# Genome Sequencing

CMSC 423
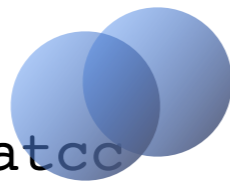Carl Kingsford

# Genome Sequencing

ACCGTCCAATTGG...
TGGCAGGTTAACC...

E.g. human: 3 billion bases
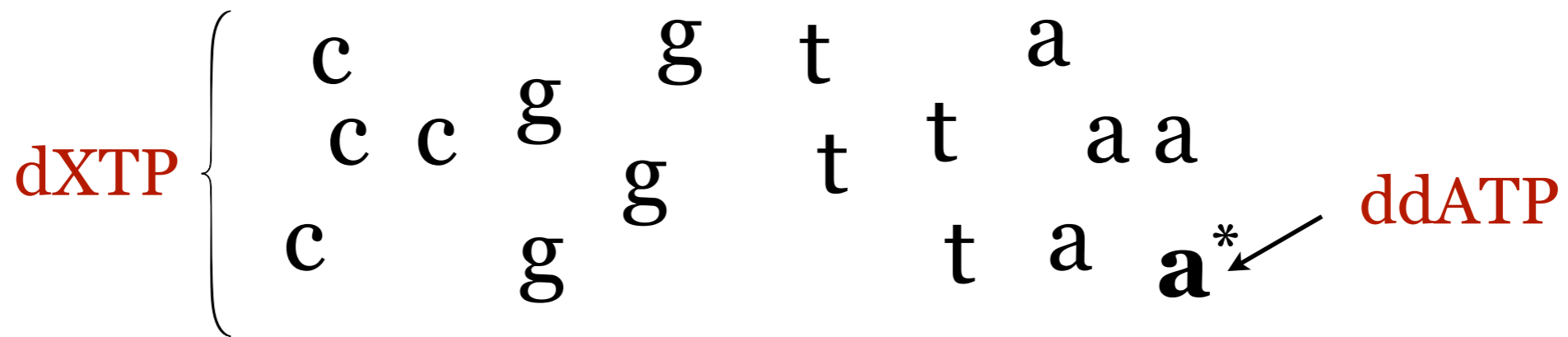split into 23 chromosomes

Main tool of traditional sequencing: DNA Synthesis

*DNA polymerase*: enzyme that will
grow a complementary DNA strand.

gacgatcggtttatcc
ctgctagccaaataggctaatactacgga

# Sanger Sequencing: Finding the As

$$\text{dXTP} \left\{ \begin{array}{c} \text{c} \quad\quad \text{g} \quad\quad \text{t} \quad\quad\quad \text{a} \\ \text{c} \quad \text{c} \quad \text{g} \quad\quad\quad\quad \text{t} \quad \text{a} \; \text{a} \\ \text{g} \quad \text{g} \quad\quad \text{t} \\ \text{c} \quad\quad\quad \text{g} \quad\quad\quad\quad \text{t} \quad \text{a} \quad \mathbf{a}^* \end{array} \right.$$

ddATP

gacgatcgg ttt**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**tt**A**tg**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**tt**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**tt**A**tg**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**tt**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

gacgatcgg ttt**A**tccg**A**\*
ctgctagcc aaa**T**aggc**T**aa**T**ac**T**acgga

# Size → Sequence

gacgatcggtttt**A***

gacgatcggtttt**A***

gacgatcggtttt**A**ccg**A***

gacgatcggtttt**A**ccg**A***

gacgatcggtttt**A**ccg**A**tt**A***

gacgatcggtttt**A**ccg**A**tt**A***

gacgatcggtttt**A**ccg**A**tt**A**tg**A***

gacgatcggtttt**A**ccg**A**tt**A**tg**A***


gacgatcggtttatccgattat**G***

gacgatcggtttatcc**G***


gacgatcggtttat**C***

gacgatcggtttatc**C***

# Size → Sequence

Single lane: ddXTP
that fluoresce
different colors

gacgatcggtttt**A**\*

gacgatcggtttt**A**\*

gacgatcggtttt**A**tccg**A**\*

gacgatcggtttt**A**tccg**A**\*

gacgatcggtttt**A**tccg**A**tt**A**\*

gacgatcggtttt**A**tccg**A**tt**A**\*

gacgatcggtttt**A**tccg**A**tt**A**tg**A**\*

gacgatcggtttt**A**tccg**A**tt**A**tg**A**\*

gacgatcggtttatccgattat**G**\*

gacgatcggtttatcc**G**\*

gacgatcggtttat**C**\*

gacgatcggtttatc**C**\*



Size

A  C  G  T

# Size → Sequence

Single lane: ddXTP that fluoresce different colors

gacgatcggtttt**A***

gacgatcggtttt**A***

gacgatcggttt**A**ccg**A***

gacgatcggttt**A**ccg**A***

gacgatcggttt**A**ccg**A**tt**A***

gacgatcggttt**A**ccg**A**tt**A***

gacgatcggttt**A**ccg**A**tt**A**tg**A***

gacgatcggttt**A**ccg**A**tt**A**tg**A***

gacgatcggtttatccgattat**G***

gacgatcggtttatcc**G***

gacgatcggtttat**C***

gacgatcggtttatc**C***

Size
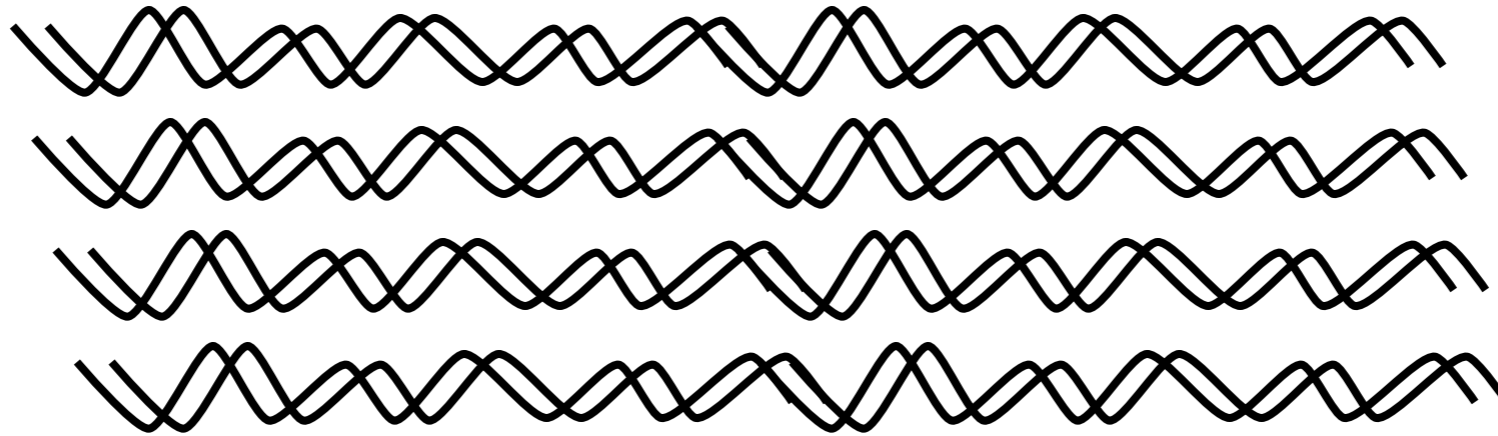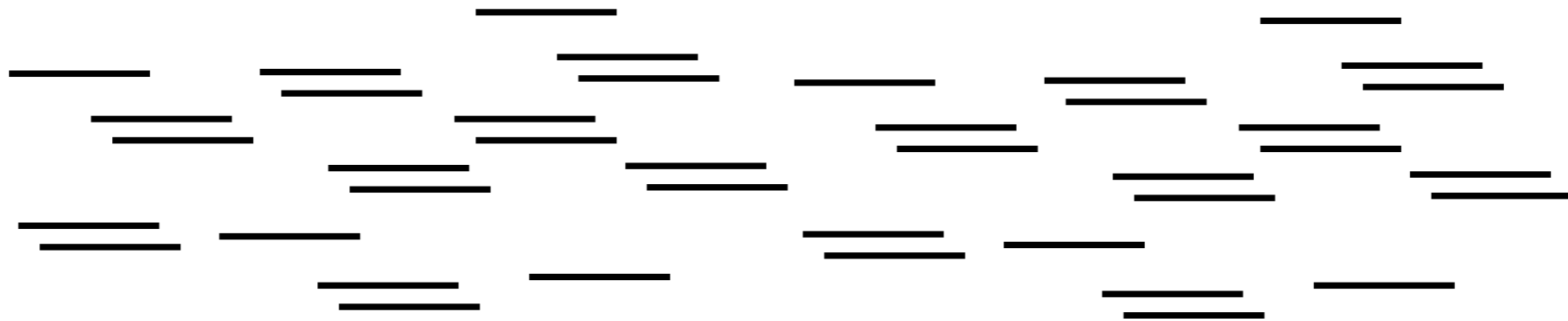
A   C   G   T

Main problem: larger fragments take a long time to be sorted correctly (or don't sort correctly ever) → 800-1000 letter maximum
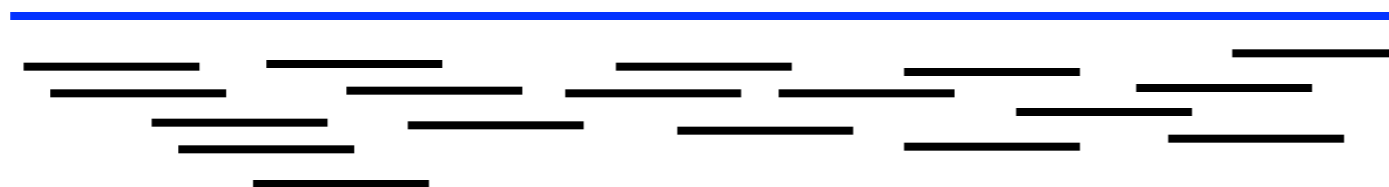
# Shotgun Sequencing

Many copies
of the DNA

Shear it, randomly breaking them into many small pieces,
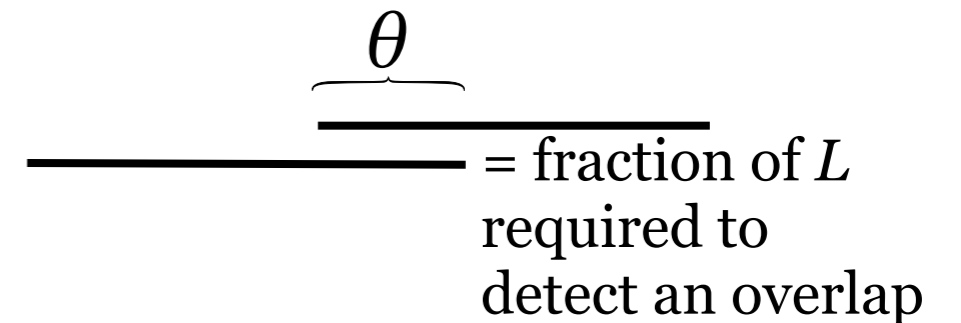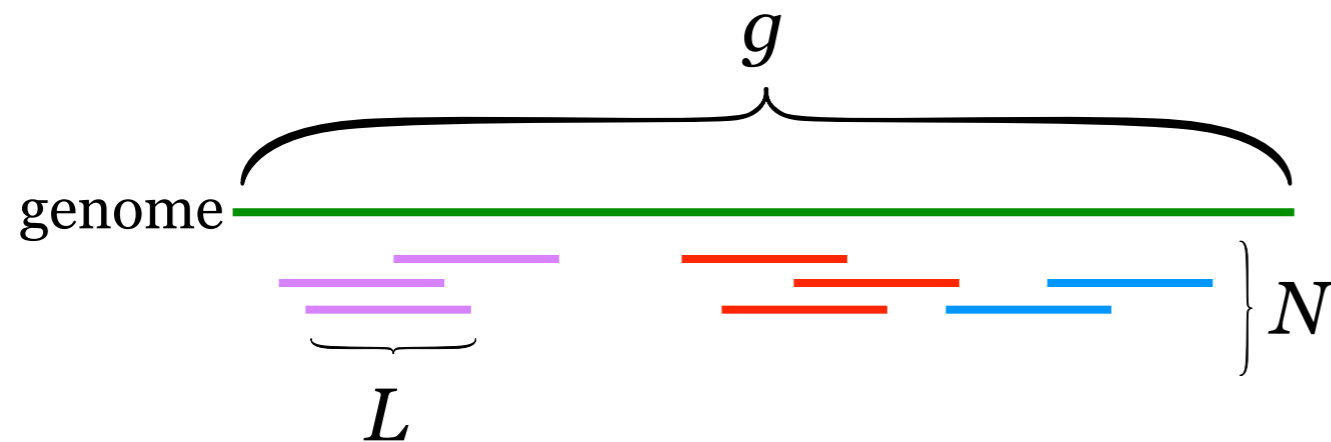read ends of each:

Assemble into original genome:

# Lander-Waterman Statistics

How many reads to we need to be sure we cover the whole genome?



An ***island*** is a contiguous group of reads that are connected by overlaps of length $\geq \theta L$.
(Various colors above)

Want: Expression for expected # of islands given $N, g, L, \theta$.

# Expected # of Islands

$\lambda := N/g$ = probability a read starts at a given position
(assuming random sampling)

Pr($k$ reads start in an interval of length $x$)
  $x$ trials, want $k$ "successes," small probability $\lambda$ of success
  Expected # of successes = $\lambda x$
  Poisson approximation to binomial distribution:

$$\mathrm{Pr}(k \text{ reads in length } x) = e^{-\lambda x}\frac{(\lambda x)^k}{k!}$$

Expected # of islands = $N \times$ Pr(read is at rightmost end of island)

$\underline{\phantom{xxx}(1\text{-}\theta)\text{L}\phantom{xxx}}|\phantom{x}\theta\text{L}$    $= N \times$ Pr(0 reads start in $(1\text{-}\theta)L$)

$$= Ne^{-\lambda(1-\theta)L}\frac{\lambda^0}{0!} \quad \text{(from above)}$$

$$= Ne^{-\lambda(1-\theta)L}$$

$$= Ne^{-(1-\theta)LN/g} \quad \leftarrow LN/g \text{ is called the } \textbf{coverage } c.$$
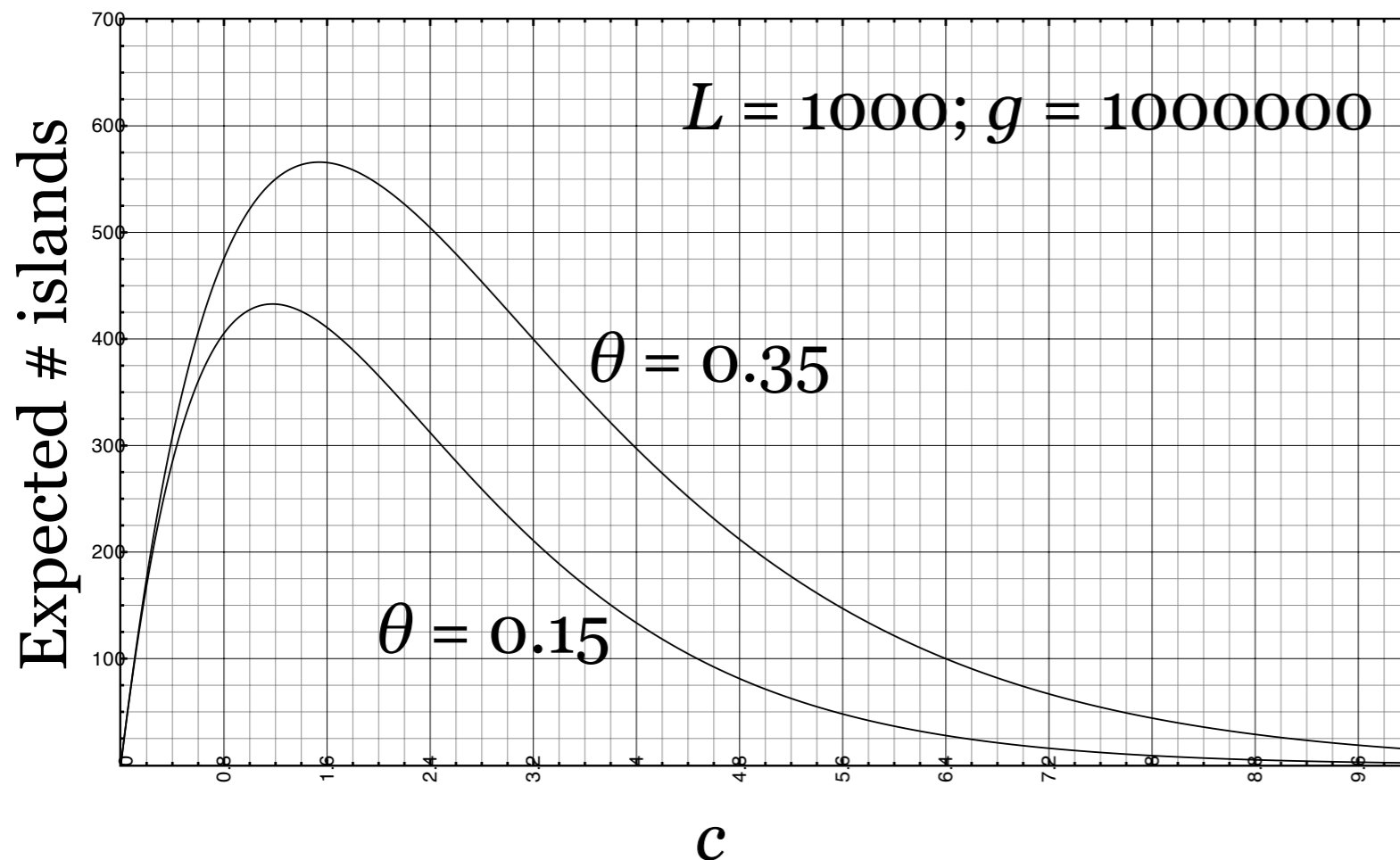
# Expected # of Islands, 2

Rewrite to depend more directly on the things we can control: c and $\theta$

$$\text{Expected \# of islands} = Ne^{-(1-\theta)LN/g}$$

$$= Ne^{-(1-\theta)c}$$

$$= \frac{L/g}{L/g}Ne^{-(1-\theta)c}$$

$$= \frac{g}{L}ce^{-(1-\theta)c}$$



L = 1000; g = 1000000

$\theta = 0.35$

$\theta = 0.15$

# Summary

- "Sanger" sequencing widely used up through 2006 or 2007, including for the human genome project.

- Won Sanger his second Nobel prize.

- Lander-Waterman statistics estimate the number of islands you will get for a given coverage.

  - Used as a way to guess how much sequencing you need to do for a given technology and genome size.

  - Often hard in practice to guess the genome size g before you've sequenced it.