UNIVERSITY OF
EXETER

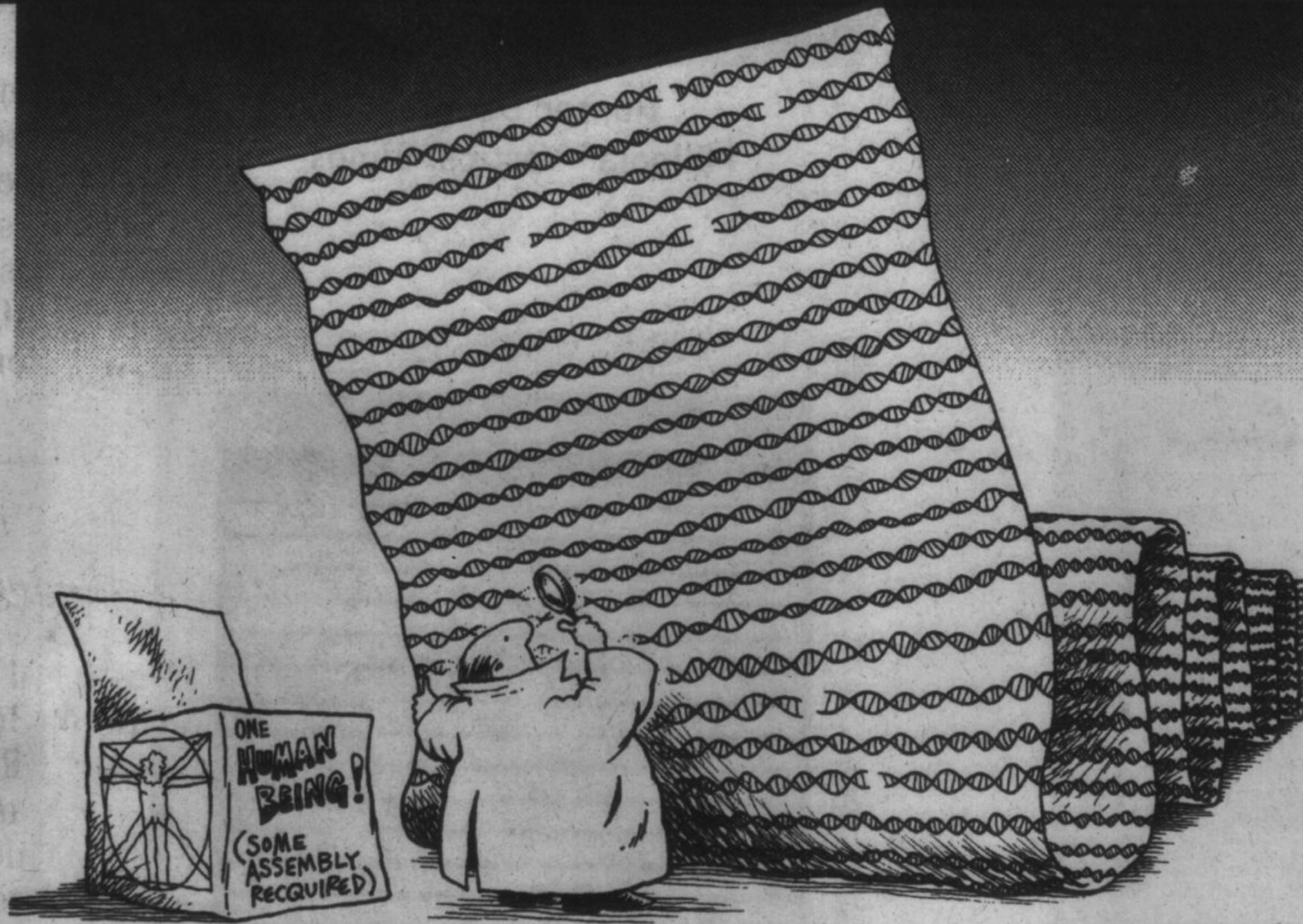# Assembly and annotation using Galaxy

Konrad Paszkiewicz

Sequencing Service, University of Exeter, UK.

25th May 2011

ONE
HUMAN
BEING!
(SOME
ASSEMBLY
REQUIRED)

BY AUTH FOR THE PHILADELPHIA INQUIRER

# Overview

- Why de-novo assembly?
- What is de-novo assembly?
- Types of assemblers
- Annotation
- A toy example in Galaxy
- Future developments

UNIVERSITY OF
EXETER

# Sequencing - 2007

### PRODUCTION

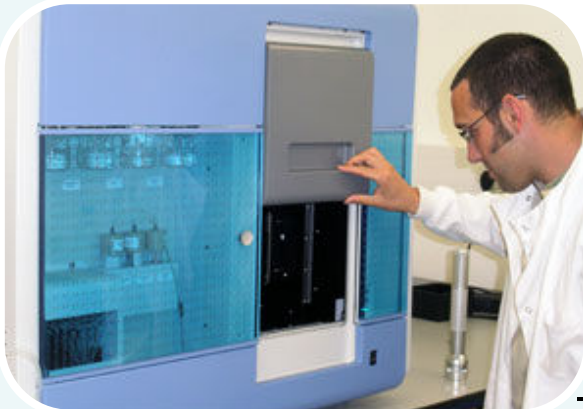Rooms of equipment
Subcloning > picking > prepping
35 FTEs
3-4 weeks

### SEQUENCING

74x Capillary Sequencers
10 FTEs
15-40 runs per day
**1-2Mb per instrument per day**
**120Mb total capacity per day**

UNIVERSITY OF
EXETER

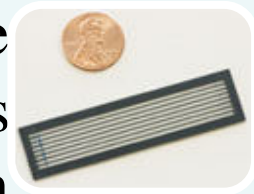# 2<sup>nd</sup> generation sequencing - Today

**PRODUCTION**

1x Cluster Station
1 FTE
1 day

**SEQUENCING**

1x Genome Analyzer
Same FTE as above
1 run per 3-10 days
- 90Gb per instrument per run

UNIVERSITY OF
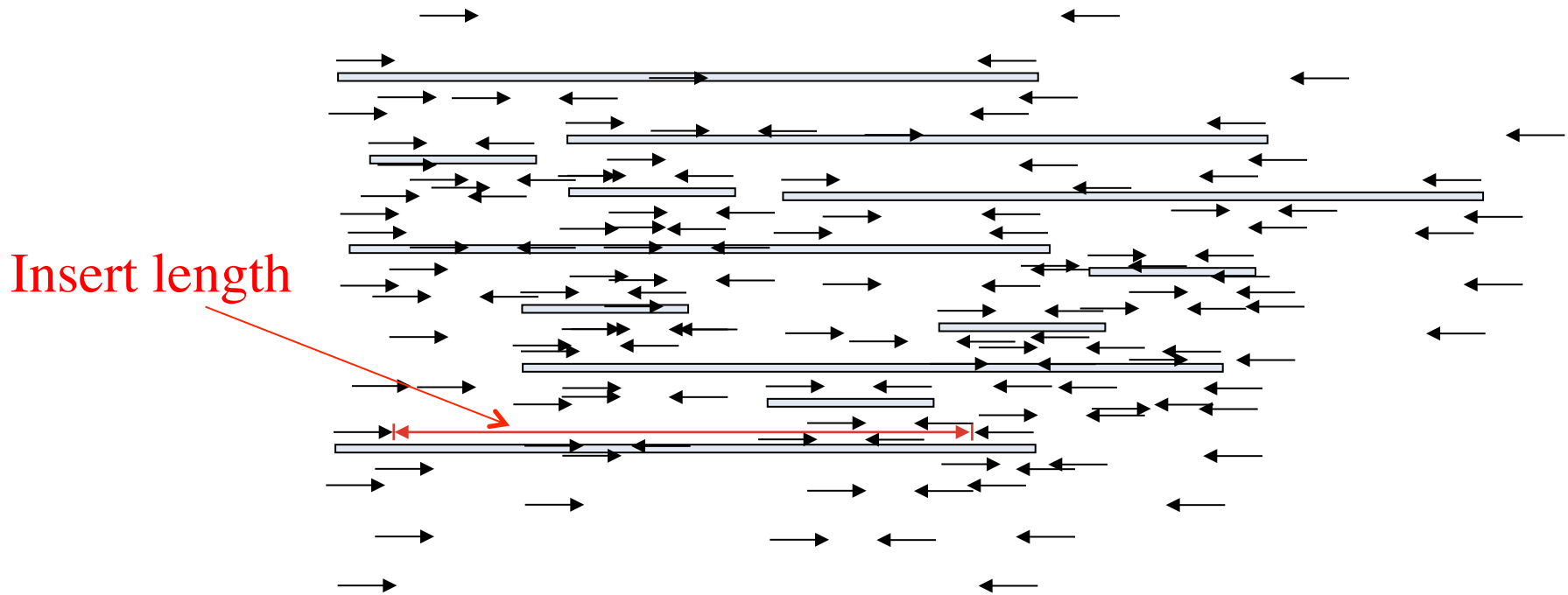EXETER

# Why de-novo assembly?

# Why is de-novo sequence assembly useful?

- No reference genome available
- What is the most suitable reference genome?
  (e.g.    species definition problem in bacteria)

  – What's new in a genome?

    • Remapping will not tell you what is new in a genome (e.g. plasmids, novel genes, novel chromosomes)

  – What's really missing from a genome?

    • Remapping may fail to detect homologous regions
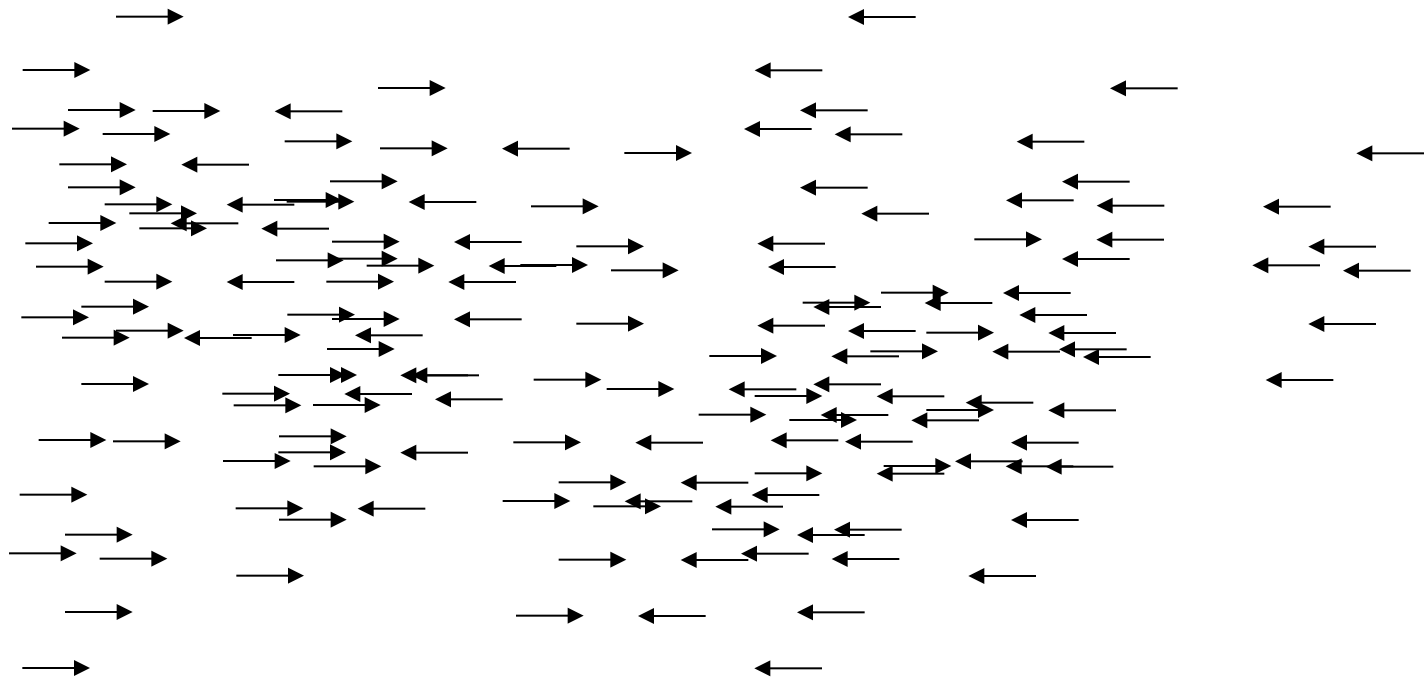
# What is de-novo assembly?

# De-novo sequence assembly

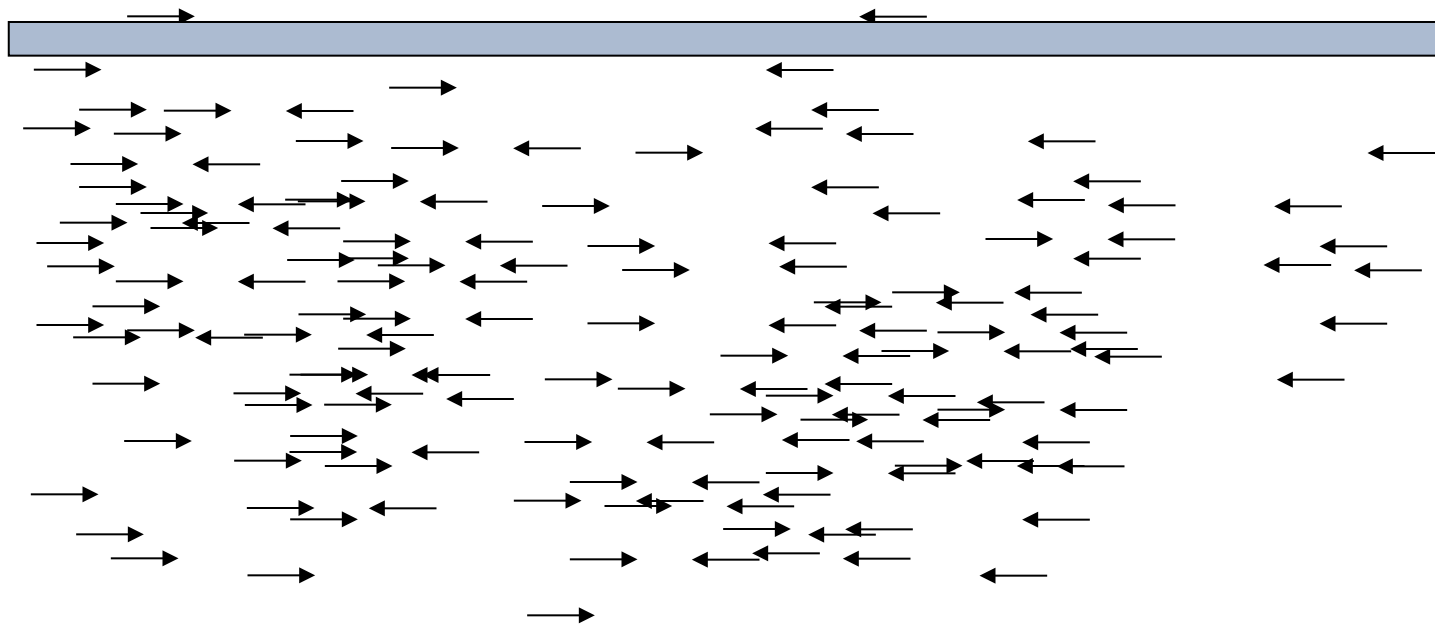1. Sequence DNA fragments from each end



Insert length

# De-novo Sequence Assembly

1. Sequence DNA fragments from each end

2. Reads aligned to generate contigs

# De-novo Sequence Assembly

1. Sequence DNA fragment from each end
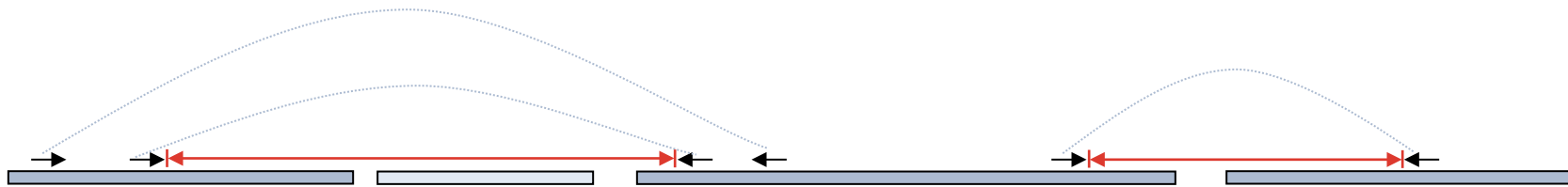
2. Reads aligned to generate contigs

# De-novo Sequence Assembly

1. Sequence clones from each end

2. Reads aligned to generate contigs

3. Supercontigs derived from paired reads on different contigs
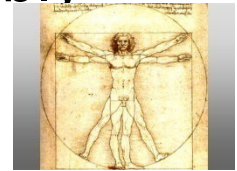
# De-novo Sequence Assembly

1. Sequence reads from each end
2. Reads aligned to generate contigs
3. Supercontigs derived from paired reads on different contigs



4. Ordering of contigs is determined
5. Different insert lengths and read lengths can resolve ambiguities

# De-novo assembly:
# It's not just for genomes.

1. Traditional single homogenous genome assembly

2. Single organism transcriptomes without a reference
   - Estimates of expression

3. Genomic/transcriptomic assembly of symbionts and metagenomes

# Metagenomics

# Denovo Sequence Assembly

- **Caveats**
  - No assembly is perfect
  - Assemblies from 2nd generation tend to be worse in a number of ways than Sanger based-assemblies
    - \+ Easier to generate data
    - \+ Easier to generate lots of assemblies
    - \- Shorter reads/higher error rates
    - \- Man/brainpower is more thinly spread
    - \- Harder to evaluate assemblies
    - \- Harder to annotate and compare between samples

UNIVERSITY OF
EXETER

# Types of assemblers

# Types of assemblers

- 4 categories, many variations

- Each tends to have its own niche

- Memory and hardware requirements can differ substantially

- Galaxy has support either in-built or via Galaxy Tool-shed for Velvet, MIRA, AbySS, Phrap Newbler

- **Typically a parameter scan is need to get the 'best' assembly**

| Name | Read Type | Algorithm | Reference |
|---|---|---|---|
| SUTTA | long & short | B&B | (Narzisi and Mishra [25], 2010) |
| ARACHNE | long | OLC | (Batzoglou al. [14], 2002) |
| CABOG | long & short | OLC | (Miller et al. [13], 2008) |
| Celera | long | OLC | (Myers et al. [12], 2000) |
| Edena | short | OLC | (Hernandez et al. [16], 2008) |
| Minimus (AMOS) | long | OLC | (Sommer et al. [15], 2007) |
| Newbler | long | OLC | 454/Roche |
| CAP3 | long | Greedy | (Huang and Madan [7], 1999) |
| PCAP | long | Greedy | (Huang et al. [8], 2003) |
| Phrap | long | Greedy | (Green [6], 1996) |
| Phusion | long | Greedy | (Mullikin and Ning [9], 2003) |
| TIGR | long | Greedy | (Sutton et al. [5], 1995) |
| ABySS | short | SBH | (Simpson et al. [19], 2009) |
| ALLPATHS | short | SBH | (Butler et al. [46,47], 2008/2011) |
| Euler | long | SBH | (Pevzner et al. [17], 2001) |
| Euler-SR | short | SBH | (Chaisson and Pevzner [35], 2008) |
| Ray | long & short | SBH | (Boisvert et al. [48], 2010) |
| SOAPdenovo | short | SBH | (Li et al. [20], 2010) |
| Velvet | long & short | SBH | (Zerbino and Birney [18,49], 2008/2009) |
| PE-Assembler | short | Seed-and-Extend | (Ariyaratne and Sung [50], 2011) |
| QSRA | short | Seed-and-Extend | (Bryant et al. [23], 2009) |
| SHARCGS | short | Seed-and-Extend | (Dohm et al. [21], 2007) |
| SHORTY | short | Seed-and-Extend | (Hossain et al. [51], 2009) |
| SSAKE | short | Seed-and-Extend | (Warren et al. [22], 2007) |
| Taipan | short | Seed-and-Extend | (Schmidt et al. [24], 2009) |
| VCAKE | short | Seed-and-Extend | (Jeck et al. [52], 2007) |

Reads are defined as "long" if produced by Sanger technology and "short" if produced by Illumina technology . Note that Velvet was designed for micro-reads (e.g. Illumina) but long reads can be given in input as additional data to resolve repeats in a greedy fashion.
doi:10.1371/journal.pone.0019175.t001

Narzisi G, Mishra B, Comparing De Novo Genome Assembly:
The Long and Short of It. 2011  PLoS ONE 6(4):

De novo assembly of short sequence reads
Paszkiewicz, K. Studholme, D.
Briefings in Bioinformatics
August 2010 11(5): 457-472

UNIVERSITY OF
EXETER

# Annotation

# Annotation

Identification of

genes
exons
promoters
signal peptides
regulatory regions
alleles
non-coding RNAs
repeats...

2 broad categories of annotation methodology:
Sequence homology-based (e.g. Blast)
Profile/HMM-based (e.g. PFAM, TMHMM, SignalP)

UNIVERSITY OF
EXETER

# Annotation

To do this effectively it is often necessary to gather additional data:

e.g.

ChIP-Seq
RNA-seq

# Annotation

Exon structure
Transcription start sites



Annotated gene structure

# A toy example in Galaxy

# Denovo sequencing project

A new beta-proteobacterium which secretes elemental metal

60% GC content

Approximately 8 Mb genome

**Method:** 1 lane Illumina sequencing
Mass spectrometry

**Aim:** Which genes(s) are responsible for translocation?

# Process

1. Uploading files from Illumina sequencing

2. Filtering reads

3. De-novo assembly

4. Annotation

5. Locating secretion protein using mass-spectrometry information

# 1. Uploading files

# 2nd generation sequencing output formats

Illumina

SoLID/ABI-Life

Roche 454

Ion Torrent

FASTQ (various flavours)　　Colourspace FASTA　　SFF　　SFF or FASTQ

UNIVERSITY OF
EXETER

# Uploading FASTQ files



Or (maybe) via Get data from UISR A/ENA

# Uploading FASTQ files

# 2. Filtering reads

# All platforms have errors and artefacts



Illumina          SoLID/ABI-Life          Roche 454          Ion Torrent

1. Removal of low quality bases
2. Removal of adaptor sequences
3. Platform specific artefacts (e.g homopolymers)

# Illumina artefacts

## Sequence-specific error profile of Illumina sequencers

Kensuke Nakamura[1,*], Taku Oshima[2], Takuya Morimoto[2,3], Shun Ikeda[1], Hirofumi Yoshikawa[4,5], Yuh Shiwa[5], Shu Ishikawa[2], Margaret C. Linak[6], Aki Hirai[1], Hiroki Takahashi[1], Md. Altaf-Ul-Amin[1], Naotake Ogasawara[2] and Shigehiko Kanaya[1]

[1]Graduate School of Information Science, [2]Graduate School of Biological Sciences, Nara Institute of Science and Technology, 8916-5 Takayama-cho, Ikoma, Nara 630-0192, Japan, [3]Biological Science Laboratories, Kao Corporation, 2606 Akabane, Ichikai, Haga, Tochigi 321-3497, [4]Department of Bioscience, Tokyo University of Agriculture, [5]Genome Research Center, NODAI Research Institute, Tokyo University of Agriculture, 1-1-1 Sakuragaoka Setagaya-ku, Tokyo, 156-8502, Japan and [6]Department of Chemical Engineering and Material Science, University of Minnesota, 223 Amundson Hall, 421 Washington Avenue S.E., Minneapolis, MN 55455, USA

### ABSTRACT

We identified the sequence-specific starting positions of consecutive miscalls in the mapping of reads obtained from the Illumina Genome Analyser (GA). Detailed analysis of the miscall pattern indicated that the underlying mechanism involves sequence-specific interference of the base elongation process during sequencing. The two major sequence patterns that trigger this sequence-specific error (SSE) are: (i) inverted repeats and (ii) GGC sequences. We speculate that these sequences favor dephasing by inhibiting single-base

platforms [Illumina/Solexa Genome Analyser (4), Life Technologies/ABI SOLiD System (5) and Roche/454 Genome Sequencer FLX (6)], the Illumina Genome Analyser (GA) is, at the moment, the most popular choice for the analysis of genomic information (7). The Illumina/Solexa sequencers are characterized by: (i) solid-phase amplification and (ii) a cyclic reversible termination (CRT) process, also termed sequencing-by-synthesis (SBS) technology (8). The sequencer can generate hundreds of millions of relatively short (30–100 bp) read sequences per run.

The application of data obtained from this NGS technology can be roughly categorized into the following three

Nakamura, K. et al Sequence-specific error profile of Illumina sequencers
*Nucl. Acids Res. (2011) May 16, 2011*

UNIVERSITY OF
EXETER

# Illumina artefacts

1. GC rich regions are under represented
     a. PCR
     b. Sequencing
2. Substitutions more common than insertions
3. GGC/GCC motif is associated with low quality and mismatches
4. Filtering low quality reads exacerbates low coverage of GC regions

*Assembly and/or filtering software should account for this technology specific bias but doesn't yet*

# Quality controlling workflow

# Quality controlling workflow

# Quality controlling workflow

# Quality visual summaries

# 3. De-novo Assembly

# Assembly workflow

# Velvet optimiser for genomic de-novo assembly

- De-bruijn graph assembler
- Runs a selection of k-mer lengths and parameters
- Selects optimum assembly based on contig length and N50 size (adjustable)
- Originally written by Simon Gladman, CSIRO
- Available at the Galaxy Tool Shed

# However...

- We need a method of benchmarking the assembly using biological knowledge

- GC value

- Genome size ~ Total number of bp in contigs?

- Fraction of genes fully assembled
  - Measured against closely related genome

- Manual finishing, gap closure only if really necessary

- Most assemblies only need to be 'good-enough'... whatever that means...

# Assembly results

# Assembly statistics

# Assembly statistics

# Taxonomy of contigs

# 4. Annotation

# Annotation workflow

# Still to be included

- De-novo gene prediction
- EST and other evidence needs to be included
- tRNAs
- RepeatMasker
- Non-coding features
- Other annotation software pipelines

# Can we incorporate these?

# Can we incorporate these?

# Do we want to incorporate these?

## Is the service sustainable if it becomes really popular?

## If so, locally? really Web services?

# Denovo sequencing project

A new bacterium which secretes elemental metal

60% GC content

Approximately 8 Mb genome

**Aim:** Which genes(s) are responsible for translocation?

# 5. Where is the secretory protein?

# Mass spectrometry evidence



MTITASQSRTEVVVRSA..

# Locate peptide within contigs ORFs using BlastP

MTITASQSRTEVVVRSA....



Contig 204 ORF 17

# Check with annotation tools

• SignalP predicts a signal peptide using both NN and HMM

• TMHMM also predicts that the peptide is external

• PFAM reports a DUF (Domain of Unknown Function)

• BlastP NR reports Hypothetical proteins

ORF located and characterised as coding for a novel metal export factor

# Summary

- Filtered and formatted raw data

- Assembled a draft 8 Mb genome – no finishing

- Evaluated metrics and taxonomy of contigs

- Called ORFs bacterial codon usage table

- Basic annotation with BlastP against NCBI NR

- PFAM, SignalP, TMHMM

- Identified peptide within contigs

- No hits in PFAM, NCBI NR. Signal peptide present

- Time frame < 1 day

UNIVERSITY OF
EXETER

# Other assemblers

- Minimus2 (Galaxy wrapper by Edward Kirton)
  - Merge contigs from different assemblies

- MIRA (Galaxy wrapper by Peter Cock, SCRI)
  - Recent upgrades for PacBio and Ion Torrent

- AbySS (Galaxy wrapper by Edward Kirton)

- Newbler (Galaxy wrapper by Edward Kirton)
  - Roche/454 proprietary assembler and remapper

- Phrap (Galaxy wrapper by Edward Kirton)
  - Sanger read assembly

- String Graph Assembler (Jared Simpson, Sanger)
  - Useful for large (> human) genomes with short reads

# Available at Galaxy Toolshed

# Other applications

# Oases optimiser for de-novo RNA-seq

- Sister program of Velvet
- Runs a selection of kmer lengths
- Combines all results
- Uses these as a scaffold to assemble transcripts at shortest kmer length

# Galaxy denovo RNA-seq Pipeline

# Future developments

# Community to-do/wish list

- Adding tools dedicated to evaluating assembly quality (e.g. Using EST sequences or related sequences)

- Tools to aid in finishing assemblies

- AFG or other assembly-format visualisation

- Collating and formatting annotation (e.g. GFF files)

- Metagenomics/transcriptomics (e.g. MetaVelvet)

- Gene prediction software

- Blast2Go

- Comparison of GO or PFAM terms between samples

- Enabling workflows of workflows

- AMOS tools (Amos validate etc), web-services

# Future developments

A single Illumina GAIIx run can produce data for ~ 100 bacterial genomes in less than a week.

Cost: ~10,000 Euro

**Question**:    How do we deal with 100s of
comparisons between datasets in Galaxy?

Do we want to?
Do we have a choice?

# Rapid Pneumococcal Evolution in Response to Clinical Interventions

Nicholas J. Croucher,[1] Simon R. Harris,[1] Christophe Fraser,[2] Michael A. Quail,[1] John Burton,[1] Mark van der Linden,[3] Lesley McGee,[4] Anne von Gottberg,[5] Jae Hoon Song,[6] Kwan Soo Ko,[7] Bruno Pichon,[8] Stephen Baker,[9] Christopher M. Parry,[9] Lotte M. Lambertsen,[10] Dea Shahinas,[11] Dylan R. Pillai,[11] Timothy J. Mitchell,[12] Gordon Dougan,[1] Alexander Tomasz,[13] Keith P. Klugman,[4,5,14] Julian Parkhill,[1] William P. Hanage,[2,15] Stephen D. Bentley[1]*

Epidemiological studies of the naturally transformable bacterial pathogen *Streptococcus pneumoniae* have previously been confounded by high rates of recombination. Sequencing 240 isolates of the PMEN1 (Spain[23F]-1) multidrug-resistant lineage enabled base substitutions to be distinguished from polymorphisms arising through horizontal sequence transfer. More than 700 recombinations were detected, with genes encoding major antigens frequently affected. Among these were 10 capsule-switching events, one of which accompanied a population shift as vaccine-escape serotype 19A isolates emerged in the USA after the introduction of the conjugate polysaccharide vaccine. The evolution of resistance to fluoroquinolones, rifampicin, and macrolides was observed to occur on multiple occasions. This study details how genomic plasticity within lineages of recombinogenic bacteria can permit adaptation to clinical interventions over remarkably short time scales.

# DNA sequencing generations

| Then + Now | Now | Now + anticipated | Anticipated |
|---|---|---|---|
| **1st Gen** Sanger | **2nd Gen** -parallised | **3rd Gen** -single mol or electronic | **Next** -single mol AND electronic |
| •Low throughput<br>•High cost<br>•Accurate<br>•Broad user base | •Optical<br>•Amplification needed<br>•Highly parallel<br>•Improved cost and Throughput<br>•New applications | •Optical<br>•Single-molecule<br>•Highly parallel<br>•Cost similar<br>•New applications<br><br>•Or electronic, clonal | •Direct electrical (no optics)<br>•Single-molecule, highly parallel<br>•Transformation of workflow<br>•Designed to broaden user base, deliver step change in cost, power<br>•New applications |
| Sanger | GAII (Solexa/Illumina)<br>SOLiD (ABI/LIFE)<br>454 FLX (454/Roche) | Helicos<br>Pacific Biosciences<br>Ion Torrent<br>(LIFE Starlight) | Nanopores |

Estimated cost of a human genome using these technologies

| | | |
|---|---|---|
| $70M | $200k --- $50k ---- $20k --- 15k--- | ?$5k - $1k? |

# Questions?

**Konrad Paszkiewicz**

**k.h.paszkiewicz@exeter.ac.uk**

"We need to start thinking about how to train people, both health-care professionals and scientists, to be facile in bioinformatics. We need to foster development of professionals who have expertise analyzing large data sets of the size that biologists haven't had to think about. We need to entice smart people into genomics."

*Eric Green,*

*Director National Human Genome Research Institute*

# Acknowledgements

University of Exeter

- Murray Grant
- Karen Moore
- Alex Moorhouse

UNIVERSITY OF
EXETER