# Graphs and Genomes

## Michael Schatz

CSH

# Dynamic Programming Matrix

Compute the optimal alignment of ABC…XY..N and DEF…UV…M

|   | 0 | A | B | C | ... | X | Y | ... | N |
|---|---|---|---|---|-----|---|---|-----|---|
| **0** |   |   |   |   |   |   |   |   |   |
| **D** |   |   |   |   |   |   |   |   |   |
| **E** |   |   |   |   |   |   |   |   |   |
| **F** |   |   |   |   |   |   |   |   |   |
| **...** |   |   |   |   |   |   |   |   |   |
| **U** |   |   |   |   |   |   |   |   |   |
| **V** |   |   |   |   |   |   |   |   |   |
| **...** |   |   |   |   |   |   |   |   |   |
| **M** |   |   |   |   |   |   |   |   |   |

# Dynamic Programming Matrix

Compute the optimal alignment of ABC…XY..N and DEF…UV…M

|   | 0 | A | B | C | … | X | Y | … | N |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 1 | 2 | 3 |   | X | X+1 |   | N |
| **D** | 1 |   |   |   |   |   |   |   |   |
| **E** | 2 |   |   |   |   |   |   |   |   |
| **F** | 3 |   |   |   |   |   |   |   |   |
| **…** |   |   |   |   |   |   |   |   |   |
| **U** | U |   |   |   |   |   |   |   |   |
| **V** | U+1 |   |   |   |   |   |   |   |   |
| **…** |   |   |   |   |   |   |   |   |   |
| **M** | M |   |   |   |   |   |   |   |   |

Top row and first column are easy: it takes L-edits to transform and empty string into a length L string

# Dynamic Programming Matrix

Compute the optimal alignment of "ABC…XY..N" and "DEF…UV…M"

|   | **0** | **A** | **B** | **C** | **...** | **X** | **Y** | **...** | **N** |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 1 | 2 | 3 | | X | X+1 | | N |
| **D** | 1 | | | | | | | | |
| **E** | 2 | | | | | | | | |
| **F** | 3 | | | | | | | | |
| **...** | | | | | | | | | |
| **U** | U | | | | | $\gamma$ | $\alpha$ | | |
| **V** | U+1 | | | | | $\beta$ | $\Omega$ | | |
| **...** | | | | | | | | | |
| **M** | M | | | | | | | | |

$$\Omega = \min \begin{cases} \text{"Up" + 1} & \alpha+1 \\ \text{"Left" + 1} & \beta+1 \\ \text{"Diagonal" + 0/1} & \gamma+1 \end{cases}$$

| Up | Left | Diagonal |
|---|---|---|
| ABC...XY- | ABC....X**Y** | ABC...X**Y** |
| DEF....U**V** | DEF...UV- | DEF...U**V** |
| $\alpha$ | $\beta$ | $\gamma$ |

# Biological Networks



Figure 5 Putative regulatory elements shared between groups of correlated and anticorrelated genes

Vanessa M. Brown et al. Genome Res. 2002; 12: 868-884

tRNA Synthetase, tRNA, ATP, Roger Sayle with RasMol, Glaxo Wellcome, 1995

BANANA$

EUKARYA

Stramenopiles
Alveolates
Plantae
Red algae
Animalia
Slime molds
Fungi
Entamoebae
Heteroloboea
Physarum
Kinetoplastids
Euglenoids
Microsporidians
Trichomonads
Diplomonads

BACTERIA
Plant Chloroplasts
Cyanobacteria
Mycoplasma
Agrobacterium
Plant Mitochondria
Enterobacteria
Sulfolobus
Thermoplasma
Halobacteria
Methanobacteria
ARCHAEA

# Graphs



- Nodes
  - People, Proteins, Genes, Neurons, Sequences, Numbers, …

- Edges
  - A is connected to B
  - A is related to B
  - A regulates B
  - A precedes B
  - A interacts with B
  - A activates B
  - …

# Graph Types



List

Tree

Directed Acyclic Graph

Cycle

Complete

# Representing Graphs

## Adjacency Matrix
Good for dense graphs
Fast, Fixed storage: $N^2$ bits

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   | I | I | I |   |   |
| B |   |   |   | I | I |   |   |
| C |   |   |   |   |   | I | I |
| D |   |   |   |   |   | I |   |
| E |   |   |   |   |   | I |   |
| F |   |   |   |   |   |   | I |
| G |   |   |   |   |   |   |   |

## Adjacency List
Good for sparse graphs
Compact storage: 4 bytes/edge

A: C, D, E          D: F
B: D, E             E: F
C: F, G             G:

## Edge List
Easy, good if you (mostly) need
to iterate through the edges
8 bytes / edge

A,C          B,C          C,F
A,D          B,D          C,G
A,E          B,E          D,F
        E,F          F,G

**_Tools_**
**_Matlab:_** http://www.mathworks.com/
**_Graphviz_**:  http://www.graphviz.org/
**_Gephi_**: https://gephi.org/
**_Cytoscape_**: http://www.cytoscape.org/

```
digraph G {
      A->B
      B->C
      A->C
}
dot –Tpdf -og.pdf g.dot
```

# Network Characteristics

| | *C. elegans* | *D. melanogaster* | *S. cerevisiae* |
|---|---|---|---|
| # Nodes | 2646 | 7464 | 4965 |
| # Edges | 4037 | 22831 | 17536 |
| Avg. / Max Degree | 3.0 / 187 | 6.1 / 178 | 7.0 / 283 |
| # Components | 109 | 66 | 32 |
| Largest Component | 2386 | 7335 | 4906 |
| Diameter | 14 | 12 | 11 |
| Avg. Shortest Path | 4.8 | 4.4 | 4.1 |
| Data Sources | 2H | 2x2H, TAP-MS | 8x2H, 2xTAP, SUS |
| Degree Distributions |  |  |  |

***Small World***: Avg. Shortest Path between nodes is small
***Scale Free***: Power law distribution of degree – preferential attachment

# Network Motifs

- ## Network Motif
  - Simple graph of connections
  - Exhaustively enumerate all possible 1, 2, 3, … k node motifs

- ## Statistical Significance
  - Compare frequency of a particular network motif in a real network as compared to a randomized network

- ## Certain motifs are "characteristic features" of the network



**Network Motifs: Simple Building Blocks of Complex Networks**
Milo et al (2002) *Science. 298:824-827*

# Modularity

- Community structure
  - Densely connected groups of vertices, with only sparser connections between groups
  - Reveals the structure of large-scale network data sets



- Modularity
  - The number of edges falling within groups minus the expected number in an equivalent network with edges placed at random
  - Larger positive values => Stronger community structure
  - Optimal assignment determined by computing the eigenvector of the modularity matrix

$$Q = \frac{1}{4m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) (s_i s_j + 1)$$

Normalization factor    Adjacency matrix    Indicates same group

Random Prob. (product of degrees)

**Modularity and community structure in networks.**
Newman ME (2006) *PNAS. 103(23) 8577-8582*

# Kevin Bacon and Bipartite Graphs

Find the ***shortest***
path from
Kevin Bacon
to
Jason Lee

72

60

35

Breadth First Search:
4 hops

31

Bacon Distance:
2

45

# BFS

**BFS(start, stop)**
// initialize all nodes dist = -1
start.dist = 0
list.addEnd(start)
while (!list.empty())
  ***cur = list.begin()***
  if (cur == stop)
    print cur.dist;
  else
    foreach child in cur.children
      if (child.dist == -1)
        child.dist = cur.dist+1
        ***list.addEnd(child)***

0

A,B,C
B,C,D,E
C,D,E,F,L

D,E,F,L,G,H
E,F,L,G,H,I
F,L,G,H,I,J
L,G,H,I,J,X
G,H,I,J,X,O
H,I,J,X,O

I,J,X,O,M
J,X,O,M
X,O,M,N
O,M,N
M,N

N

[How many nodes will it visit?]

[What's the running time?]

[What happens for disconnected components?]

D:2 — I:3
A:1 — E:2 — J:3 — N:4
0 — B:1 — F:2 — X:3
C:1 — G:2 — L:2 — O:3
H:2 — M:3

# BFS

**BFS(start, stop)**
// initialize all nodes dist = -1
start.dist = 0
list.addEnd(start)
while (!list.empty())
  *cur = list.begin()*
  if (cur == stop)
    print cur.dist;
  else
    foreach child in cur.children
      if (child.dist == -1)
        child.dist = cur.dist+1
        *list.addEnd(child)*

0

A,B,C
B,C,D,E
C,D,E,F,L

D,E,F,L,G,H
E,F,L,G,H,I
F,L,G,H,I,J
L,G,H,I,J,X
G,H,I,J,X,O
H,I,J,X,O

I,J,X,O,M
J,X,O,M
X,O,M,N
O,M,N
M,N

N

# DFS

**DFS(start, stop)**
// initialize all nodes dist = -1
start.dist = 0
list.addEnd(start)
while (!list.empty())
  *cur = list.end()*
  if (cur == stop)
    print cur.dist;
  else
    foreach child in cur.children
      if (child.dist == -1)
        child.dist = cur.dist+1
        *list.addEnd(child)*

0

A,B,C

A,B,G,H
A,B,G,M

A,B,G
A,B,L
A,B,O
A,B,N
A,B,J
A,B,E,F
A,B,E,K
A,B,E

A,B

A
D
I

# BFS and TSP

- BFS computes the shortest path between a pair of nodes in $O(|E|) = O(|N|^2)$

- What if we wanted to compute the shortest path visiting every node once?
  - Traveling Salesman Problem

ABDCA: 4+2+5+3 = 14

ACDBA: 3+5+2+4 = 14*

ABCDA: 4+1+5+1 = 11

ADCBA: 1+5+1+4 = 11*

ACBDA: 3+1+2+1 = 7

ADBCA: 1+2+1+3= 7 *

# Greedy Search

# Greedy Search

***Greedy Search***

cur=graph.randNode()

while (!done)

   next=cur.getNextClosest()



Greedy:   ABDCA = 5+8+10+50= 73

Optimal:   ACBDA = 5+11+10+12 = 38

Greedy finds the global optimum only when

1. Greedy Choice: Local is correct without reconsideration
2. Optimal Substructure: Problem can be split into subproblems

Optimal Greedy: Making change with the fewest number of coins

# TSP Complexity

- ## No fast solution
  - Knowing optimal tour through n cities doesn't seem to help much for n+1 cities

  [How many possible tours for n cities?]

- ## Extensive searching is the only provably correct algorithm
  - Brute Force: $O(n!)$
    - ~20 cities max
    - $20! = 2.4 \times 10^{18}$

# Branch-and-Bound

- Abort on suboptimal solutions as soon as possible
  - ADBECA = 1+2+2+2+3 = 10
  - ABDE = 4+2+30 > 10
  - ADE = 1+30 > 10
  - AED = 1+30 > 10
  - …

- Performance Heuristic
  - Always gives the optimal answer
  - Doesn't always help performance, but often does
  - Current TSP record holder:
    - 85,900 cities
    - $85900! = 10^{386526}$

[When not?]

# TSP and NP-complete

- TSP is one of many extremely hard problems of the class NP-complete

  - Extensive searching is the only way to find an exact solution

  - Often have to settle for approx. solution

- WARNING: Many biological problems are in this class

  - Find a tour the visits every node once (Genome Assembly)

  - Find the smallest set of vertices covering the edges (Essential Genes)

  - Find the largest clique in the graph (Protein Complexes)

  - Find the highest mutual information encoding scheme (Neurobiology)

  - Find the best set of moves in tetris

  - …

  - http://en.wikipedia.org/wiki/List_of_NP-complete_problems

# Break

# What is your genome?

Like Dickens, we must computationally reconstruct a genome from short fragments

# Sequencing a Genome

1. Shear & Sequence DNA

2. Construct assembly graph from overlapping reads

...AGCCTAGGGATGCGCGACACGT

GGATGCGCGACACGTCGCATATCCGGTTTGGTCAACCTCGGACGGAC

CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

# Assembly Complexity

# Assembly Complexity



**The advantages of SMRT sequencing**
Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology.* 14:405

# Milestones in Genome Assembly



1977. Sanger *et al.*
1st Complete Organism
5375 bp



1995. Fleischmann *et al.*
1st Free Living Organism
TIGR Assembler. 1.8Mbp



1998. C.elegans SC
1st Multicellular Organism
BAC-by-BAC Phrap. 97Mbp



2000. Myers *et al.*
1st Large WGS Assembly.
Celera Assembler. 116 Mbp



2001. Venter *et al.*, IHGSC
Human Genome
Celera Assembler/GigaAssembler. 2.9 Gbp



2010. Li *et al.*
1st Large SGS Assembly.
SOAPdenovo 2.2 Gbp

# Assembly Applications

- Novel genomes

- Metagenomes

- Sequencing assays
  - Structural variations
  - Transcript assembly
  - …

# Ingredients for a good assembly

## Coverage



### High coverage is required
- Oversample the genome to ensure every base is sequenced with long overlaps between reads
- Biased coverage will also fragment assembly

## Read Length



### Reads & mates must be longer than the repeats
- Short reads will have *false overlaps* forming hairball assembly graphs
- With long enough reads, assemble entire chromosomes into contigs

## Quality



### Errors obscure overlaps
- Reads are assembled by finding kmers shared in pair of reads
- High error rate requires very short seeds, increasing complexity and forming assembly hairballs

**Current challenges in *de novo* plant genome sequencing and assembly**
Schatz MC, Witkowski, McCombie, WR (2012) *Genome Biology*. 12:243

# Typical sequencing coverage



Imagine raindrops on a sidewalk

We want to cover the entire sidewalk but each drop costs $1

# 1x sequencing



Balls in Bins
Total balls: 1000

Histogram of balls in each bin
Total balls: 1000  Empty bins: 361

# 8x sequencing



Balls in Bins
Total balls: 8000

Histogram of balls in each bin
Total balls: 8000  Empty bins: 1

# Poisson Distribution

The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

Formulation comes from the limit of the binomial equation

Resembles a normal distribution, but over the positive values, and with only a single parameter.

*Key property:*
- *The standard deviation is the square root of the mean.*

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

# de Bruijn Graph Construction

- $D_k = (V,E)$
  - $V$ = All length-k subfragments (k < l)
  - E = Directed edges between consecutive subfragments
    - Nodes overlap by k-1 words

Original Fragment

| It was the best of |

Directed Edge

| It was the best | → | was the best of |

- Locally constructed graph reveals the global sequence structure
  - Overlaps between sequences implicitly computed

de Bruijn, 1946
Idury and Waterman, 1995
Pevzner, Tang, Waterman, 2001

# de Bruijn Graph Assembly

It was the best

was the best of

the best of times,

best of times, it

of times, it was

times, it was the

it was the worst

was the worst of

the worst of times,

worst of times, it

it was the age

was the age of

the age of foolishness

the age of wisdom,

age of wisdom, it

of wisdom, it was

wisdom, it was the

After graph construction, try to simplify the graph as much as possible

# Repetitive regions

| Repeat Type | Definition / Example | Prevalence |
|---|---|---|
| Low-complexity DNA / Microsatellites | $(b_1 b_2 \ldots b_k)^N$ where $1 \le k \le 6$ <br> CACACACACACACACACA | 2% |
| SINEs (Short Interspersed Nuclear Elements) | *Alu* sequence (~280 bp) <br> Mariner elements (~80 bp) | 13% |
| LINEs (Long Interspersed Nuclear Elements) | ~500 – 5,000 bp | 21% |
| LTR (long terminal repeat) retrotransposons | Ty1-copia, Ty3-gypsy, Pao-BEL (~100 – 5,000 bp) | 8% |
| Other DNA transposons | | 3% |
| Gene families & segmental duplications | | 4% |

- Over 50% of mammalian genomes are repetitive
  - Large plant genomes tend to be even worse
  - Wheat: 16 Gbp; Pine: 24 Gbp

# Repeats and Coverage Statistics

A      $R_1$      B      $R_2$      C         $R_1 + R_2$

- If $n$ reads are a uniform random sample of the genome of length $G$, we expect $k = n\Delta/G$ reads to start in a region of length $\Delta$.

  - If we see many more reads than k (if the arrival rate is > A) , it is likely to be a collapsed repeat

$$\Pr(X - copy) = \binom{n}{k}\left(\frac{X\Delta}{G}\right)^k\left(\frac{G - X\Delta}{G}\right)^{n-k}$$

$$A(\Delta, k) = \ln\left(\frac{\Pr(1 - copy)}{\Pr(2 - copy)}\right) = \ln\left(\frac{\dfrac{(\Delta n/G)^k}{k!}e^{\frac{-\Delta n}{G}}}{\dfrac{(2\Delta n/G)^k}{k!}e^{\frac{-2\Delta n}{G}}}\right) = \frac{n\Delta}{G} - k\ln 2$$

**The fragment assembly string graph**
Myers, EW (2005) Bioinformatics. 21(suppl 2): ii79-85.

# Paired-end and Mate-pairs

*Paired-end sequencing*

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation

300bp

*Mate-pair sequencing*

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads

10kbp

10kbp circle

2x100 @ ~10kbp (outies)

2x100 @ 300bp (innies)

# Scaffolding

- Initial contigs (*aka* unipaths, unitigs) terminate at
  - *Coverage gaps*: especially extreme GC
  - *Conflicts*: errors, repeat boundaries

- Use mate-pairs to resolve correct order through assembly graph
  - Place sequence to satisfy the mate constraints
  - Mates through repeat nodes are tangled

- Final scaffold may have internal gaps called sequencing gaps
  - We know the order, orientation, and spacing, but just not the bases. Fill with Ns instead

# N50 size

Def: 50% of the genome is in contigs as large as the N50 value

Example:  1 Mbp genome

50%

| 1000 |

| 300 | 100 | 45 | 45 | 30 | 20 | 15 | 15 | 10 | . | . | . | . | . |

N50 size = 30 kbp
   (300k+100k+45k+45k+30k = 520k >= 500kbp)

Note:
   N50 values are only meaningful to compare when base genome
   size is the same in all cases

# Whole Genome Alignment with MUMmer

Slides Courtesy of Adam M. Phillippy

University of Maryland

# Goal of WGA

- For two genomes, *A* and *B*, find a mapping from each position in *A* to its corresponding position in *B*



CCGGTAGGCTATTAAACGGGGTGAGGAGCGTTGGCATAGCA

CCGGTAGGCTATTAAACGGGGTGAGGAGCGTTGGCATAGCA

# Not so fast...

- Genome *A* may have insertions, deletions, translocations, inversions, duplications or SNPs with respect to *B* (sometimes all of the above)



CCGGTAGGATATTAAACGGGGTGAGGAGCGTTGGCATAGCA

CCGCTAGGCTATTAAACCCCGGAGGAG....GGCTGAGCA

# WGA visualization

- How can we visualize *whole* genome alignments?

- With an alignment dot plot
  - *N* x *M* matrix
    - Let $i$ = position in genome *A*
    - Let $j$ = position in genome *B*
    - Fill cell *(i,j)* if $A_i$ shows similarity to $B_j$

  - A perfect alignment between *A* and *B* would completely fill the positive diagonal

**Alignment of 2 strains of *Y. pestis***
http://mummer.sourceforge.net/manual/

# 3rd Gen Long Read Sequencing



PacBio RS II

CSHL/PacBio

Oxford Nanopore

CSHL/ONT

# PacBio SMRT Sequencing

Imaging of fluorescently phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



http://www.pacificbiosciences.com/assets/files/pacbio_technology_backgrounder.pdf

# SMRT Sequencing Data



| Match | 83.7% |
|---|---|
| Insertions | 11.5% |
| Deletions | 3.4% |
| Mismatch | 1.4% |

```
TTGTAAGCAGTTGAAAACTATGTGTGGATTTAGAATAAAGAACATGAAAG
||||||||||||||||||||||||| ||||||| ||||||||||||| |||
TTGTAAGCAGTTGAAAACTATGTGT-GATTTAG-ATAAAGAACATGGAAG

ATTATAAA-CAGTTGATCCATT-AGAAGA-AAACGCAAAAGGCGGCTAGG
| ||||||| ||||||||||||| |||| | |||||| |||||| ||||||
A-TATAAATCAGTTGATCCATTAAGAA-AGAAACGC-AAAGGC-GCTAGG

CAACCTTGAATGTAATCGCACTTGAAGAACAAGATTTTATTCCGCGCCCG
| ||||||| |||| || ||||||||||||||||||||||||||||||||
C-ACCTTG-ATGT-AT--CACTTGAAGAACAAGATTTTATTCCGCGCCCG

TAACGAATCAAGATTCTGAAAACACAT-ATAACAACCTCCAAAA-CACAA
| ||||||| |||||||||||||| || || |||||||||| |||||
T-ACGAATC-AGATTCTGAAAACA-ATGAT----ACCTCCAAAAGCACAA

-AGGAGGGGAAAGGGGGGAATATCT-ATAAAAGATTACAAATTAGA-TGA
 ||||||    ||    |||||||| || |||||||||||||| || |||
GAGGAGG---AA-----GAATATCTGAT-AAAGATTACAAATT-GAGTGA

ACT-AATTCACAATA-AATAACACTTTTA-ACAGAATTGAT-GGAA-GTT
||| ||||||||| | |||||||||||| ||| |||||||| |||| |||
ACTAAATTCACAA-ATAATAACACTTTTAGACAAAATTGATGGGAAGGTT

TCGGAGAGATCCAAAACAATGGGC-ATCGCCTTTGA-GTTAC-AATCAAA
|| |||||||||| |||||| ||| ||| |||| |||||| |||||||
TC-GAGAGATCC-AAACAAT-GGCGATCG-CTTTGACGTTACAAATCAAA

ATCCAGTGGAAAATATAATTTATGCAATCCAGGAACTTATTCACAATTAG
||||||| |||||||| |||||| ||||| |||||||||||||||||||||
ATCCAGT-GAAAATATA--TTATGC-ATCCA-GAACTTATTCACAATTAG
```

Sample of 100k reads aligned with BLASR requiring >100bp alignment

# PacBio Assembly Algorithms



| PBJelly | PacBioToCA & ECTools | HGAP & Quiver |

**Gap Filling and Assembly Upgrade**

English *et al* (2012)
*PLOS One.* 7(11): e47768

**Hybrid/PB-only Error Correction**

Koren, Schatz, *et al* (2012)
*Nature Biotechnology.* 30:693–700

**PB-only Correction & Polishing**

Chin *et al* (2013)
*Nature Methods.* 10:563–569

**< 5x**     PacBio Coverage     **> 50x**

# S. cerevisiae W303

PacBio RS II sequencing at CSHL in the McCombie Lab
- Size selection using an 7 Kb elution window on a BluePippin™ device from Sage Science



Over 175x coverage in
16 SMRTcells / 2 days
using P5-C3

Mean: 5910

83x over 10kbp

8.7x over 20kb

Max: 36,861bp

# S. cerevisiae W303

S288C Reference sequence
- 12.1Mbp; 16 chromo + mitochondria; N50: 924kbp

PacBio assembly using HGAP + Celera Assembler
- 12.4Mbp; 21 non-redundant contigs; N50: 811kbp; >99.8% id

# S. cerevisiae W303

S288C Reference sequence
- 12.1Mbp; 16 chromo + mitochondria; N50: 924kbp

PacBio assembly using HGAP + Celera Assembler
- 12.4Mbp; 21 non-redundant contigs; N50: 811kbp; >99.8% id



35kbp repeat cluster

Near-perfect assembly:
All but 1 chromosome assembled as a single contig

PacBio® Advances in Read Length

# Oxford Nanopore MinION

- Thumb drive sized sequencer powered over USB

- Capacity for 512 reads at once

- Senses DNA by measuring changes to ion flow

# Nanopore Sequencing

# Nanopore Basecalling



**Data → Events → Sequence**

ONT1  CCGACTCCGGTTACCCGCGTTGATTTGCTGGGGCAGGGCCG
      |||||||||||||||||:||||||||||||||||||||||||
REF   CCGACTCCGGTTACCAGCGTTGATTTGCTGGGGCAGGGCCG

- Hidden Markov model
- Only four options per transition
- Pore type = distinct kmer length

TAGGG → AGGGG → GGGTG
        AGGGT → GGGTT
        AGGGA → GGGTA
        AGGGC → GGGTC

- Form probabilistic path through measured states currents and transitions
  - e.g. Viterbi algorithm

Basecalling currently performed at Amazon with frequent updates to algorithm

# Nanopore Readlengths

noise

Spike-in

**Oxford Nanopore Sequencing at CSHL**
30 runs, 267k reads, 122x total coverage
Between 11 and 73k reads per run!
Mean flow cell: 50 Mbp in 2 days
Max flow cell: 446Mbp in 2 days

Mean: 5473bp

41x over 10kbp

8x over 20kb

Max: 146,992bp

# Nanopore Alignments

Mean: 6903bp

**Alignment Statistics (BLASTN)**
Mean read length at ~7kbp
Shearing targeted 10kbp
70k reads align (32%)
40x coverage

13.8x over 10kbp

1.8x over 20kb

Max: 50,900bp

# Nanopore Accuracy

**Alignment Quality (BLASTN)**
Of reads that align, average ~64% identity

# Nanopore Accuracy

**Alignment Quality (BLASTN)**
Of reads that align, average ~64% identity
"2D base-calling" improves to ~70% identity

# NanoCorr: Nanopore-Illumina Hybrid Error Correction

https://github.com/jgurtowski/nanocorr

1. BLAST Miseq reads to all raw Oxford Nanopore reads

2. Select non-repetitive alignments
   - First pass scans to remove "contained" alignments
   - Second pass uses Dynamic Programming (LIS) to select set of high-identity alignments with minimal overlaps

3. Compute consensus of each Oxford Nanopore read
   - Currently using Pacbio's pbdagcon



Post-correction %ID
Mean: ~97%

# Long Read Assembly

## S288C Reference sequence

- 12.1Mbp; 16 chromo + mitochondria
- Chromosome N50: 924kbp

### *Illumina MiSeq*
30x, 300bp PE (Flashed)

Celera Assembler
- 6953 non-redundant contigs
- N50: 59kb >99.9% id



### *Oxford Nanopore*
30x corrected reads > 6kb

NanoCorr + Celera Assembler
- 234 non-redundant contigs
- N50: 362kbp  >99.78% id



### *Pacific Biosciences*
25x corrected reads > 10kb

HGAP + Celera Assembler
- 21 non-redundant contigs
- N50: 811kb >99.8% id

# Assembly Summary

Assembly quality depends on

1. **Coverage**: low coverage is mathematically hopeless
2. **Repeat composition**: high repeat content is challenging
3. **Read length**: longer reads help resolve repeats
4. **Error rate**: errors reduce coverage, obscure true overlaps

- Assembly is a hierarchical, starting from individual reads, build high confidence contigs/unitigs, incorporate the mates to build scaffolds
  - Extensive error correction is the key to getting the best assembly possible from a given data set

- Watch out for collapsed repeats & other misassemblies
  - Globally/Locally reassemble data from scratch with better parameters & stitch the 2 assemblies together

# Thank You



http://schatzlab.cshl.edu/teaching/
@mike_schatz