Let $\lambda = np = E(x)$, so $p = \dfrac{\lambda}{n}$

$$P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$= \frac{n(n-1)\ldots(n-x+1)}{x!}\left(\frac{\lambda}{n}\right)^x \left(1-\frac{\lambda}{n}\right)^{n-x} \sim \frac{n^x}{x!}\left(\frac{\lambda}{n}\right)^x = \frac{\lambda^x}{x!} ;$$

$$\sum_x \frac{\lambda^x}{x!} = e^\lambda .$$

Normalization requires $\displaystyle\sum_x P(X = x) = 1$.

Thus $P(X = x) = \dfrac{\lambda^x}{x!} e^{-\lambda}$

# Poisson Mean & Variance

If X is a Poisson random variable, then:

- Mean: $\mu = E(X) = \lambda$ $\;= n \cdot p$
- Variance: $\sigma^2 = V(X) = \lambda$ $\;= n \cdot p$ (it was $pp(1-p)$ for binomial)
- Standard deviation: $\sigma = \lambda^{1/2}$

Note: Variance = Mean

Note: Standard deviation/Mean = $\lambda^{-1/2}$
        decreases with $\lambda$

# Poisson Distribution in Genome Assembly

# Poisson Example: Genome Assembly

- Goal: DNA sequence of the entire genome of an organism

- Problem: Sequencers generate short reads of random portions of a genome

- Solution: assemble genome from short reads using computers

- Whole Genome Shogun Assembly pioneered by Craig Venter in 1990s

- The human genome was jointly announced in 2001 by the Human Genome Project (public) and Celera Genomics (Craig Venter's company)

# Short Reads assemble into Contigs



Figure 5.1.

# Promise of Genomics



Drew Sheneman, New Jersey -- The Newark Star Ledger, E-mail Drew.

I think I found the corner piece!

Cost per Raw Megabase of DNA Sequence

# Current sequencing technologies

| Technology | Read Length | Error Rate | Cost per Gbase |
|---|---|---|---|
| Illumina NovaSeq | 75-500 bp | ~0.1% | $5-$150 |
| BGI DNBSEQ | 35-300 bp | ~0.1% | $5-$120 |
| Ion Torrent | 200-600 bp | ~0.5% | $70-$1000 |
| PacBio | 10,000-25,000 bp | 13% | $7-$40 |
| Oxford Nanopore | 10,000-100,000+ bp | 3-10% | $30-$60 |

MinION, a palm-sized gene sequencer made by
UK-based Oxford Nanopore Technologies

# How many short reads do we need?

**Input**

**Output**

**Low coverage:**

A few pieces to
assemble

many contigs,
many gaps

**High coverage:**

many pieces to
assemble

a few contigs, a
few gaps

# Genome Assembly

Whole-genome "shotgun" sequencing starts by copying and fragmenting the DNA

("Shotgun" refers to the random fragmentation of the whole genome; like it was fired from a shotgun)

Input: GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
35bp

Copy    GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
by      GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
PCR:    GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
        GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

Fragment:   GGCGTCTA    TATCTCGG    CTCTAGGCCCTC    ATTTTTT
            GGC    GTCTATAT    CTCGGCTCTAGGCCCTCA    TTTTTT
            GGCGTC  TATATCT    CGGCTCTAGGCCCT       CATTTTTT
            GGCGTCTAT    ATCTCGGCTCTAG    GCCCTCA    TTTTTT

Courtesy of Ben Langmead. Used with permission.

# Assembly

Assume sequencing produces such a large # fragments that almost all genome positions are *covered* by many fragments...

...but we don't know what came from where

Reconstruct
this

CTAGGCCCTCAATTTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG

From these

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

Courtesy of Ben Langmead. Used with permission.

# Assembly

Overlaps between short reads help to put them together

<div align="center">

CTAGGCCCTCAATTTTT

CTCTAGGCCCTCAATTTTT

GGCTCTAGGCCCTCATTTTTT

CTCGGCTCTAGCCCCTCATTTT

TATCTCGACTCTAGGCCCTCA            177 nucleotides

TATCTCGACTCTAGGCC

TCTATATCTCGGCTCTAGG

GGCGTCTATATCTCG

GGCGTCGATATCT

GGCGTCTATATCT

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT     35 nucleotides

</div>

# Where is the Poisson?

- G - genome length (in bp)
- L - short read average length
- N – number of short read sequenced
- λ – sequencing coverage redundancy = LN/G
- x- number of short reads covering a given site on the genome

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

Poisson as a limit of Binomial: For a given site on the genome for each short read Prob(site covered): p=L/G is very small. Number of attempts (short reads): N is very large. Their product (sequencing redundancy): λ = NL/G is O(1).



Ewens, Grant, Chapter 5.1

# What fraction of the genome is missing?

# What fraction of genome is covered?

- Coverage: $\lambda = NL/G$,
  
  *X – random variable equal to the number of times a given site is covered by short reads.*
  
  *Poisson: $P(X=x) = \lambda^x \exp(-\lambda)/x!$*
  
  *$P(X=0) = \exp(-\lambda)$, $P(X>0) = 1 - \exp(-\lambda)$*

- *Total length covered: $G*[1 - \exp(-\lambda)]$*

| $\lambda$ | 2 | 4 | 6 | 8 | 10 | 12 |
|---|---|---|---|---|---|---|
| Mean proportion of genome covered | .864665 | .981684 | .997521 | .999665 | .999955 | .999994 |

Table 5.1. The mean proportion of the genome covered for different values of $\lambda$

# How long should be the length $L_{ov}$ of the overlap to connect two short reads into a contig?



If DNA was a random chain with $p_A = p_C = p_G = p_T = 1/4$

$L_{ov} \sim 16\text{-}20$ would be enough

$2 \cdot G \cdot 4^{-L_{ov}} = 2 \cdot 3\text{x}10^9 \cdot 4^{-16} = 1.4$

$2 \cdot 3\text{x}10^9 \cdot 4^{-20} = 0.0055 << 1$

# How many contigs?



$$\text{P(short read can be extended by another short read)} = \frac{L - L_o}{G} = p$$

$$\text{P(short read cannot be extended by any short reads)} = e^{-pN} \approx Ne^{-\lambda}$$

$$\text{number of contigs} = Ne^{-pN} \approx Ne^{-\lambda}$$

# How many contigs?

- A given short read is the
right end of a contig if and only if
no left ends of other short reads fall within it.
- The left end of another short read has the probability
$p=(L-1)/G$ to fall within a given read. There are
$N-1$ other reads. Hence the expected number of left
ends inside a given shot read is
$p \cdot (N-1)=(N-1) \cdot (L-1)/G \approx \lambda$
- If significant overlap required to merge two short reads
is $L_{ov}$, modified $\lambda$ is given by $(N-1) \cdot (L- L_{ov})/G$
- Probability that no left ends fall inside a short read is
$exp(- \lambda)$. Thus the Number of contigs is $N_{contigs}=Ne^{- \lambda}$:

| $\lambda$ | 0.5 | 0.75 | 1 | 1.5 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|
| Mean number of contigs | 60.7 | 70.8 | 73.6 | 66.9 | 54.1 | 29.9 | 14.7 | 6.7 | 3.0 | 1.3 |

Table 5.2. The mean number of contigs for different levels of coverage, with $G = 100,000$ and $L = 500$.

# Average length of a contig?

- Length of a genome covered:
  $G_{covered} = G \cdot P(X>0) = G \cdot (1 - exp(-\lambda))$

- *Number of contigs* $N_{contigs} = N \cdot e^{-\lambda}$

- Average length of a contig =

$<L> = \sum_i L_i / N_{contigs} = G_{covered} / N_{contigs} =$

$G \cdot (1 - exp(-\lambda)) / N \cdot e^{-\lambda} = L \cdot (1 - exp(-\lambda)) / \lambda \cdot e^{-\lambda}$

| $\lambda$ | 2 | 4 | 6 | 8 | 10 |
|---|---|---|---|---|---|
| Mean contig size | 1,600 | 6,700 | 33,500 | 186,000 | 1,100,000 |

Table 5.3. The mean contig size for different values of $a$ for the case $L = 500$.

# Matlab exercise: Poisson distribution

- Generate a sample of size 100,000 for Poisson-distributed random variable X with λ =2, 6, 20

- Plot the approximation to the Probability Mass Functions based on these samples. Combine them in the same figure.

- Calculate the mean and variance of this sample and compare it to theoretical calculations:
  E[X]= λ and V[X]=λ

# Matlab exercise: Poisson distribution

- **Stats=100000; lambda=2;**
- **r2=random('Poisson',lambda,Stats,1);**
- **mu_p=sum(r2)./Stats;**
- **disp(mu_p);**
- **var_p=sum((r2-mu_p).^2)./Stats;**
- **disp(var_p);**
- **std_p=sqrt(var_p)**
- **[a,b]=hist(r2, 0:max(r2));**
- **p_p=a./sum(a);**
- **figure; stem(b,p_p);**
- **figure; semilogy(b,p_p,'ko-');**

# Estimate

- Human genome is $3 \times 10^9$ bp long
- Chromosome 1 is about G=$0.25 \times 10^9$ bp
- Illumina generates short reads L=100 bp long
- What number of reads *N* are needed to completely assemble the 1st chromosome*?*
- The formula to use is: $1=N_{contigs}=Ne^{-\lambda}=Ne^{-NL/G}$
- Answer: N=$4.4 \times 10^7$ short (100bp) reads
  Test: 4.4e7*exp(-4.4e7*100/0.25e9)=0.9997
- What coverage redundancy $\lambda$ will it be?
  Answer: $\lambda =NL/G$=17.6 coverage redundancy

# How much would it cost to assemble human genome now?

- Human Genome Project: $2.7 billion in 1991 dollars.

- Now a de novo full assembly of the whole human genome would now cost $3 \times 10^9$ x 17.6 /$10^9$ x 10$/GBase =$ 530

- 2nd genome (and after) would be even cheaper as we would already have a reference genome to which we can map short reads. (Puzzle: picture on the box)

- But this is a naïve estimate. In reality, there are complications. See the next slides:

# What spoils these estimates?

```
>gi|224514922|ref|NT_024477.14| Homo sapiens chromosome 12 genomic
contig, GRCh37.p13 Primary Assembly (displaying 3' end)
CGGGAAATCAAAAGCCCCTCTGAATCCTGCGCACCGAGATTCTCCCCAGCCAAGGTGAGGCGGCAGCAGT
GGGAGATCCACACCGTAGCATTGGAACACAAATGCAGCATTACAAATGCAGACATGACACCGAAAATATA
ACACACCCCATTGCTCATGTAACAAGCACCTGTAATGCTAATGCACTGCCTCAAAACAAAATATTAATAT
AAGATCGGCAATCCGCACACTGCCGTGCAGTGCTAAGACAGCAATGAAAATAGTCAACATAATAACCCTA
ATAGTGTTAGGGTTAGGGTCAGGGTCCCGGTCCGGGTCGGGGTCCGGGTCCGGGGTCCGGGTCAGGGTGA
GGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGT
TAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
GTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTA
GGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGT
TAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
GTTAGGGTTAGGGTTAGGGTTAG
```

**FIGURE 8.11**   A BLASTN search of the human genome (all assemblies) database was performed at the NCBI website using TTAGGGTTAGGGTTAGGG as query (i.e., three TTAGGG repeats). There were matches to hundreds of genomic scaffolds. This figure shows an example (NT_024477.14) assigned to the telomere of chromosome 12q having many dozens of TTAGGG repeats. These occurred at the 3' end of the genomic contig sequence.

There were 100s of matches while one expects << 1 match:

$$2 \cdot 3\times10^9 \cdot 4^{-18} = 0.08 << 1$$

DNA repeats make assembly difficult

# Repeats are like sky puzzle pieces

# How many repeats are in eukaryotic genomes?



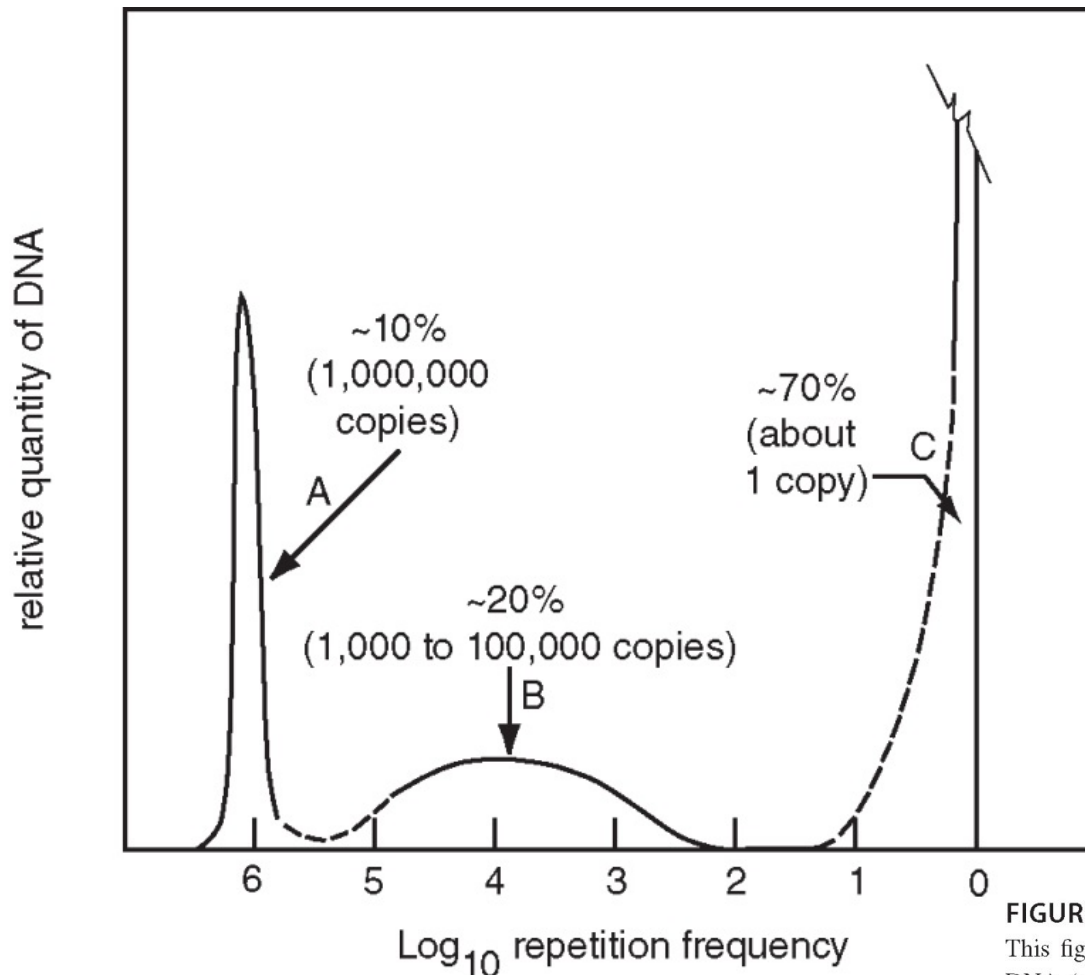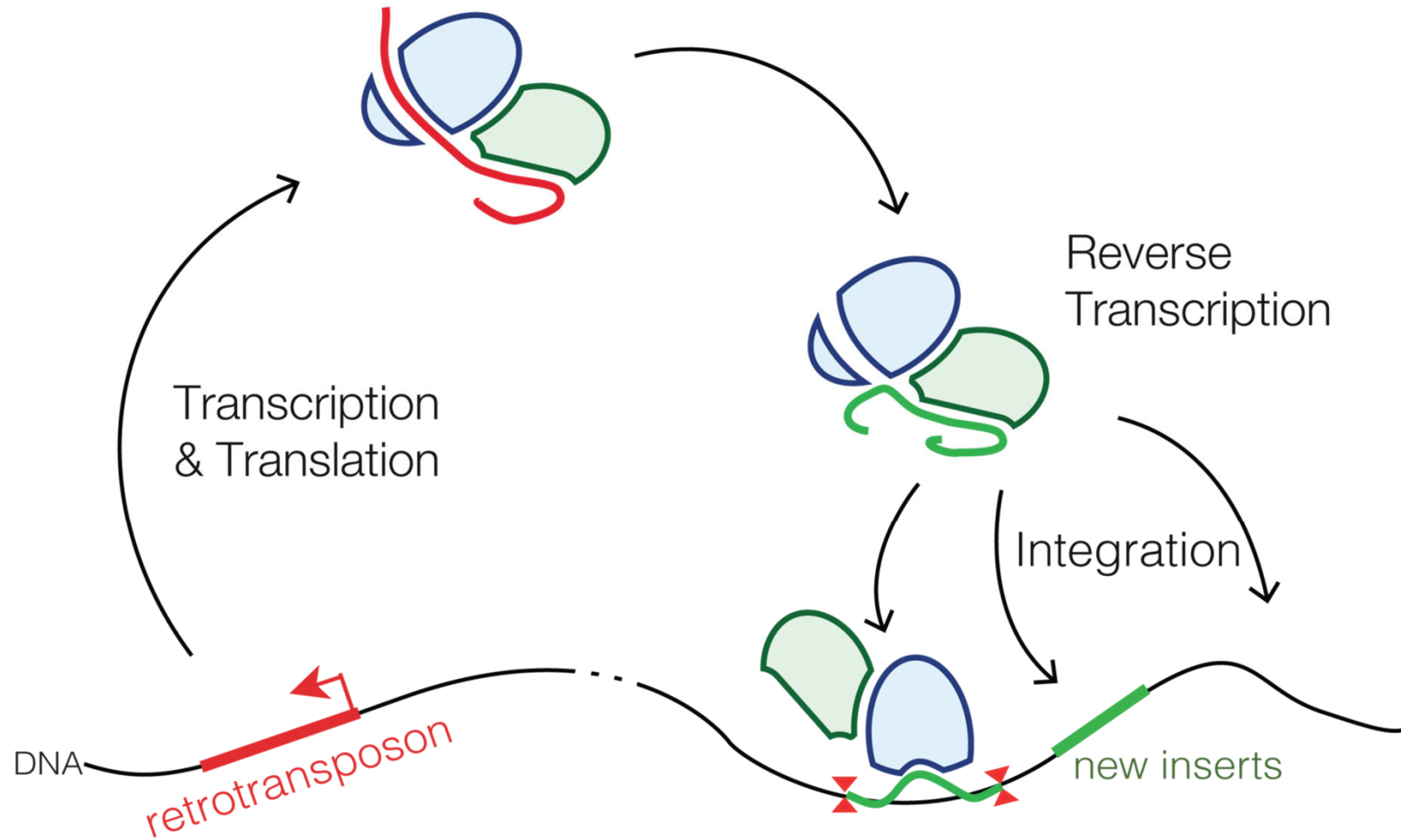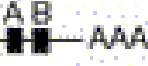Data for mouse genome obtained in 1961 (sic!) using DNA denaturation and renaturation curves

**FIGURE 8.6** The complexity of genomic DNA can be estimated by denaturing then renaturing DNA. This figure (redrawn from Britten and Kohne, 1968) depicts the relative quantity of mouse genomic DNA (y axis) versus the logarithm of the frequency with which the DNA is repeated. The data are derived from a $C_0 t_{1/2}$ curve, which describes the percent of genomic DNA that reassociates at particular times and DNA concentrations. A large $C_0 t_{1/2}$ value implies a slower reassociation reaction. Three classes are apparent. The fast component accounts for 10% of mouse genomic DNA (arrow A), and represents highly repetitive satellite DNA. An intermediate component accounts for about 20% of mouse genomic DNA and contains repeats having from 1000 to 100,000 copies. The slowly reassociating component, comprising 70% of the mouse genome, corresponds to unique, single-copy DNA. Britten and Kohne (1968) obtained similar profiles from other eukaryotes, although distinct differences were evident between species. Used with permission.

Formation of
Ribonucleoprotein complexes

Reverse
Transcription

Transcription
& Translation

Integration

DNA

retrotransposon

new inserts

# Almost all transposable elements in mammals fall into one of four classes

Classes of interspersed repeat in the human genome



| | | | Length | Copy number | Fraction of genome |
|---|---|---|---|---|---|
| LINEs | Autonomous | ORF1   ORF2 (pol)   AAA | 6–8 kb | 850,000 | 21% |
| SINEs | Non-autonomous | A B   AAA | 100–300 bp | 1,500,000 | 13% |
| Retrovirus-like elements | Autonomous | gag   pol   (env) | 6–11 kb | 450,000 | 8% |
| | Non-autonomous | (gag) | 1.5–3 kb | | |
| DNA transposon fossils | Autonomous | transposase | 2–3 kb | 300,000 | 3% |
| | Non-autonomous | | 80–3,000 bp | | |

Slide by Ross Hardison, Penn State U.