# Introduction to Bioinformatics
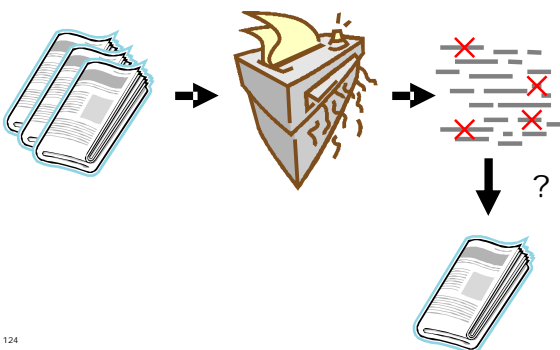
Genome sequencing & assembly

---

## Genome sequencing & assembly

- p DNA sequencing
  - n How do we obtain DNA sequence information from organisms?
- p Genome assembly
  - n What is needed to put together DNA sequence information from sequencing?
- p First statement of sequence assembly problem (according to G. Myers):
  - n Peltola, Söderlund, Tarhio, Ukkonen: Algorithms for some string matching problems arising in molecular genetics. Proc. 9th IFIP World Computer Congress, 1983

123

---

## Recovery of shredded newspaper



124

---

## DNA sequencing

- p DNA sequencing: resolving a nucleotide sequence (whole-genome or less)
- p Many different methods developed
  - n Maxam-Gilbert method (1977)
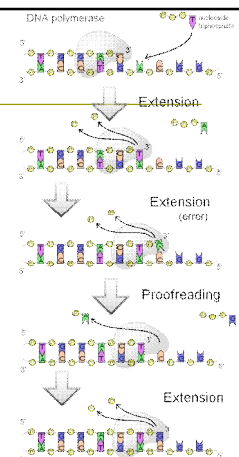  - n Sanger method (1977)
  - n High-throughput methods

125

---

## Sanger sequencing: sequencing by synthesis

- p A sequencing technique developed by Fred Sanger
- p Also called *dideoxy sequencing*

126

---

## DNA polymerase

- p A *DNA polymerase* is an enzyme that catalyzes DNA synthesis
- p DNA polymerase needs a *primer*
  - n Synthesis proceeds always in 5'->3' direction



127

## Dideoxy sequencing

p In Sanger sequencing, chain-terminating dideoxynucleoside triphosphates (ddXTPs) are employed

  n ddATP, ddCTP, ddGTP, ddTTP lack the 3'-OH tail of dXTPs

p A mixture of dXTPs with small amount of ddXTPs is given to DNA polymerase with DNA template and primer
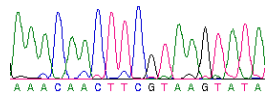
p ddXTPs are given fluorescent labels

## Dideoxy sequencing

p When DNA polymerase encounters a ddXTP, the synthesis cannot proceed

p The process yields copied sequences of different lengths

p Each sequence is terminated by a labeled ddXTP

## Determining the sequence

p Sequences are sorted according to length by capillary electrophoresis

p Fluorescent signals corresponding to labels are registered

p *Base calling*: identifying which base corresponds to each position in a read

  n Non-trivial problem!

AAACAACTTCGTAAGTATA

Output sequences from base calling are called reads

## Reads are short!

p Modern Sanger sequencers can produce quality reads up to ~750 bases[1]

  n Instruments provide you with a quality file for bases in reads, in addition to actual sequence data

p Compare the read length against the size of the human genome ($2.9 \times 10^9$ bases)

p Reads have to be assembled!

[1] *Nature Methods* - 5, 16 - 18 (2008)

## Problems with sequencing

p Sanger sequencing error rate per base varies from 1% to 3%[1]

p Repeats in DNA

  n For example, ~300 base *Alu* sequence repeated is over million times in human genome

  n Repeats occur in different scales

p What happens if repeat length is longer than read length?

  n We will get back to this problem later

[1] Jones, Pevzner (2004)

## Shortest superstring problem

p Find the shortest string that "explains" the reads

p *Given a set of strings (reads), find a shortest string that contains all of them*

## Example: Shortest superstring

Set of strings: {000, 001, 010, 011, 100, 101, 110, 111}

Concetenation of strings: 000001010011100101110111

```
                    010
                    110
                    011
                    000
Shortest superstring: 0001110100
                    001
                     111
                      101
                       100
```

## Shortest superstrings: issues

p NP-complete problem: unlike to have an efficient (exact) algorithm

p Reads may be from either strand of DNA

p Is the shortest string necessarily the correct assembly?

p What about errors in reads?

p Low *coverage* -> gaps in assembly
  n Coverage: average number of times each base occurs in the set of reads (e.g., 5x coverage)

## Sequence assembly and combination locks

p What is common with sequence assembly and opening keypad locks?

## Whole-genome shotgun sequence

p *Whole-genome shotgun sequence assembly* starts with a large sample of genomic DNA

1. Sample is randomly partitioned into *inserts* of length > 500 bases
2. Inserts are multiplied by cloning them into *a vector* which is used to infect bacteria
3. DNA is collected from bacteria and sequenced
4. Reads are assembled

## Assembly of reads with Overlap-Layout-Consensus algorithm

p Overlap
  n Finding potentially overlapping reads

p Layout
  n Finding the order of reads along DNA

p Consensus (Multiple alignment)
  n Deriving the DNA sequence from the layout

p Next, the method is described at a very abstract level, skipping a lot of details

Kececioglu, J.D. and E.W. Myers. 1995. Combinatorial algorithms for DNA sequence assembly. *Algorithmica* 13: 7-51.

## Finding overlaps

p First, pairwise overlap alignment of reads is resolved

p Reads can be from either DNA strand: The *reverse complement* r* of each read r has to be considered

acggagtcc
agtccgcgctt

$r_1$: tgagt, $r_1$*: actca
$r_2$: tccac, $r_2$*: gtgga

## Example sequence to assemble

5' – CAGCGCGCTGCGTGACGAGTCTGACAAAGACGGTATGCGCATCG TGATTGAAGTGAAACGCGATGCGGTCGGTCGGTGAAGTTGTGCT - 3'

p 20 reads:

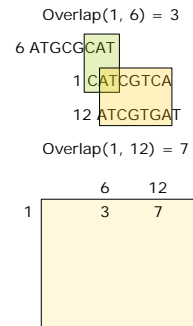| # | Read | Read* | # | Read | Read* |
|---|------|-------|---|------|-------|
| 1 | CATCGTCA | TCACGATG | 11 | GGTCGGTG | CACCGACC |
| 2 | CGGTGAAG | CTTCACCG | 12 | ATCGTGAT | ATCACGAT |
| 3 | TATGCGCA | TGCGCATA | 13 | GCGCTGCG | CGCAGCGC |
| 4 | GACGAGTC | GACTCGTC | 14 | GCATCGTG | CACGATGC |
| 5 | CTGACAAA | TTTGTCAG | 15 | AGCGCGCT | AGCGCGCT |
| 6 | ATGCGCAT | ATGCGCAT | 16 | GAAGTTGT | ACAACTTC |
| 7 | ATGCGGTC | GACCGCAT | 17 | AGTGAAAC | GTTTCACT |
| 8 | CTGCGTGA | TCACGCAG | 18 | ACGCGATG | CATCGCGT |
| 9 | GCGTGACG | CGTCACGC | 19 | GCGCATCG | CGATGCGC |
| 10 | GTCGGTGA | TCACCGAC | 20 | AAGTGAAA | TTTCACTT |

140

## Finding overlaps

p Overlap between two reads can be found with a *dynamic programming* algorithm
   n Errors can be taken into account
p Dynamic programming will be discussed more on next lecture
p Overlap scores stored into the overlap matrix
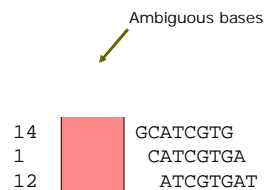   n Entries (i, j) below the diagonal denote overlap of read $r_i$ and $r_j^*$

Overlap(1, 6) = 3

6 ATGCGCAT
1 CATCGTCA
12 ATCGTGAT

Overlap(1, 12) = 7

| | 6 | 12 |
|---|---|---|
| 1 | 3 | 7 |

141

## Finding layout & consensus

p Method extends the assembly *greedily* by choosing the best overlaps
p Both orientations are considered
p Sequence is extended as far as possible

Ambiguous bases

```
14       GCATCGTG
1        CATCGTGA
12        ATCGTGAT
```

Consen

142

## Finding layout & consensus

p We move on to next best overlaps and extend the sequence from there
p The method stops when there are no more overlaps to consider
p A number of contigs is produced
p Contig stands for contiguous sequence, resulting from merging reads

```
2              CGGTGAAG
10             GTCGGTGA
11             GGTCGGTG
7     ATGCGGTC
---------------------
      ATGCGGTCGGTGAAG
```
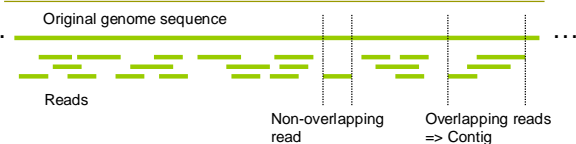
143

## Whole-genome shotgun sequencing: summary

Original genome sequence

Reads

Non-overlapping read

Overlapping reads => Contig

p Ordering of the reads is initially unknown
p Overlaps resolved by aligning the reads
p In a $3 \times 10^9$ bp genome with 500 bp reads and 5x coverage, there are $\sim 10^7$ reads and $\sim 10^7(10^7-1)/2 = \sim 5 \times 10^{13}$ pairwise sequence comparisons
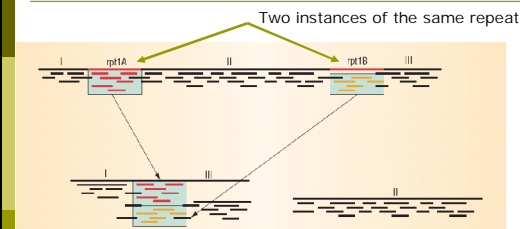
144

## Repeats in DNA and genome assembly

Two instances of the same repeat



Figure 2. Repeat sequence. The top represents the correct layout of three DNA sequences. The bottom shows a repeat collapsed in a misassembly.
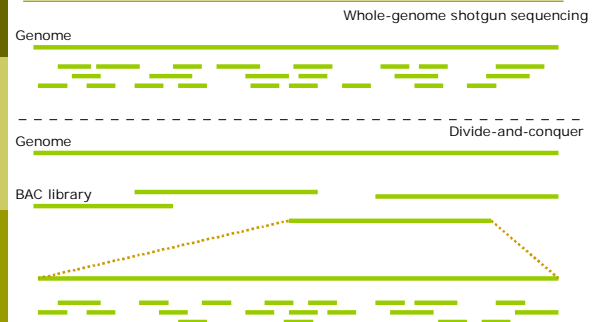
*Pop, Salzberg, Shumway (2002)*

145

## Repeats in DNA cause problems in sequence assembly

- p  Recap: if repeat length exceeds read length, we might not get the correct assembly
- p  This is a problem especially in eukaryotes
  - n  ~3.1% of genome consists of repeats in Drosophila, ~45% in human
- p  Possible solutions
  1. Increase read length – feasible?
  2. Divide genome into smaller parts, with known order, and sequence parts individually

146

## "Divide and conquer" sequencing approaches: BAC-by-BAC



Whole-genome shotgun sequencing
Genome

Genome   Divide-and-conquer

BAC library

147

## BAC-by-BAC sequencing

- p  Each BAC (Bacterial Artificial Chromosome) is about 150 kbp
- p  Covering the human genome requires ~30000 BACs
- p  BACs shotgun-sequenced separately
  - n  Number of repeats in each BAC is significantly smaller than in the whole genome…
  - n  …needs much more manual work compared to whole-genome shotgun sequencing

148

## Hybrid method

- p  Divide-and-conquer and whole-genome shotgun approaches can be combined
  - n  Obtain high coverage from whole-genome shotgun sequencing for short contigs
  - n  Generate of a set of BAC contigs with low coverage
  - n  Use BAC contigs to "bin" short contigs to correct places
- p  This approach was used to sequence the brown Norway rat genome in 2004

149

## Paired end sequencing

- p  *Paired end* (or *mate-pair*) sequencing is technique where
  - n  both ends of an insert are sequenced
  - n  For each insert, we get two reads
  - n  We know the distance between reads, and that they are in opposite orientation



Read 1   k   Read 2

  - n  Typically read length < insert length

150

## Paired end sequencing

- p  The key idea of paired end sequencing:
  - n  Both reads from an insert are unlikely to be in repeat regions
  - n  If we know where the first read is, we know also second's location



Repeat region

Read 1   k   Read 2

- p  This technique helps to WGSS higher organisms

151

5

## First whole-genome shotgun sequencing project: Drosophila melanogaster



- p Fruit fly is a common *model organism* in biological studies
- p Whole-genome assembly reported in Eugene Myers, *et al.*, A Whole-Genome Assembly of *Drosophila*, *Science* 24, 2000
- p Genome size 120 Mbp

http://en.wikipedia.org/wiki/Drosophila_melanogaster

152

## Sequencing of the Human Genome

- p The (draft) human genome was published in 2001
- p Two efforts:
  - n Human Genome Project (public consortium)
  - n Celera (private company)
- p HGP: BAC-by-BAC approach
- p Celera: whole-genome shotgun sequencing



HGP: Nature 15 February 2001 Vol 409 Number 6822

Celera: Science 16 February 2001 Vol 291, Issue 5507

153

## Genome assembly software

- p phrap (Phil's revised assembly program)
- p AMOS (A Modular, Open-Source whole-genome assembler)
- p CAP3 / PCAP
- p TIGR assembler

154

## Next generation sequencing techniques

- p Sanger sequencing is the prominent first-generation sequencing method
- p Many new sequencing methods are emerging

- p See Lars Paulin's slides (course web page) for details

155

## Next-gen sequencing: 454

- p Genome Sequencer FLX (454 Life Science / Roche)
  - n >100 Mb / 7.5 h run
  - n Read length 250-300 bp
  - n >99.5% accuracy / base in a single run
  - n >99.99% accuracy / base in consensus

156

## Next-gen sequencing: Illumina Solexa

- p Illumina / Solexa Genome Analyzer
  - n Read length 35 - 50 bp
  - n 1-2 Gb / 3-6 day run
  - n > 98.5% accuracy / base in a single run
  - n 99.99% accuracy / consensus with 3x coverage

157

## Next-gen sequencing: SOLiD

- SOLiD
  - Read length 25-30 bp
  - 1-2 Gb / 5-10 day run
  - >99.94% accuracy / base
  - >99.999% accuracy / consensus with 15x coverage

158

## Next-gen sequencing: Helicos

- Helicos: Single Molecule Sequencer
  - No amplification of sequences needed
  - Read length up to 55 bp
    - Accuracy does not decrease when read length is increased
    - Instead, throughput goes down
  - 25-90 Mb / h
  - >2 Gb / day

159

## Next-gen sequencing: Pacific Biosciences

- Pacific Biosciences
  - Single-Molecule Real-Time (SMRT) DNA sequencing technology
  - Read length "thousands of nucleotides"
    - Should overcome most problems with repeats
  - Throughput estimate: 100 Gb / hour
  - First instruments in 2010?

160