# DISCUSSION CLUB: CURRENT RESEARCH

Nick Waters

December 6, 2016

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

# INTRODUCTION

My project: Comparitive Genomics of soil-Adapted E. coli

Given our 155 sequenced soil-adapted isolates, what can we learn about E. coli genomics?

- Phylogeny

## My project: Comparitive Genomics of soil-Adapted E. coli

Given our 155 sequenced soil-adapted isolates, what can we learn about E. coli genomics?

- Phylogeny
- Genomic Restructuring

My project: Comparitive Genomics of soil-Adapted E. coli

Given our 155 sequenced soil-adapted isolates, what can we learn about E. coli genomics?

- Phylogeny
- Genomic Restructuring
- Virulence/AMR

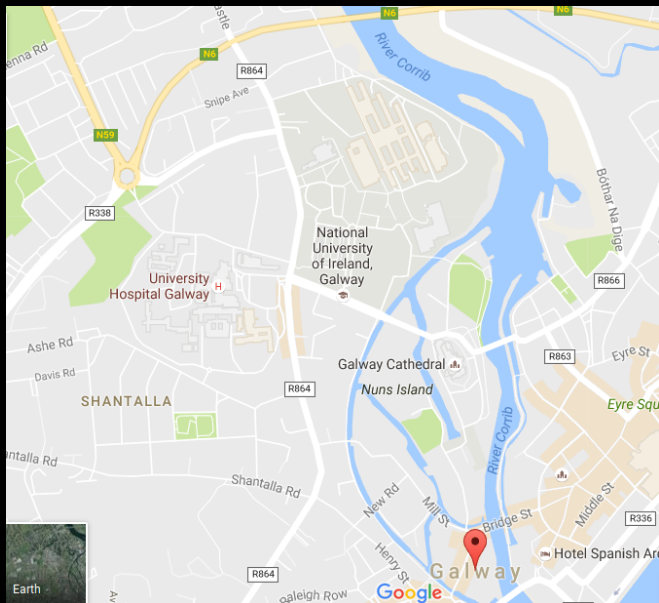My project: Comparitive Genomics of soil-Adapted E. coli

Given our 155 sequenced soil-adapted isolates, what can we learn about E. coli genomics?

- Phylogeny
- Genomic Restructuring
- Virulence/AMR
- Detection

My project: Comparitive Genomics of soil-Adapted E. coli

Given our 155 sequenced soil-adapted isolates, what can we learn about E. coli genomics?

- Phylogeny
- Genomic Restructuring
- Virulence/AMR
- Detection

# SHORT READY ASSEMBLY: BACKGROUND

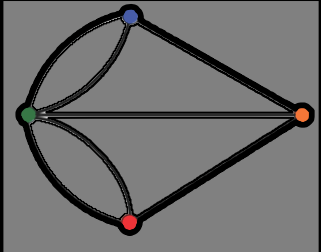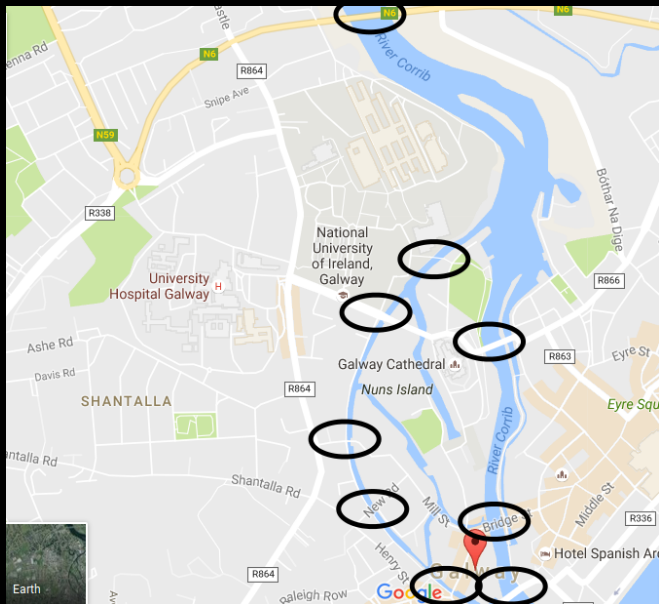Source[Chaisson et al., 2015]

Source[Chaisson et al., 2015]

Source[Compeau et al., 2011]

a

Sequence coverage gaps

b

Segmental duplication-associated gaps

c

Satellite-associated gaps

d

Muted gaps

Nature Reviews | Genetics

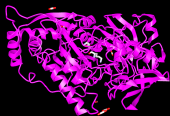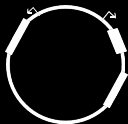Source[Chaisson et al., 2015]

Source: T. Seemann

Repeated regions cannot be resolved with kmers shorter than the repeat!

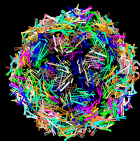# Repeated regions cannot be resolved with kmers shorter than the repeat!



Transporters          Ω Plasmids          Prophages          Ribosomes

# IS IT HOPELESS?

| method | benefits | drawbacks |
|---|---|---|
| PCR + Sanger | it works | its difficult |
| re-sequencing | improve coverage | issues with repeats |
| long reads | solves repeats | cost, availibility |
| reference assisted | easy to perform | not reliable |

LAW OF THE PROBABILITY LEVER: Slight changes can make highly improbable events almost certain

1. Within a taxanomic group, GC content is largely conserved (kmer strain typing, etc).

Source: David Hard

LAW OF THE PROBABILITY LEVER: Slight changes can make highly improbable events almost certain

1. Within a taxanomic group, GC content is largely conserved (kmer strain typing, etc).
2. Within a taxonomic group, genome size is largely conserved.

Source: David Hard

Law of the Probability Lever: Slight changes can make highly improbable events almost certain

1. Within a taxanomic group, GC content is largely conserved (kmer strain typing, etc).
2. Within a taxonomic group, genome size is largely conserved.
3. Bacterial genomes are dense.

Source: David Hard

LAW OF THE PROBABILITY LEVER: Slight changes can make highly improbable events almost certain

1. Within a taxanomic group, GC content is largely conserved (kmer strain typing, etc).
2. Within a taxonomic group, genome size is largely conserved.
3. Bacterial genomes are dense.
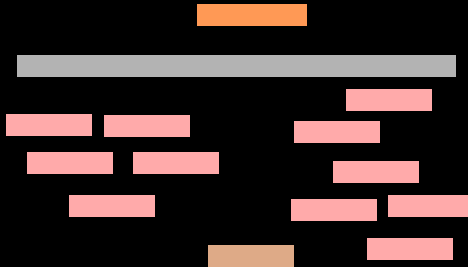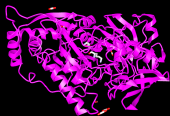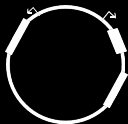4. Nucleotide order is not random.

Source: David Hard

Figure: Bridge Reconstruction. Pink fragments are reads. Grey shows the gene of interest with interupted coverage. Orange fragemnt is a pseudoread generated from this situation under the hypothesis that the beige fragment exists but is underrepresented
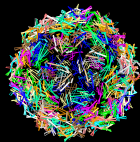
# Repeated regions cannot be resolved with kmers shorter than the repeat!
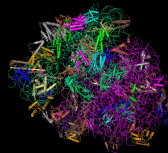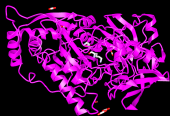


Transporters   Ω Plasmids   Prophages   Ribosomes

# Repeated regions cannot be resolved with kmers shorter than the repeat!



Transporters    Ω Plasmids    Prophages    Ribosomes

rDNA: ribosomal DNA operon

- Prokaryotes: 16S, 23S, 5S
- Conserved within taxa
- Repeated within the genome (1x to >14x)

1. Since the rDNA structure is conserved within taxa, rDNA flanking regions may be conserved
2. Regions flanking the rDNA region will be unique within genomes
3. If flanking regions are unique, they can be used to build "long reads"

# HYPOTHESIS 1: RIBOSOMAL OPERONS

# rDNA flanking regions are conserved conserved

# rDNA flanking regions are conserved

# HYPOTHESIS 2: FLANKING UNIQUENESS

Shannon Entropy by Position
scannedScaffolds

16S ribosomal RNA

23S ribosomal RNA

5S ribosomal RNA

Shannon Entropy

Consensus Coverage

Position (bp)

# HYPOTHESIS 3: LONG READ CONSTRUCTION

- Automated method for constructing select "long reads" from Illumina data
- Written in python3 and R, wrapping barrnap, SMALT, SPAdes, and samtools
- 5 stages:
  1. Identify rDNA clusters
  2. Extracts reads mapping to a cluster
  3. Assemble into long reads
  4. Repeat (3x default) to extend
  5. Submit rDNA long reads to de novo assembly

# DOES IT WORK?

1. synthetic reads on synthetic genome (7 E. coli Sakai rDNAs separated by 6kb random sequence)
2. synthetic reads on real genome
3. short reads from hybrid assembly
4. GAGE-B datasets

# Synthetic reads on synthetic genome

Mauve Demo

# CONCLUSIONS

1. Unpredictable
2. Single problem/solution
3. Biased by reference

⚲ The architecture of bacterial genomes can aid assembly

- The architecture of bacterial genomes can aid assembly
- rDNA flanking regions are unique within a genome

- The architecture of bacterial genomes can aid assembly
- rDNA flanking regions are unique within a genome
- riboSeed improves assemblies at best

- The architecture of bacterial genomes can aid assembly
- rDNA flanking regions are unique within a genome
- riboSeed improves assemblies at best
- riboSeed doesn't work on in all cases, but rarely introduces errors

- Benchmark against GAGE-B
- Benchmark against more hybrid assembly studies
- Find early indicator
- Apply to fungal genomic
- Apply to other conserved regions

📄 Chaisson, M. J. P., Wilson, R. K., and Eichler, E. E. (2015).
Genetic variation and the de novo assembly of human genomes.
Nature Publishing Group, 16.

📄 Compeau, P. E. C., Tesler, G., and Pevzner, P. A. (2011).
How to apply de Bruijn graphs to genome assembly.
Nature biotechnology, 29(11):987–991.

**OÉ Gaillimh**
NUI Galway

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



The James
**Hutton**
**Institute**

- Leighton Pritchard
- Ashleigh Holmnes

**OÉ Gaillimh**
NUI Galway

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

The James
**Hutton**
Institute

- Leighton Pritchard
- Ashleigh Holmnes

# Questions?