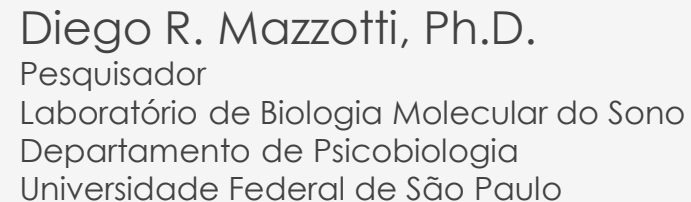


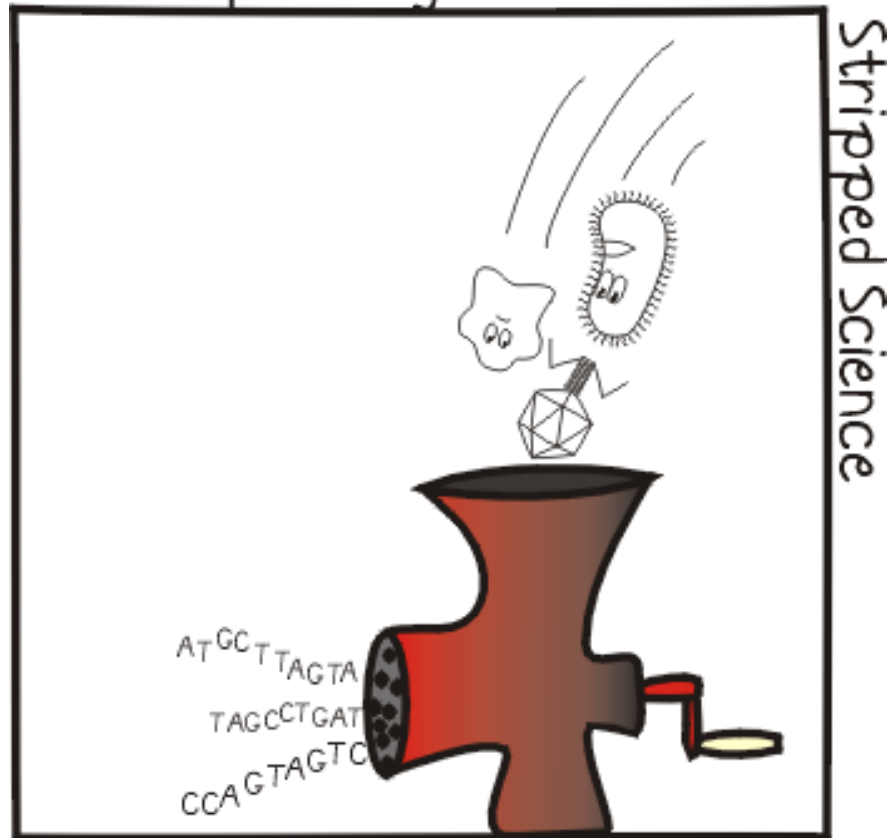
# Detecção de variantes – conceitos chave



E-mail: mazzottidr@gmail.com

# O sequenciamento de nova geração

Mass sequencing

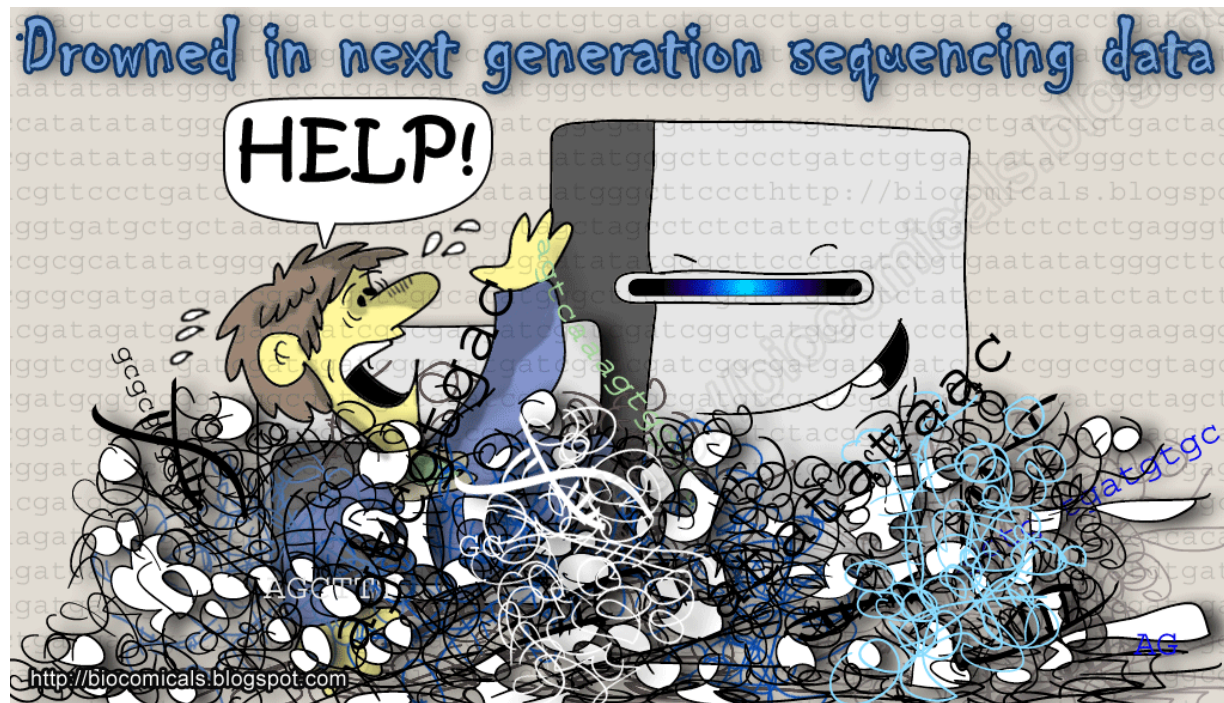


by Viktor S. Poór

Stripped Science

# Análise dos dados de sequenciamento de nova geração

- Independente da metodologia e aplicação:
  - Geração de milhões a bilhões de sequências curtas (50 – 250pb)



# Análise dos dados de sequenciamento de nova geração

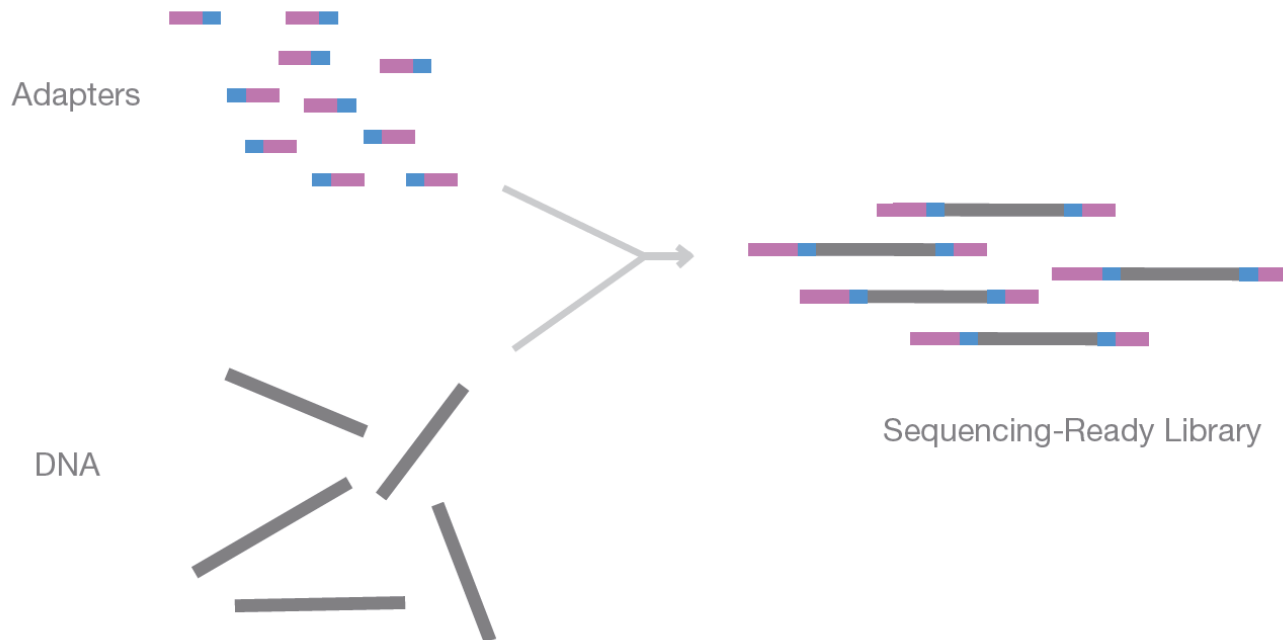
- Milhões de sequências curtas
  - Objetivos:
    - "Juntar" sequências para montar o genoma
    - Mapear as sequências em regiões conhecidas do genoma
    - Identificar o que é diferente em relação a um genoma referência (*Genome/ Exome/ Targeted Sequencing*)
    - Caracterizar níveis de expressão e identificar variantes (*Transcriptome / RNA-sequencing*)

# Análise dos dados de sequenciamento de nova geração

- Alguns conceitos chave:
  - *Read*
  - *Single-end versus Paired-End reads*
  - Alinhamento
  - Cobertura (*Coverage / Depth of Coverage*)
  - *Variant Calling*
  - Anotação

# Análise dos dados de sequenciamento de nova geração

- *Read*
  - Unidade resultante do sequenciamento de um fragmento de DNA



# Análise dos dados de sequenciamento de nova geração

- ◉ *Read*
  - ◉ O NGS é paralelo, portanto milhões de *reads* são gerados por corrida
  - ◉ *Reads* são organizados em um arquivo chamado **.fastq**
  - ◉ *Reads* nesse arquivo estão dispostos **aleatoriamente**

# Análise dos dados de sequenciamento de nova geração

## Formato **.fastq**

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*(((((***+))%%%++)(%%%) .1***-+*''))**55CCF>>>>>CCCCCCC65
```

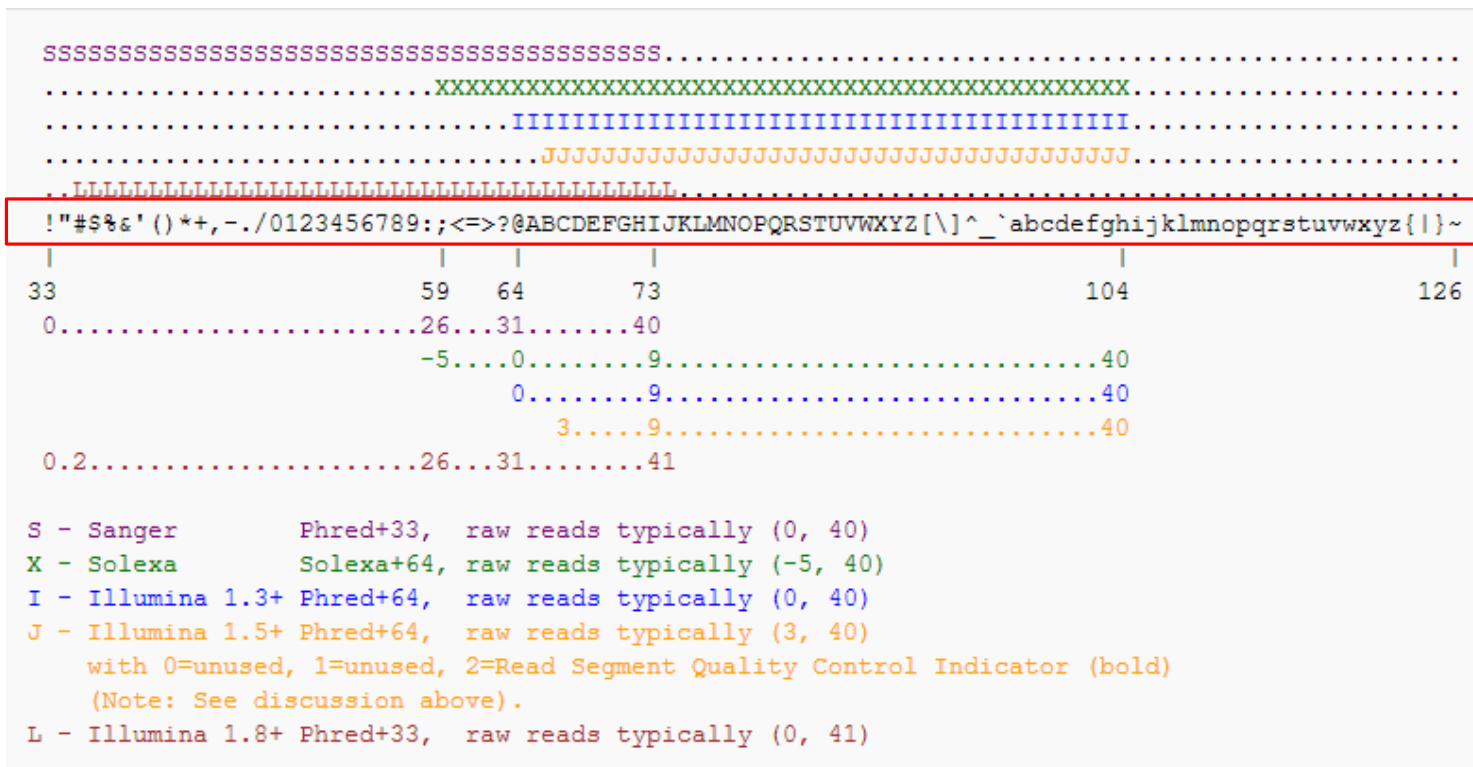
4 linhas para cada:

- 1 – Identificação do *read*
- 2 – Sequência de nucleotídeos
- 3 – Outra identificação do *read* (pouco utilizada)
- 4 – Escores de qualidade por símbolos



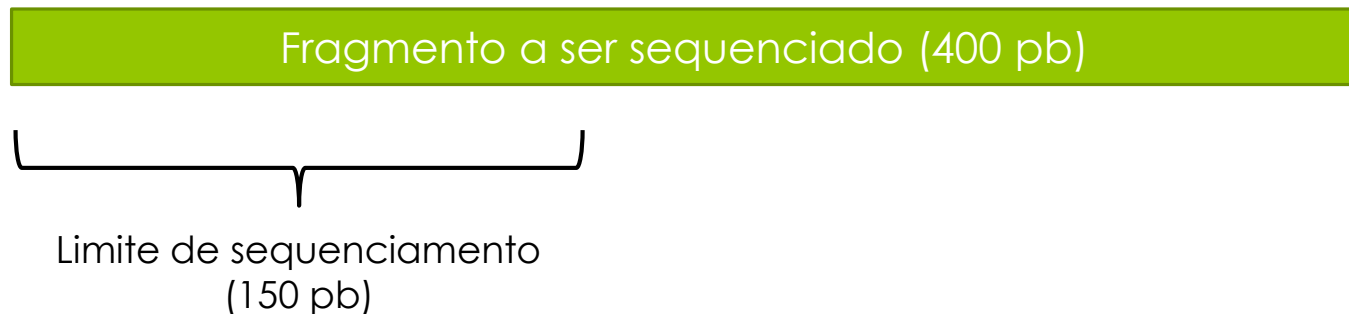
# Análise dos dados de sequenciamento de nova geração

## ● Escores de qualidade por símbolos



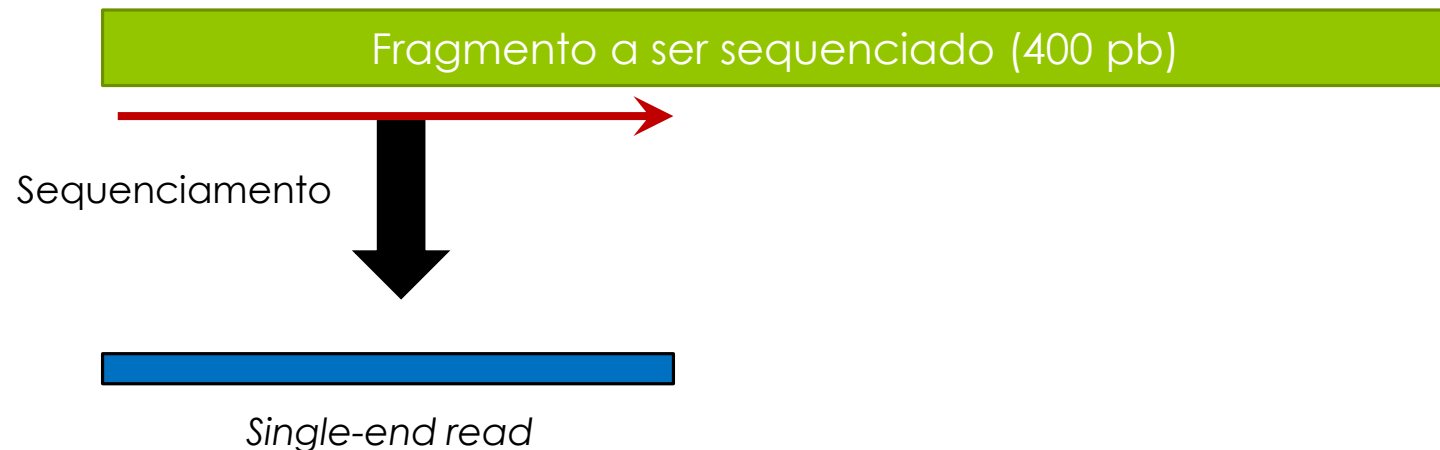
# Análise dos dados de sequenciamento de nova geração

- *Single-end versus Paired-end reads*
  - Maneira como cada fragmento de DNA é sequenciado



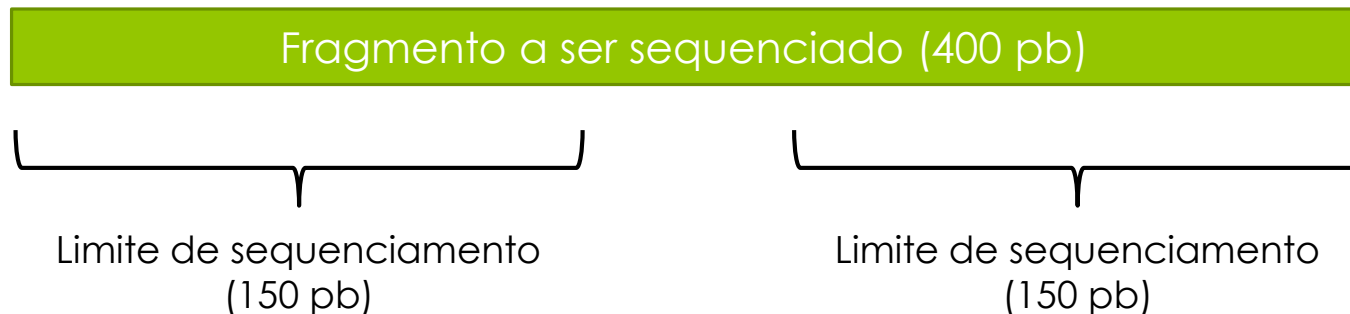
# Análise dos dados de sequenciamento de nova geração

- *Single-end versus Paired-end reads*
  - Maneira como cada fragmento de DNA é sequenciado



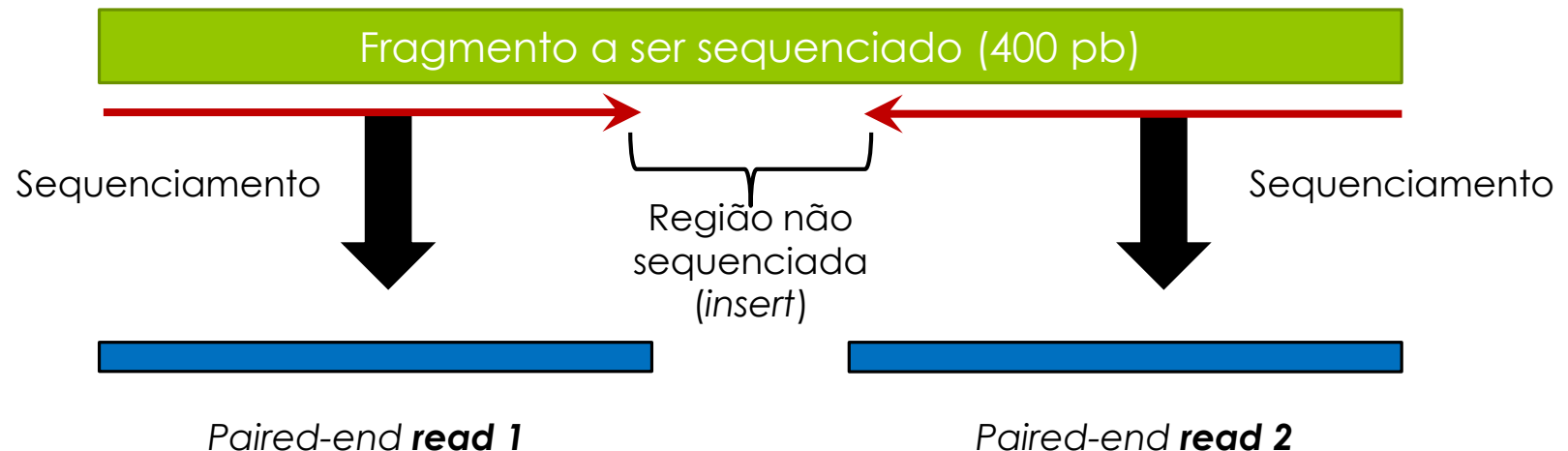
# Análise dos dados de sequenciamento de nova geração

- *Single-end versus Paired-end reads*
  - Maneira como cada fragmento de DNA é sequenciado



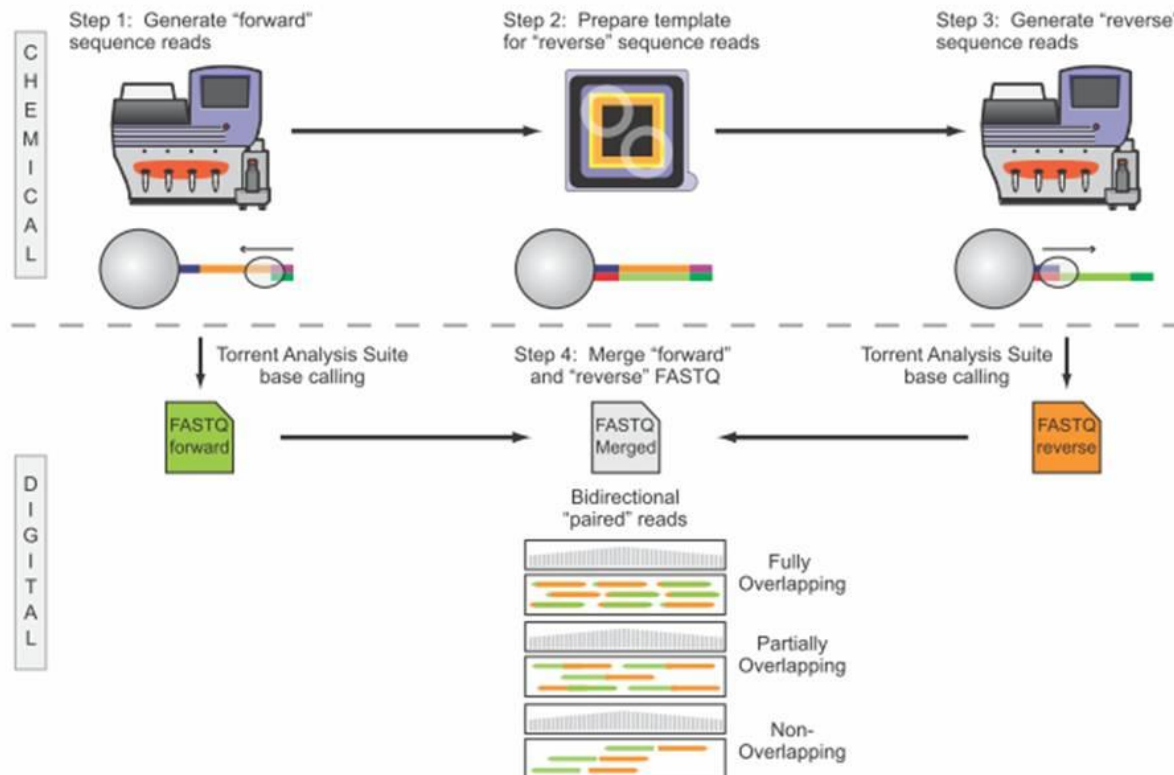
# Análise dos dados de sequenciamento de nova geração

- *Single-end versus Paired-end reads*
  - Maneira como cada fragmento de DNA é sequenciado



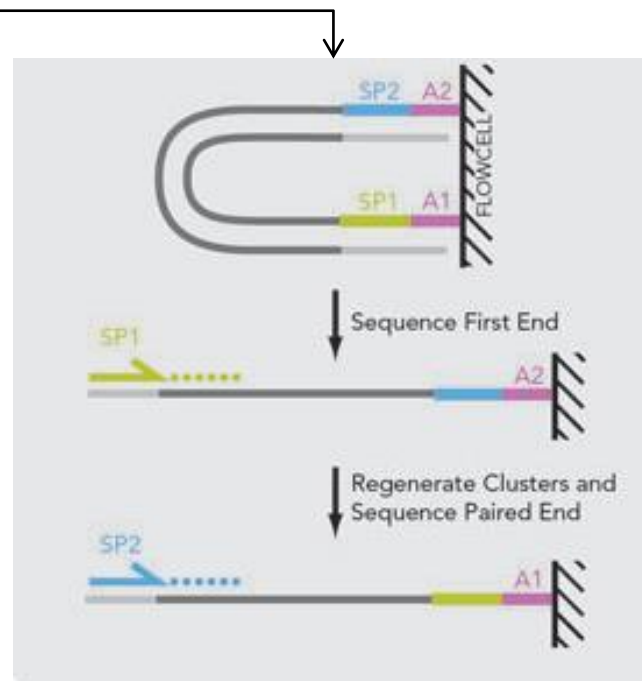
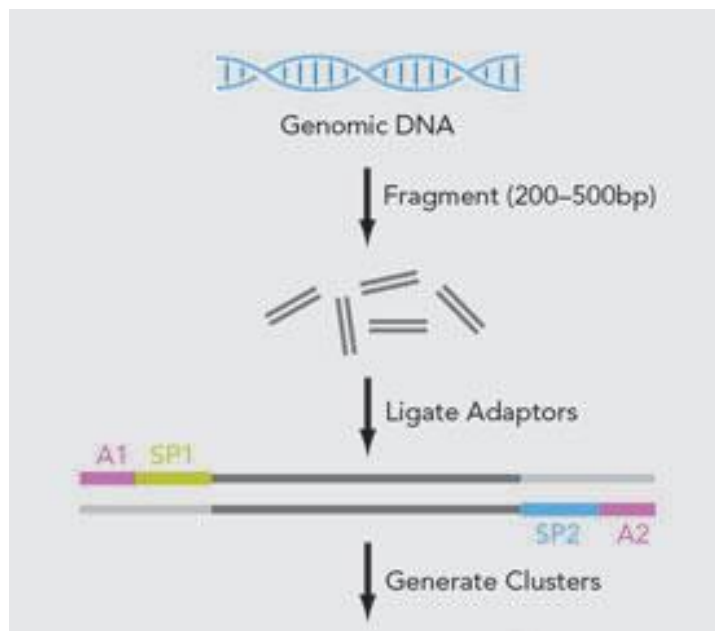
# Análise dos dados de sequenciamento de nova geração

- Single-end versus Paired-end (Ion Torrent)



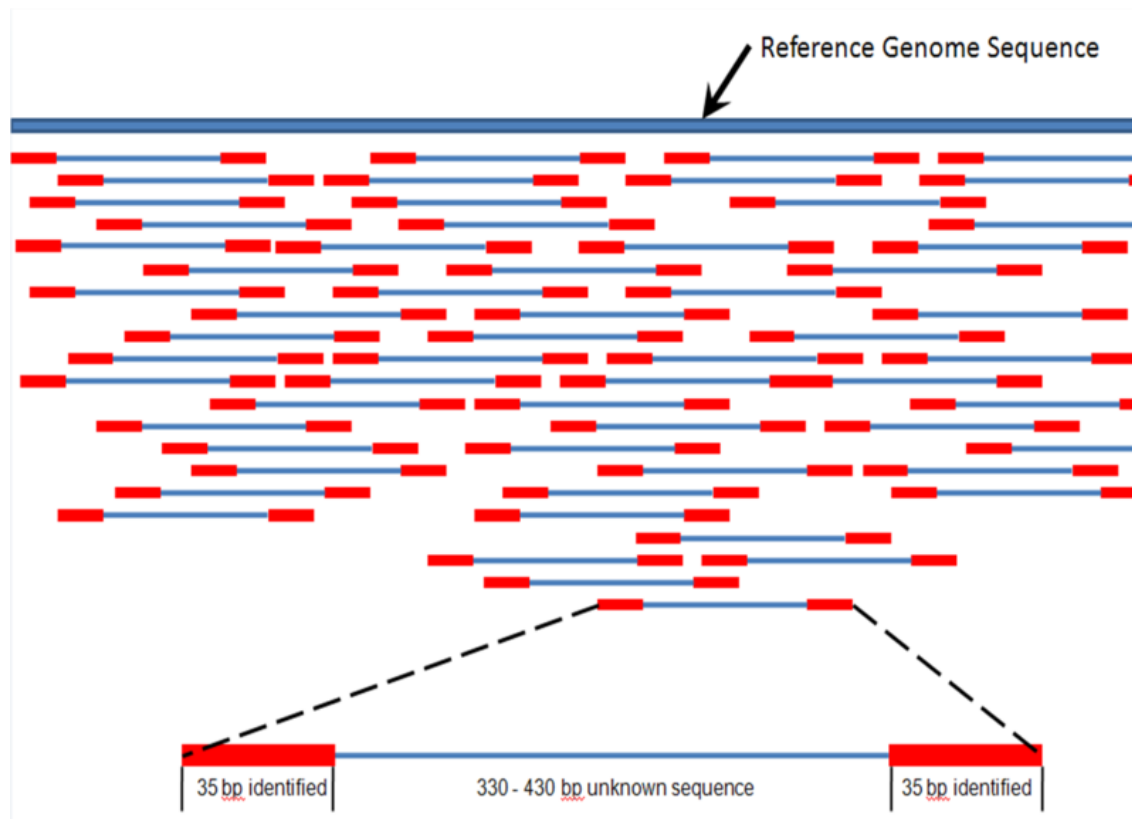
# Análise dos dados de sequenciamento de nova geração

- *Illumina paired-end*



# Análise dos dados de sequenciamento de nova geração

- *Paired-end reads*





# Análise dos dados de sequenciamento de nova geração

- Alinhamento (mapeamento)
  - Identificar a coordenada genômica de cada *read*
  - Deve ser computacionalmente rápida e permitir *mismatches*

```
...CCATAG          TATGCGCCC          CGGAAATTT          GGTATAC...
...CCAT          CTATATGCG          TCGGAAATT          CGGTATAC
...CCAT          GGCTATATG          CTATCGGAAA          GCGGTATA
...CCA          AGGCTATAT          CCTATCGGA          TTGCGGTA          C...
...CCA          AGGCTATAT          GCCCTATCG          TTTGCGGT          C...
...CC          AGGCTATAT          GCCCTATCG          AAATTTGC          ATAC...
...CC          TAGGCTATA          GCGCCCTA          AAATTTGC          GTATAC...

...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...
```

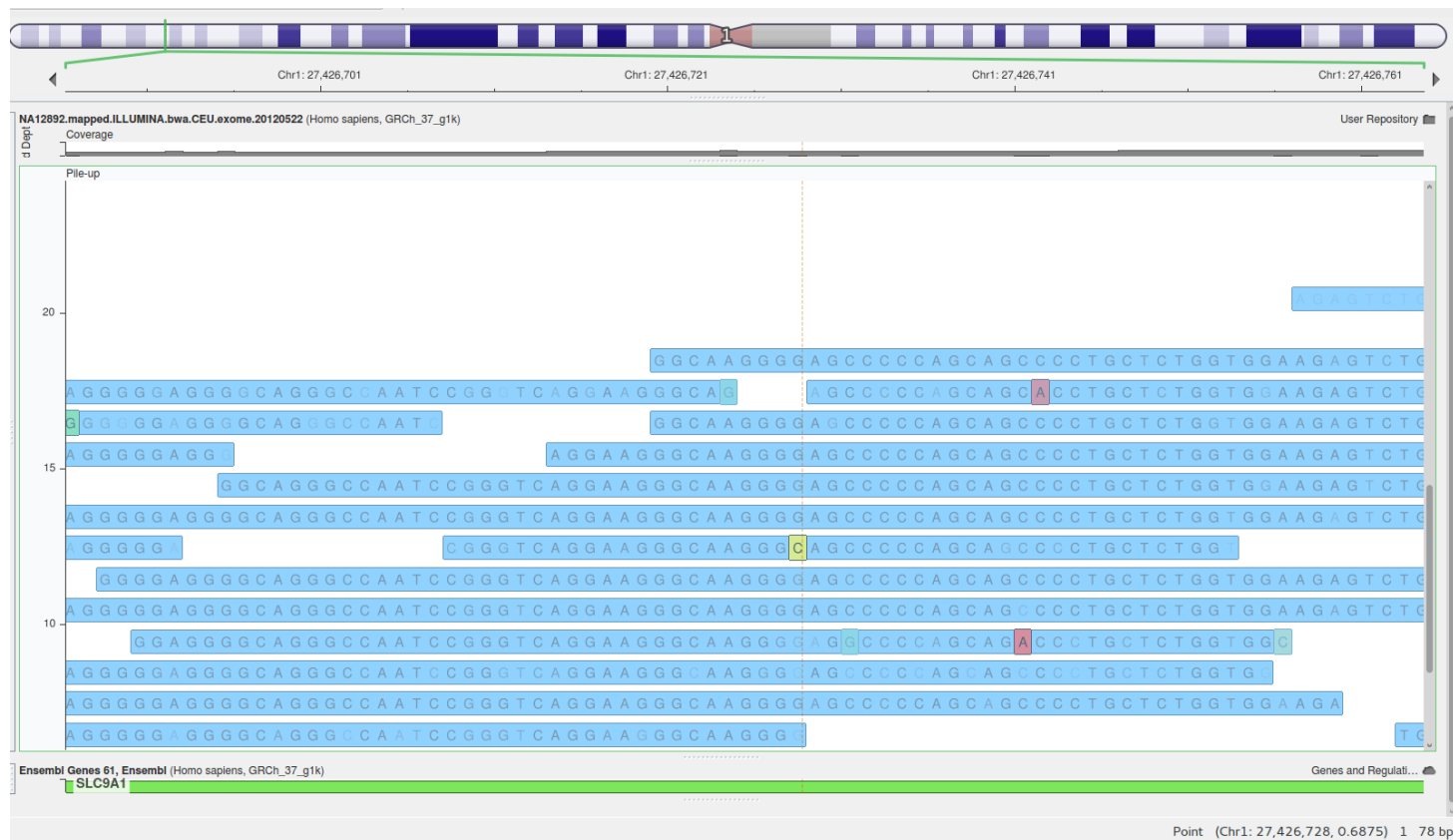
# Análise dos dados de sequenciamento de nova geração

- Exemplo de um resultado de alinhamento



# Análise dos dados de sequenciamento de nova geração

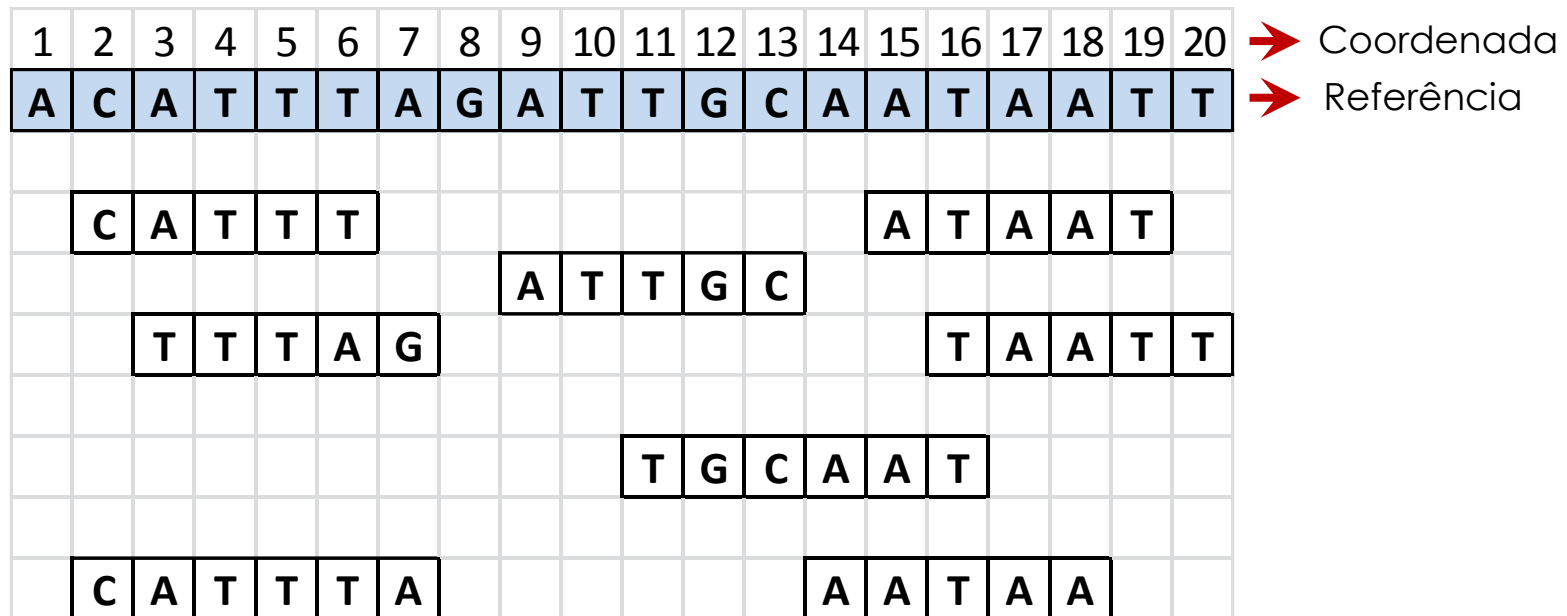
- Exemplo de um resultado de alinhamento



# Análise dos dados de sequenciamento de nova geração

- Cobertura / Coverage / Depth of Coverage
  - Três definições:
    - Quantas vezes em média cada base foi “lida” (sequenciada). Ex: a cobertura foi 30 X
    - Quantos % do genoma foi coberto pelo meu sequenciamento. Ex: 95% de cobertura
    - Quantos % do genoma foi lido pelo menos X vezes. Ex: 80% do genoma foi lido pelo menos 30 X

# Análise dos dados de sequenciamento de nova geração



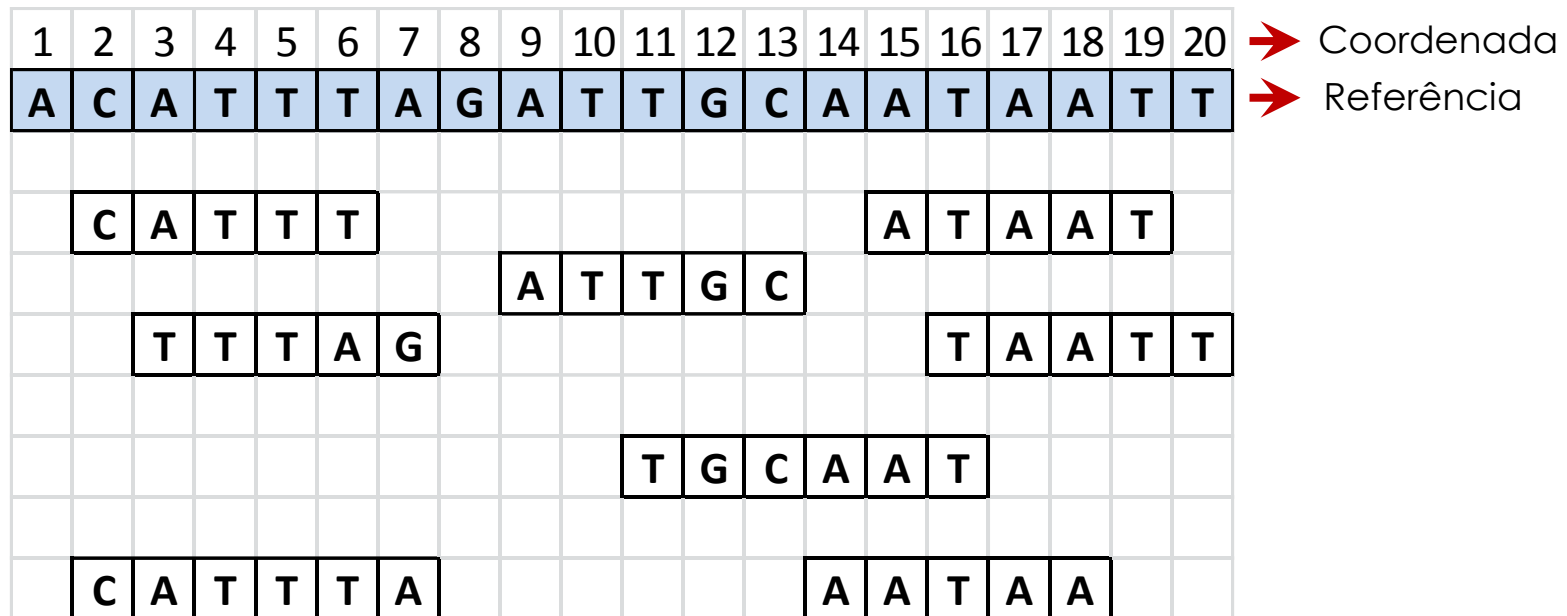
Quantos reads?

Cobertura média?

% do genoma coberto?

Coordenada e cobertura da base com maior cobertura?

# Análise dos dados de sequenciamento de nova geração



Quantos reads? → **8**

Cobertura média? →  $(0+2+3+3+3+...+3+2+1) / 20 = \mathbf{2,1\ X}$

% do genoma coberto? →  $18 / 20 = \mathbf{90\%}$

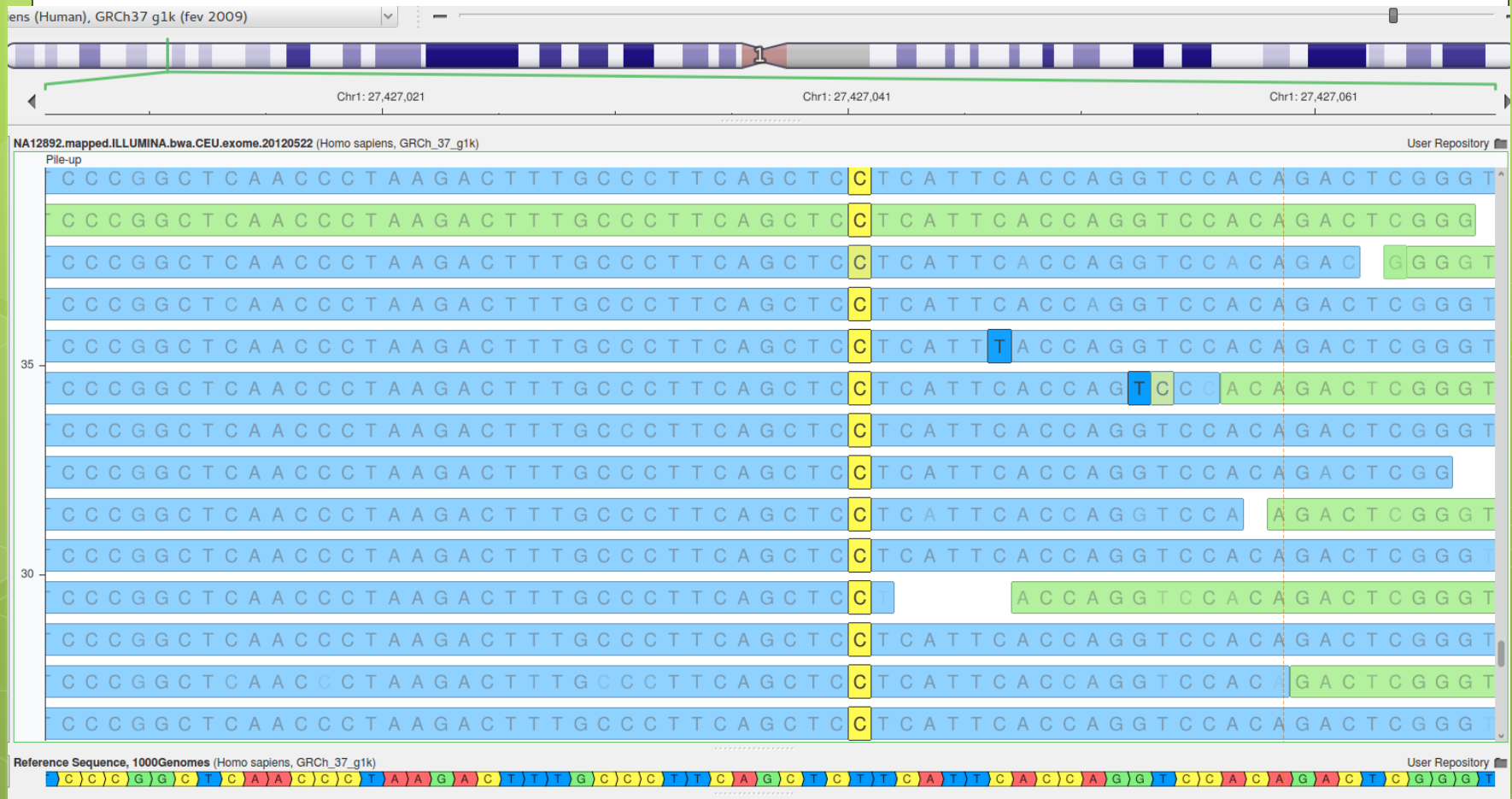
Coordenada e cobertura da base com maior cobertura? → **16 (4 X)**

# Análise dos dados de sequenciamento de nova geração

- *Variant Calling*
  - Identificar o que realmente é diferente do genoma de referência, ou seja, **identificar variantes**
  - Deve conseguir diferenciar:
    - Erro de sequenciamento versus variante real
    - Alteração em homozigose x heterozigose

# Análise dos dados de sequenciamento de nova geração

## Variant Calling

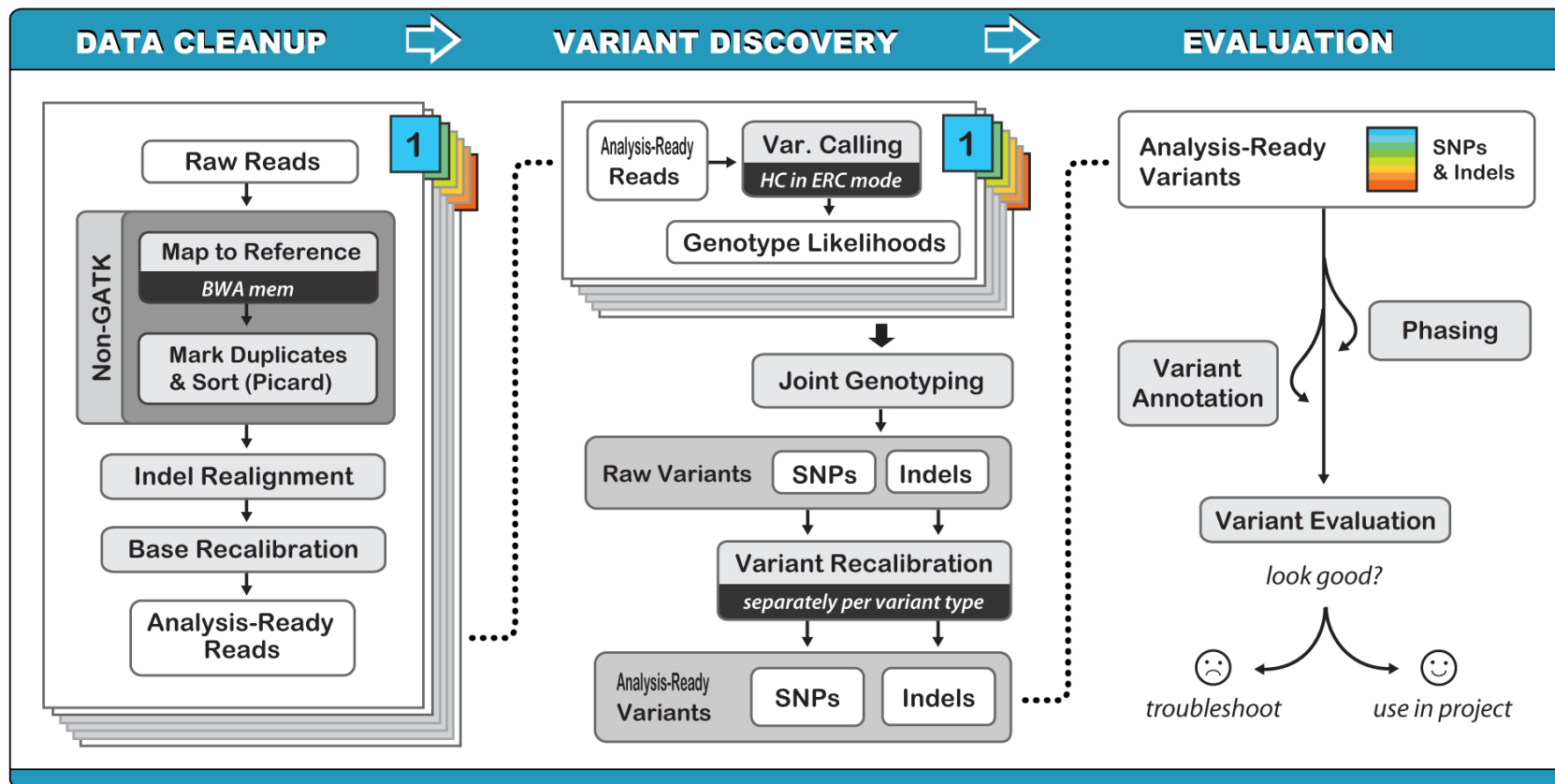




# Análise dos dados de sequenciamento de nova geração

- Anotação
  - Atribuir o máximo possível de informação sobre cada variante identificada
  - Nome do gene, troca de base, troca de aminoácido, classe funcional, escore de predição de patogenicidade, frequência em estudos populacionais, etc.

# Análise dos dados de sequenciamento de nova geração



# Análise dos dados de sequenciamento de nova geração

- Etapas da análise (tutorial):

