

# Bacterial ID by MALDI using DNA databases

SimulTOF Systems  
Marlborough, MA

Ken Parker

[Kenneth.parker@simultof.com](mailto:Kenneth.parker@simultof.com)

# History

- Mass spec has been used for bacterial typing for decades
  - starting in the **1970s**, originally via pyrolysis of lipids
- **1996**: MALDI was used to type bacteria in protein m/z range
  - ribosomal proteins correspond to many of the strongest masses
- **1997**: Pineda et al. -> matching organisms to ribosomal protein mass lists
- Since then, database look-up has largely been ignored
  - Except for selected clades (see next slide)
- Instead, **FDA approved methods** to identify bacteria via **spectral matching**
  - requires **careful curation of library spectra**
- In this talk, MALDI -> bacterial ID to DNA databases described
  - Want to extend MALDI market beyond pathogen ID
- Can easily search **10,000** strains / species
  - no limit yet.

# Discrimination of *Escherichia coli* O157, O26 and O111 from Other Serovars by MALDI-TOF MS Based on the S10-GERMS Method

Teruyo Ojima-Kato, Naomi Yamamoto, Mayumi Suzuki,  
Tomohiro Fukunaga Hiroto Tamura

PLOS ONE | [www.plosone.org](http://www.plosone.org) 1 November 2014 | Volume 9  
| Issue 11 | e113458

Identifies bacteria using the ribosomal proteins encoded together.  
Called S10 GERM for S10-spc-alpha operon gene-encoded  
ribosomal protein mass spectrum

This uses about half of the ribosomal subunits

Our method uses all of the subunits in the desired mass range

Extract ribosomal sequences keyed to organism

Remove N-terminal methionine as necessary

Calculate singly and doubly charged masses

Save database (in memory)

Find matches within tolerance

Calculate score for each organism from  
% intensity matched  
% proteins found  
Mass accuracy

Get colony

Collect spectrum

Detect peaks

Select mass range  
(2.5-10K)

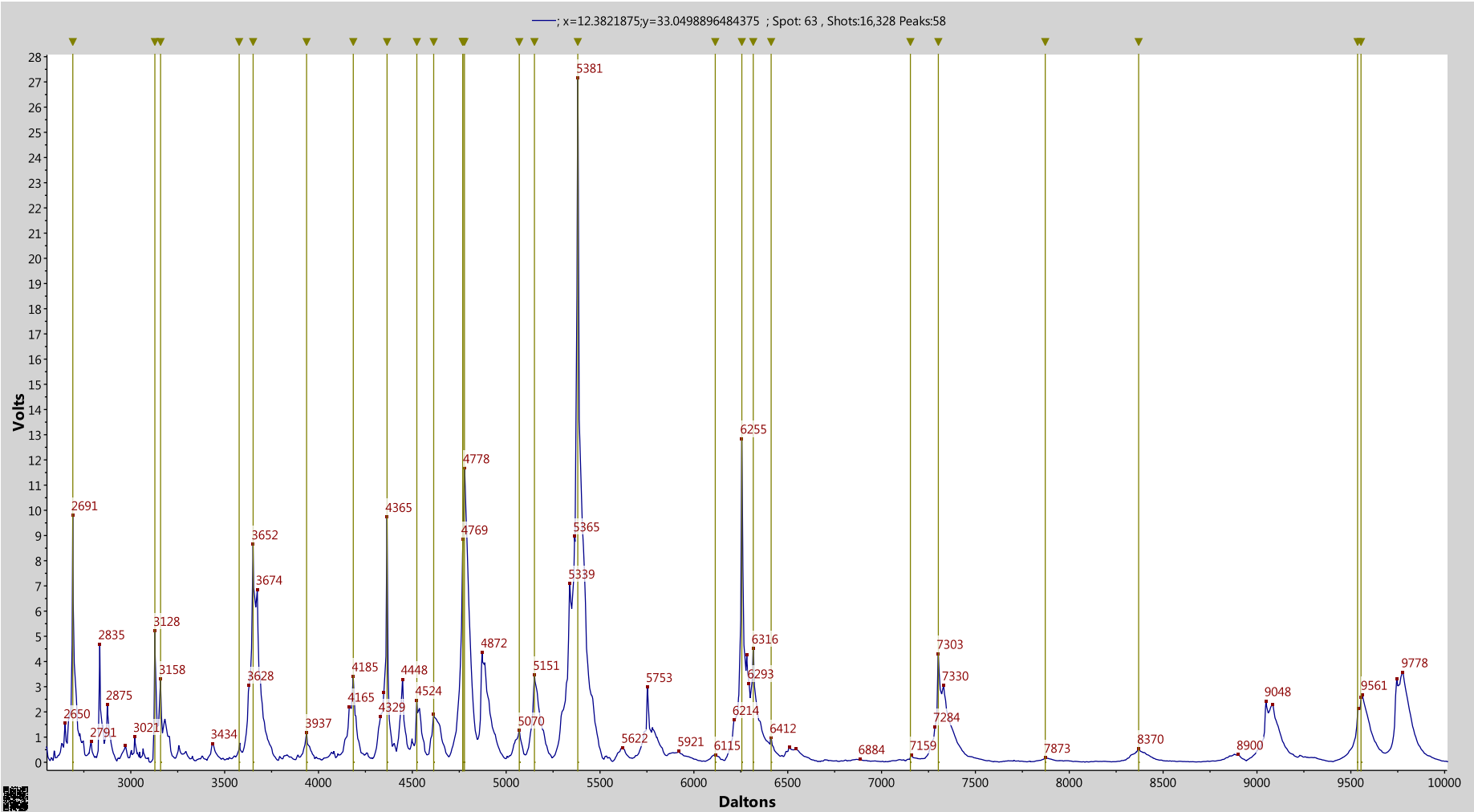
Recalibrate  
Tighten  
Constrain taxonomy

# When to search DNA databases?

- If you want to know why MALDI works (biochemistry)
- If you want to know what you are distinguishing
- When you are not happy with library searches
  - Returns scores of all organisms in the database
    - Useful for exploring relationships of closely related organisms
  - Returns sequences and names for all matches
- Can easily adjust mass tolerance or mass range to determine robustness of the answer
- Can restrict search to taxon id (saves time)
- Should work better for environmental organisms
- **Come find me with any bacterial ID problem!**

# Matches of ribosomal proteins to *S. dysenteriae* within 500 ppm

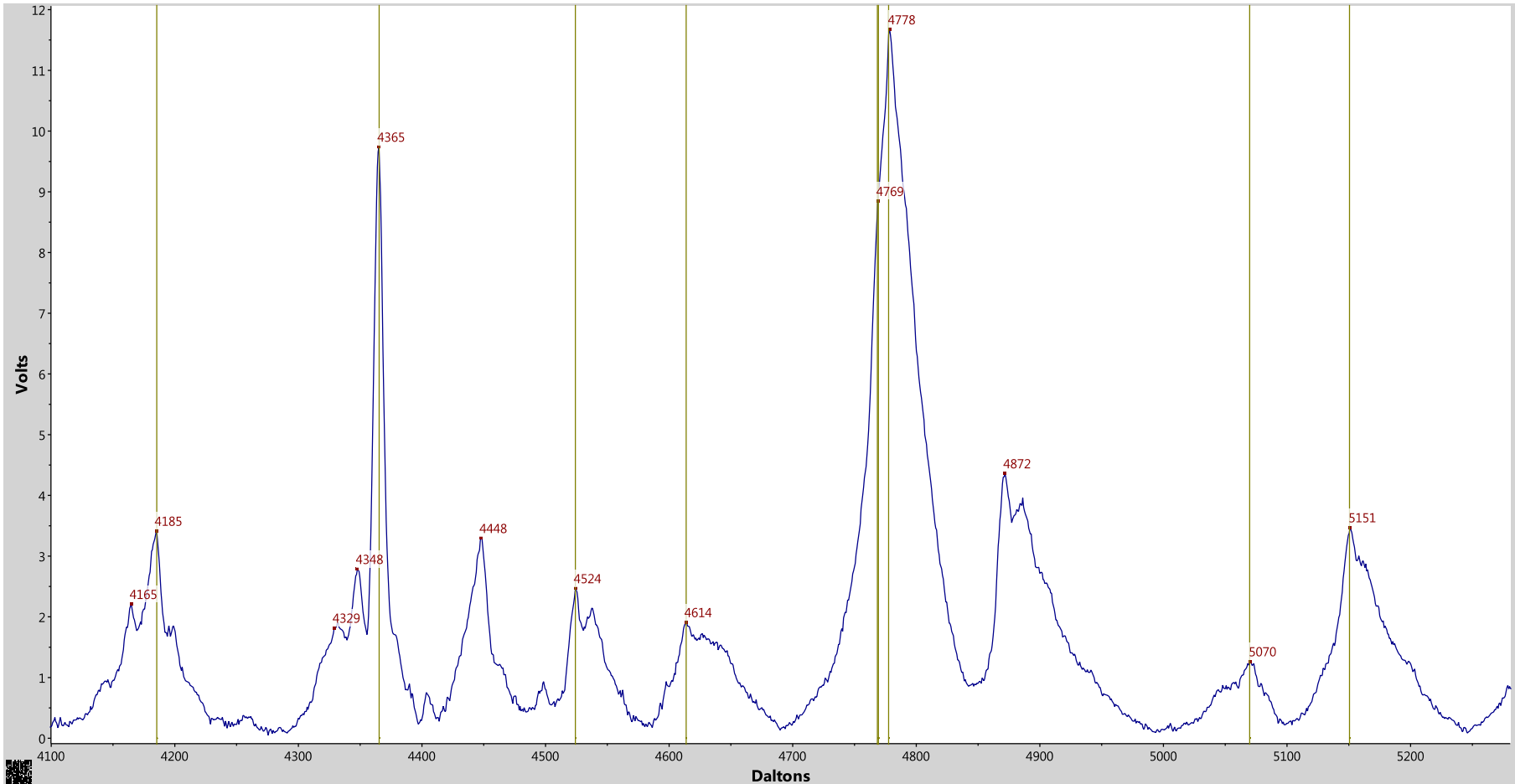
Plate csm8 spot 63



mz 2500-10000

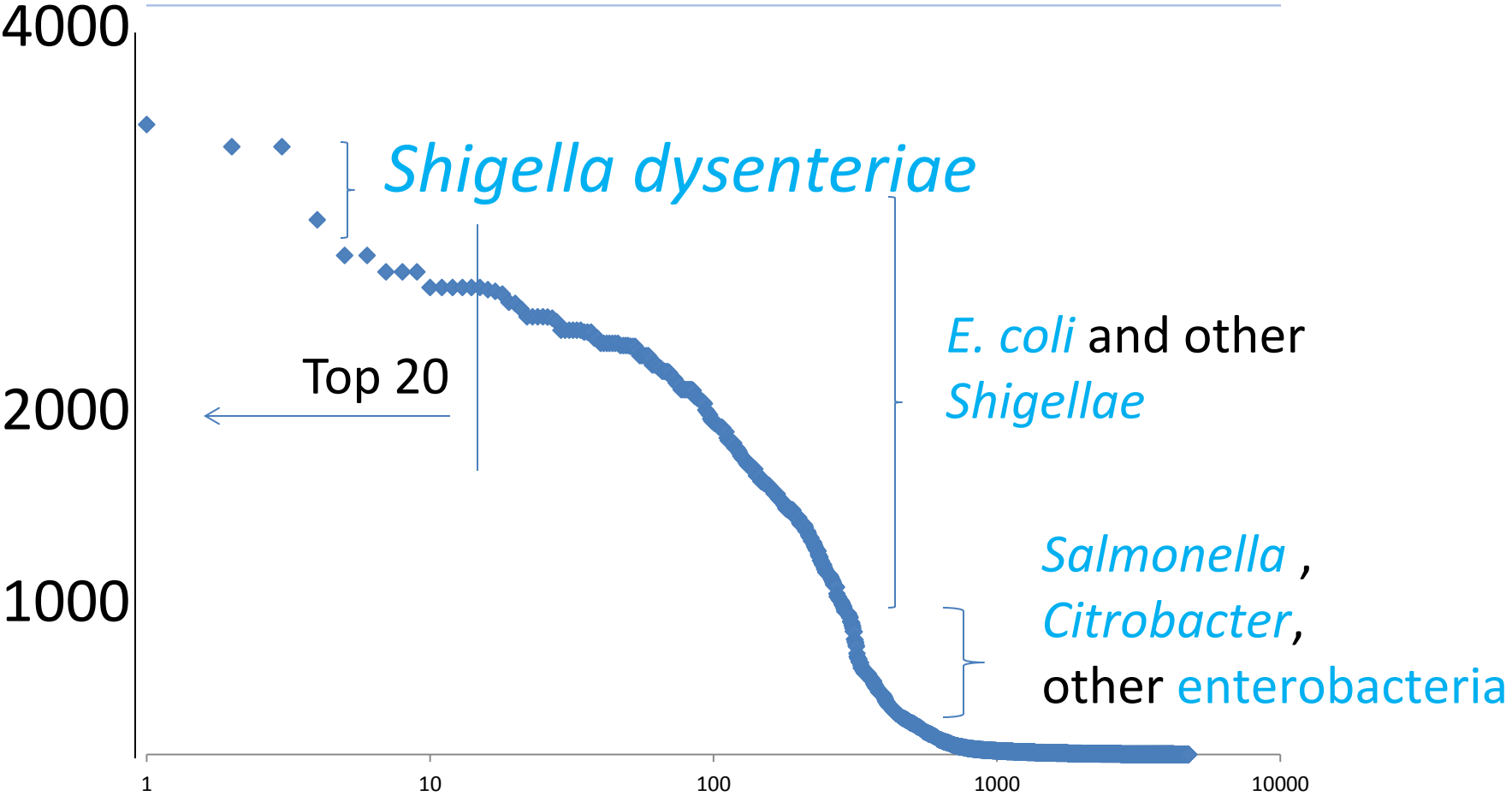
# Same spectrum

## 7 peaks matched in this region



mz 4100-5300

# Score Vs. Organism



7641 strains



# Top 20 scoring organisms; all related to *E. coli*

I	OrgID	Name	Proteins	Matches	Score	%rib	% int	ppm
1	1514	Shigella dysenteriae WRSd3	61	22	3365	36.1	64.7	80.8
2	70406	Shigella dysenteriae WRSd3	65	23	3245	35.4	64.8	80.7
3	70403	Shigella dysenteriae sd197	65	23	3245	35.4	64.8	80.7
4	6911	Escherichia coli KTE182	57	21	2856	36.8	52.6	88.7
5	8951	Escherichia coli 908573	59	21	2665	35.6	52.6	88.7
6	985	Escherichia coli KTE64	59	21	2665	35.6	52.6	88.7
7	4443	Escherichia coli MS 187-1	60	21	2577	35.0	52.6	88.7
8	1470	Escherichia coli KTE36	60	21	2577	35.0	52.6	88.7
9	1466	Escherichia coli KTE51	60	21	2577	35.0	52.6	88.7
10	2012	Escherichia coli UMEA 3240-1	61	21	2493	34.4	52.6	88.7
11	2088	Escherichia coli UMEA 3304-1	61	21	2493	34.4	52.6	88.7
12	3037	Escherichia coli 907357	61	21	2493	34.4	52.6	88.7
13	699	Escherichia coli KTE76	61	21	2493	34.4	52.6	88.7
14	1663	Escherichia coli KTE234	61	21	2493	34.4	52.6	88.7
15	1980	Escherichia coli KTE100	61	21	2493	34.4	52.6	88.7
16	70398	Shigella boydii CDC 3083-94	65	22	2482	33.8	54.2	86.3
17	15052	Escherichia coli TA280	46	17	2473	37.0	45.3	97.5
18	9460	Escherichia coli MS 79-10	58	20	2459	34.5	51.7	88.4
19	7783	Escherichia coli MS 182-1	58	20	2415	34.5	50.8	90.7
20	10135	Escherichia coli H299	59	20	2409	33.9	52.4	88.7

## Matched peaks: 9 out of top 10 explained by ribosomal proteins

N	I	theo	exper	ppm	Acc	Protein	Seq
1	1	5381.4	5381.2	-45.5	A0A090NMR8	50S ribosomal protein L34	MKRTFQ
2	2	4777.6	4778.4	158.3	A0A090NGX7	30S ribosomal protein S20(2)	ANIKSAK
3	3	3651.8	3651.5	-71.6	A0A090NXZ7	50S ribosomal protein L29(2)	MKAKEL
4	4	6255.4	6255.0	-71.4	A0A090NL18	50S ribosomal protein L33Me	AKGIREK
5	5	2691.2	2690.8	-139.6	A0A090NMR8	50S ribosomal protein L34(2)	MKRTFQ
6	6	4365.4	4365.2	-31.0	A0A090N JL4	50S ribosomal protein L36	MKVRAS
7	8	7302.5	7302.5	0.4	A0A090NXZ7	50S ribosomal protein L29	MKAKEL
8	9	5150.6	5151.0	95.2	A0A090N JM9	30S ribosomal protein S19(2)	PRSLKKG
9	10	4185.4	4185.0	-82.9	A0A090NI81	30S ribosomal protein S21(2)	PVIKVRE
10	14	3128.2	3128.3	27.5	A0A090NL18	50S ribosomal protein L33Me(2)	AKGIREK
11	19	4613.8	4613.6	-33.5	A0A090NF12	DNA-binding protein HU beta(2)	MNKSQL
12	22	6316.2	6315.7	-74.0	A0A090NP57	50S ribosomal protein L32	AVQQN
13	23	3158.6	3157.9	-210.0	A0A090NP57	50S ribosomal protein L32(2)	AVQQN
14	24	3936.6	3937.0	119.2	A0A090NC24	50S ribosomal protein L31(2)	MKKDIH
15	25	5069.8	5070.3	108.9	A0A090NNF1	30S ribosomal protein S15(2)	SLSTEAT
16	29	8369.8	8370.0	29.2	A0A090NI81	30S ribosomal protein S21	PVIKVRE
17	40	4768.5	4768.9	79.0	A0A090NIK4	DNA-binding protein HU alpha(2)	MNKQTQ
18	45	9536.0	9539.3	345.3	A0A090NIK4	DNA-binding protein HU alpha	MNKQTQ
19	48	6114.2	6115.1	149.8	A0A090NJ47	50S ribosomal protein L22(2)	METIAK
20	49	6411.6	6412.2	87.5	A0A090NXZ9	50S ribosomal protein L30	AKTIKIT
21	50	7872.1	7873.3	147.9	A0A090NC24	50S ribosomal protein L31	MKKDIH
22	56	9554.2	9554.7	50.0	A0A090NGX7	30S ribosomal protein S20	ANIKSAK

# Masses for the top 20 organism hits

Green: same mass as top hit

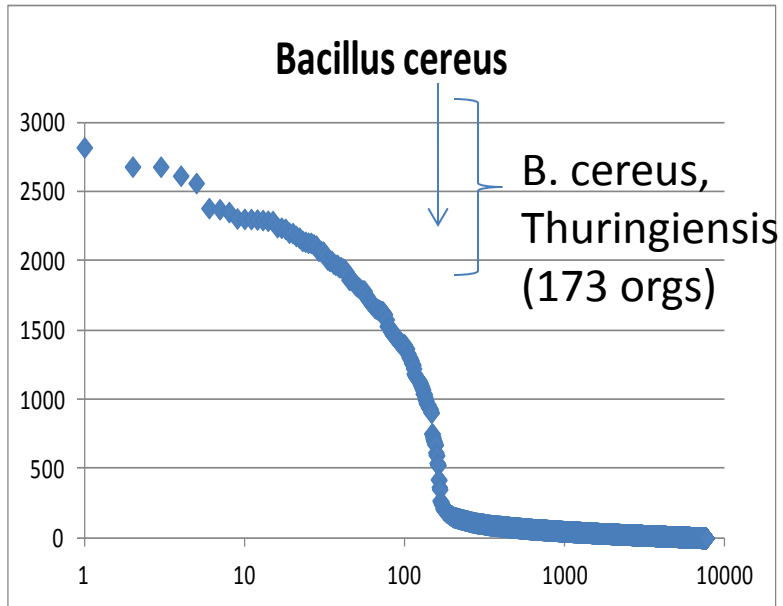
Red: subunit missing

Orange: subunit has extension

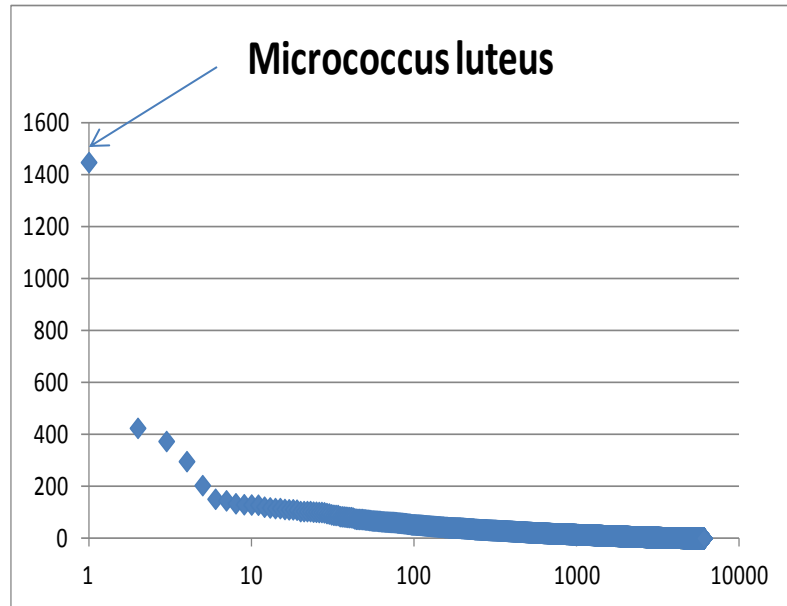
Blue: subunit has substitution (V->A)

N	I	Protein	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13	S14	S15	S16	S17	S18	S19	S20			
1	1	50S ribosomal protein L34	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381	5381		
2	2	30S ribosomal protein S20(2)	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	4778	
3	3	50S ribosomal protein L29(2)	3652	3652	3652	3638		3638	3638	3638		3638		3638	3638	3638	3638	3638	3638	3638	3638	3638	3638	3638	
4	4	50S ribosomal protein L33Me	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	6255	
5	5	50S ribosomal protein L34(2)	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	2691	
6	6	50S ribosomal protein L36	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	4365	5415	4365	4365	4365	4365	
7	8	50S ribosomal protein L29	7303	7303	7303	7274		7274	7274	7274		7274		7274	7274	7274	7274	7274	7274	7274	7274	7274	7274	7274	
8	9	30S ribosomal protein S19(2)	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	5151	
9	10	30S ribosomal protein S21(2)	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	4185	
10	14	50S ribosomal protein L33Me(2)	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	3128	
11	19	DNA-binding protein HU beta(2)	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614	4614		4614		4614		4614	
12	22	50S ribosomal protein L32	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316	6316
13	23	50S ribosomal protein L32(2)	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158	3158
14	24	50S ribosomal protein L31(2)	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937	3937
15	25	30S ribosomal protein S15(2)	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070	5070			5070	5070	5070	
16	29	30S ribosomal protein S21	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370	8370
17	40	DNA-binding protein HU alpha(2)	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769	4769
18	45	DNA-binding protein HU alpha	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539	9539
19	48	50S ribosomal protein L22(2)	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115	6115		6115	6115	6115	6115	
20	49	50S ribosomal protein L30	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412	6412		
21	50	50S ribosomal protein L31	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873	7873
22	56	30S ribosomal protein S20	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555	9555
23	20	50S ribosomal protein L25(2)	7162	5348	5348			5348	5348	5347	5348	5348	5347	5348	5347	5348	5348	5339	5347	5348	5348	5347	5348	5347	
24	53	50S ribosomal protein L35	7891	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159	7159		7159	7159	7159	7159	7159	

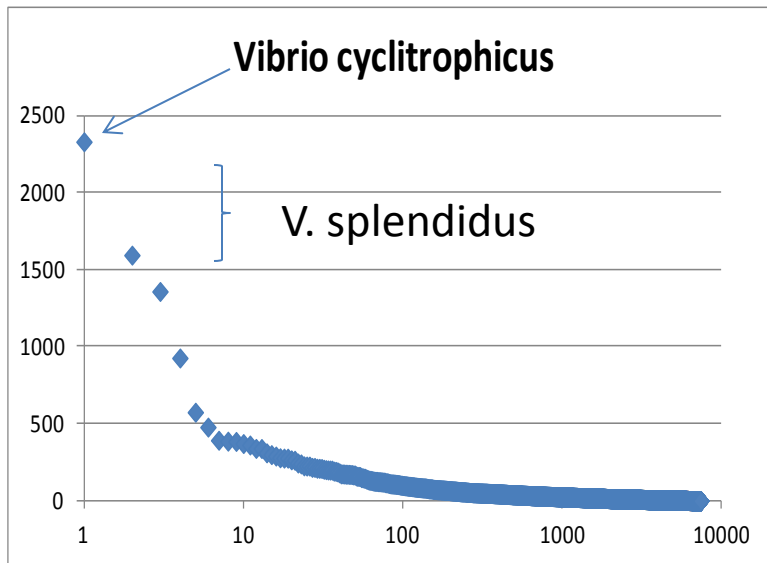
# Overgrown LB medium



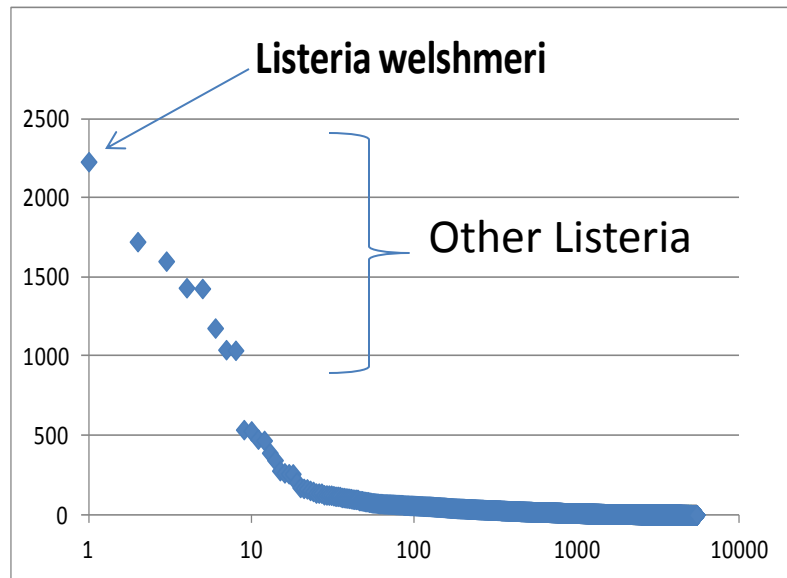
# Colorado School of Mines



# Martin Polz' lab MIT



# Colorado School of Mines



# What is a bacterial species?

- Originally:
  - disease
  - growth requirements / metabolism
  - ecological niche
- New consensus: genomic similarity -> proteomic similarity
  - **Vertical** vs Horizontal evolution
- Simplification: Substantial Ribosomal protein sequence identity
  - Highly conserved throughout prokaryotic evolution, transmitted **vertically**
  - 54 subunits typically
  - Complex nearly 1:1
- MALDI recognizes **small, abundant, basic** proteins like ribosomal subunits
- Therefore, MALDI can distinguish species
  - If organisms cannot be distinguished by MALDI, then they do not belong in different species
  - Limitations of species delineation by MALDI due to errors in species definition

# Distinguishing *Shigella* from *E. coli*

- Presumably this is clinically important
- There are 4 defined *Shigella* species
  - *Boydii* (35), *dysenteriae* (14), *flexnerii* (120), *sonnei* (583)  
vs *E. coli* (3947)
  - Known for 20 years *Shigella* not separable from *E. coli*.
  - *Shigella* 'species' not separable from one another.
  - Epidemiological differences between *Shigella* species
    - According to my primary source, Wikipedia.
  - *Shigella*'s pathogenicity due to a large plasmid that promotes invasion
    - A family of ~ 8 E3 ubiquitin ligases
    - Type III secretion system
- Carefully check ribosomal protein profile for 4 examples of *E. coli* and each *Shigella* species

# PATRIC cladogram

https://

[www.patricbrc.org/](http://www.patricbrc.org/)

## Pathosystems

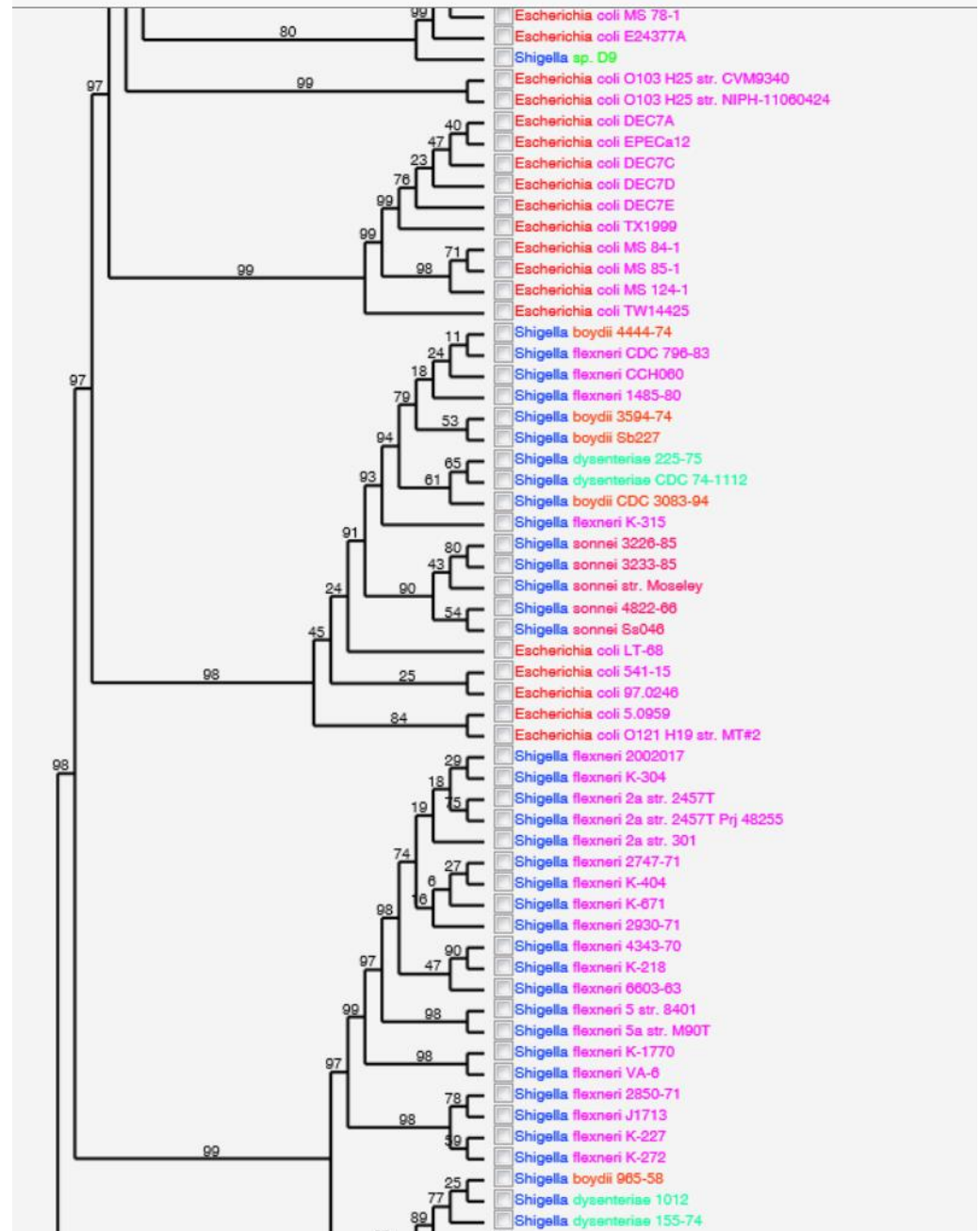
## Resource

## Integration Center

NIH/NIAID-funded project of The University of Chicago with subcontract to the Biocomplexity Institute of Virginia Tech

*Escherichia* red  
*Shigella* blue  
*dysenteriae* green

Lots of *E. coli* strains above and below, mixed with a few *Shigellae*

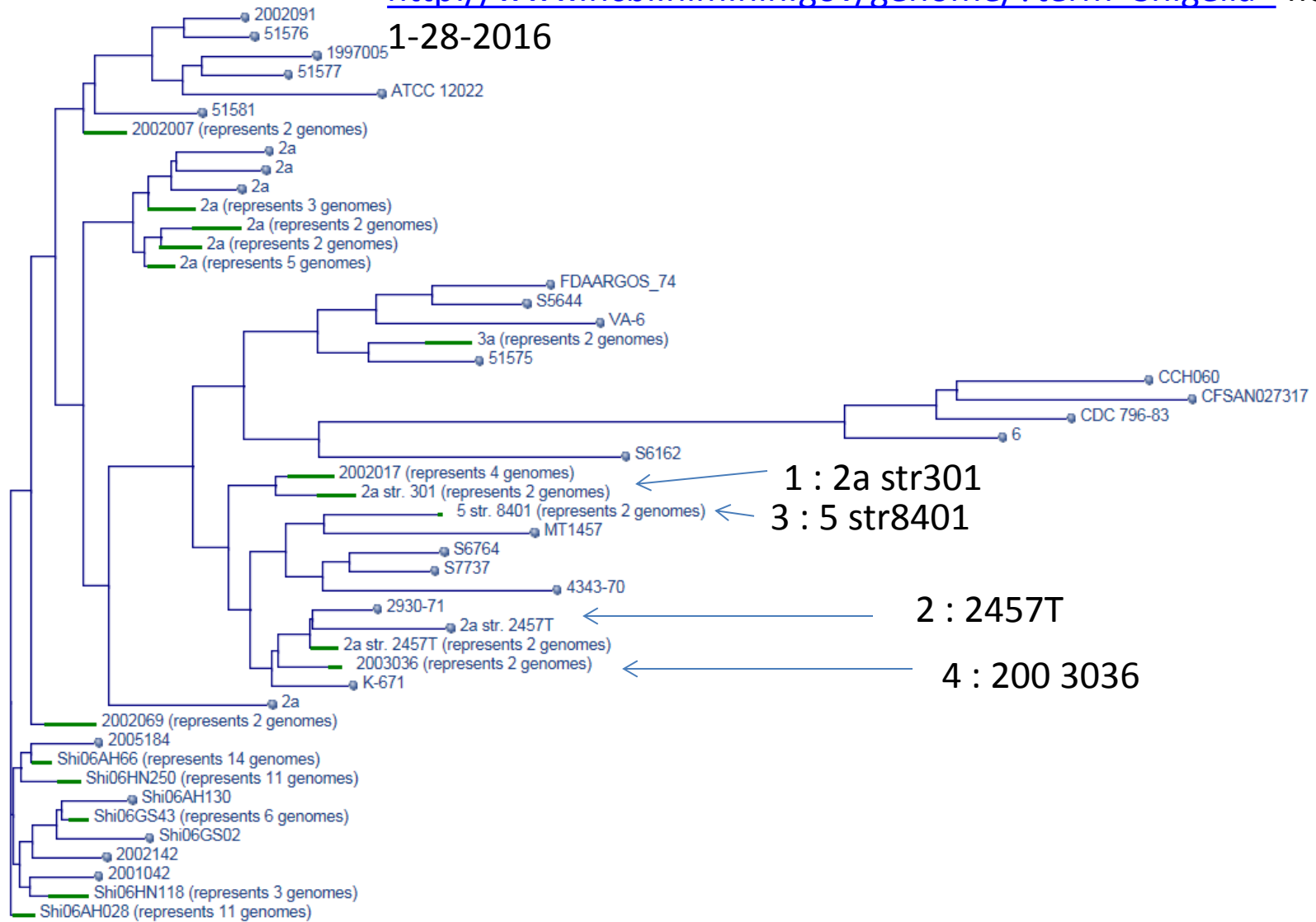


# Shigella flexneri (120)

Dendrogram (based on genomic BLAST)

[http://www.ncbi.nlm.nih.gov/genome/?term=Shigella+ flexneri](http://www.ncbi.nlm.nih.gov/genome/?term=Shigella+flexneri)

1-28-2016



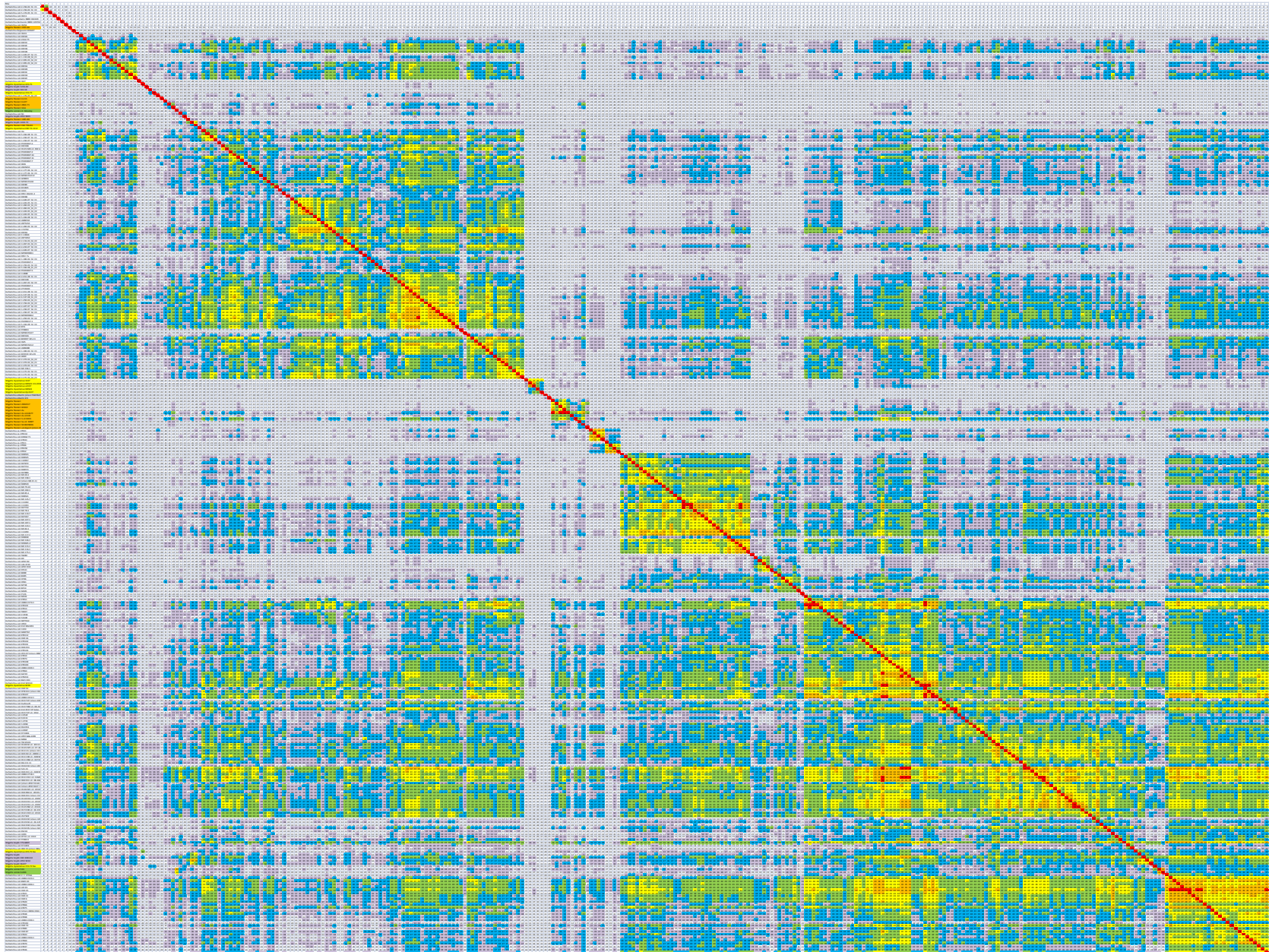




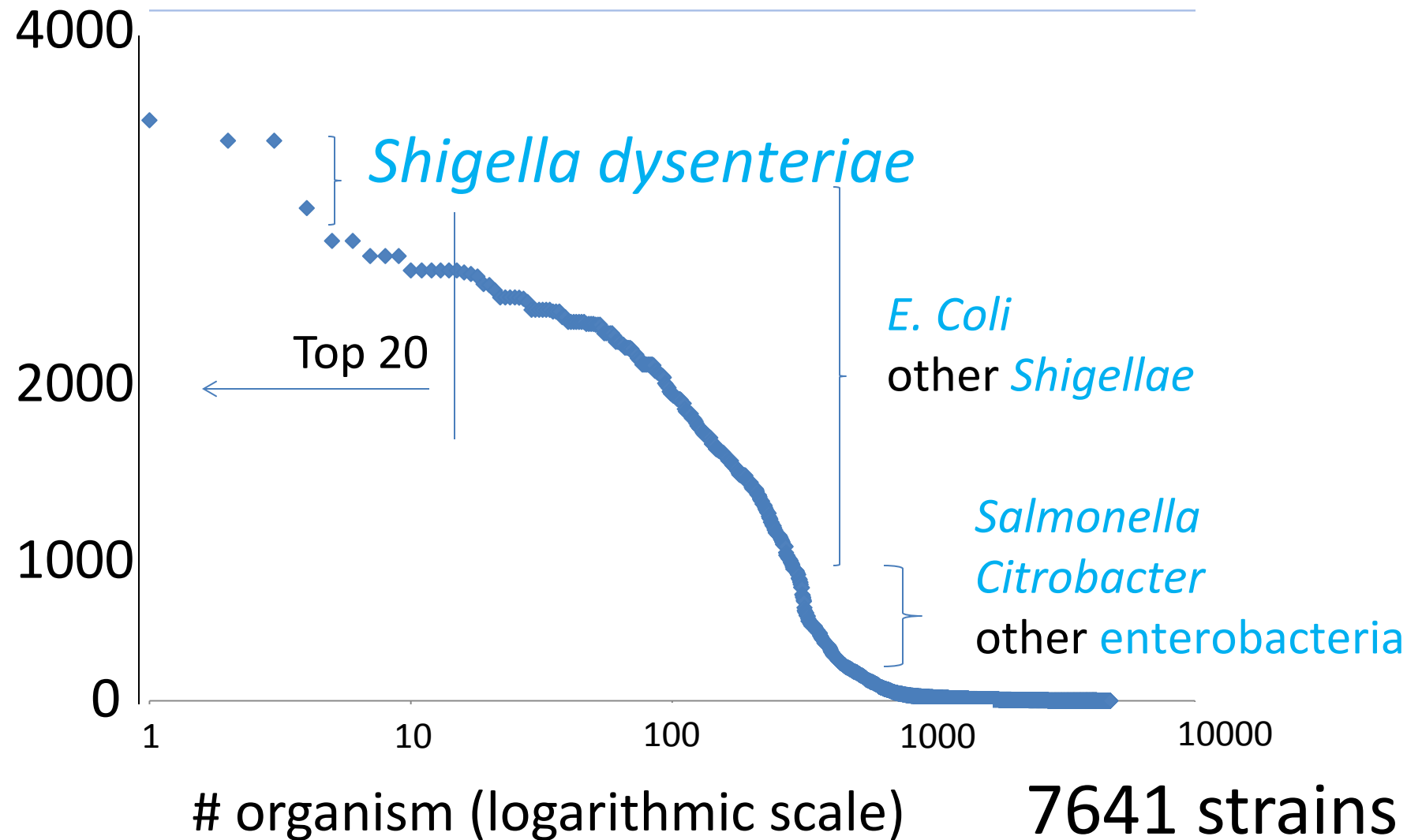




Clustering of *E. coli* / *Shigella* strains by ribosomal proteins < 200 aa long  
*Shigella* strains are colored by 'species' in the left-most column (flexneri orange)

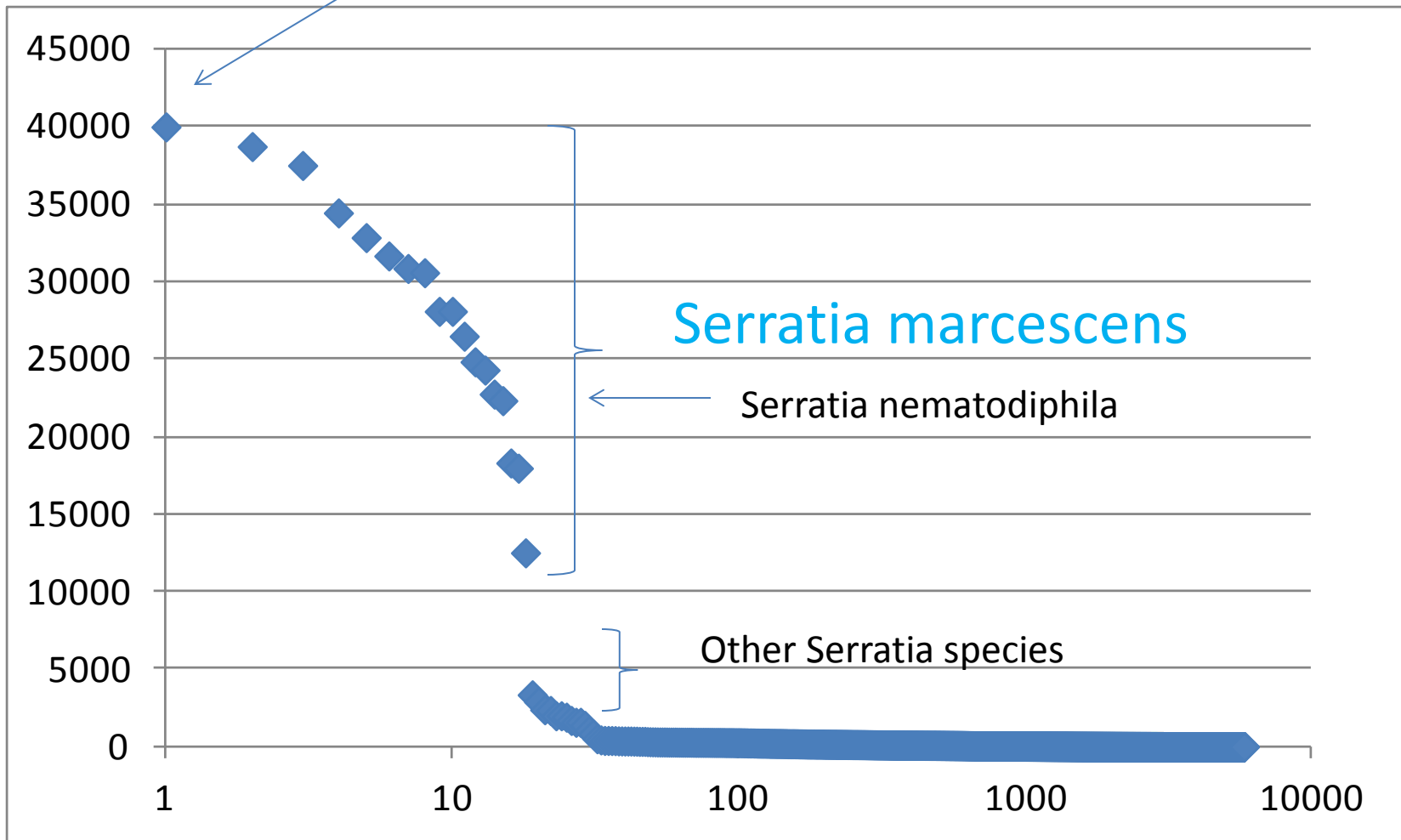


# Score Vs. Organism



In silico mapping at 500 ppm  
Ribosomal protein theoretical masses  
**786\_SSON**, thought to be *Shigella sonnei*

Self mapping - a perfect match!



# Improving public repositories

- Solution -get NCBI? to develop a ribosomal protein database
  - This was originally envisioned by Pineda et al. in 2003. (Analytical Chem. 75:3817)
  - Many more genomes available (>25000)
  - Can reduce search space by pooling genomes with identical ribosome sequences
    - Lots of identity in
      - *Staphylococcus aureus* (5584)
      - *Escherichia coli* (3947)
      - *Mycobacterium tuberculosis* (3612)

# Proving the methodology

- I have succeeded in every case library methods succeed (> 100 species)
- Both library methods and genome matching methods will fail if the organism is not in the database
- Genomic matching can be tested using any peak list (ideally with intensity)
  - Calibration is important
- Genomic matching can also be tested in silico
  - Generate a ribosomal mass list
  - Can now reduce mass tolerance to nearly zero.
- **Challenge me with a bacterial ID problem!**



# Biological confusing factors

- Some bacteria are harder to lyse so as to release ribosomal proteins
- Not all ribosomal proteins are detected
- Probably some differential degradation
- Probably at lag phase ribosomes are sometimes less prominent in the proteome

# Informatic Confounding Factors

- Some genomes have mistakes in ribosomal proteins
  - Some subunits entirely missing in some proteomes
    - In E coli and B. subtilis, only  $\sim \frac{1}{2}$  of ribosomal proteins are necessary for growth
  - N-terminal extensions or deletions in the databases
  - Some sequences are mis-annotated as ribosomal subunits
- Some ribosomal proteins are present in multiple forms, at least for some clades
  - L31 has a special zinc-deficient form
- Some proteins are modified
  - L33 -> methylated in gammaproteobacteria
- Everything works better without these problems

# Conclusions

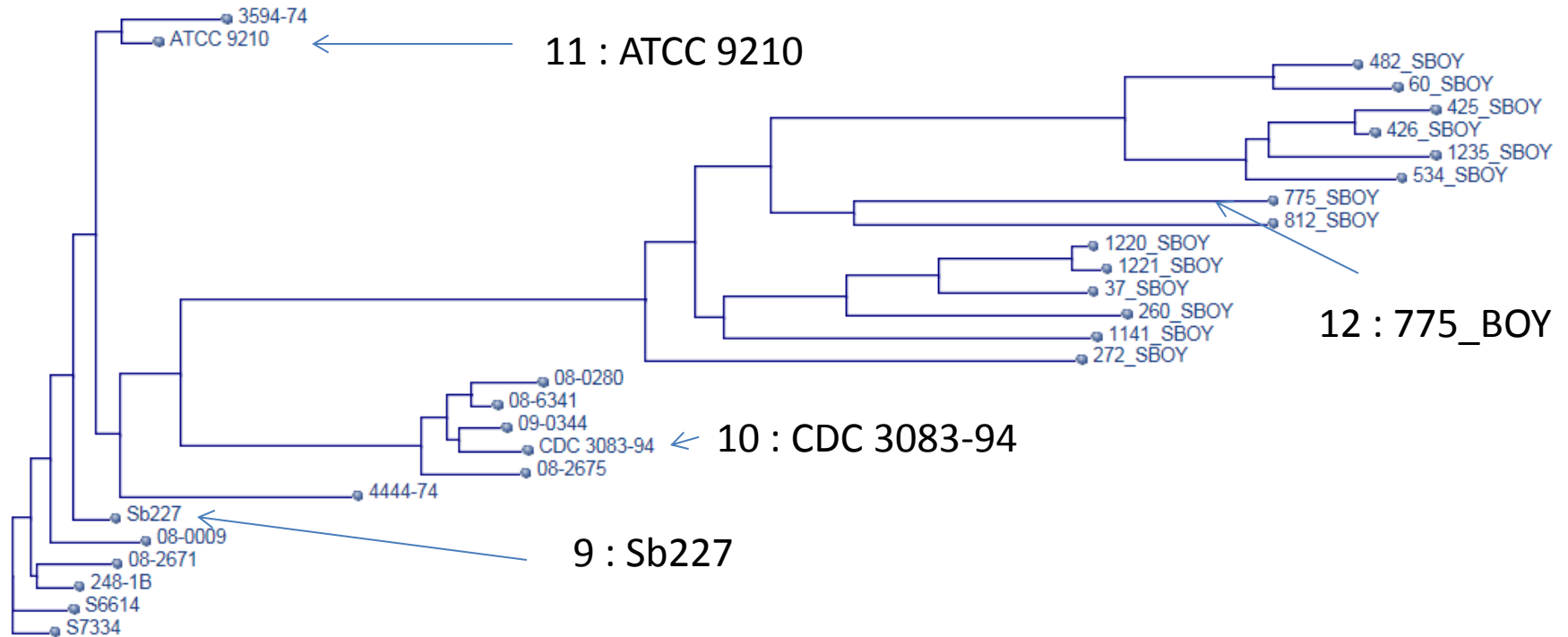
- MALDI can type to sequenced bacteria
- Same method types any sequenced bacteria to all other sequenced bacteria
- Limitations in distinguishing bacteria largely due to taxonomic confusion
  - Same principle works with entire proteome database, but is trickier
  - Can screen for modifications to ribosomal proteins
- Collaborators
  - Kent Voorhees and Chris Cox, Colorado School of Mines
  - Martin Polz, MIT
  - SimulTOF Systems

# Shigella boydii

<http://www.ncbi.nlm.nih.gov/genome/?term=Shigella+boydii>

1-28-2016

 Dendrogram (based on genomic BLAST)



1%

# Shigella dysenteriae

<http://www.ncbi.nlm.nih.gov/genome/?term=Shigella+boydii>

1-28-2016

## Reference genome: [see all organisms]

### ▣ *Shigella dysenteriae* Sd197

Submitter: Microbial Genome Center of ChMPH

Human Pathogen

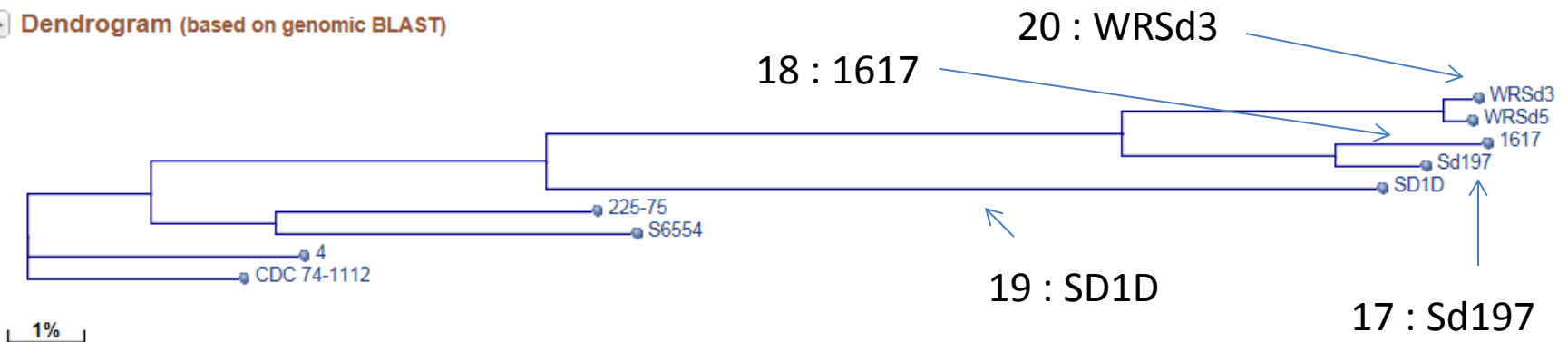
Morphology: Gram:Negative, Shape:Bacilli, Motility:No

Environment: Salinity:NonHalophilic, OxygenReq:Facultative, OptimumTemperature:37, TemperatureRange:Mesophilic, Habitat:HostAssociated

Phenotype: Disease:Dysentery

Type	Name	RefSeq	INSDC	Size (Mb)	GC%	Protein	rRNA	tRNA	Gene	Pseudogene
Chr	-	NC_007606.1	CP000034.1	4.37	51.2	4,063	22	85	4,602	427
Plsm	pSD1_197	NC_007607.1	CP000035.1	0.182726	44.8	223	-	-	224	1
Plsm	pSD197_spA	NC_009344.1	CP000640.1	0.008953	39.7	8	-	-	8	-

### ▣ Dendrogram (based on genomic BLAST)



# Shigella sonnei

<http://www.ncbi.nlm.nih.gov/genome/?term=Shigella+sonnei>

1-28-2016

Can't locate the other 3 *S. sonnei*!

Dendrogram (based on genomic BLAST)



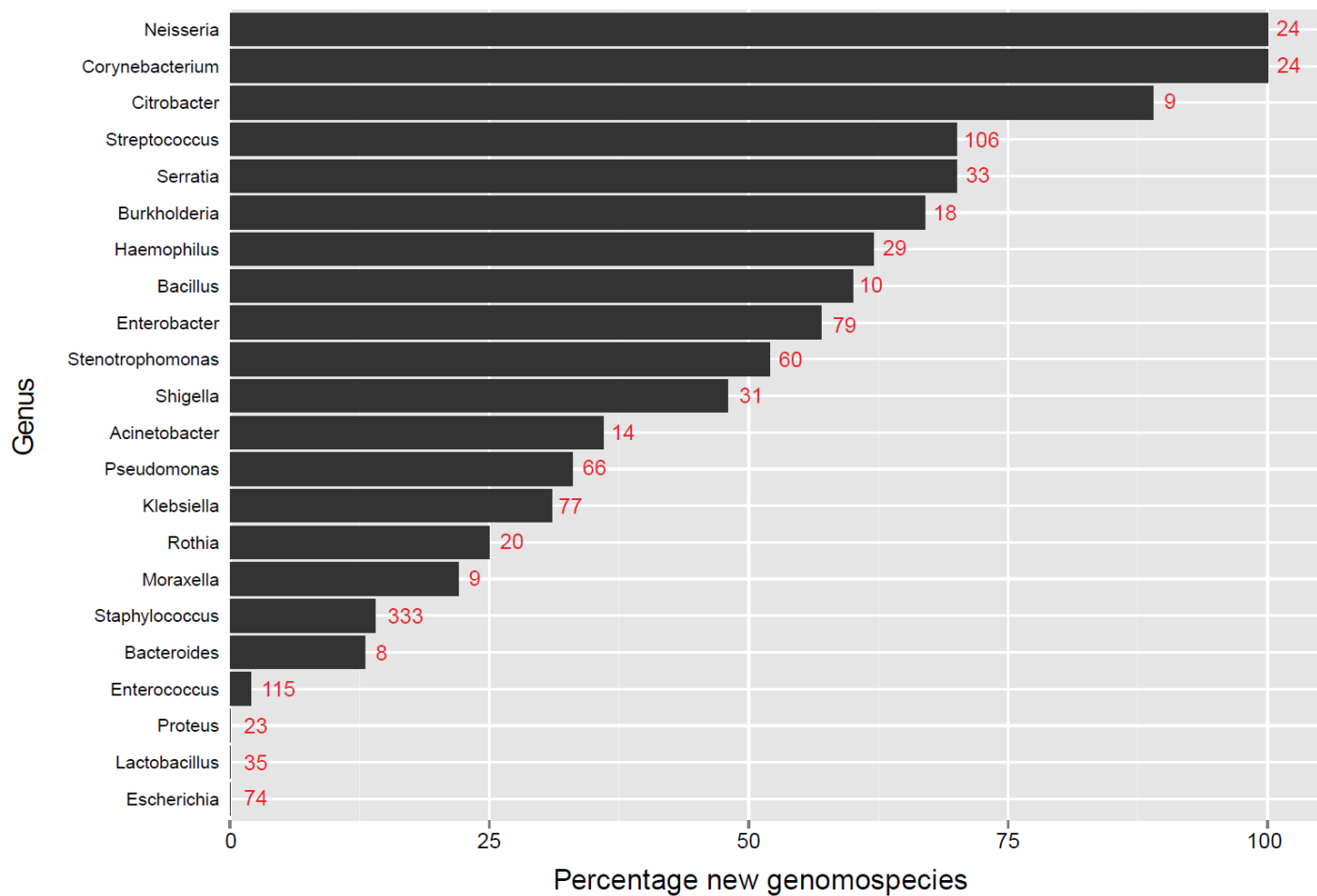
# A Year of Infection in the Intensive Care Unit: Prospective Whole Genome Sequencing of Bacterial Clinical Isolates Reveals Cryptic Transmissions and Novel Microbiota Figure S3

[David J. Roach et al; http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4521703/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4521703/)

All genera with >5 isolates sequenced are indicated.

The x-axis indicates the proportion of each genus that consisted of novel genospecies by ANIb analysis.

The red number to the right of each bar indicates the total number of isolates sampled from each genus in this study.



# Serratia marcescens 290 genomes 2-9-2016

▣ *Serratia marcescens* FGI94

Submitter: UW-Madison

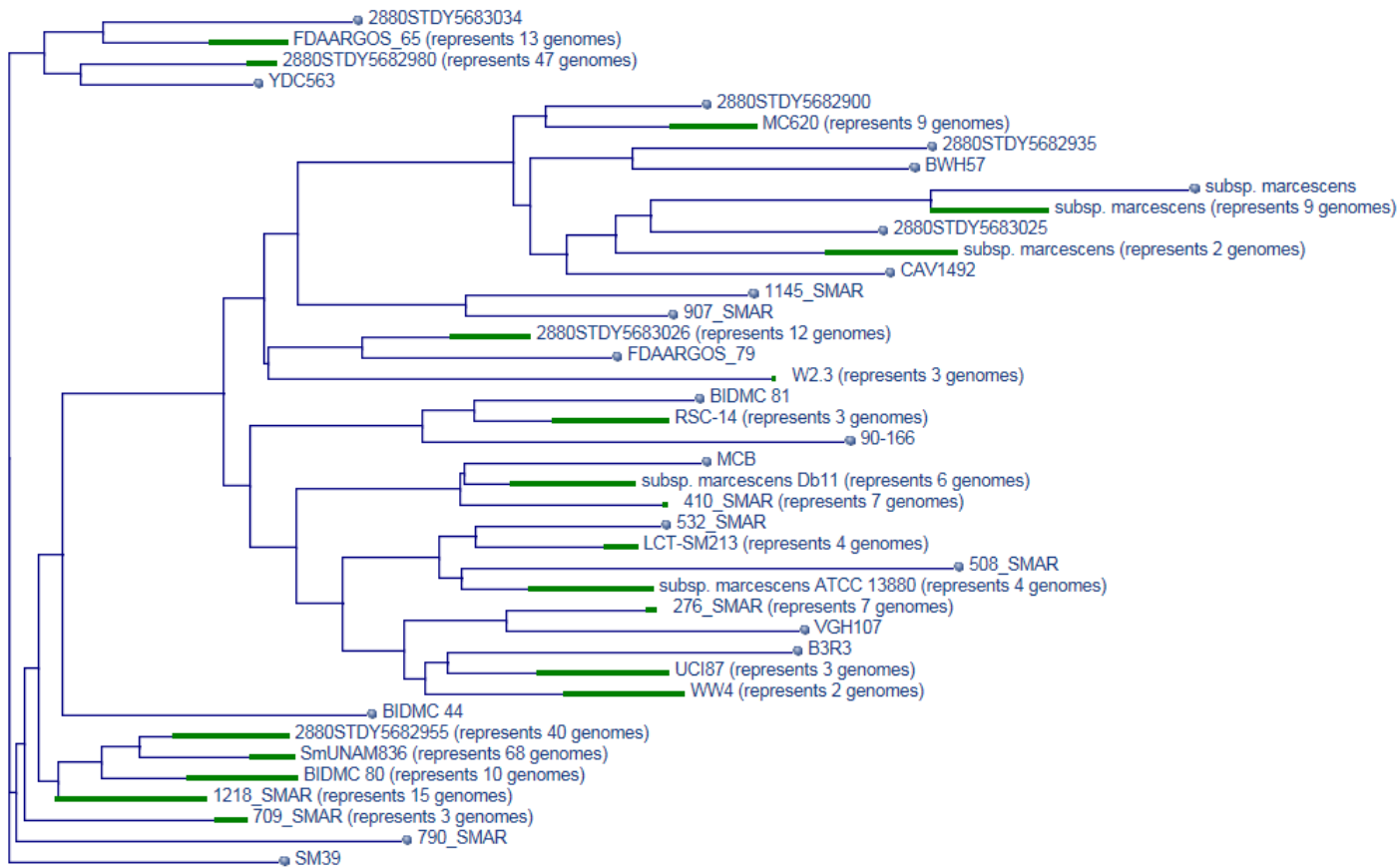
Morphology: Gram:Negative, Shape:Bacilli, Endospores:No, Motility:Yes

Environment: Salinity:NonHalophilic, OxygenReq:Facultative, TemperatureRange:Mesophilic, Habitat:HostAssociated

Phenotype: BioticRelationship:Episymbiont, TrophicLevel:Heterotroph, Disease:NA

Type	Name	RefSeq	INSDC	Size (Mb)	GC%	Protein	rRNA	tRNA	Other RNA	Gene	Pseudogene
	-	NC_020064.1	CP003942.1	4.86	58.9	4,290	22	83	3	4,436	38

▣ Dendrogram (based on genomic BLAST)



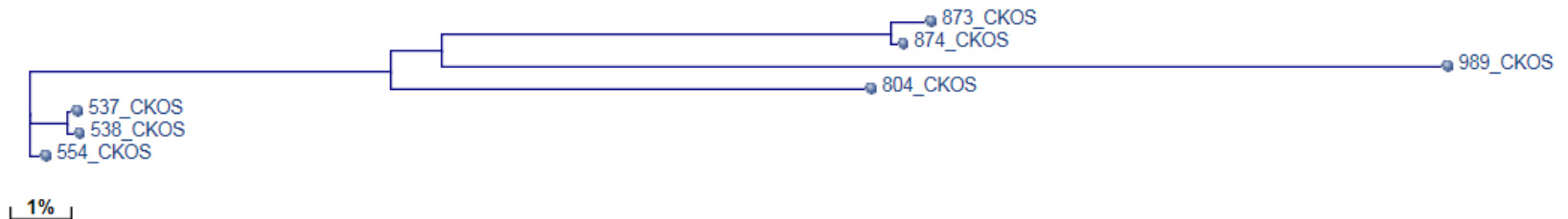


# Citrobacter koseri 14 genomes 2-9-2016

## Replicon Info

Type	Name	RefSeq	INSDC	Size (Mb)	GC%	Protein	rRNA	tRNA	Gene	Pseudogene
-	-	<a href="#">NC_009792.1</a>	<a href="#">CP000822.1</a>	4.72	53.8	4,222	22	83	4,377	50
pCKO3	-	<a href="#">NC_009793.1</a>	<a href="#">CP000823.1</a>	0.01	55.2	11	-	-	11	-
pCKO2	-	<a href="#">NC_009794.1</a>	<a href="#">CP000824.1</a>	0.01	51.3	8	-	-	8	-

## Dendrogram (based on genomic BLAST)



# Citrobacter freundii 52 genomes 2-9-2016

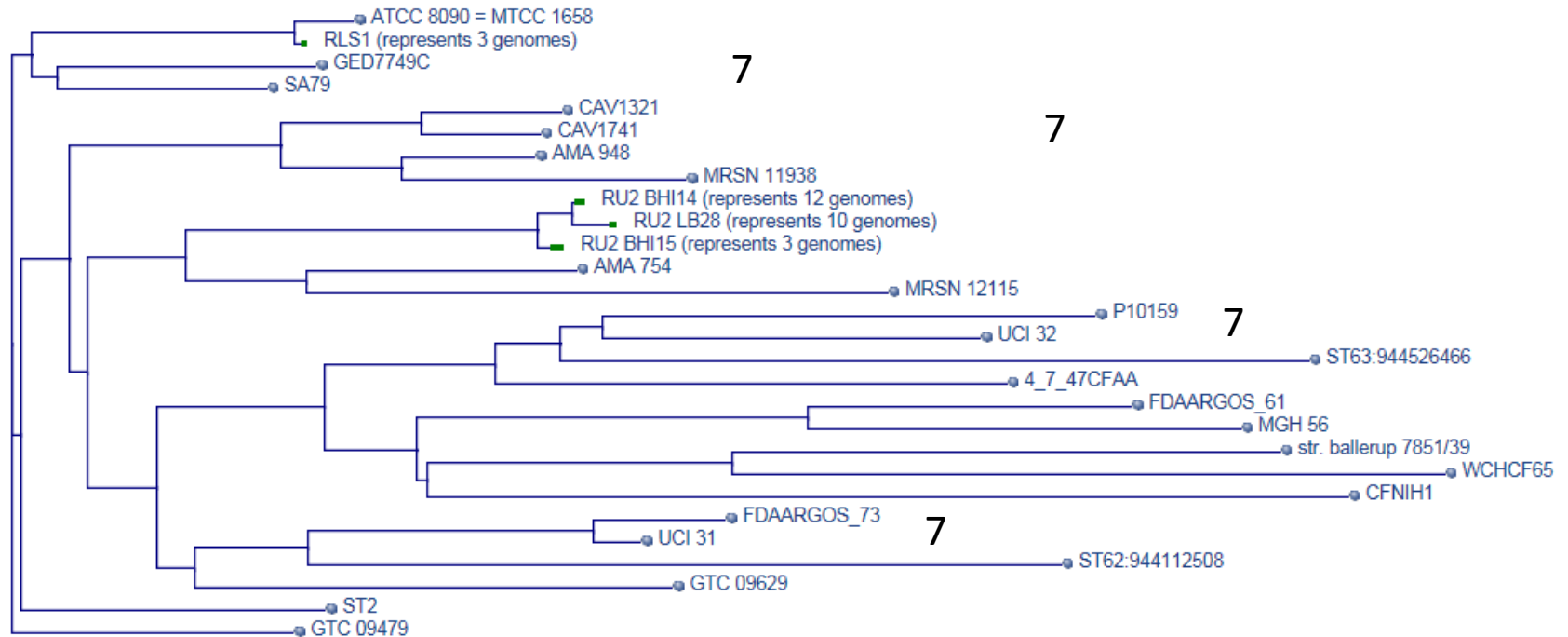
**Representative genome:** [\[see all organisms\]](#)

▣ *Citrobacter freundii* CFNIH1

Submitter: NISC

Type	Name	RefSeq	INSDC	Size (Mb)	GC%	Protein	rRNA	tRNA	Other RNA	Gene	Pseudogene
Chr	-	NZ_CP007557.1	CP007557.1	5.1	52.2	4,672	25	83	1	4,847	66
Plasm	pKEC-a3c	NZ_CP007558.1	CP007558.1	0.272297	52.7	298	-	-	-	309	11

▣ **Dendrogram (based on genomic BLAST)**



1%

10





