

# GRC MEETING: 2016 - 2017

## Comparative Genomics of Soil-adapted *E. coli*

---

Nick Waters

March 27, 2017

Department of Microbiology  
School of Natural Sciences  
National University of Ireland, Galway

# Overview

---



- Genome Assembly
- Phylogenetics, MLST, and Collection Overview
- Virulence
- Plans

# Genome Assembly

---



**Table:** Bacterial Genome Completion as of 4/1/17

Total	Complete genome	Chromosome	Contig	Scaffold
85799	6255	1143	38429	39972

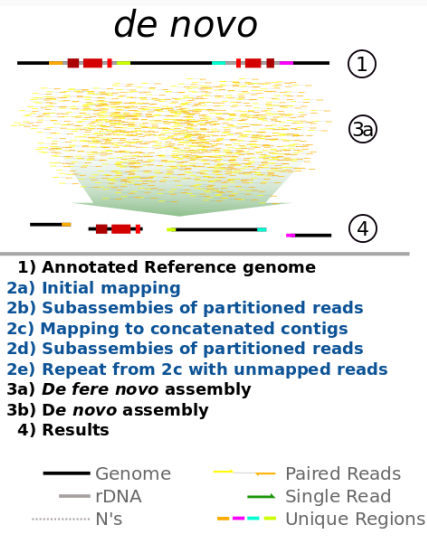
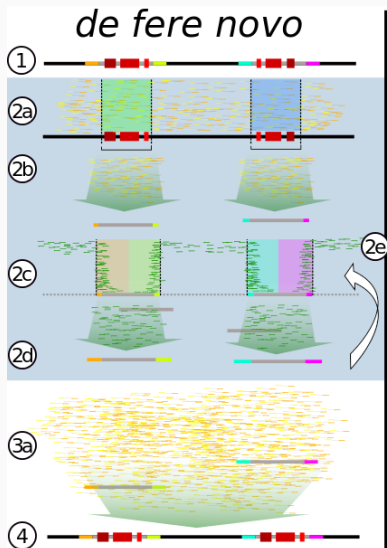


**Table:** Hits per Search Term in NCBI's SRA

Search term	hits	percentage
'illumina'	2242225	(94.27)
'pacbio'	21131	(0.89)
'ion'	30560	(1.28)
'roche'	42445	(1.78)
'oxford'	12301	(0.52)
'solid'	29791	(1.25)
Total	2378453	(100)



1. Within a taxonomic group, GC content is largely conserved.
2. Within a taxonomic group, genome size is largely conserved.
3. Bacterial genomes are dense.
4. Nucleotide order is not random.

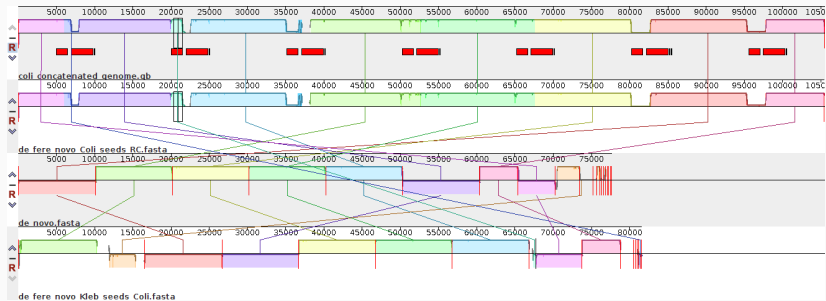






- Simulated Reads on a Simulated Genome
- Simulated Reads on Real Genomes
- Hybrid Assembly Validation
- GAGE-B Dataset

# Benchmarking: Simulated Genome





**Table:** riboSeed on simulated Sakai reads

Coverage	Ref. rDNAs	De novo(skip, miss)	De fere novo (skip, miss)
10	7	0 (7, 0)	3 (4, 0)
20	7	0 (7, 0)	6 (1, 0)
50	7	0 (7, 0)	6 (1, 0)
100	7	0 (7, 0)	6 (1, 0)



**Table:** riboSeed on *Pseudomonas aeruginosa* strain BAMCPA07-48

Coverage	Re. rDNA	De novo (skip, miss)	De fere novo (skip, miss)
200	4	1(3, 0)	4(0, 0)



- Finish benchmarking
- Applicability to fungal and archaeal genomes
- Improve Signal-to-noise ratio

# Aggregating Metadata

---



Problem: lack of single repository for the data related to the collection

- Sample Isolation data
- Phenotypic tests
- Phylotyping



- Library preparataion and sequencing QC
- Average nucleotide identity
- In silico Clermont PCR
- MLST





After eliminating isolates falling beneath the 95% ANI threshold and those with poor sequencing data, the collection consists of 153 isolates.

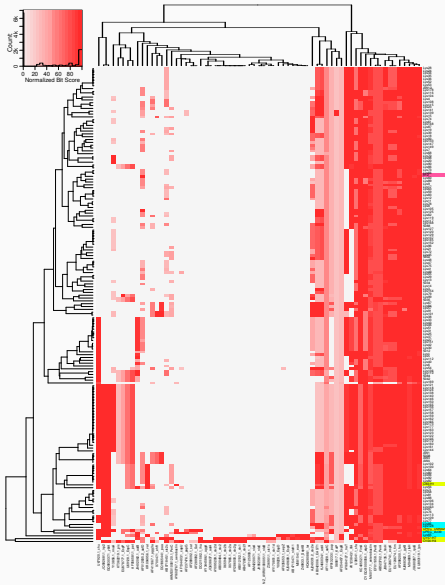
- Prevents duplication of work
- Aids automation
- Allows investigation of meta-variables

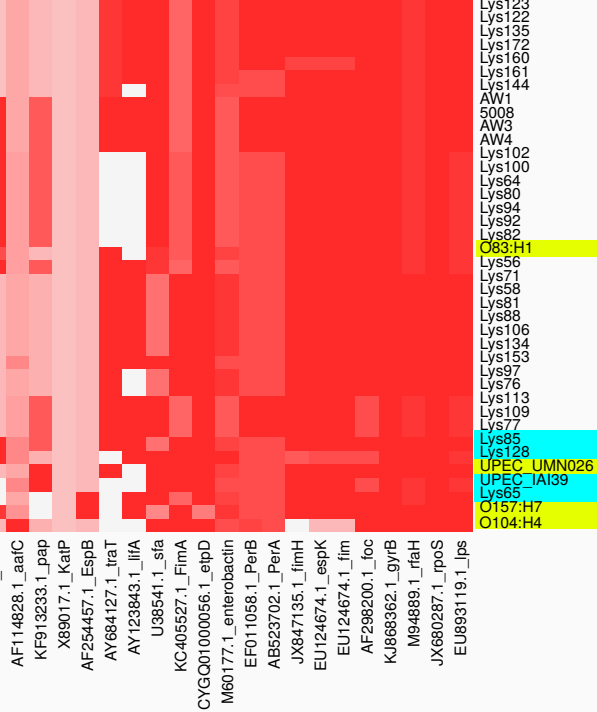
# Preliminary Virulence Profiling

---



- Search literature for genes implicated in virulence
- Select representative sequences for 50 virulence factors
- Use reciprocal translated blast to find occurrences
- Filter results, visualize







- Compare with recent tools (ARIBA, VirulenceFinder, etc)
- Experimentally assess phenotypes as needed

# Plans for 2017-2018

---



riboSeed development:

- Improve progressive QC to identify problems early
- Investigate characteristics of rDNAs that may be predictive of riboSeed success

Exploring the E. coli pangenome:

- Attempt to isolate genomic trends indicative of soil-adaptation
- Establish pangenomic context of the the collection
- Correlate metadata with pangenome

Curli loss:

- Utilize phenotypic data from Y. Soronin and others to determine additional causes of curli loss





**OÉ Gaillimh**  
NUI Galway

## NUIG Microbiology

- Dr. Fiona Brennan
- Dr. Florence Abram
- Matthias Waibel
- Stephen Nolan
- Camilla Thorn



## James Hutton Institute, Dundee

- Dr. Leighton Pritchard
- Dr. Ashleigh Holmes

QUESTIONS?