



Plataforma distribuída con Spark para el análisis de secuencias Next-generation

HPC

MSc. Vicente Machaca Arceda

November 19, 2021

Overview



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones

Table of Contents



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones



El análisis de secuencias de ADN obtenidas por *Next-generation sequencing* es una tarea vital en los estudios de Bioinformática. Pero lamentablemente, las herramientas tradicionales no son lentas y no aprovechan el poder computacional de un sistema distribuido



Desarrollar una herramienta para el análisis de secuencias obtenidas de *Next-generation sequencing*.

En esta versión, se implementó la funcionalidad para analizar la calidad de las lecturas.

Table of Contents



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones

Secuenciamiento de ADN

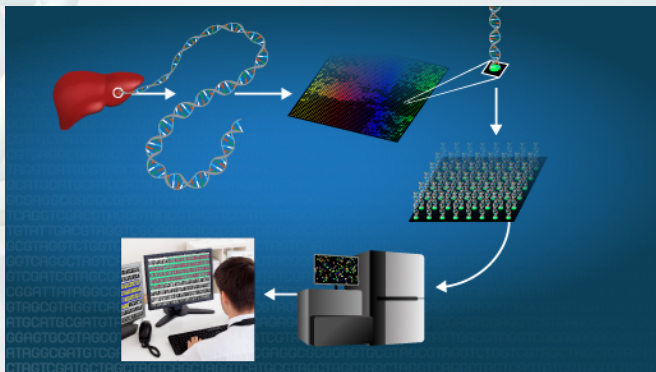


Figure: Secuenciamiento de ADN

Next-generation sequencing



Figure: Next-generation sequencing

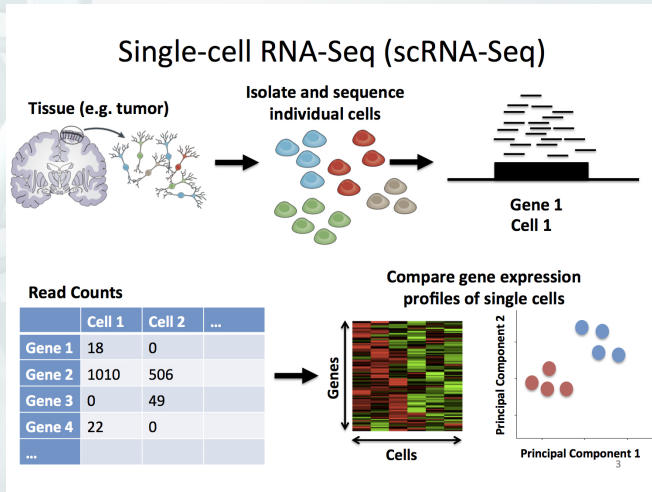


Figure: Single cell RNAseq

Table of Contents



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones

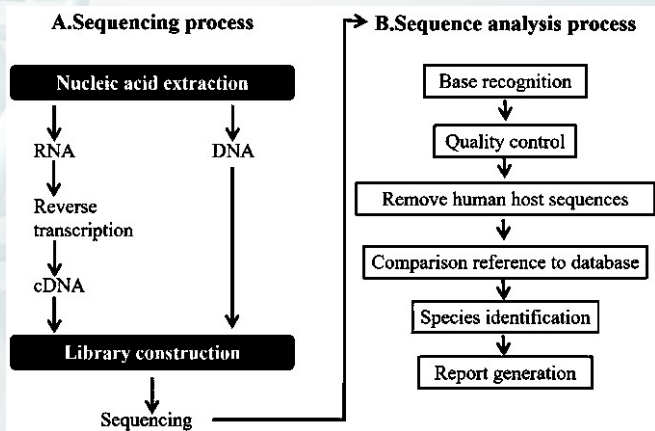


Figure: Phases comunes realizadas en un análisis de secuencias de ADN (Next-generation).



Table: Herramientas utilizadas para el proyecto

Herramienta	Version
Spark	3.2.0
Python	3.8.10
Pyspark	3.2.0



Table: Computadoras utilizadas en el sistema distribuído

Nombre de PC	Especificaciones
Desktop Asus	Procesador i7 de séptima generación y 8GB de memoria RAM. Sistema operativo Linux.
Laptop Asus	Procesador i5 de quinta generación y 4GB de memoria RAM. Sistema operativo Linux.

Propuesta

Hardware





- ▶ Conteo de la cantidad de secuencias.
- ▶ Conteo total de las bases nitrogenadas.
- ▶ Computo de la longitud de todas las secuencias.
- ▶ Computo del promedio de las longitudes de las secuencias.
- ▶ Computo de la ocurrencia de cada base nitrogenada.
- ▶ Análisis de contenido por base.

Table of Contents



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones



Para evaluar el desempeño de la propuesta se evaluó las secuencias con el código **ERR3014700**, estas fueron descargadas de NCBI.



```
1  {
2    "bases": 1843156,
3    "total_seqs": 462393,
4    "seqs_len": [
5      523,      600,      599,      600,      599,
6      600,      600,      529,      600,      538
7    ],
8    "seqs_len_mean": 564,
9    "bases_ocurrence": [
10     [      "C",      71397200      ],
11     [      "N",      4054      ],
12     [      "D",      3004      ],
13     [      "G",      70502420      ],
14     [      "A",      59888213      ],
15     [      "F",      16160      ],
16     [      "T",      59010438      ],
17     [      "B",      203      ],
18     [      "E",      4144      ]
19   ]
20 }
```

Figure: Estadísticas de las lecturas.

Resultados

Análisis de contenido por base

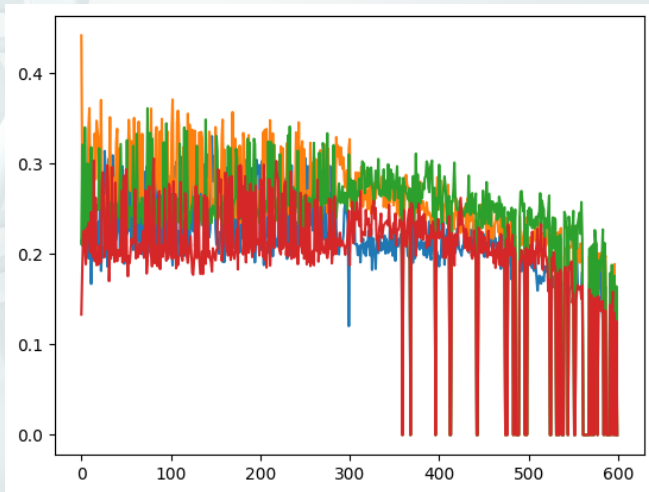


Figure: Análisis de contenido por base.

Table of Contents



Introducción
Problema
Objetivos

Marco teórico

Propuesta

Resultados

Conclusiones



- ▶ En este proyecto se ha desarrollado una herramienta distribuída que permite hacer el análisis masivo de grandes cantidades de lecturas de ADN (Next-generation sequencing).
- ▶ El proyecto se enfoco en el análisis de calidad de las lecturas de ADN, este es un paso crucial en cualquier experimento de Bioinformática. Como resultado, la propuesta obtiene estadísticas referentes a las ocurrencias de las bases nitrogenadas y además se desarrollo un análisis de contenido por base.

References I



