



# ChIP-seq introduction

Huitian Diao (Yolanda)

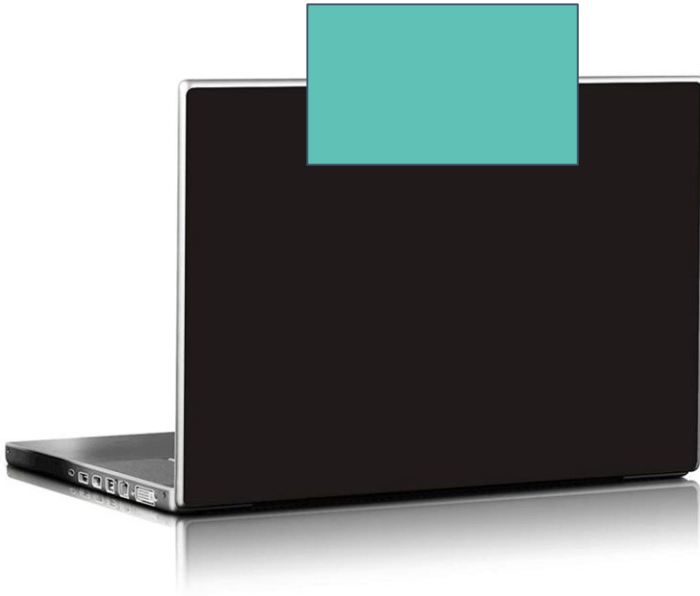
References:

ENCODE ChIP-seq pipeline: <https://www.encodeproject.org/pipelines/ENCPL138KID/>

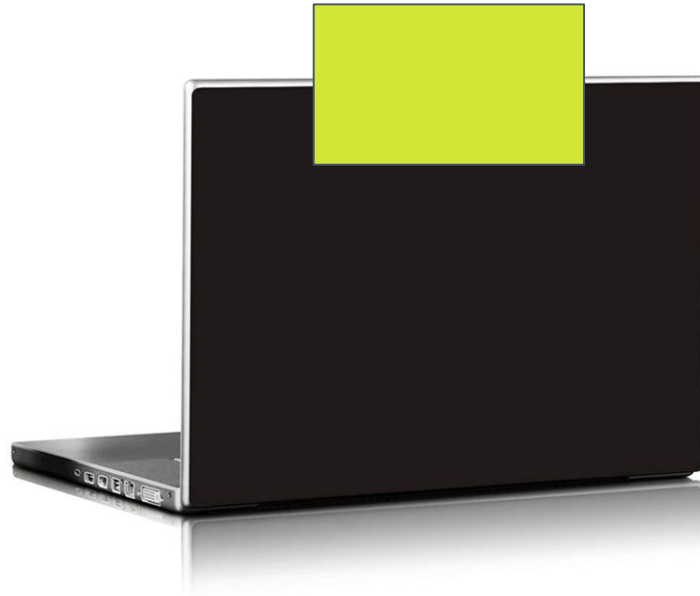


# Using sticky notes for feedback

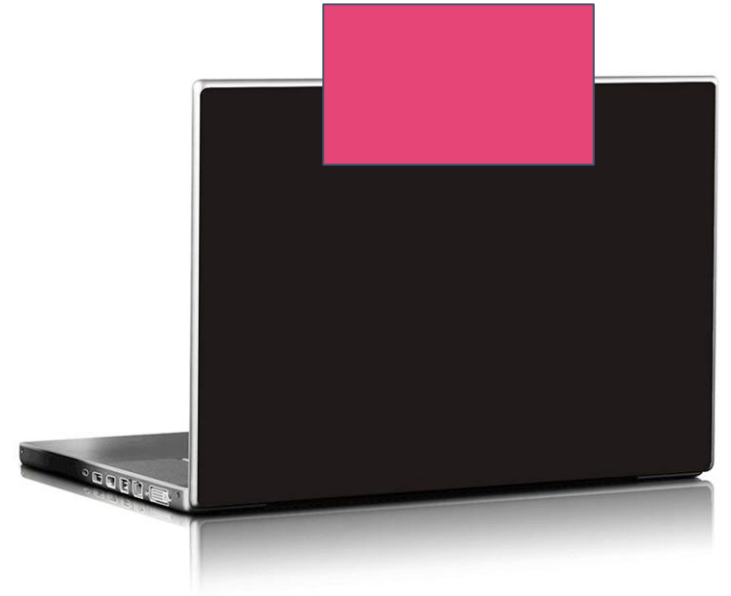
---



“I’ve got a good handle  
on things...”



“I think I understand  
but I’m still working  
through things...”

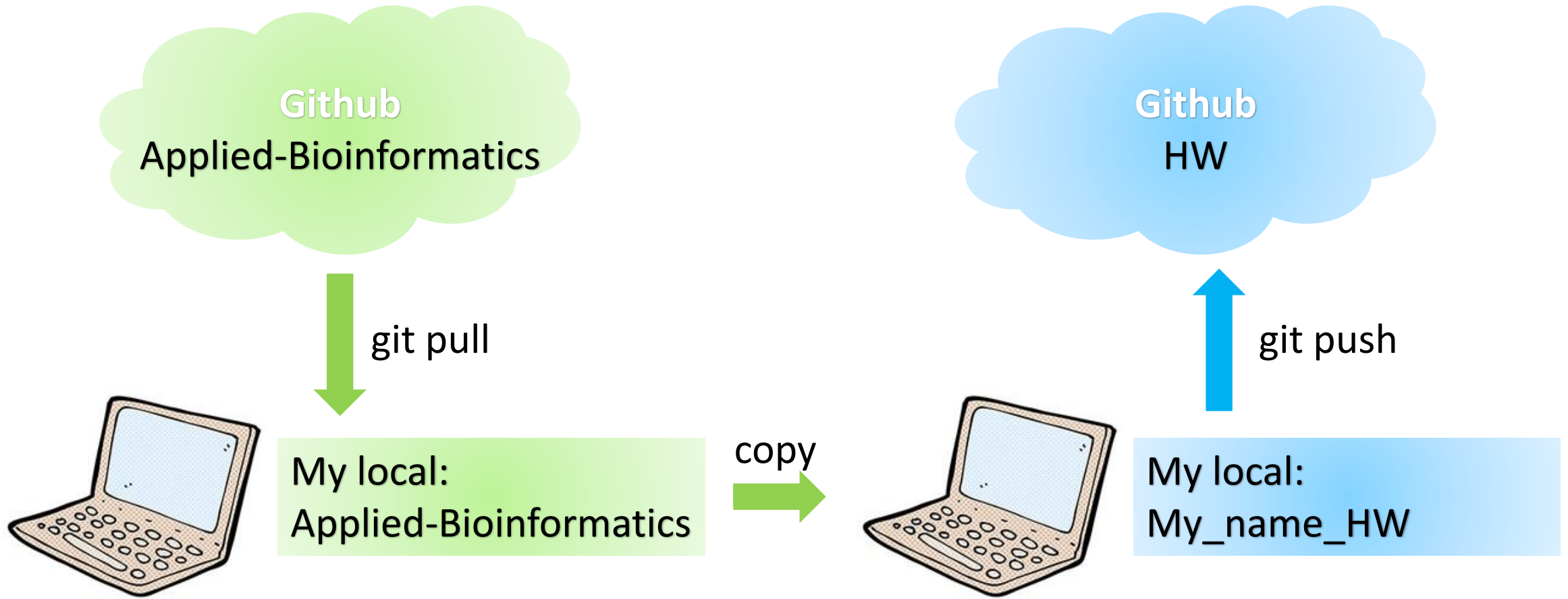


“I’m not understanding  
and I’m a little lost right  
now...”



# Get course material

---





# Q&A.1

- Why when using some packages we need to include the full path but for others we don't have to?

## Step 1. Where are the executables?

```
%%bash
```

```
which grep
```

```
which fastqc
```

```
/usr/bin/grep
```

```
/Users/yolandatiao/anaconda/bin/fastqc
```

## Step 1. How the computer find the executables?

```
%%bash
```

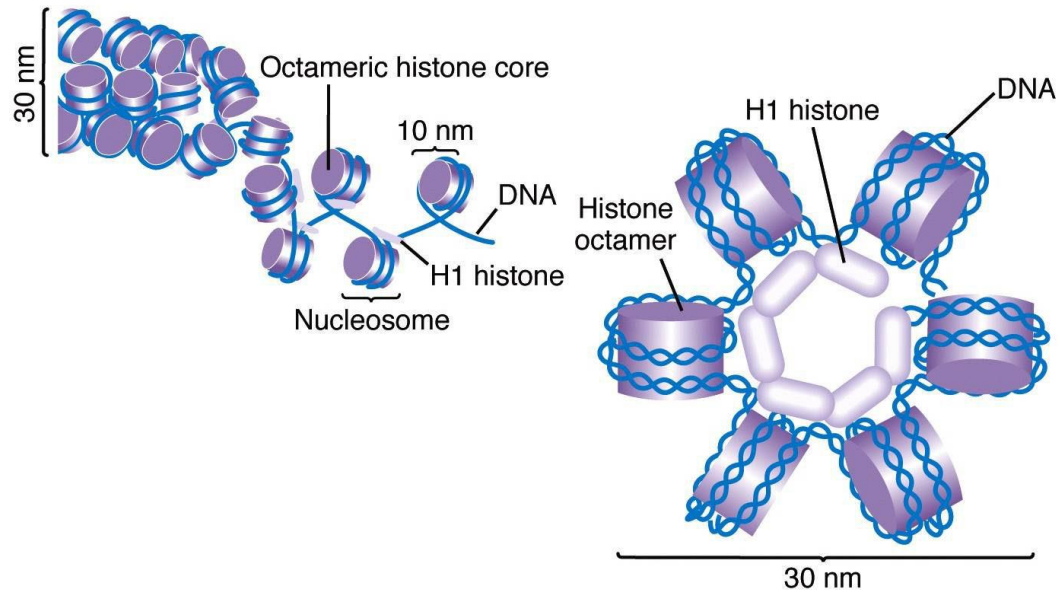
```
printenv | grep "PATH"
```

```
PATH=/Users/yolandatiao/anaconda/envs/HDpython3/bin:/Users/yolandatiao/anaconda/envs/HDpython3/bin:/Users/yolandatiao/Documents/0_Bioinformatics2017/201804_Cousera_Unix/Code/newCmd:/Users/yolandatiao/Documents/0_Bioinformatics2017/201804_Cousera_Unix/Code/Commands:/Users/yolandatiao/anaconda/bin:/Users/yolandatiao/anaconda/bin:/Users/yolandatiao/anaconda2/bin:/usr/local/bin:/usr/bin:/usr/sbin:/sbin:/Users/yolandatiao/.local/bin:/Library/Frameworks/Mono.framework/Versions/Current/Commands
```

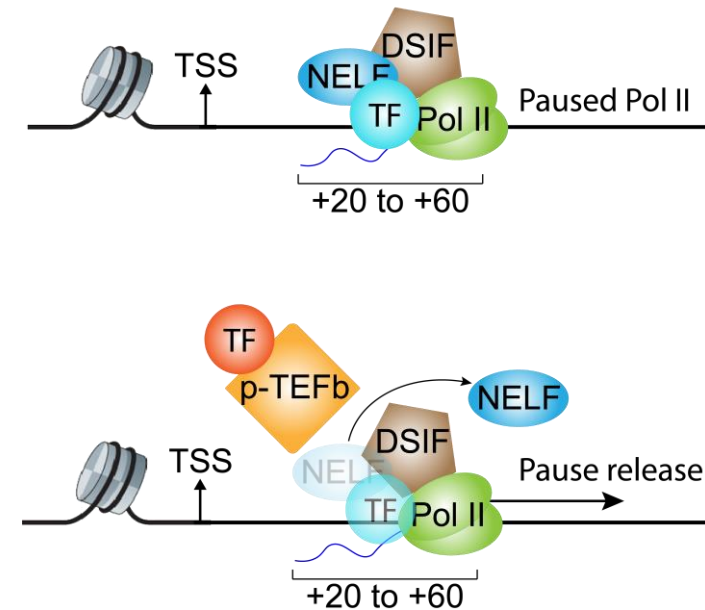


# Why ChIP-seq?

## DNA in the nucleus is highly condensed



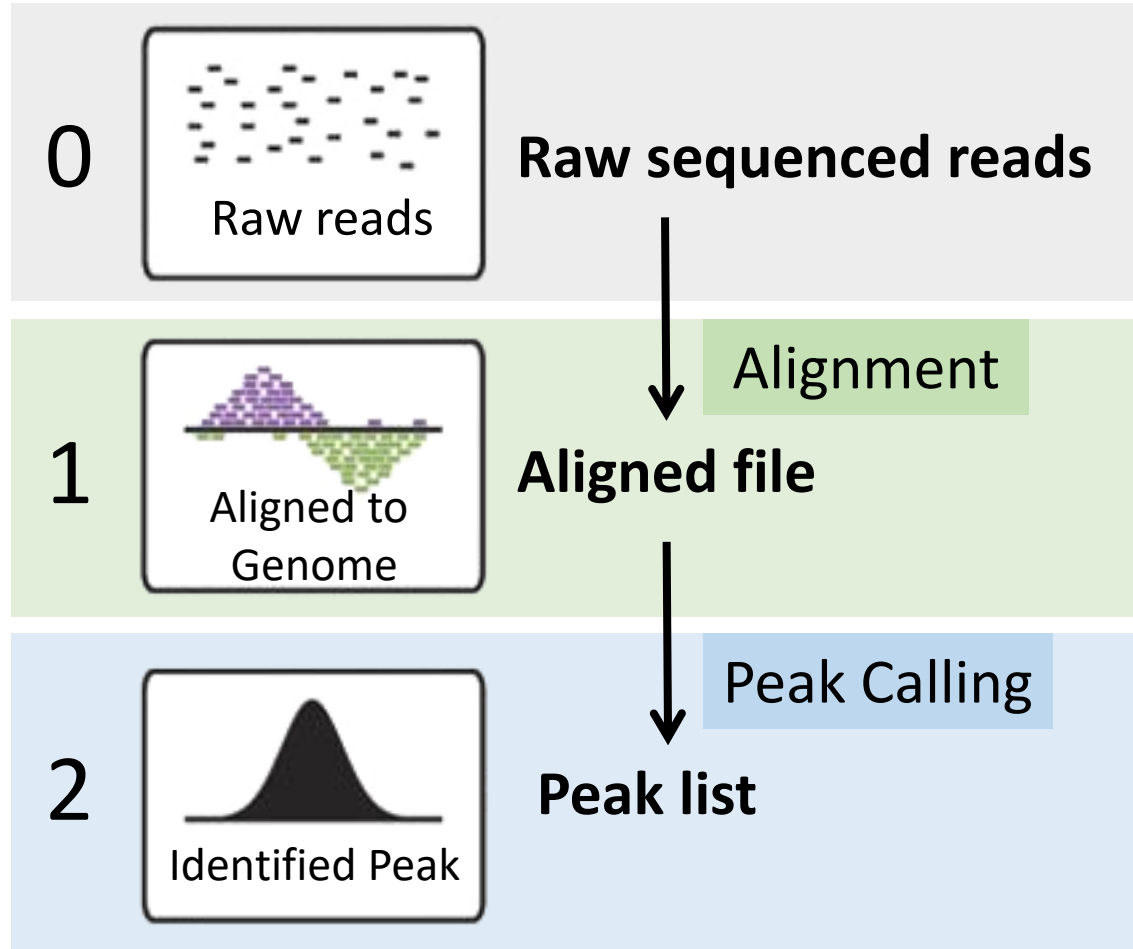
## Protein-DNA interaction is crucial For transcription regulation



**ChIP-seq:** identify DNA binding sites for proteins



# ChIP-seq analysis steps



0.0 | Acquire raw reads | fastq-dump



# Practice 4.1 (Dump-fastqc)

## Tips:

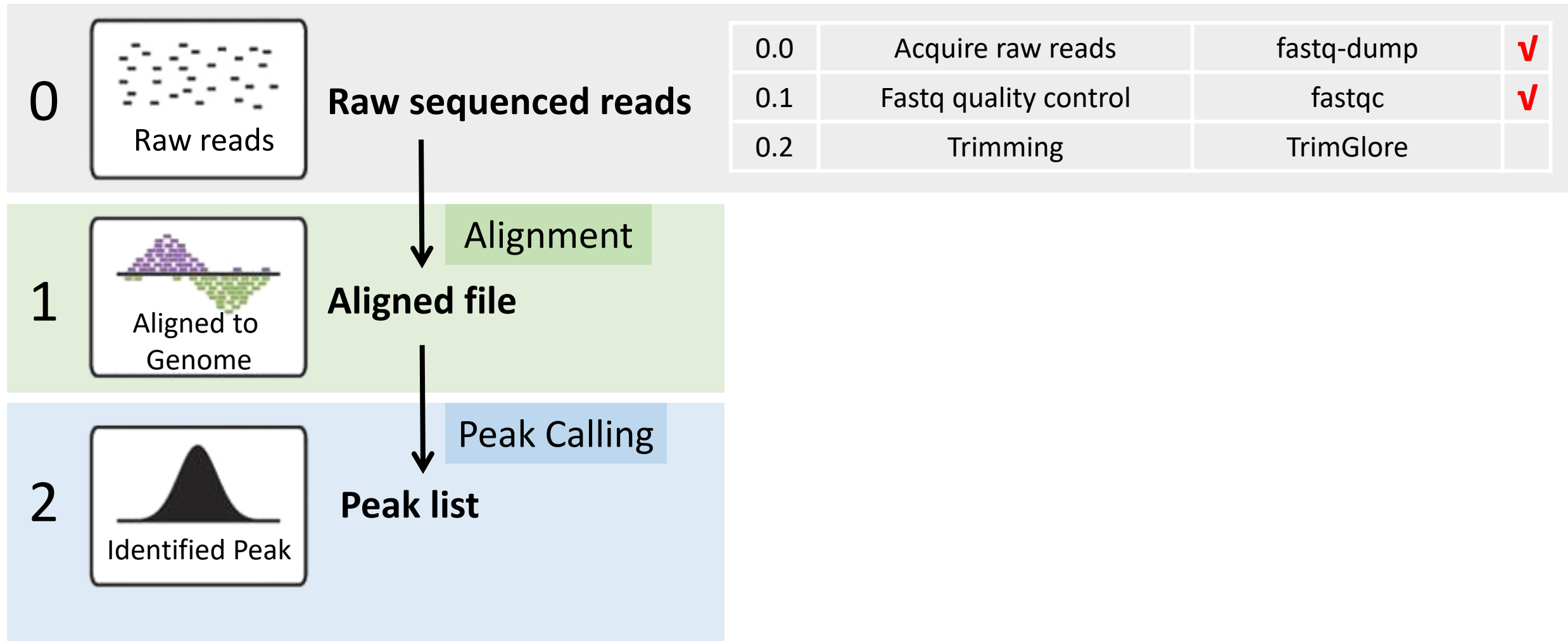
- SRA run selector:  
<https://www.ncbi.nlm.nih.gov/Traces/study/?go=home>
- fastq-dump
- fastqc

## Objectives:

1. **Find Run numbers for ChIP-seq data from this paper:**  
[B. H., Immunity, 2016](#)
  - 1.1 Search for GEO accession number from paper
  - 1.2 Find SubSeries for ChIP-seq data
  - 1.3 Search for ChIP-seq data accession number in SRA Run Selector
  - 1.4 Save **RunInfo Table**
2. **For the first file in the RunInfoTable:**  
Save the first 50000 spots into a file named **test\_50000.fastq**
3. **Run fastqc for test\_50000.fastq**  
Check the quality of the sequence



# ChIP-seq analysis steps







# Contents

---

## **1. Intro to ChIP-seq**

- 1.1 mechanism of ChIP-seq
- 1.2 ChIP-seq analysis intro
- 1.3 Fastqc

## **2. Alignment and Quality control**

- 2.1 Trim\_galore
- 2.2 Alignment and filter

## **3. Peak calling**

- 3.1 Intro to peak calling algorithm
- 3.2 MACS2 peak calling
- 3.3 ChIP-QC

## **4. Data visualization**

- 4.1 Data preparation for visualization
- 4.2 UCSC genome browser

## **5. Peak annotation and pathway analysis**

- 5.1 Differential analysis with DiffBind
- 5.2 Peak annotation with ChIPseeker

## **6. Downstream analysis**

- 6.3 Pathway analysis with ClusterProfiler



## Practice 4.2 (Dump-multiqc)

---

### Tips:

#### Q2.2: Self-help

**Google:** generate array of number in shell

**Google:** shell seq output scientific notation

**Google:** shell convert stdout to array

#### Q3: Self-help

**Google:** multiqc documentation

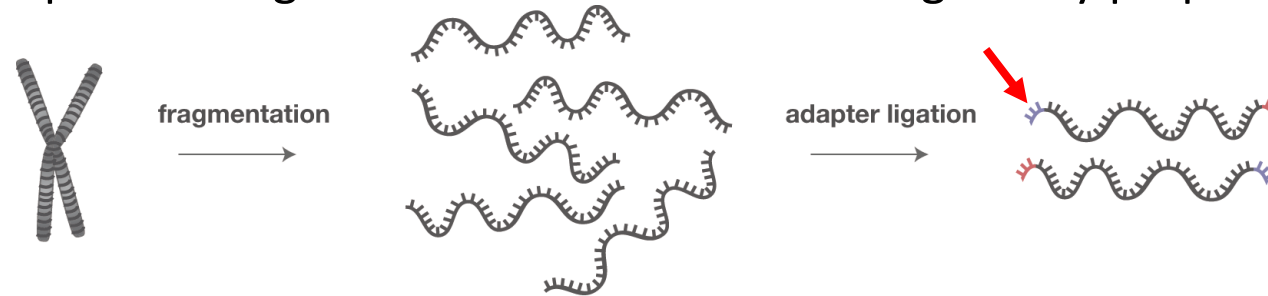
### Questions:

1. **Find Run numbers for ChIP-seq data from this paper:**  
[B. H., Immunity, 2016](#)
  - 1.1 Search for GEO accession number from paper
  - 1.2 Find SubSeries for ChIP-seq data
  - 1.3 Search for ChIP-seq data accession number in SRA Run Selector
  - 1.4 Save **RunInfo Table**
2. **Build a folder named Sample\_fastq**
  - 2.1 Redirect to the new folder
  - 2.2 Save the first 500 spots of each file into **SRRXXXXX.500.fastq** with a **for loop**
  - 2.3 Run **fastqc** on all the files
3. **Run multiqc on the fastqc outputs**  
Read fastqc and multiqc outputs. What did you find?

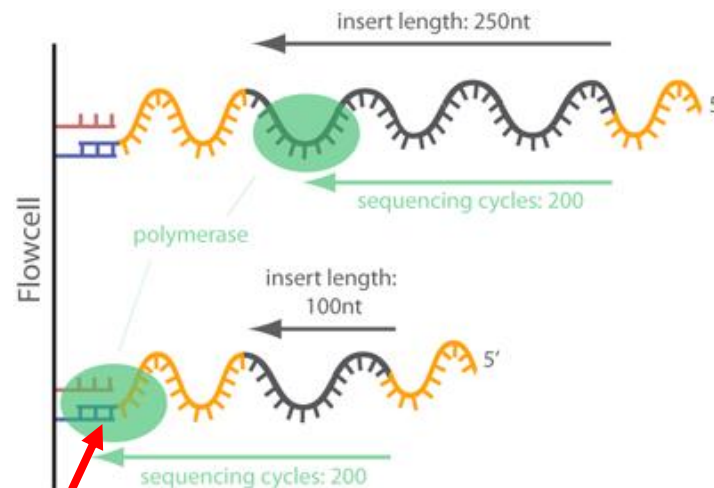


# Why trim Adapter -- Adapter contamination leads to lower alignment rate

Adapters are ligated to DNA molecules during library preparation

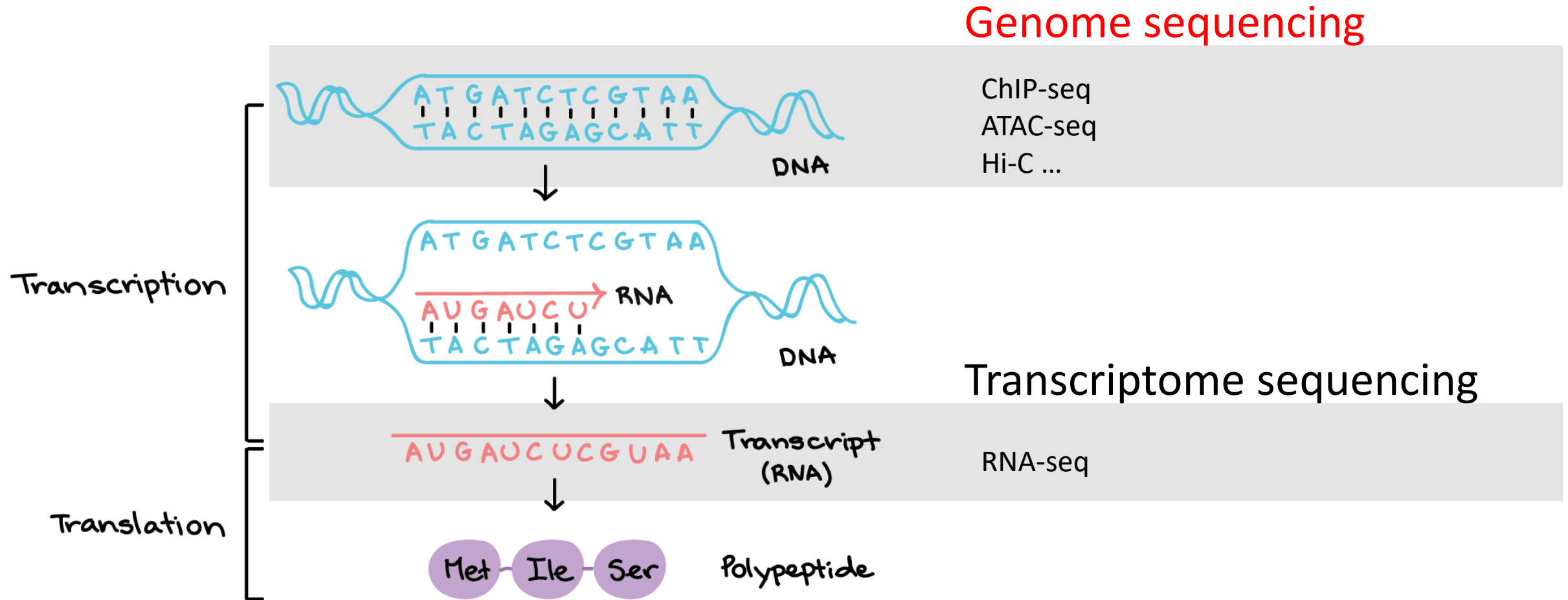


In Illumina sequencing, 3' end adapter would be sequenced if DNA insert is shorter than number of sequencing cycles





# Alignment: ChIP-seq v.s. RNA-seq





# File formats – fastq, sam

## Fastq File

Sequence identifier

Sequence

Quality score identifier (+)

Quality score

```
@SIM:1:FCX:1:15:6329:1045 1:N:0:2
TCGCACTCAACGCCCTGCATATGACAAGACAGAATC
+
<>;##=><9=AAAAAAAAAAAA9#:<#<;<<<????#=#
```

## Sam File

Header

```
@HD VN:1.3 SO:coordinate
@SQ SN:contigA LN:443
@SQ SN:contigB LN:1493
@SQ SN:contigC LN:328
```

Alignment info



```
readID43GYAX15:7:1:1202:19894/1 256 contig43 613960 1 65M * 0 0
CCAGCGCGAACGAAATCCGCATGCGTCTGGTCGTTGCACGGAACGGCGGGCGGTGTGATGCAC
GGC EDDEEDEE=EE?DE??DDDBADEBEFFFD BEFFEBBCBC=?BEEEE@=:?:?7?:8-
6?7?@??# AS:i:0 XS:i:0 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:i:0 MD:Z:65 YT:Z:UU
```

Name

Sam Flag

Ref  
Name

Position

MAPQ  
quality

CIGAR  
string

Mate  
name

Mate Pos

Temp  
length

Read  
Sequence

Read  
Quality

Add. Info

\* Columns separated by Tab (/t)

All you want to know about file formats is here: <https://genome.ucsc.edu/FAQ/FAQformat.html>