# Next Generation Sequencing (NGS): Reads Mapping

## 北京大学生物信息学中心 高歌

**Ge Gao, Ph.D.**

**Center for Bioinformatics, Peking University**

# Unit 2:
# NGS: Reads Mapping

## 北京大学生物信息学中心 高歌
## Ge Gao, Ph.D.
## Center for Bioinformatics, Peking University
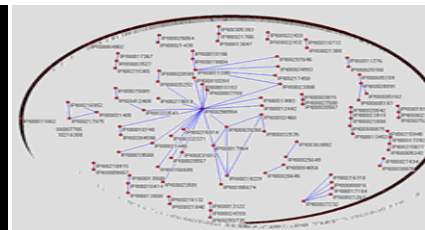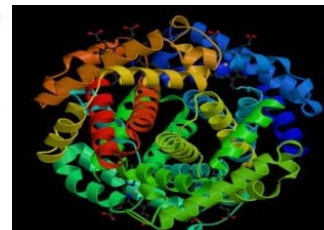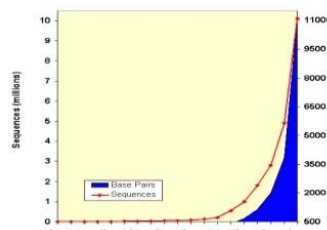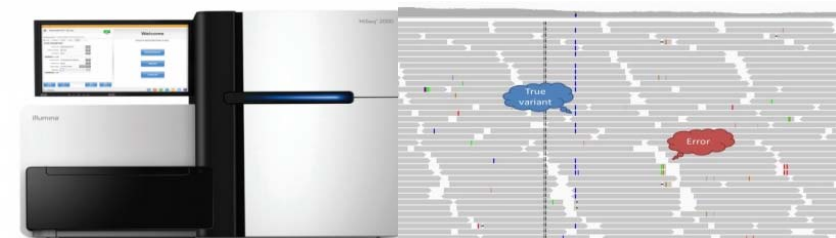
# Reads Mapping



Map-Making / Cartography:
Establish relationship between locations
(http://en.wikipedia.org/wiki/Cartography)

<u>Technological</u>: Reads is usual too short to be used/assembled *de novo*

<u>Scientific</u>: Taking full usage of existing annotation/knowledge

mapping

Reference Genome

Mapped alignment

Calling Genetic Variants

SNP

Measuring Abundance: RNA-Seq, ChIP-Seq, etc.

Copyright ©Peking University

# Mapping: Input Data

mapping

Reference Genome

- Reference Genome
  - Nucleotide
  - **Length**: Hundreds of Mb *per* chromosome
  - ~3 Gb in total (for human genome)

- Reads
  - Nucleotide, with **various qualities** (relatively **high error rate**: 1e-2 ~ 1e-5)
  - **Length**: 36~80 bp *per* read
  - Hundreds of Gbs *per* run

# "Embedded" Alignment



Genome

Read

One sequence is *"embedded"* in the other sequence (NGS Reads, PCR primer, *etc.*)

What we need here is actually a hybrid "global-local" alignment

- ✓ "Global" for short sequence (i.e. NGS Read)
- ✓ But "Local" for long sequence (i.e. Reference Genome)
- ✓ In particular, the surrounding "overhang" gaps should be not penalized.

| | δ | Gap open |
| | ε | Gap Extension |

| M | *Match* |
| X | *Insert* at sequence X (delete at sequence Y) |
| Y | *Insert* at sequence Y (delete at sequence X) |

| | M | X | Y |
|---|---|---|---|
| M | 1-2δ | δ | δ |
| X | 1-ε | ε | 0 |
| Y | 1-ε | 0 | 0 |

Genomic chromosome: m = hundreds of Mb

Sequencing Read: n = 36~80bp

Most of paths will just fail eventually!

In real world, the speed will be a BIG problem!

# BLAST Ideas: Seeding-and-extending

1. Find matches (seed) between the query and subject
2. Extend seed into High Scoring Segment Pairs (HSPs)
   - Run Smith-Waterman algorithm on the specified region only.
3. Assess the reliability of the alignment.

# Alphabetical Index

Copyright © Peking University

**Keys**

**Index**

Index Function

Data

Data Block 1

Data Block 2

Data Block *i*

Data Block *n*

Copyright © Peking University

# Hash #

Hash function maps (partial) data into (hashed) keys for following-up indexing

| Data | | Key | Hash function | (Hashed) keys |
|------|------|-----|---------------|---------------|

**HBS: A naive hash function**

Let's assume: A = 1, C = 2, G = 4, T = 8, then: $HBS(S) = \sum_i HBS(S_i)$, e.g:
$HBS(\text{AAAAA}) = 1 + 1 + 1 + 1 + 1 = 5$
$HBS(\text{GTACG}) = 4 + 8 + 1 + 2 + 4 = 19$
...

123456789012345678901234567890
TAACCTAACCCTAACCCAACCCTAACCC

Reference Genome

CCTAA

HBS

2+2+8+4+4
=20

**Index Table**

| ... | |
| --- | --- |
| 20 | |
| ... | |

**Address Table**
(CCTAA, 11)
...

# Pigeonhole principle (抽屉原理)

"In mathematics, the pigeonhole principle states that if $n$ items are put into $m$ pigeonholes with $n > m$, then at least one pigeonhole must contain more than one item."

After splitting the read into $n$ (non-overlapped) blocks, there will be at least $n$-$m$ perfectly-matched blocks (i.e. without any mismatch with in the block) by allowing up-to-$m$ mismatches.



One mismatch

Two mismatches

ELAND
MAQ
SOAP1
...

# Prefix Tree



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

# Suffix Tree



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

# Burrows–Wheeler transform (BWT)



Pos

X = googol$

(Li H, et. al, Bioinformatics, 2009)

i S(i)    B[i]

lo$oogg

(6,3,0,5,2,4,1)

BOWTIE
BWA
SOAP3
...

One of candidate sequence

Query sequence



(Source: Bedell *et al.* 2003)

Score

Length of extension

$X = S$

*Trim to max*

$$F\left(0,0\right) = 0$$

$$F\left(i,j\right) = \max \begin{cases} F\left(i-1,j-1\right) + s\left(x_i, y_j\right) \\ F\left(i-1,j\right) + d \\ F\left(i,j-1\right) + d \\ 0 \end{cases}$$

**Quality**: Given $p$ = the probability of a base calling is *wrong*, its Quality Score can be written as

$$Q = -10 * log_{10}(p)$$

| p | Q |
|---|---|
| 0.1 | 10 |
| 0.01 | 20 |
| 0.001 | 30 |
| 0.0001 | 40 |

```
0          10            20            30            40
|          |             |             |             |
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
|          |             |             |             |
0          10            20            30            40
```

**Mapping Quality**

Given reference sequence *z* (length *L*), a read sequence *x* (length *l*), *u* is the alignment position of *x* on *z*, the probability that *z* actually coming from the position *u* is ***p(z|x,u)***

(*Genome Res*. 2008 Nov;18(11):1851.)

$$p(z \mid x,u) = \prod_{mismatch} p(z_i) \qquad SQ(u) = \log(p(z \mid x,u)) = \sum_{mismatch} p(z_i) = \sum_{mismatch} Q(z_i)$$

```
Read:  ACGT  (Quality: 30 30 25 20)
Ref:   ACGTACGGA
       ACGT          0+ 0+ 0+ 0    SQ(0)
        ACGT        30+30+25+20    SQ(1)
         ACGT       30+30+25+20    SQ(2)
          ACGT      30+30+25+20    SQ(3)
           ACGT      0+ 0+ 0+20    SQ(4)
            ACGT    30+30+ 0+20    SQ(5)
```

# Mapping Quality

If we assume that a uniform NULL model, i.e. the read can randomly come from all possible positions with equal probability, then the error of mapping to a specified position $u$ could be written as

$$E(u) = \frac{SQ(u)}{\sum_i SQ(i)}$$

(*Genome Res*. 2008 Nov;18(11):1851.)

```
Read: ACGT (Quality: 30 30 25 20)
Ref:  ACGTACGGA        SQ(u)         E(u)
      ACGT        0+ 0+ 0+ 0        0/415
       ACGT      30+30+25+20      105/415
        ACGT     30+30+25+20      105/415
         ACGT    30+30+25+20      105/415
          ACGT    0+ 0+ 0+20       20/415
           ACGT  30+30+ 0+20       80/415
```

# Genetic Variants

- SNV: Single Nucleotide Variant
  - Substitution (SNP)
  - Indel: insertion/deletion

- Structural Variation (SV)
  - Large-scale insertion/deletion
  - Inversion
  - Translocation
  - Copy Number Variation (CNV)



Nature Reviews | Genetics

# SNP Calling is NOT Genotyping

- "SNP calling aims to determine in which positions there are polymorphisms or in which positions at least one of the bases differs from a reference sequence"

- "Genotype calling is the process of determining the genotype for each individual and is typically only done for positions in which a SNP or a 'variant' has already been called."

(Source: *Nature Reviews Genetics* 12, 443-451)

# Counting: an intuitive (and naïve) approach



- Counting  high-confident , non-reference allele (i.e. Quality >= 20)
    - Freq <20% or > 80%: homozygous genotype
    - Otherwise: heterozygous

- Works well for "deeply sequenced regions" (DSR), i.e. depth > 25x
    - But suffer from under-calling of heterozygous genotypes for low-coverage regions
    - And can't give an objective measurement for reliability

(Source: *Nature Reviews Genetics* 12, 443-451)

# A Simple Probabilistic Model for Genotyping

1. For a diploid genome, there will be at most two different alleles (A and a) observed at a given site:
   - 3 possible genotypes: <A,A>, <A,a>, <a,a>
   - <u>Number of A</u>: k; <u>Number of a</u>: n-k


2. Then, the probability for each genotypes is
   - P(D|<A,A>) = the probability that we have (n-k) sequencing errors at this site $\prod_{n-k} P(x_i)$
   - Similarly, we can see the P(D|<a,a>) = $\prod_k P(x_i)$
   - P(D|<A,a>) = 1 − (P(D|<A,A>) + P(D|<a,a>))


3. Bayes Formula can be further employed to calculate posterior probabilities, i.e. P(<A,A>|D), P(<a,a>|D), and P(<A,a>|D) if we can estimate the prior probabilities P(<A,A>), P(<a,a>) and P(<A,a>)

# Genome Analysis ToolKit (GATK)



Phase 1: nGS data processing — Typically by lane

Input — Raw reads → Mapping → Local realignment → Duplicate marking → Base quality recalibration → Output: Analysis-ready reads

Phase 2: variant discovery and genotyping — Typically multiple samples simultaneously but can be single sample alone

Input: Sample 1 reads · · · Sample N reads → SNPs, Indels, Structural variation (SV) → Raw variants

Phase 3: integrative analysis

Raw indels, Raw SNPs, Raw SVs → External data → Pedigrees, Known variation, Population structure, Known genotypes → Variant quality recalibration ↔ Genotype refinement → Analysis-ready variants

# 生物信息学：导论与方法
## Bioinformatics: Introduction and Methods



https://www.coursera.org/course/pkubioinfo