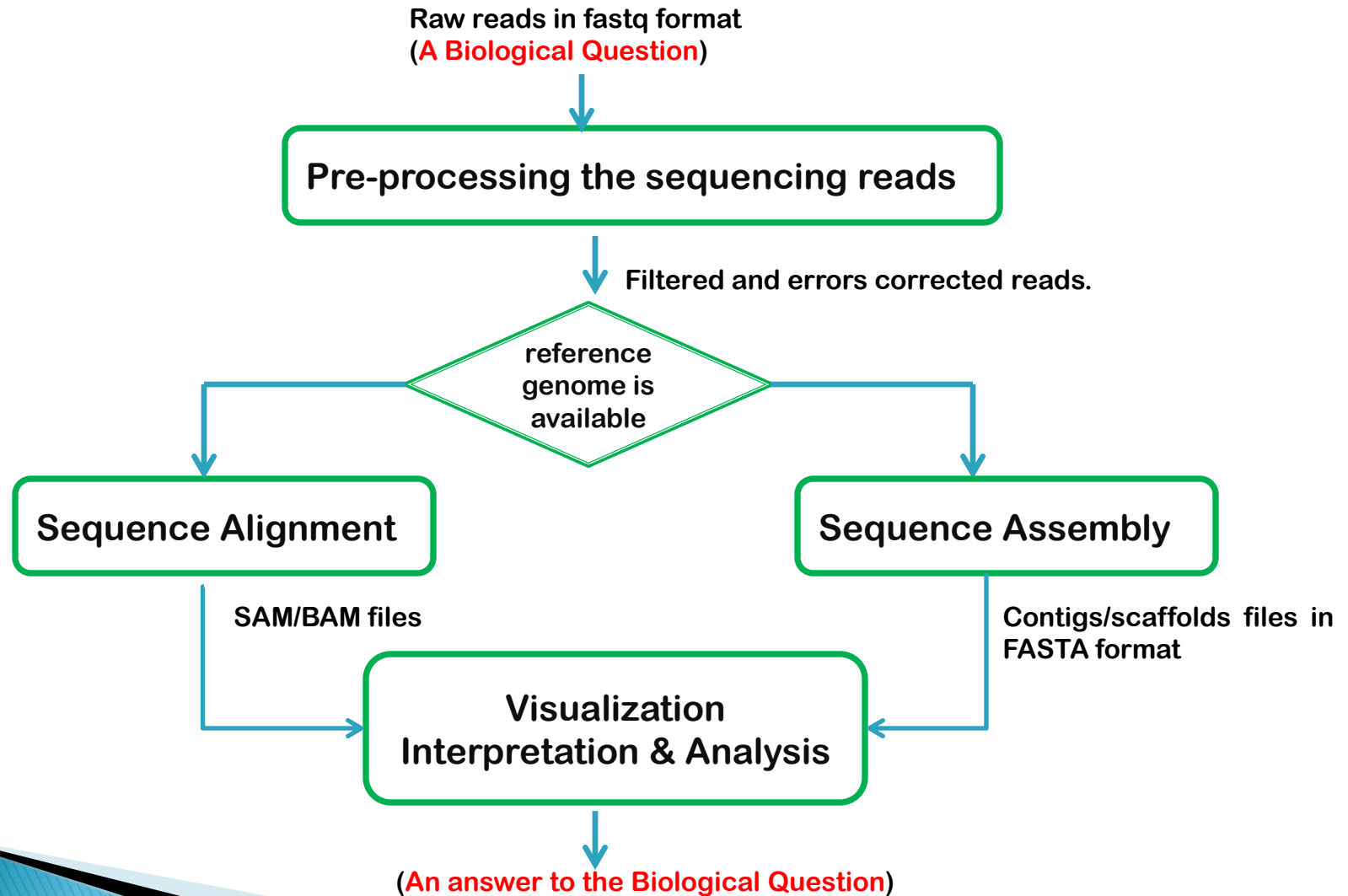# Informatics on High-throughput Sequencing Data

**(Summer Course 2020 )**
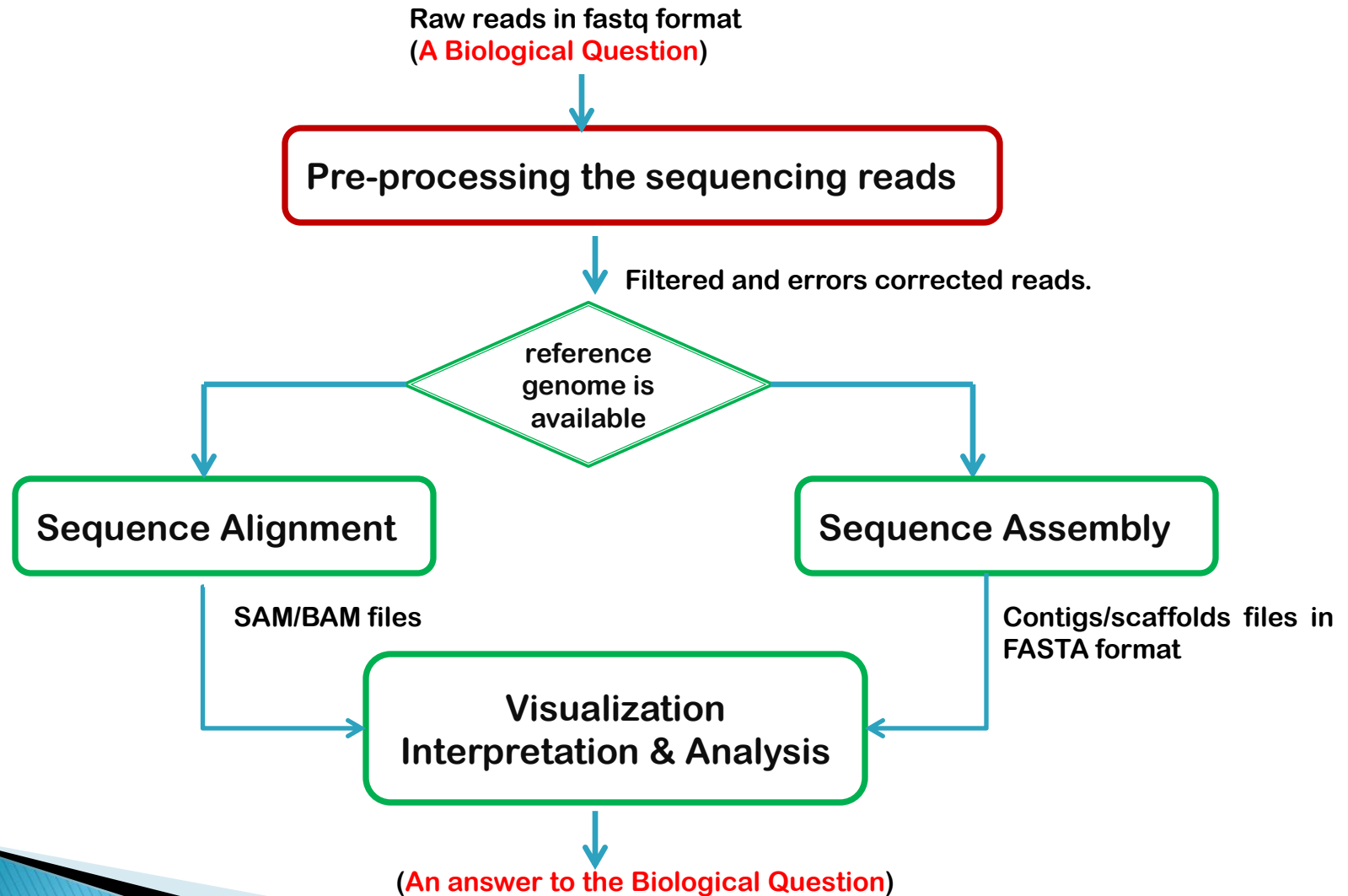
**Day 13**

# A typical Data Analysis pipeline
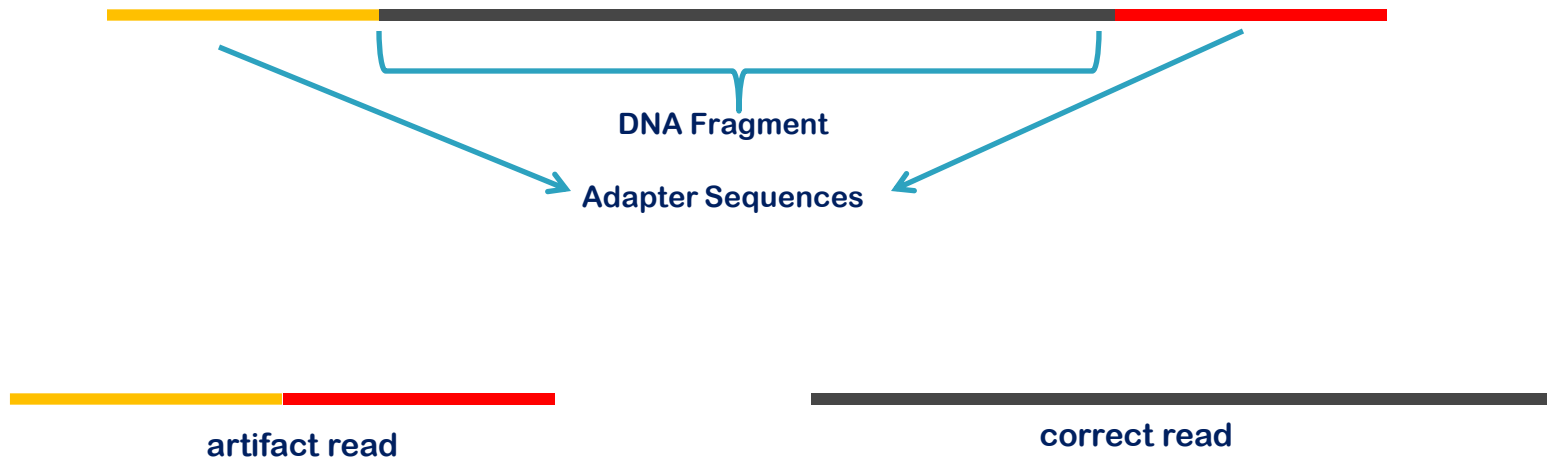
Raw reads in fastq format
(**A Biological Question**)

↓

**Pre-processing the sequencing reads**

↓ Filtered and errors corrected reads.

reference genome is available

**Sequence Alignment**

SAM/BAM files

**Sequence Assembly**

Contigs/scaffolds files in FASTA format

**Visualization Interpretation & Analysis**

↓

(**An answer to the Biological Question**)

# Pre-processing reads

➢ **Filter out garbage reads (Reads Trimming)**
  ✓ **Reads with low quality base calls.**
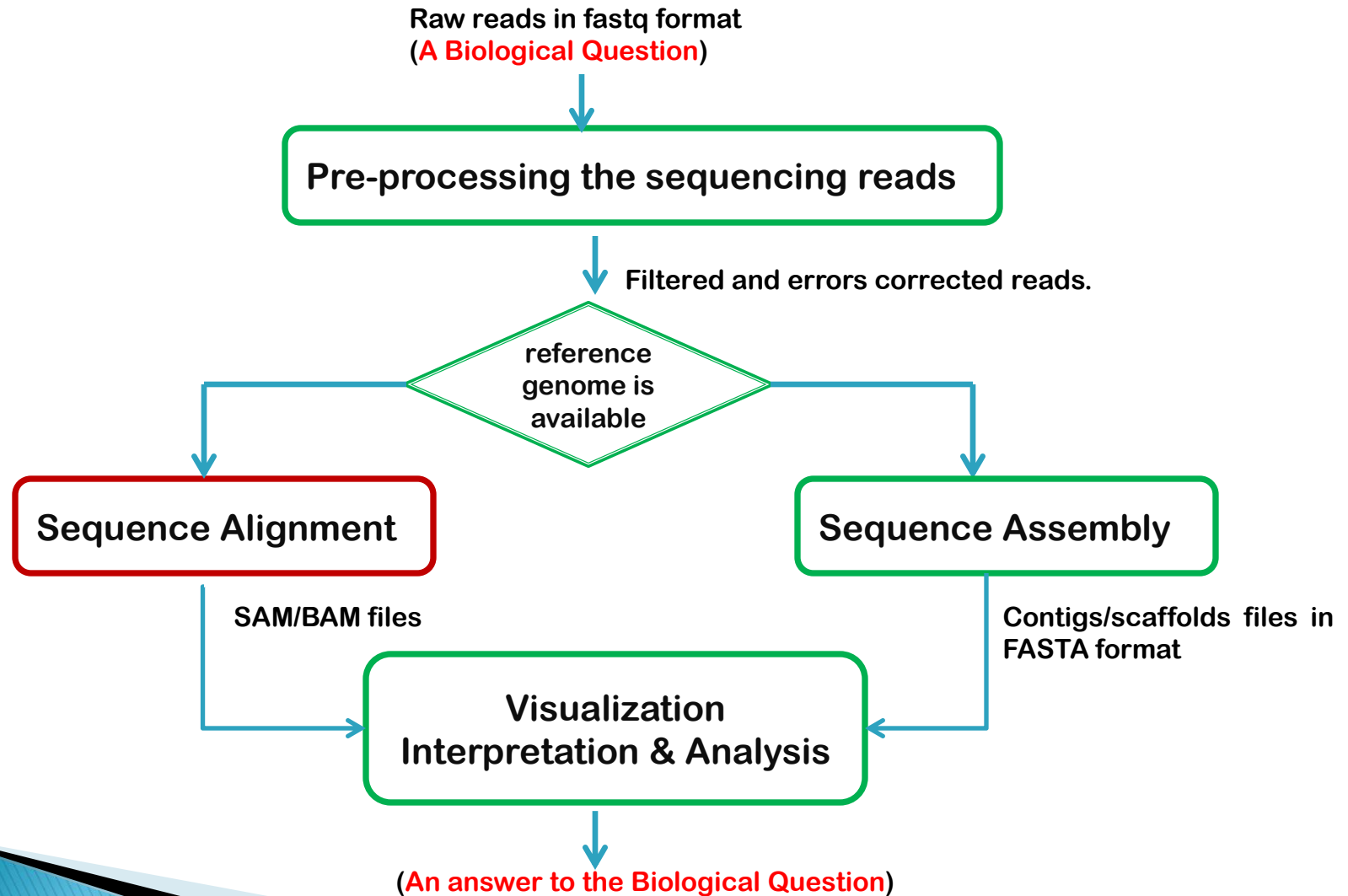  ✓ **Reads that are clearly artifacts with chemistry.**

**A typical read for DNA sequencing process**

**DNA Fragment**

**Adapter Sequences**

**artifact read**                    **correct read**

# Pre-processing reads

➢ **Filter out garbage reads (Reads Trimming)**
- ✓ **FASTX**
- ✓ **AfterQC**
- ✓ **Trimmomatic**

➢ **Errors detection and correction**
- ✓ **Quake**
- ✓ **Lighter**
- ✓ **Musket**

# A typical Data Analysis pipeline

Raw reads in fastq format
(**A Biological Question**)

**Pre-processing the sequencing reads**

Filtered and errors corrected reads.

reference genome is available

**Sequence Alignment**

**Sequence Assembly**

SAM/BAM files

Contigs/scaffolds files in FASTA format

**Visualization Interpretation & Analysis**

(**An answer to the Biological Question**)

# Sequence alignment

**Sequencing Reads**



**Reference Genome**

```
GATCACAGGTCTATCACCCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTTT
CGTCTGGGGGGTATGCACGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCCTATGTC
GCAGTATCTGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT
ACAGGCGAACATACTTACTAAAGTGTGTTAATTAATTAATGCTTGTAGGACATAATAATA
ACAATTGAATGTCTGCACAGCCACTTTCCACACAGACATCATAACAAAAAATTTCCACCA
AACCCCCCCTCCCCCGCTTCTGGCCACAGCACTTAAACACATCTCTGCCAAACCCCAAAA
ACAAAGAACCCTAACACCAGCCTAACCAGATTTCAAATTTTATCTTTTGGCGGTATGCAC
TTTTAACAGTCACCCCCCAACTAACACATTATTTTCCCCTCCCACTCCCATACTACTAAT
CTCATCAATACAACCCCCGCCCATCCTACCCAGCACACACACACCGCTGCTAACCCCATA
CCCCGAACCAACCAAACCCCAAAGACACCCCCCACAGTTTATGTAGCTTACCTCCTCAAA
GCAATACACTGACCCGCTCAAACTCCTGGATTTTGGATCCACCCAGCGCCTTGGCCTAAA
CTAGCCTTTCTATTAGCTCTTAGTAAGATTACACATGCAAGCATCCCCGTTCCAGTGAGT
TCACCCTCTAAATCACCACGATCAAAAGGAACAAGCATCAAGCACGCAGCAATGCAGCTC
AAAACGCTTAGCCTAGCCACACCCCCACGGGAAACAGCAGTGATTAACCTTTAGCAATAA
ACGAAAGTTTAACTAAGCTATACTAACCCCAGGGTTGGTCAATTTCGTGCCAGCCACCGC
GGTCACACGATTAACCCAAGTCAATAGAAGCCGGCGTAAAGAGTGTTTTAGATCACCCCC
TCCCCAATAAAGCTAAAACTCACCTGAGTTGTAAAAAACTCCAGTTGACACAAAATAGAC
TACGAAAGTGGCTTTAACATATCTGAACACACAATAGCTAAGACCCAAACTGGGATTAGA
TACCCCACTATGCTTAGCCCTAAACCTCAACAGTTAAATCAACAAAACTGCTCGCCAGAA
CACTACGAGCCACAGCTTAAAACTCAAAGGACCTGGCGGTGCTTCATATCCCTCTAGAGG
AGCCTGTTCTGTAATCGATAAACCCCGATCAACCTCACCACCTCTTGCTCAGCCTATATA
```

# Sequence alignment software

| Aligner | Approach | Applications | Availability |
|---------|----------|--------------|--------------|
| BWA-mem | Burrows-Wheeler | DNA, SE, PE, SV | open-source |
| Bowtie2 | Burrows-Wheeler | DNA, SE, PE, SV | open-source |
| Novoalign | hash-based | DNA, SE, PE | free for academic use |
| TopHat | Burrows-Wheeler | RNA-seq | open-source |
| STAR | hash-based (reads) | RNA-seq | open-source |
| GSNAP | hash-based (reads) | RNA-seq | open-source |

# Genome Indexing

## INDEX

# Genome Indexing

Step1: hash/index the genome

CATGGTCATTGGTTCC

| Kmer/Hash | Genome Positions |
|-----------|------------------|
| CAT | 1,7 |
| ATG | 2 |
| TGG | 3,10 |
| GGT | 4,11 |
| GTC | 5 |
| TCA | 6 |
| ATT | 8 |
| TTG | 9 |
| GTT | 12 |
| TTC | 13 |
| TCC | 14 |

Slides curiosity from Aaron Quinlan: https://github.com/quinlan-lab/applied-computational-genomics

# Genome Indexing

Step2: use the index to find reads locations

Toy genome    **CATGGTCATTGGTTCC**

| Kmer/Hash | Genome Positions |
|-----------|------------------|
| CAT | 1,7 |
| ATG | 2 |
| TGG | 3,10 |
| GGT | 4,11 |
| GTC | 5 |
| TCA | 6 |
| ATT | 8 |
| TTG | 9 |
| GTT | 12 |
| TTC | 13 |
| TCC | 14 |

Read    **TGGTCA**

# BWA-MEM



Slides curiosity from Aaron Quinlan: https://github.com/quinlan-lab/applied-computational-genomics

# BWA-MEM workflow

This takes a long time, but you do it *once*

Create BWT of reference genome.

`$ bwa index grch38.fa`

Output is in SAM format. Use multiple threads if you have a computer with multiple CPUs.

Align paired-end FASTQ to BWT index.

`$ bwa mem -t 16 grch38.fa 1.fq 2.fq > sample.sam`

# BWA-MEM

# BWA-MEM

# BWA-MEM

- bzip2 –d bwa-0.7.17.tar.bz2

- tar xvf bwa-0.7.17.tar

- make

- ./bwa index  wu_0.v7.fas
  - wu_0.v7.fas.amb
  - wu_0.v7.fas.ann
  - wu_0.v7.fas.bwt
  - wu_0.v7.fas.pac
  - wu_0.v7.fas.sa

    .amb is text file, to record appearance of N (or other non-ATGC) in the ref fasta.
    .ann is text file, to record ref sequences, name, length, etc.
    .bwt is binary, the Burrows-Wheeler transformed sequence.
    .pac is binary, packaged sequence (four base pairs encode one byte).
    .sa is binary, suffix array index.

- ./bwa mem –t 16 wu_0.v7.fas wu_0_A_wgs.fastq > results.sam

# BWA-MEM

```
@SQ      SN:Chr1 LN:29923332
@SQ      SN:Chr2 LN:19386101
@SQ      SN:Chr3 LN:23042017
@SQ      SN:Chr4 LN:18307997
@SQ      SN:Chr5 LN:26567293
@SQ      SN:chloroplast  LN:154478
@SQ      SN:mitochondria LN:366924
@PG      ID:bwa  PN:bwa  VN:0.7.17-r1198-dirty   CL:./bwa mem -t 16 wu_0.v7.fas wu_0_A_wgs.fastq
```

# SAM

| Field | Meaning |
|---|---|
| GAII05_0002:1:2:12086:1654 | Read ID |
| 16 | Flag |
| Chr2 | Chr |
| 1694072 | start |
| 0 | MAPQ |
| 51M | CIGAR |
| * | Mate Chr |
| 0 | Mate start |
| 0 | Mate dis |
| CCTTGTAAAATCATTATTAATGTTTTTAAACCCCTTTTAAAAATCCTTGTA | read |
| CCCCCCCCCCCCCCCCCCCBBCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC | qual |
| NM:i:1<br>MD:Z:20C30<br>AS:i:46<br>XS:i:46 | Tag-Type-Value |

# Thanks!
## // | ?