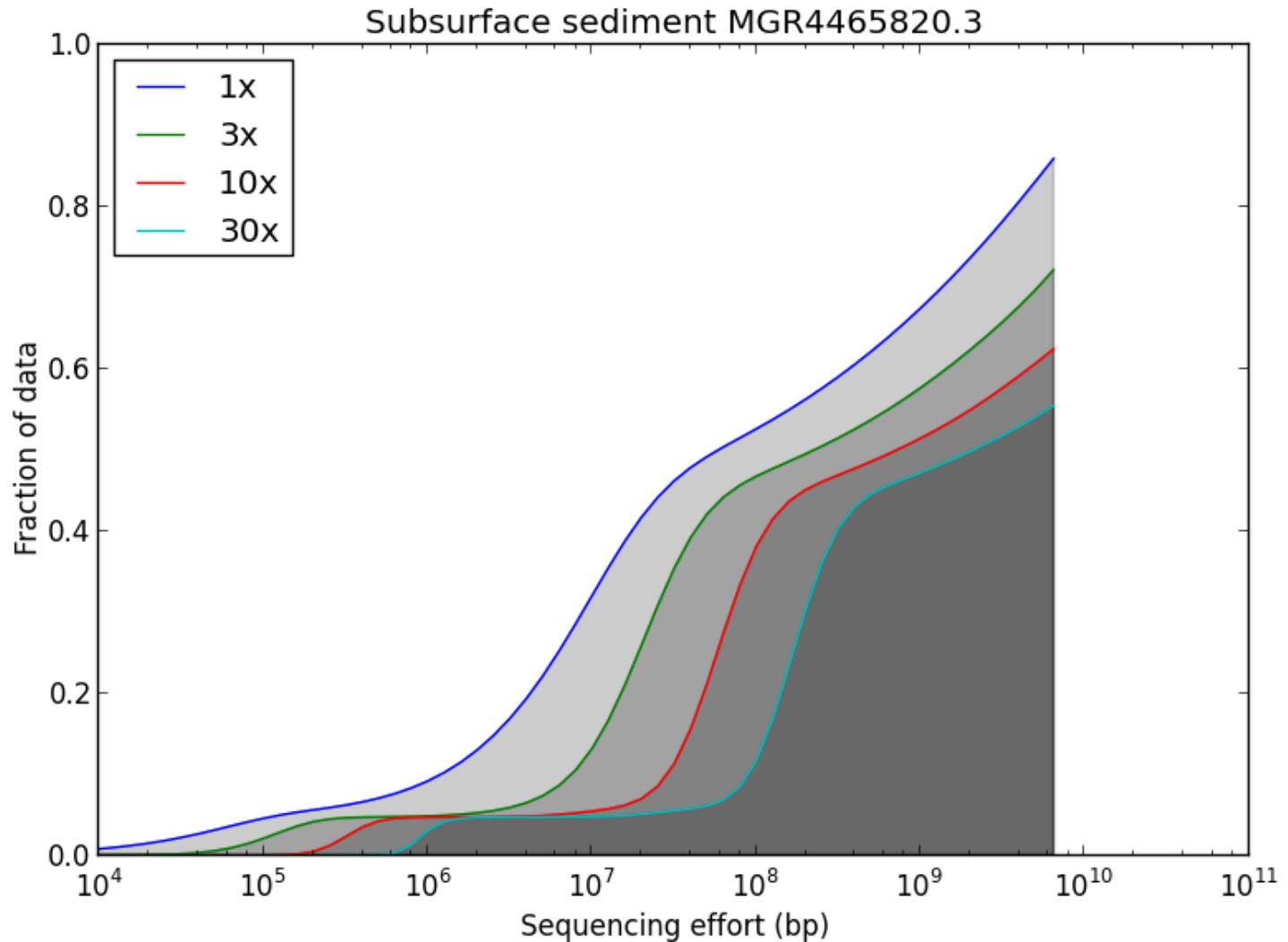
A soil sample



A microbiome sample



A subsurface sample





Introduction

Sequence quality

Assembly

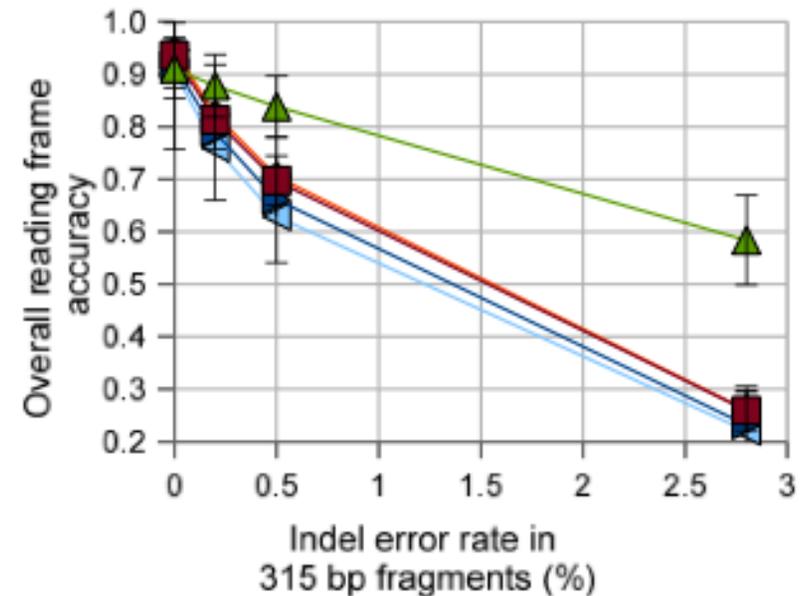
Gene prediction

Functional annotation

Current MG-RAST work

Not all gene callers created equal

- Question:
 - What happens if I vary the tool chain?
- Existing approaches rely on:
 - Compare results of different studies (ie multiple pipelines)
- Here we study 5 different popular gene finding tools for metagenome on **simulated data**
- **Effects are dramatic**
 - Accuracy goes drops dramatically with moderate error
- Comparison of data requires identical tool chain



▲ FGS
■ MGA
◆ MGM
▶ PRODG
◀ OPH



From: Trimble et al, BMC Bioinformatics, 2012



Introduction

Sequence quality

Assembly

Gene prediction

Functional annotation

Current MG-RAST work

Experiment: Does metagenomics work?

- Simple test, pick 4 different sets of samples
- **Red**: Human microbiome, different body parts (sizes from 1-16GBp each)
- **Blue**: Feces from pregnant women in Finland (~12GBp each)
- **Green**: Soil in Illinois (~12GBp each)
- **Black**: Bioreactor (Beer-to-Caproate) 4 points in a time series (3GBp each)

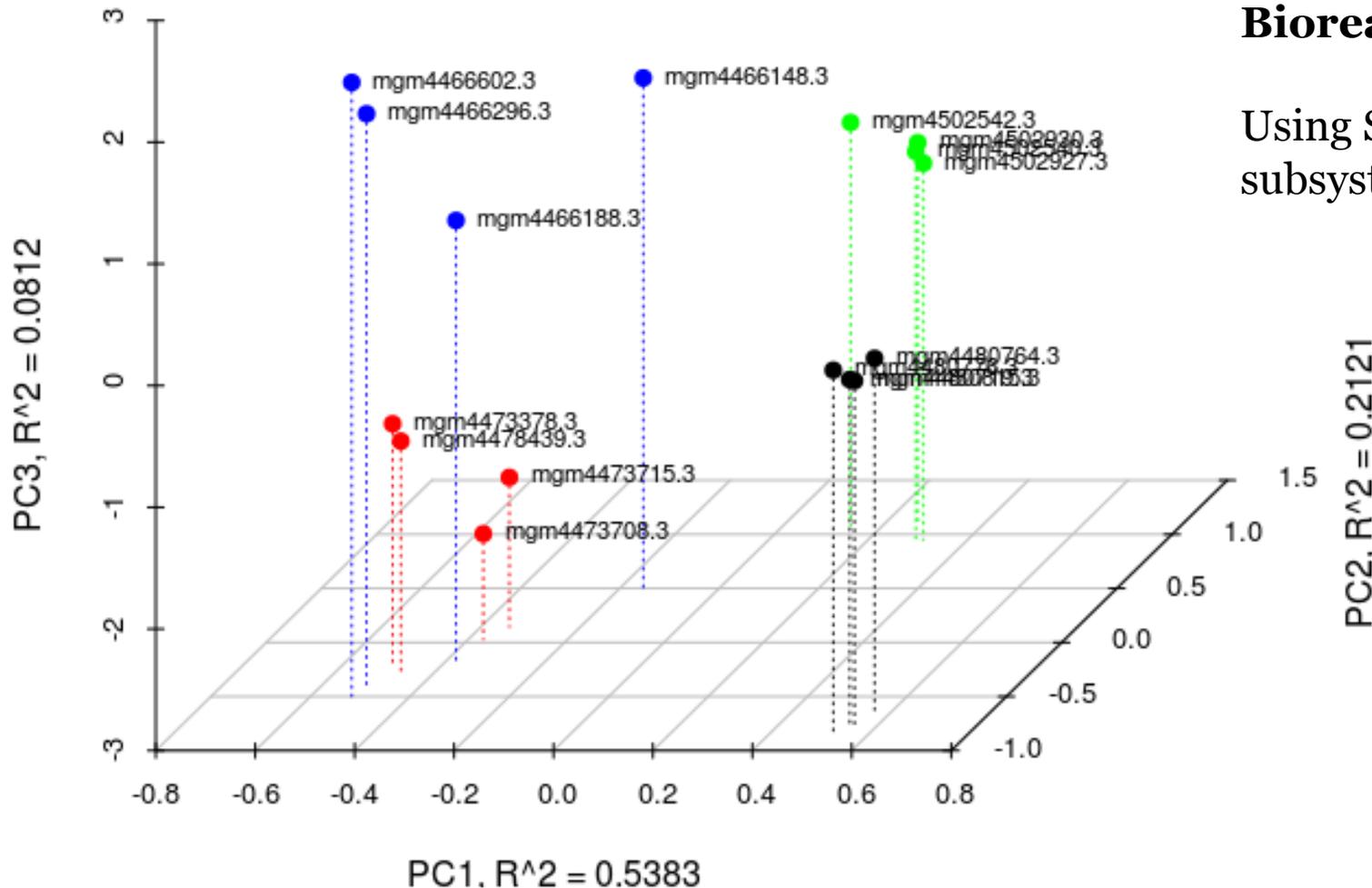
PCoA - four projects



Shotgun metagenomics works → Yeah!

Human microbiome
Feces
Soil
Bioreactor

Using SEED
subsystems Level 3



Or does it?

Human microbiome

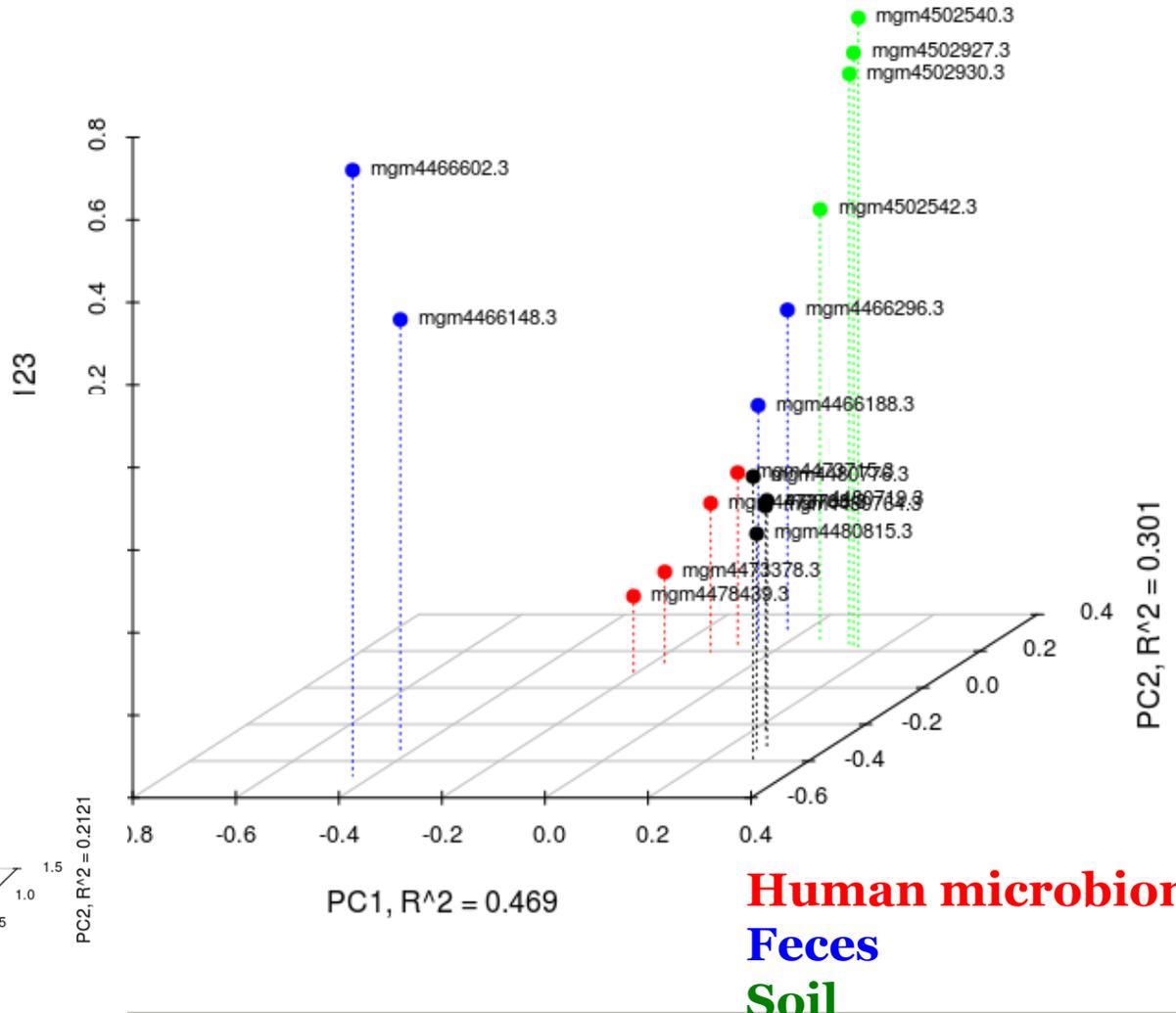
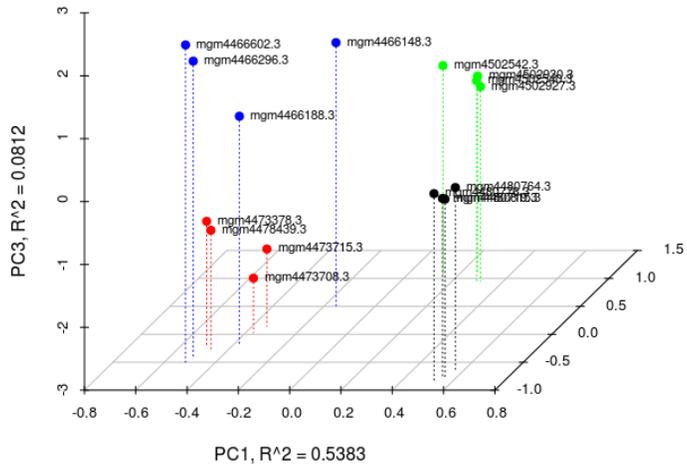
Feces

Soil

Bioreactor

KEGG instead of SEED

PCoA - four projects



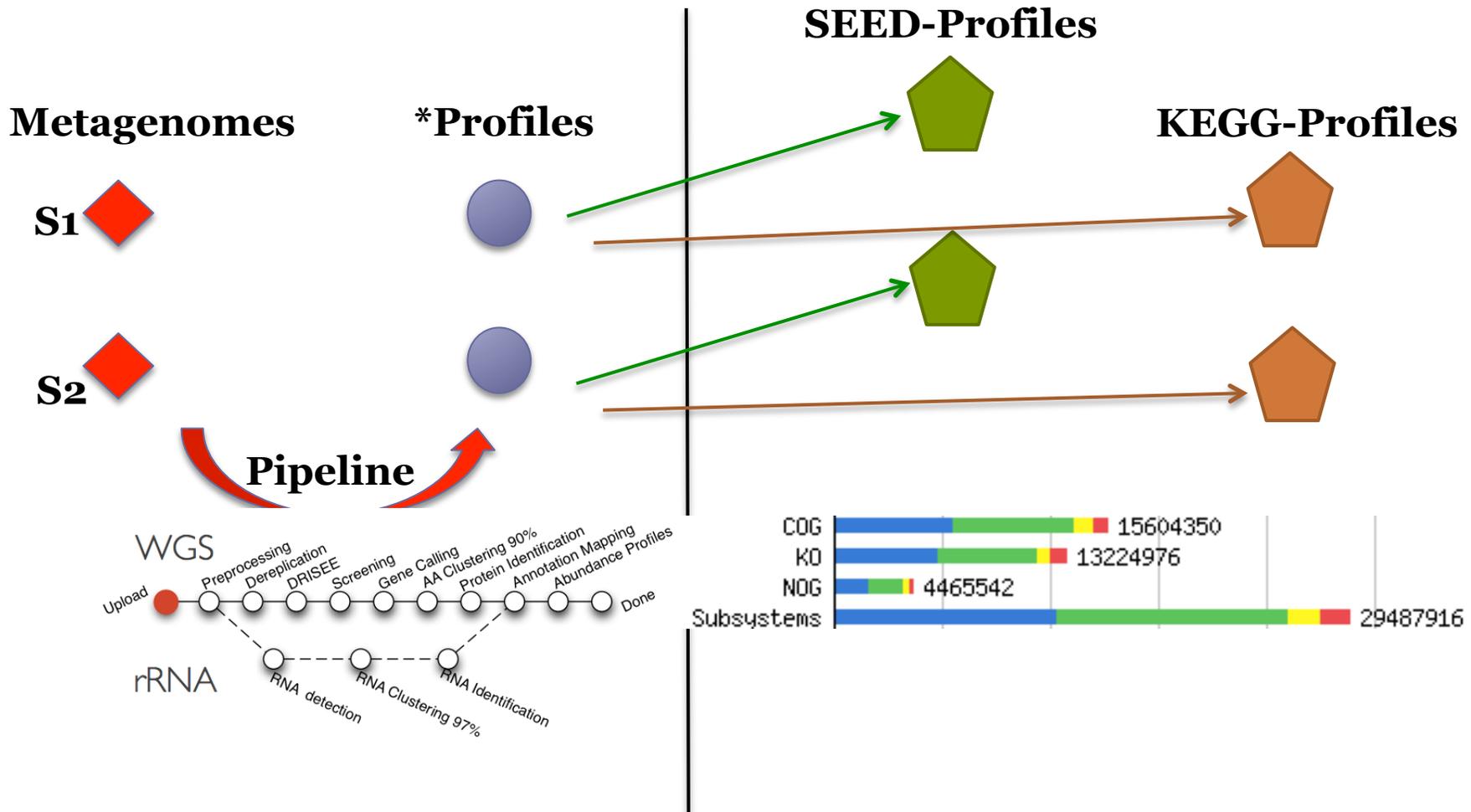
Human microbiome

Feces

Soil

Bioreactor

What has happened?





Introduction

Sequence quality

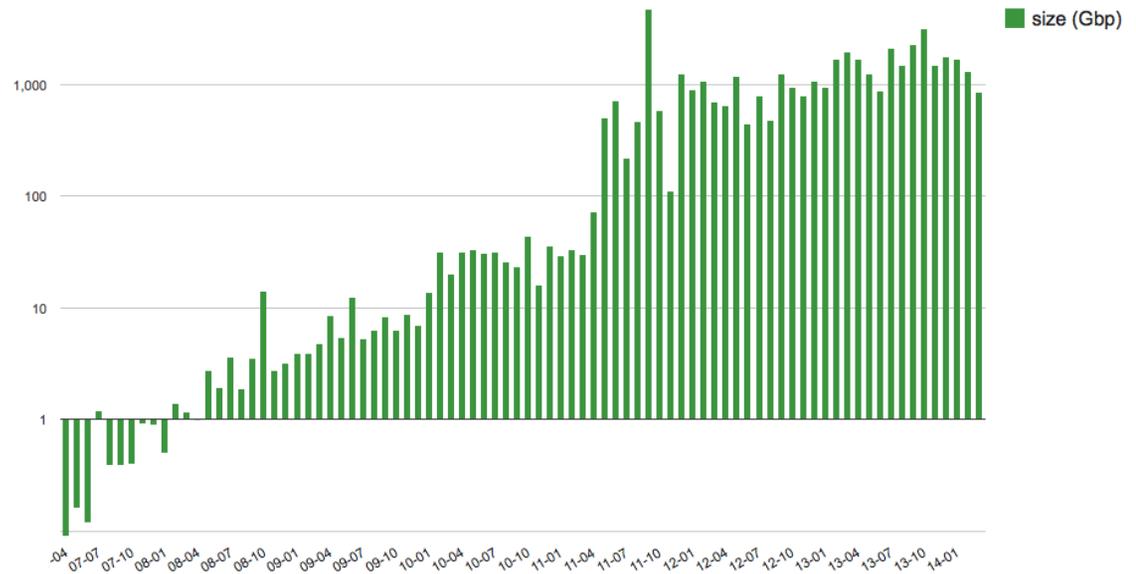
Assembly

Gene prediction

Functional annotation

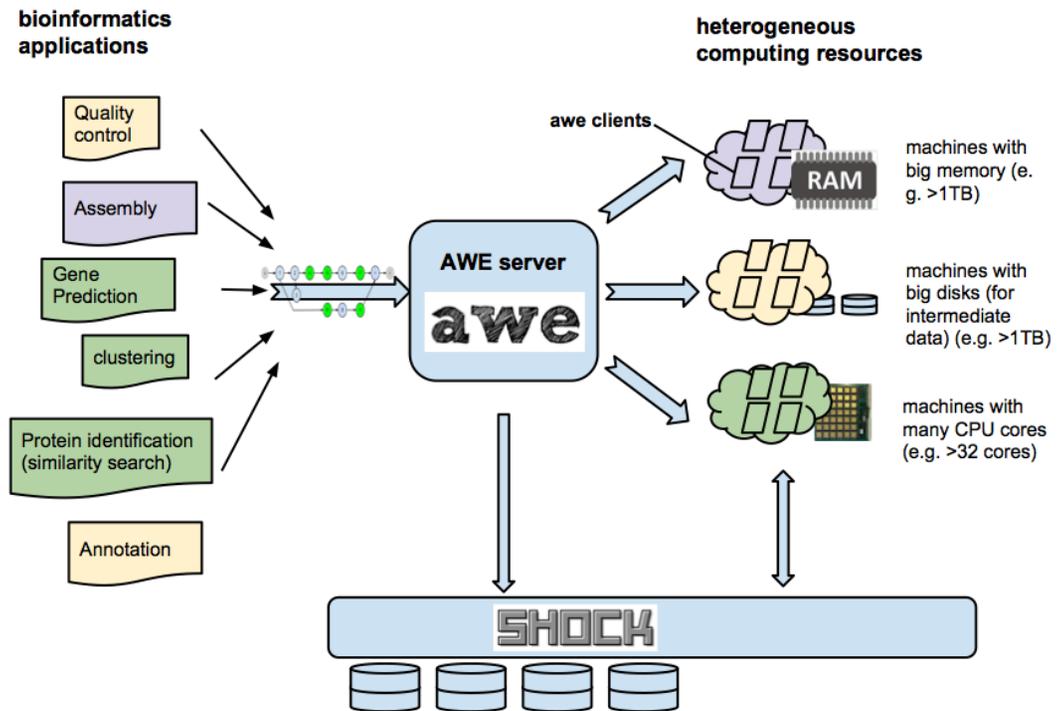
Current MG-RAST work

Today



From a capability problem to a capacity problem

built engineering tools accelerate platform



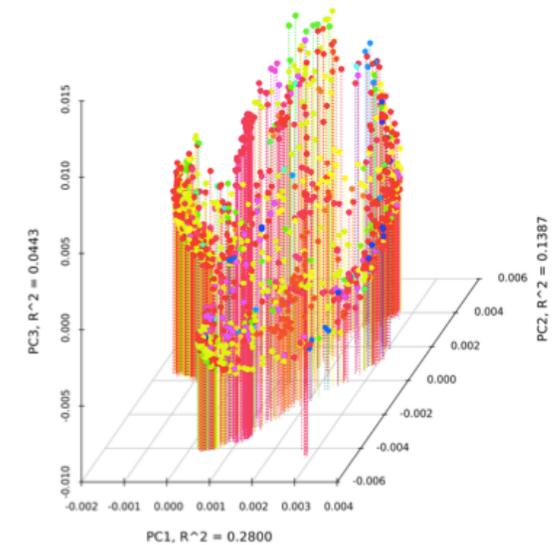
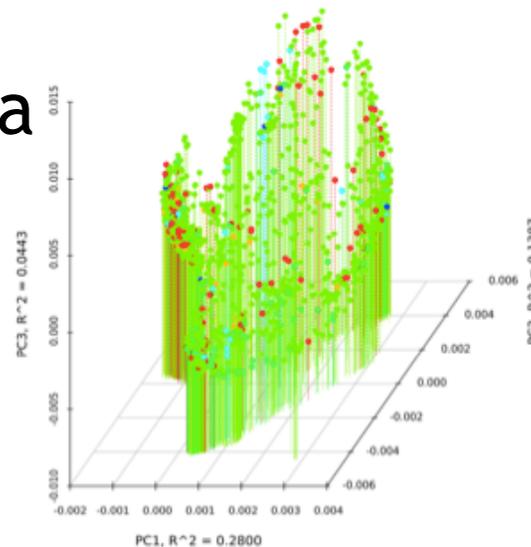
Common analysis tasks



- Assign taxonomy and function at scale
 - User sets parameters for annotation transfer at analysis time
- analyze a single sample
 - Taxonomy, rarefaction, function (if WGS)
 - K-mer profile, functional categories
 - QC tools
- Compare many samples wrt function and taxonomy
 - Extract subsets e.g. functions for set of species or species for a set of functions
 - Use subsets for comparison
- Compute normalizations, heatmaps, PCoA, tables

Example: large scale comparison

- Analysis of 2075 soil WGS data sets
- clustered by normalized subsystem abundance
- Painted with metadata



Example: scaling

- Use of computational appliances for metagenomic analysis
 - Appliances are self contained virtual machines (VMs)
 - Throughput **over 100GBp/day**, can scale much higher
 - Using KBase and non-KBase computers
- Have **third party users contribute computing** from their institution to the analysis of their data
 - e.g. UOregon (Jessica Green's group), OSU, Hudson Alpha, University of Oklahoma

- 
- Provide a unifying API
 - Move to Kmer based analysis whenever possible
 - Re-Build in Kbase using novel technology
 - Apply for big data analytics

Unifying API

- REST and JSON-RPC API
- Supports multiple levels of access
- R, Perl, Python, Java, ...

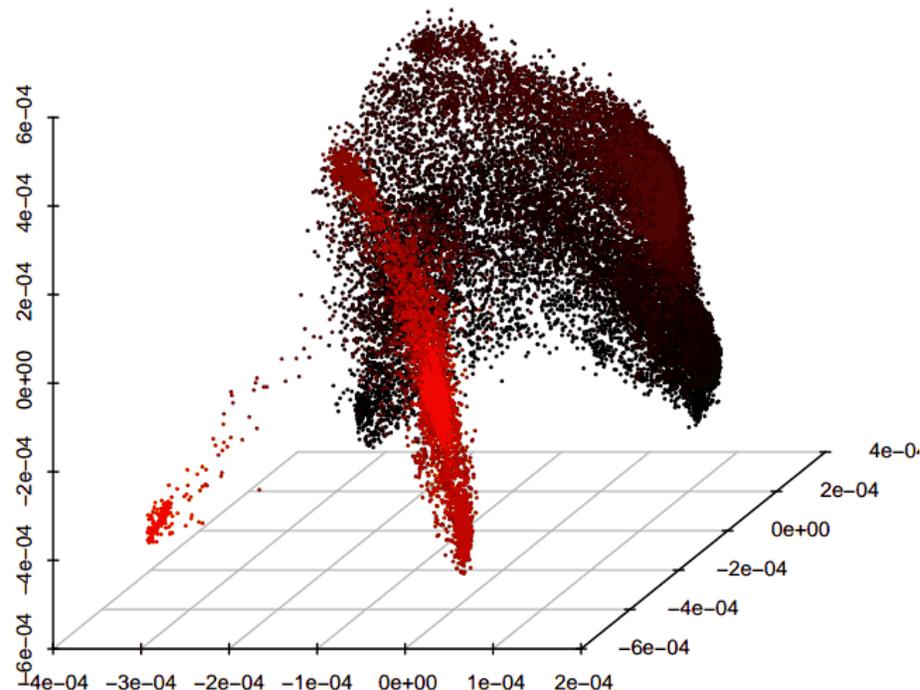
- Re-building Web-UI using API currently

- Opening up formerly closed application is liberating

Big data analytics in Bio:: compare many environmental samples

28k shotgun metagenomes
Normalized Subsystem level 3
Red=foreground

Using matR
R language client



Comparing metabolic models of a microbial community using different classes of sequence data and computational methods.

- Metabolic models can help understand interactions and dependencies of community members
- Identification of cultivation conditions for “unculturable” members enables use of additional research techniques, e.g. wet lab
- Deeper understanding of microbial communities may allow to predict the effect of modifications in culture conditions
- Possible applications in industry and medicine (e.g. optimization of the production of biofuels and biopharmaceuticals)
- Different various methods produce differing results

→ VERY MUCH WORK IN PROGRESS ←

From reads to taxa / functions

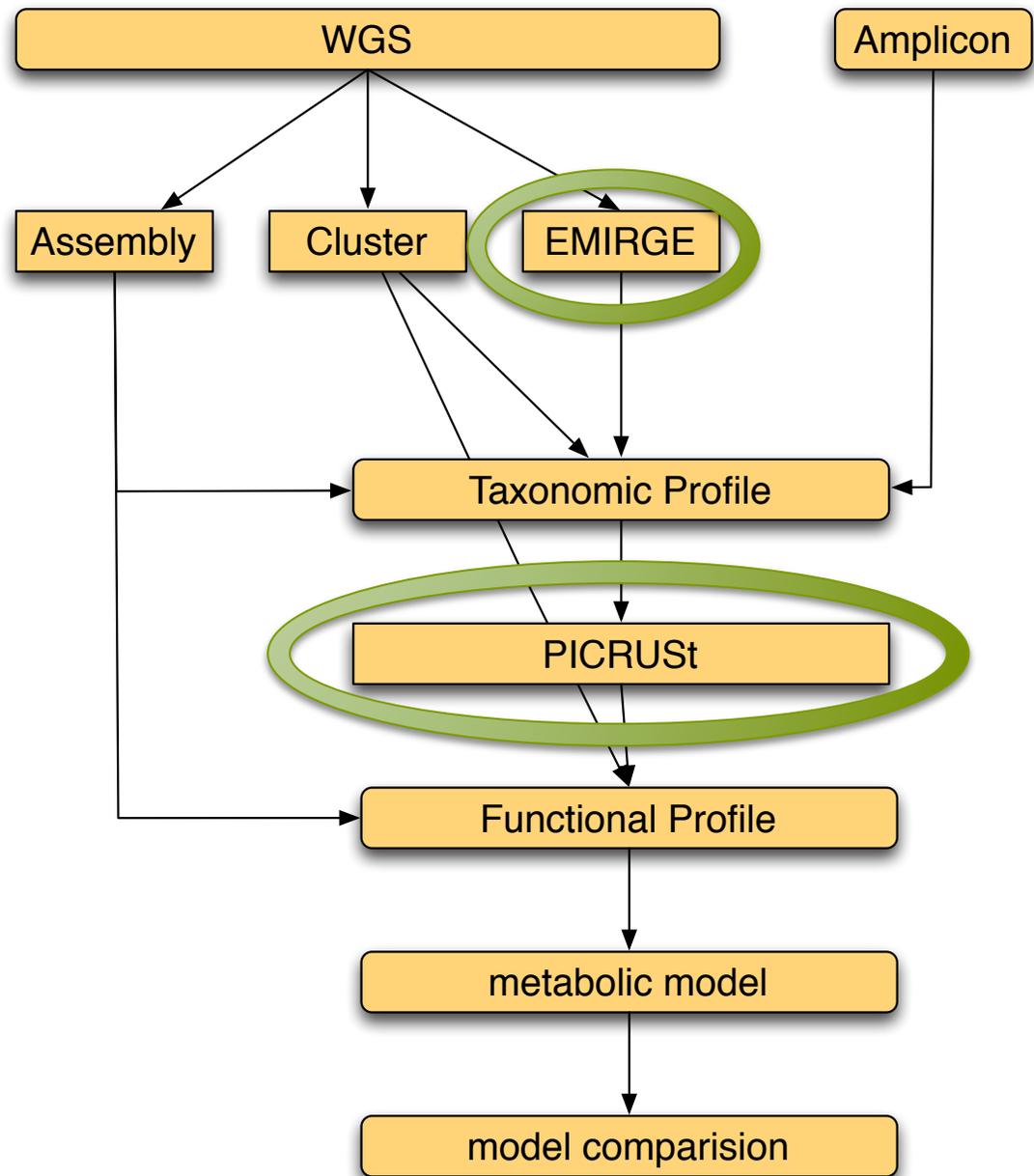
- Functional genes
 - From assembled contigs (e.g. velvet, meta-velvet, KiKi, IDBA-UD,)
 - From assembled contigs and singletons
 - From clustering
- 16s ribosomal genes
 - From assembled contigs (see above)
 - From assembled contigs and singletons
 - From clustering
 - From EMIRGE (Miller et al, Genome Biology, 2011)

From taxa / reads to models

- Functions derived from observed genes
 - using Kbase Communities pipeline (based on Meyer *et al*, BMC Bioinf. 2008)
- Functions inferred via 16s proxy
 - using PICRUSt (Langille, *et al* Nature Biotechnology, 2013)
- ➔ Next ModelSEED to create model and fill gaps
 - (Henry, *et al* Nature Biotechnology, 2010)

For all steps: lengthy, expensive computational tasks

- **Input:** environmental shotgun sequences
- **Output:** several models of community metabolism



Taxonomic reconstructions on the same data

treatment	GBp	sequences	# 16S	# OTUs-97%	ID
16S amplicon (Meyer et al, BMC Bioinformatics 2008)	0.005	20,913	3,650	2269	mgm4506684.3
Clustering (Meyer et al, BMC Bioinformatics 2008)	4.763	47,165,484	630,700	15999	mgm4509402.3
EMIRGE (Miller et al, Genome Biology, 2011)	4.763	47,165,484	20	4	mgm4509402.3
KiKi (unpublished)	0.604	3,400,711	429,985	1021	mgm4512893.3
Khmer (Pell, et al PNAS 2012)	1.208	3,652,356	14	n/A	mgm4510006.3

Creates mixed bag metabolic model

View Metabolic Model Details - Output 14:50:58, 2/10/2014

mgm4506694.3.otu.picrust.kegg2ss.anno.raw.001
1917 reactions, 1556 compounds, 1 gapfill runs

[Reactions](#)
[Compounds](#)
[Compartment](#)
[Biomass](#)
[Gapfill](#)
[Gapgen](#)

Show entries Search all:

Id (compartment)	Name	Equation	Genome Features Mapped to this Reaction	
rxn00001 (c0)	Pyrophosphate phosphohydrolase	$\text{H}_2\text{O}[\text{c0}] + \text{PPi}[\text{c0}] \Rightarrow (2)\text{H}^+[\text{c0}] + (2)\text{Phosphate}[\text{c0}]$	kb mganno.167.tail.0.g.0.peg.6652 kb mganno.167.tail.0.g.0.peg.4886 kb mganno.167.tail.0.g.0.peg.6194	rxn00001 (c0)
rxn00002 (c0)	Urea-1-carboxylate amidohydrolase	$(3)\text{H}^+[\text{c0}] + \text{Allophanate}[\text{c0}] + \text{H}_2\text{O}[\text{c0}] \Rightarrow (2)\text{NH}_3[\text{c0}] + (2)\text{CO}_2[\text{c0}]$	kb mganno.167.tail.0.g.0.peg.603 kb mganno.167.tail.0.g.0.peg.177 kb mganno.167.tail.0.g.0.peg.5928	rxn00002 (c0)
rxn00006 (c0)	hydrogen-peroxide:hydrogen-peroxide oxidoreductase	$(2)\text{H}_2\text{O}_2[\text{c0}] \Rightarrow (2)\text{H}_2\text{O}[\text{c0}] + \text{O}_2[\text{c0}]$	kb mganno.167.tail.0.g.0.peg.4256	rxn00006 (c0)
rxn00007 (c0)	alpha,alpha-Trehalose glucosyltransferase	$\text{H}_2\text{O}[\text{c0}] + \text{TRHL}[\text{c0}] \Leftrightarrow (2)\text{D-Glucose}[\text{c0}]$	kb mganno.167.tail.0.g.0.peg.348 kb mganno.167.tail.0.g.0.peg.3887	rxn00007 (c0)
rxn00010 (c0)	Glyoxylate carboxy-lyase (dimerizing)	$\text{H}^+[\text{c0}] + (2)\text{Glyoxalate}[\text{c0}] \Rightarrow \text{Tartronate semialdehyde}[\text{c0}] + \text{CO}_2[\text{c0}]$	kb mganno.167.tail.0.g.0.peg.5640	rxn00010 (c0)

Showing 1 to 5 of 1,917 entries
[First](#)
[Previous](#)
[1](#)
[2](#)
[3](#)
[4](#)
[5](#)
[Next](#)
[Last](#)

KBase has UIs for: EMIRGE, PICRUST, ModelSEED, ...

e.g.
Gapfilling the models

Gapfill Metabolic Model ... Last run: 10:27:27, 2/11/2014

Fill in missing core metabolism functions in a draft model.

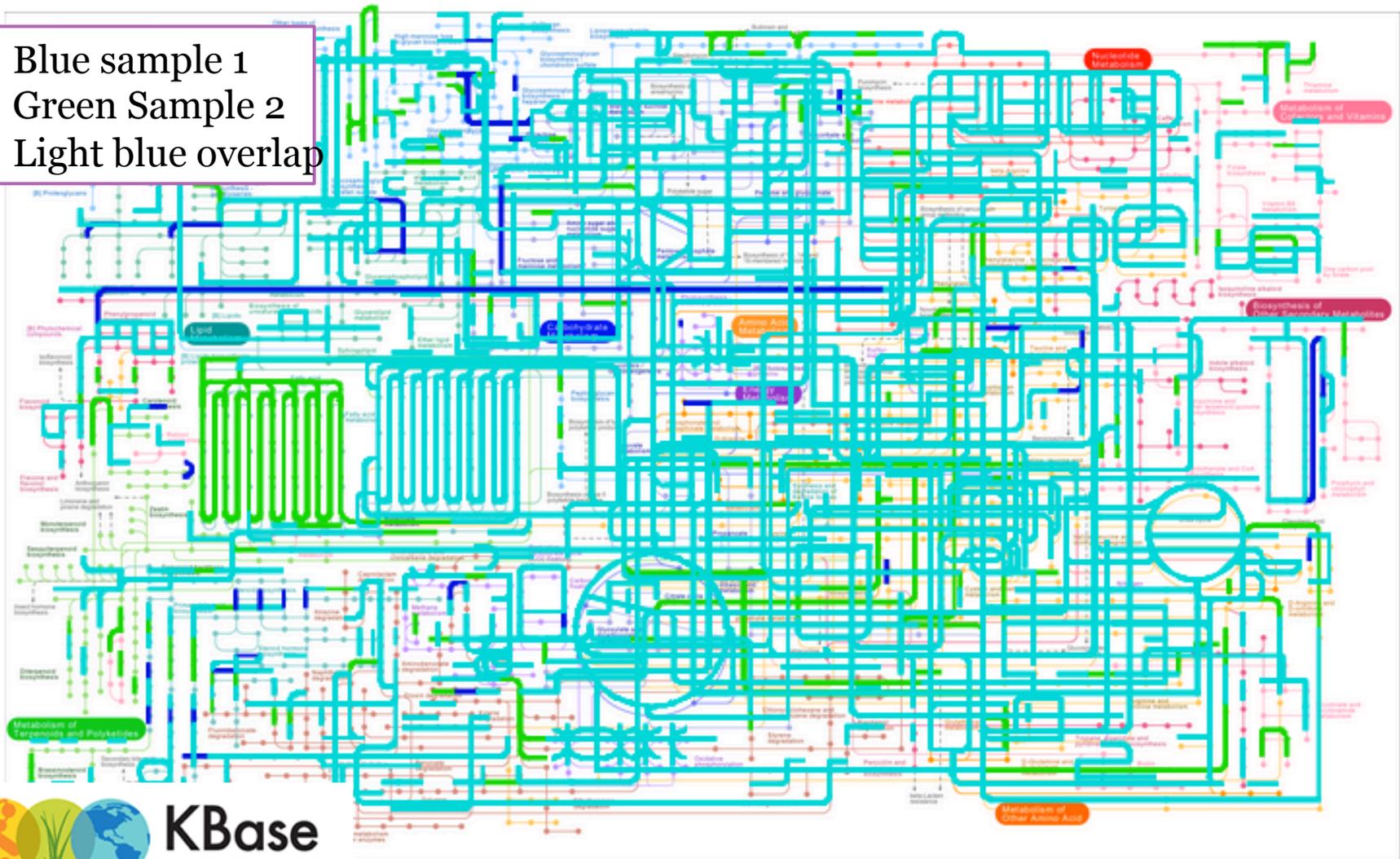
Workspace	<input type="text"/>	name of workspace, default is current
Model Name	<input type="text" value="mgm4506684.3.otu.picrust.kegg2ss.anno.raw.001"/>	workspace ID of model

 Delete Run

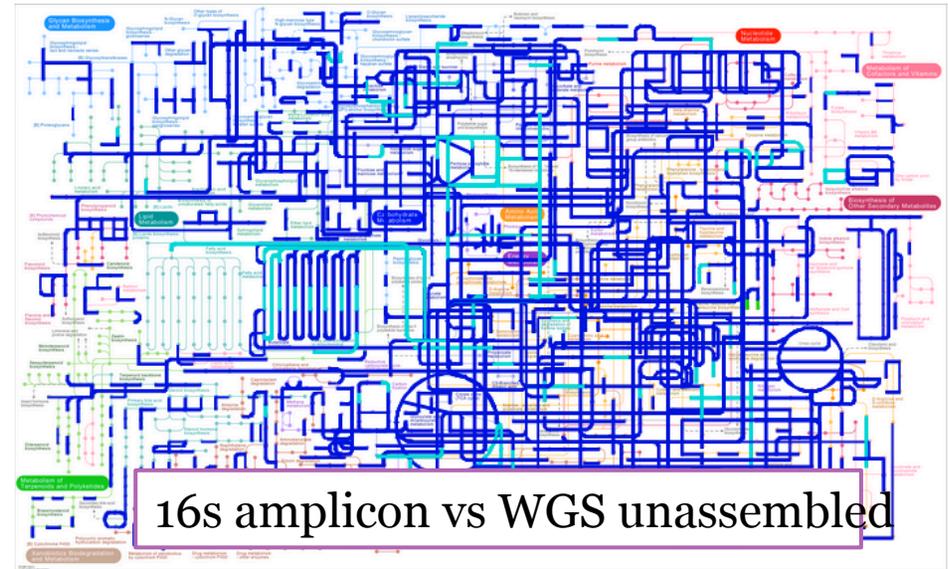
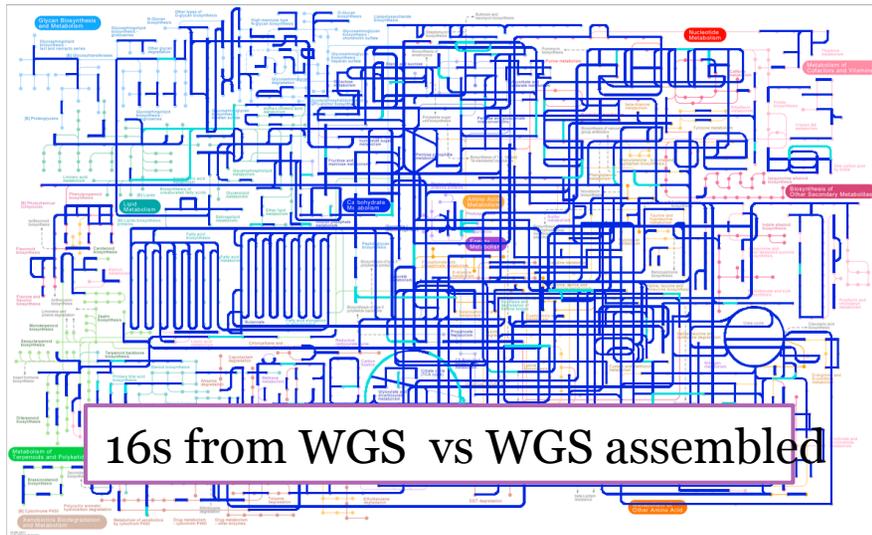
Step 2 / 2: Gapfill Model Starting

Visual Comparison two different samples

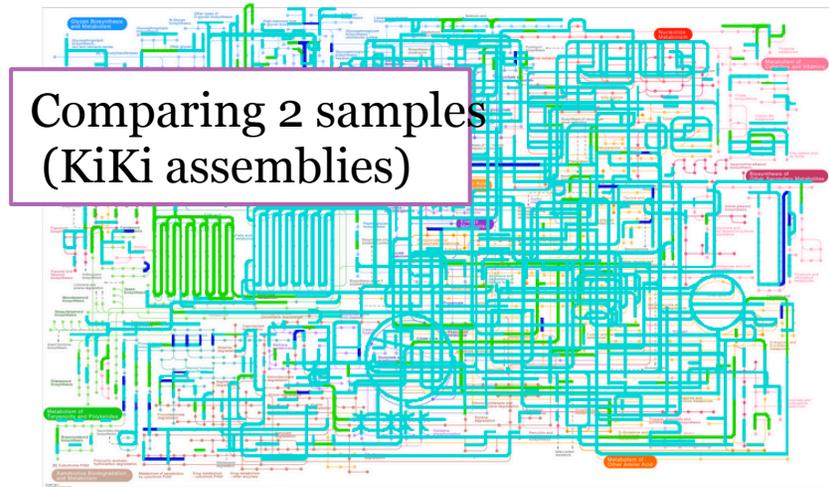
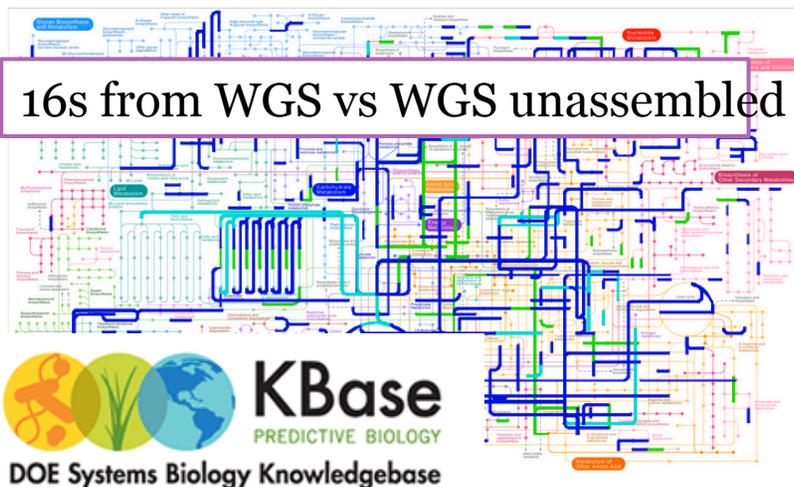
Blue sample 1
Green Sample 2
Light blue overlap



Comparing multiple technologies



Data from Steve Allison, Princeton,



Acknowledgements

Team:

- Jared Bischof
- Dan Braithewaite
- Adina Chuang-Howe
- Narayan Desai
- Mark d' Souza
- Katya Drybinski
- Elizabeth M. Glass
- **Wolfgang Gerlach**
- Travis Harrison
- **Kevin Keegan**
- Tobias Paczian
- Hunter Matthews
- **Wei Tang**
- **Will Trimble**
- **Andreas Wilke**
- Jared Wilkening

Collaborators:

- D. Antonopoulos (Argonne)
- A. Arkin (Berkeley)
- E. Chang (UChicago)
- Dawn Field (Oxford)
- F.-O. Glöckner (MPI Bremen)
- Jack Gilbert (Argonne)
- Jeff Grethe (CalIT2, CAMERA)
- Sarah Hunter / Guy Cochrane (EBI)
- Ken Kemner (Argonne)
- Rob Knight (Colorado)
- Nikos Kyrpides (DOE JGI)
- J. Tiedje (MSU)
- Owen White (UMaryland, HMP DACC)

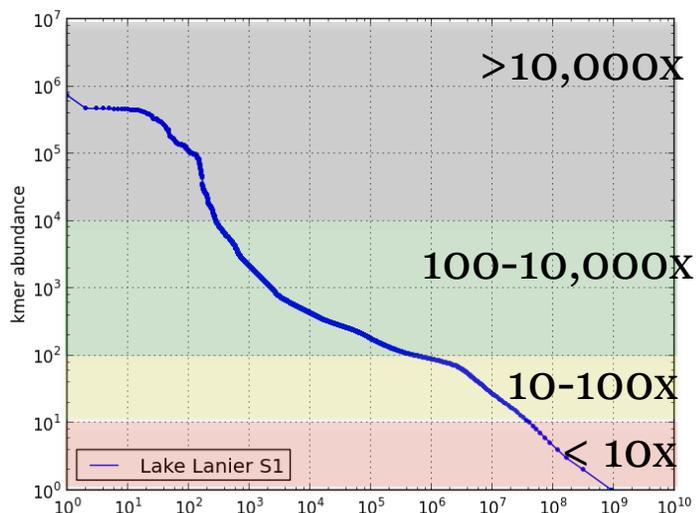


Thank you





Stratifying kmers - heterogeneous population of underlying reads

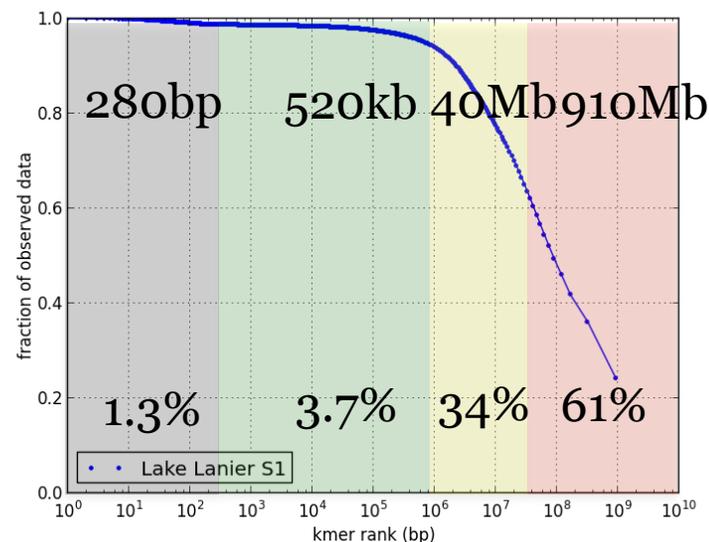
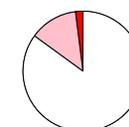
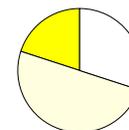


Example: Lake Lanier S-1

expect *but can't prove* mostly complete

expect mostly fragmented / incomplete

expect assembly is hopeless



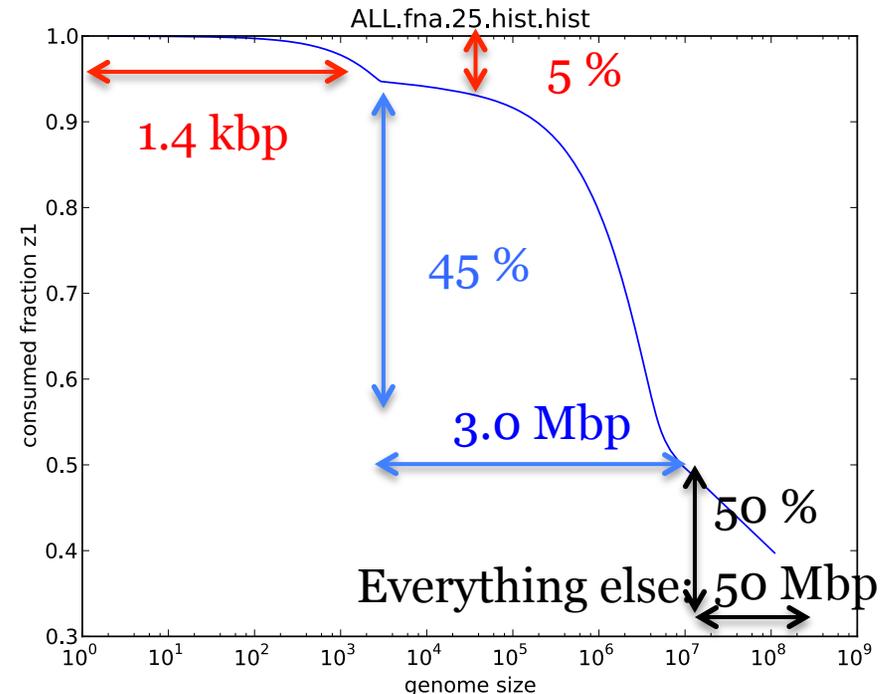
Average completeness must be some sort of weighted average of the completeness of the components.

I can deliver estimates of approximate total genome size(s) in each coverage band.

Decadal coverage bands (1-9, 10-99, 100-999) look go

Predicting replicons before assembly

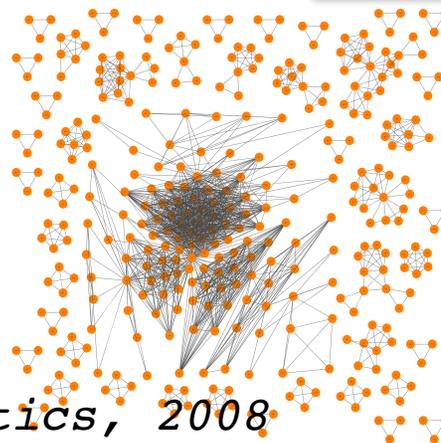
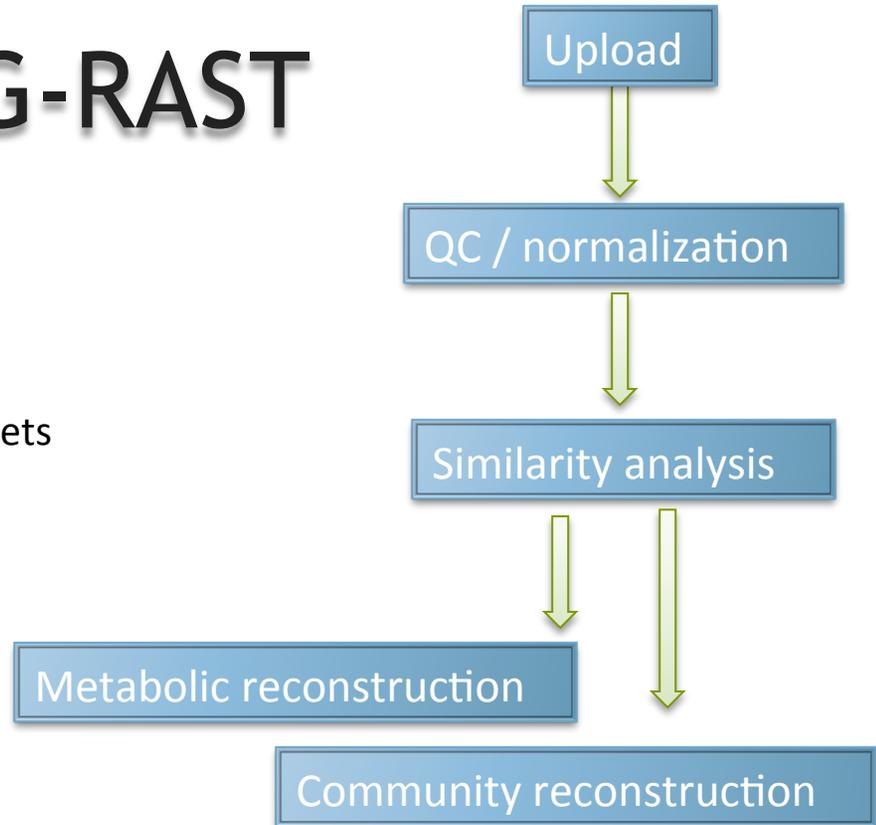
- Using **K-mer spectra** to predict (pan-) genome size
 - K-mer= unique word, easily computed
- In addition to alpha diversity
 - 300 OTU data set
- Using k-mer size 25
- **Red** and **blue** replicons were missing in assembly
 - Allows adjustment of parameters



From: Williams et al, BMC Genomics, 2013

Brief history of MG-RAST

- December 2007 (v1)
 - 100+ groups and ~250 data submitters
 - 100+ data sets, ~10+ GBp total size
- October 2009 (v2)
 - **Pre-publication sharing available**
 - ~1500 data submitters, ~300 public data sets
 - 6000+ data sets
 - 200+ GBp total data set size
 - About ~30 GBp/month throughput
- March 2011 (release v3)
 - 2500+ data submitters
 - ~2000 public data sets
 - 25,000 data sets total
 - Throughput:
 - 47GBp in 24h
 - 3000 submissions in 24h
- March 2012 (v 3.1.2)
 - 12 Terabasepairs (10^{12} bp)
- May 2012 (3.20)
 - 13.8 TBp (10^{12} bp)
 - 128 billion sequences



Metabolic Reassembly 1.0	ID: 441208.3	ID: 441162.3	ID: 441148.3	ID: 441078.3
Amino Acids and Derivatives	1000	1000	1000	1000
Carbohydrates	1000	1000	1000	1000
Cell Division and Cell Cycle	1000	1000	1000	1000
Cell Wall and Capsule	1000	1000	1000	1000
Clustered Gene Neighborhoods	1000	1000	1000	1000
Co-factors, Vitamins, Prosthetic Groups, Signaling	1000	1000	1000	1000
DNA Metabolism	1000	1000	1000	1000
Fatty Acids and Lipids	1000	1000	1000	1000
Macromolecular Synthesis	1000	1000	1000	1000
Metabolism	1000	1000	1000	1000
Metabolism of Aromatic Compounds	1000	1000	1000	1000
Metabolism Transport	1000	1000	1000	1000
Motility and Chemotaxis	1000	1000	1000	1000
Message Metabolism	1000	1000	1000	1000
Nucleosides and Nucleotides	1000	1000	1000	1000
Protein Metabolism	1000	1000	1000	1000
Protein Folding	1000	1000	1000	1000
Protein Synthesis	1000	1000	1000	1000
Protein Transport	1000	1000	1000	1000
Regulation and Cell Signaling	1000	1000	1000	1000
Secondary Metabolism	1000	1000	1000	1000
Storage	1000	1000	1000	1000
Structure	1000	1000	1000	1000
Transcription	1000	1000	1000	1000
Translation	1000	1000	1000	1000
Unclassified	1000	1000	1000	1000
Vitamins	1000	1000	1000	1000

Meyer et al., BMC Bioinformatics, 2008

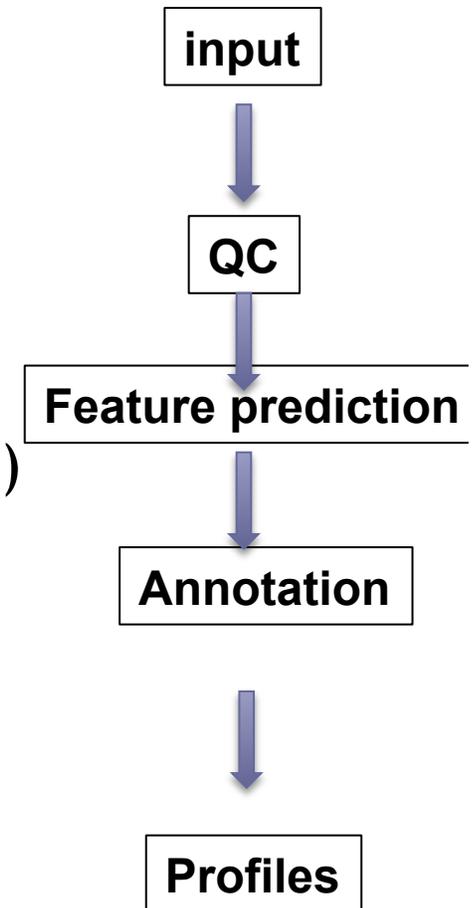
simplified

MG-RAST in a slide

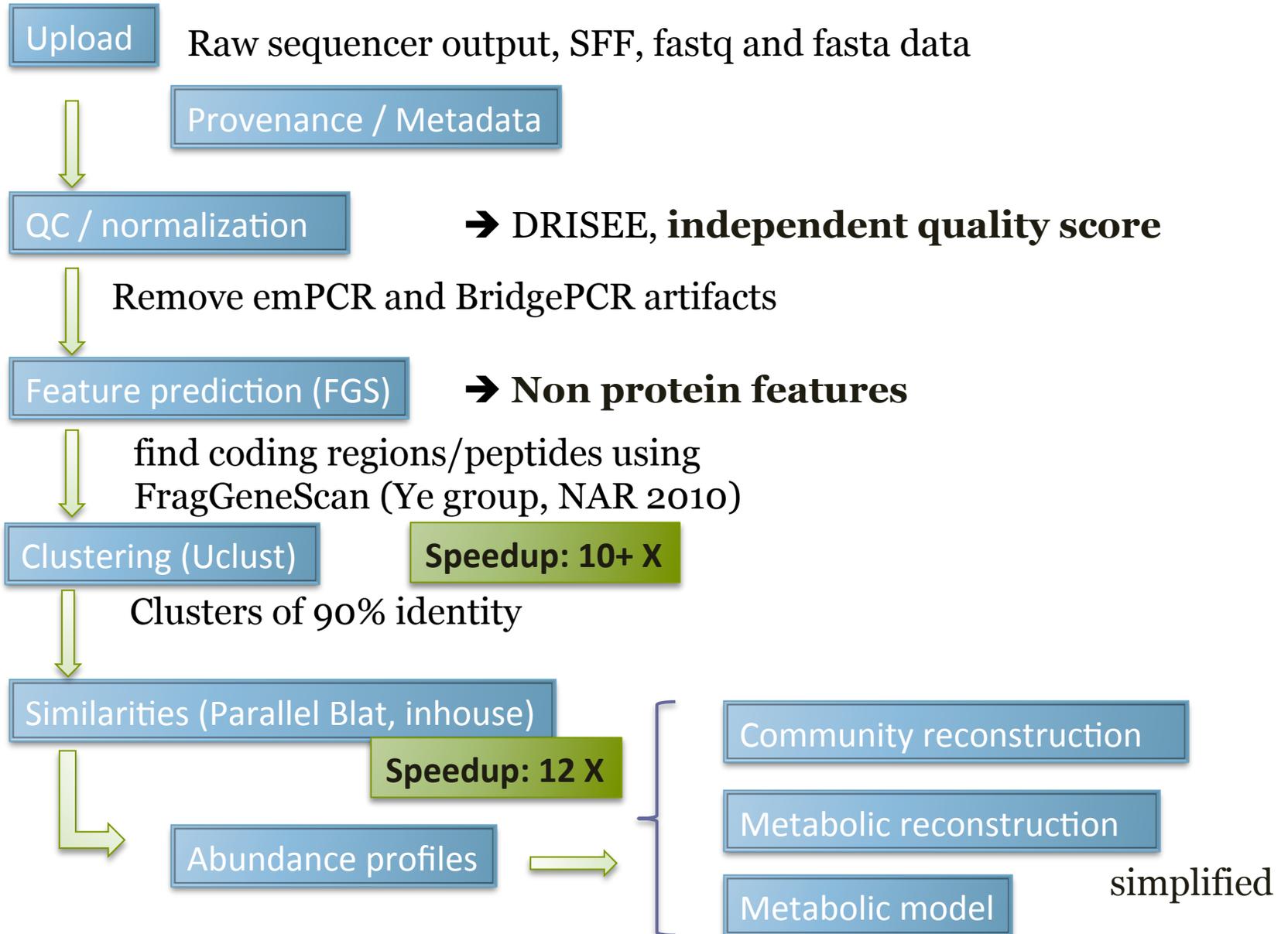
- 80,000 data sets, 10,000+ public
 - 27TBp analyzed, 240 billion sequences
- Normalized data from many groups
- Intensive QC (innovative QC methods)
- Mapping against known annotations
 - Using many sources (SEED, KEGG, COG, NOGs, ...)
- Allow comparison and meta-analysis
 - Require metadata → partnership with GSC
- Creates abundance profile
 - Counting occurrence of taxa, phyla, ...

and

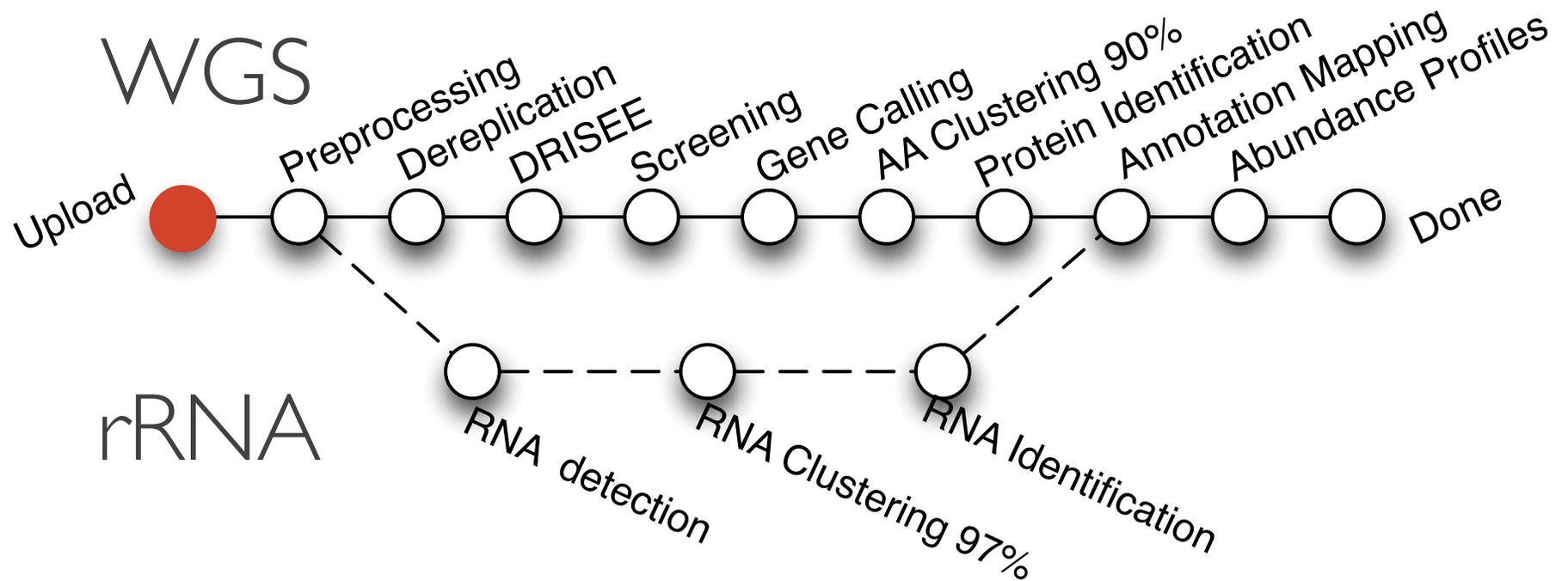
 - Count functions (e.g. subsystems, KEGG pathways, etc.)



Improved Algorithm



Pipeline Overview



MG-RAST is an automated pipeline for annotating individual NGS-sequencing reads, transforming raw sequences into *meaningful* output











Access to MG-RAST

- via web-UI (<http://metagenomics.anl.gov>)
- ftp/http downloads
- Web services API [BETA]
 - <http://api.metagenomics.anl.gov>
 - C, C++, Java, Perl, Python, Ruby, Javascript, ...
 - RESTful, JSON objects
- R client (matR) [BETA]

Pipeline changes?

- MGA vs FGS?
 - MGA misses 70+%
- BLAST vs BLAT?
 - BLAST is more sensitive
- What is the effect of a different protein database?
 - E.g. TrEMBL
- What is the effect of using protein motifs?
 - E.g. interpro?
- **Assembly is the biggest unknown...**

Lessons learned from GOS

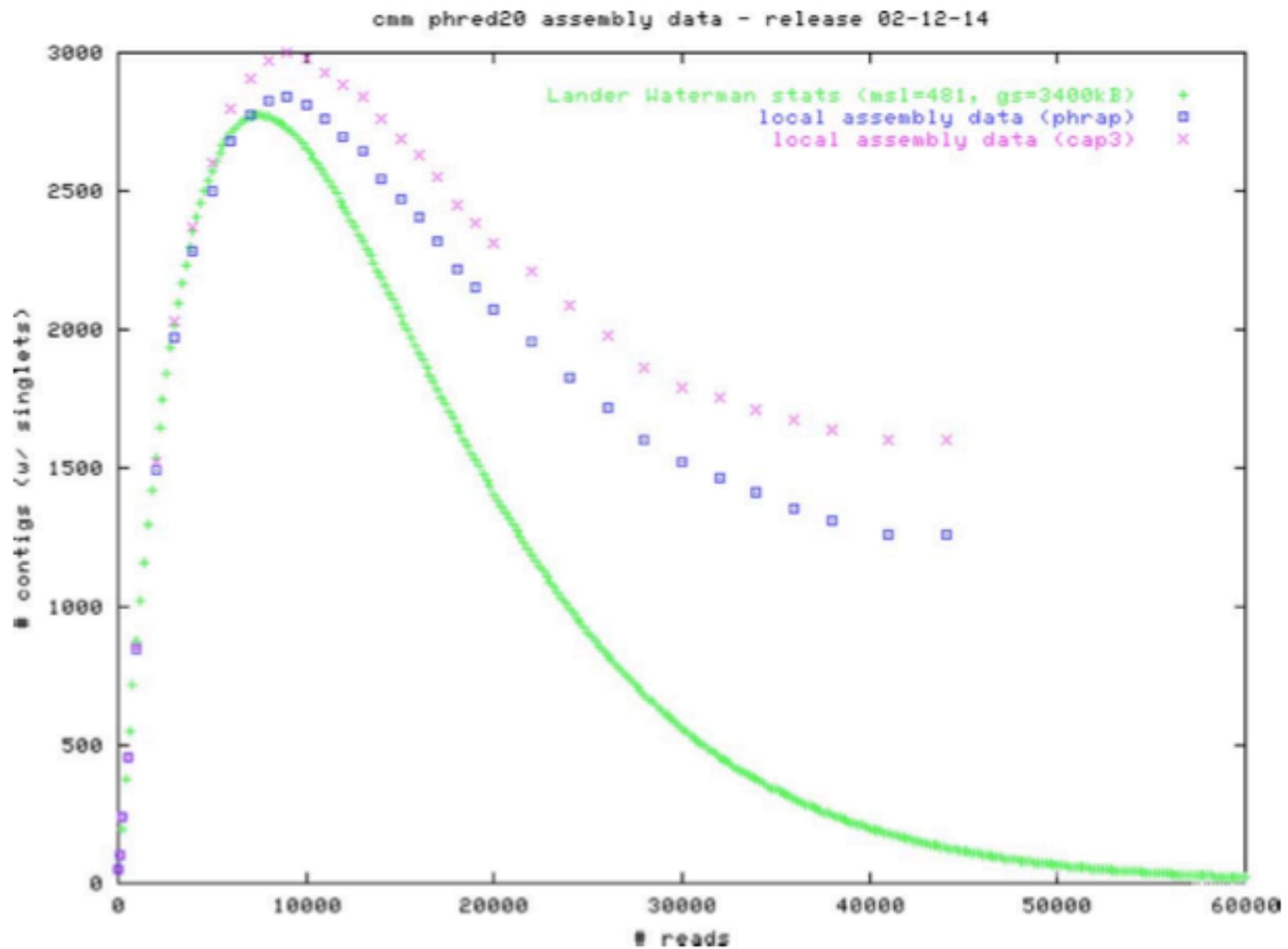
- Old style:
 - “Lets sequence as much as we can afford”
 - “Metagenomics is like genomics”
- Today:
 - Often 16s amplicon study first
 - replicates (biological and technical)
 - “design for statistics”
 - “replicate or lie” (Jim Prosser)
 - metadata
 - Genomics Standards Consortium provides tools
 - Provide good QC
 - Identify signal vs. noise ratio
 - Throw away bad data when needed (!)
 - Identify appropriate analysis workflow
 - Perform assembly?

→ Design for statistics

**→ Metadata
(r)evolution**

→ Data hygiene

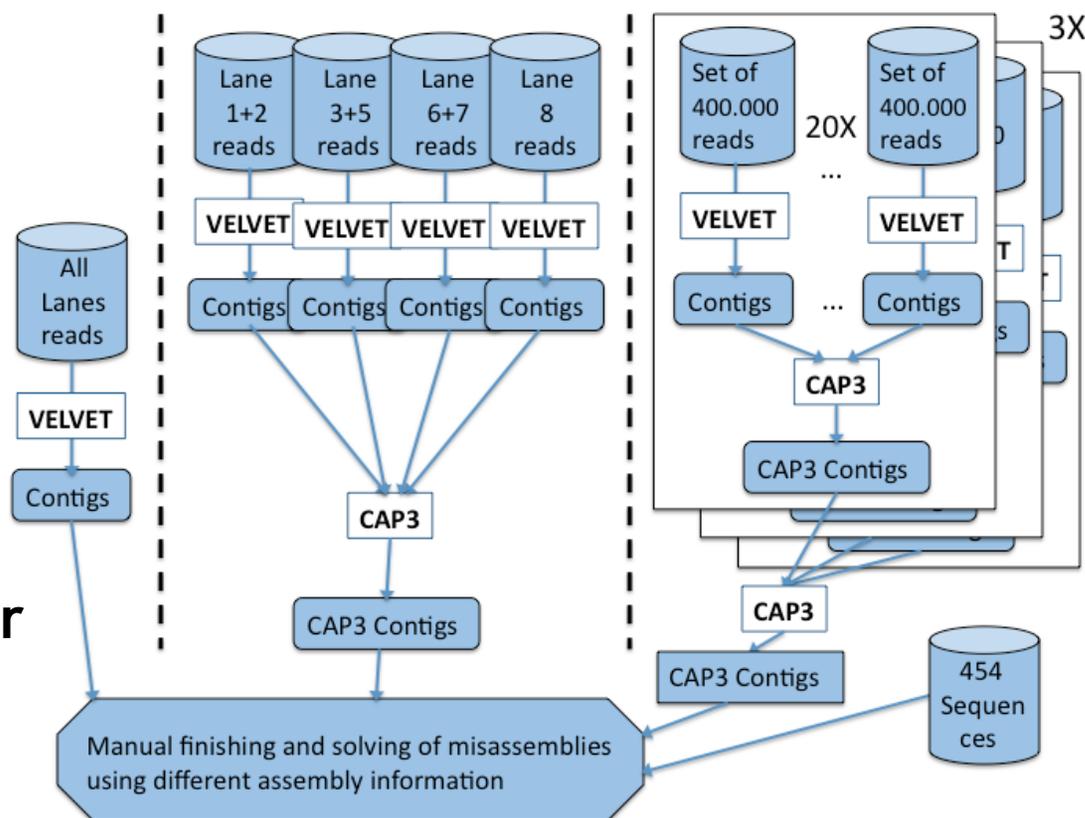
→ Tool chain matters



Kaiser et al, J. Biotech, 2003

Our subsampling pipeline

- Build multiple sets of reads
 - All
 - 2 lanes at a time
 - Sets of 400k reads
- Run velvet on each set
 - K-mer= 51
 - K-mer=17
- Initially used cap3 to integrate between sets
- replaced cap3 with **newbler**
 - added 50% to contig yield!
- Class of small repeats
100-200bp requires longer reads (using 454)
 - Manually integrated

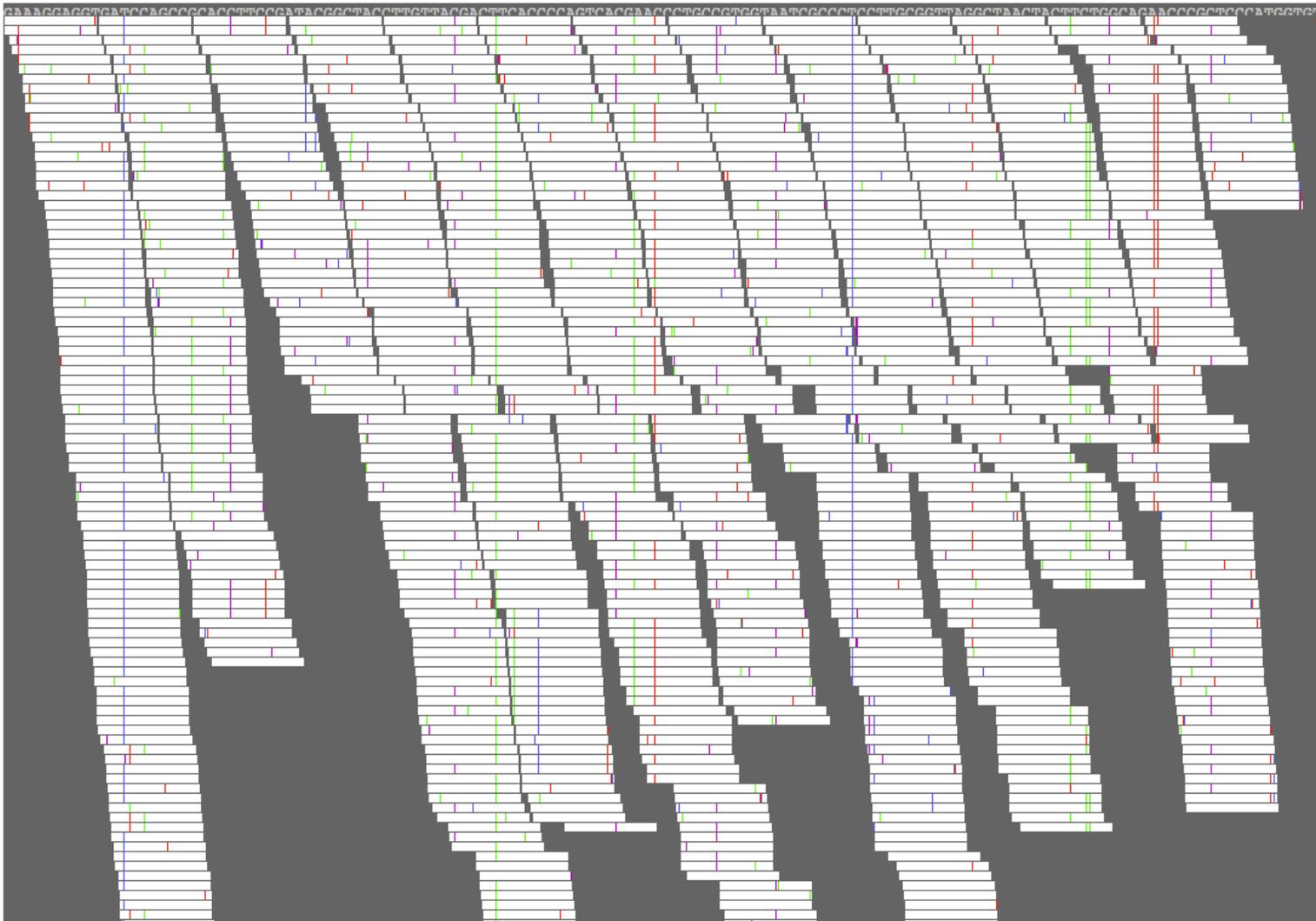




What does having 45 million reads give us?

- GOS assumed nature was *pseudo* clonal
- In GOS assembly reduced diversity to consensus sequences
- Now...

There are many strains in the data

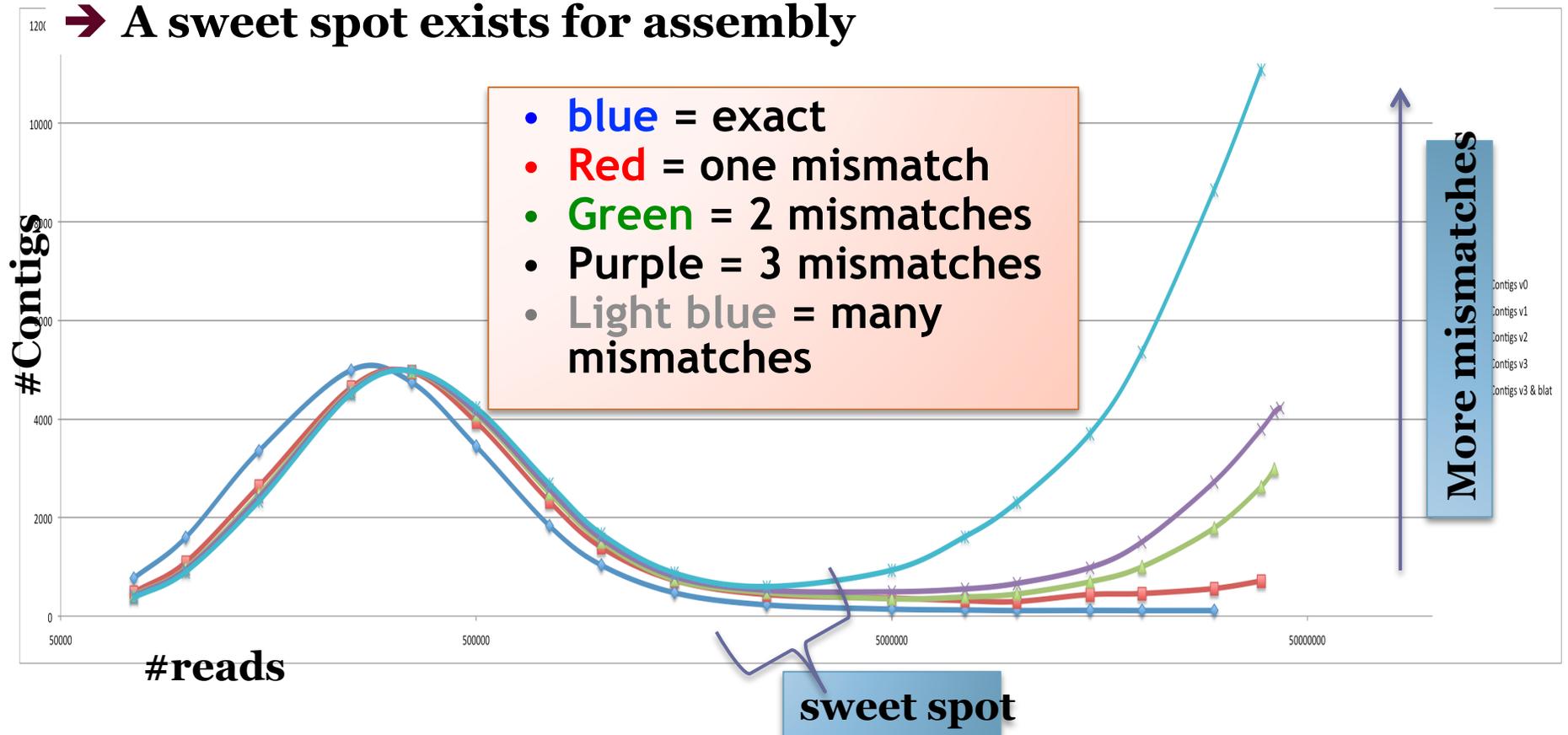


Effect of strain variation on assembly

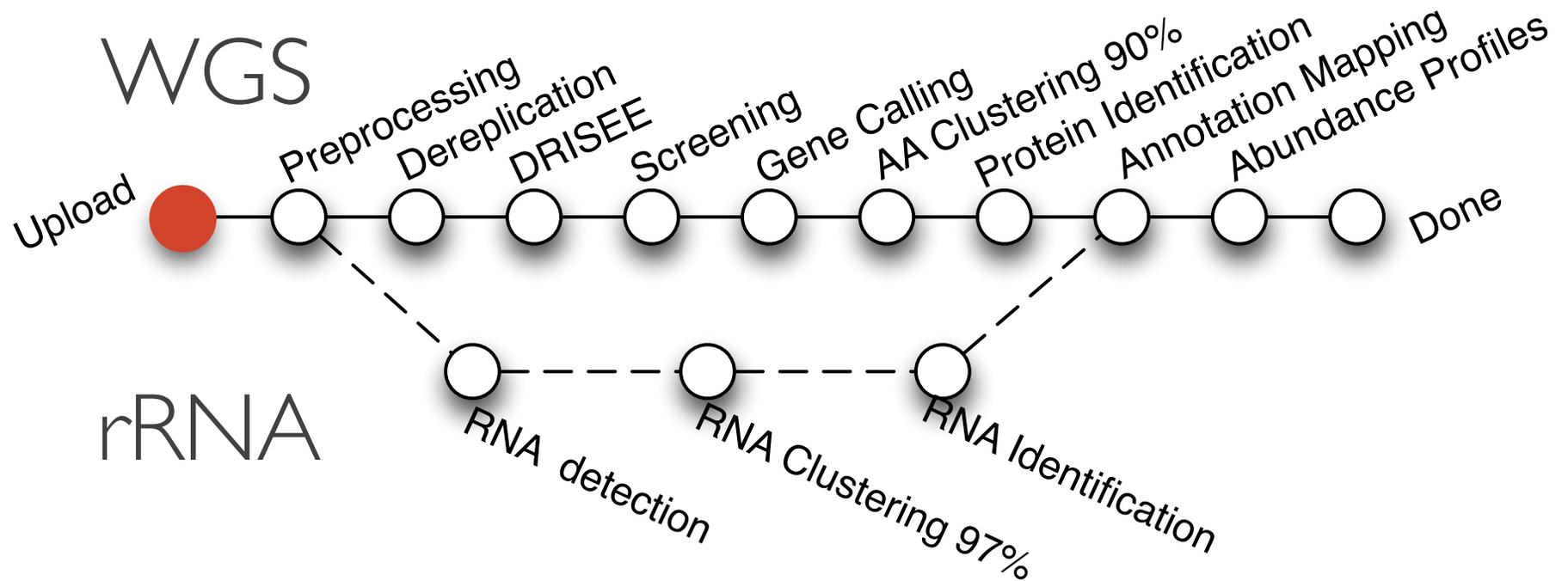
- we created several artificial subsets using the genome as reference
 - allowing more and more variation (repeated with bootstrapping)

➔ **Confirms strain variation breaks assembly tools**

➔ **A sweet spot exists for assembly**



Pipeline Overview



MG-RAST is an automated pipeline for annotating individual NGS-sequencing reads, transforming raw sequences into *meaningful* output