



OPTIMIZACIÓN DE LA BÚSQUEDA DE PATRONES GENÉTICOS MEDIANTE ÁRBOLES

Integrantes: Constanza Araya ([GitHub/Cnstnza](#))

Profesor: Christian Vasquez ([GitHub/chvasquezr](#))

14 de noviembre de 2025

Índice

1	Introduccion	2
2	Objetivos	3
3	Estructuras de datos	4
3.1	Lista enlazada	4
3.1.1	Nodo	4
3.2	Arboles	6
3.2.1	Arbol 4-ario (Trie)	8
4	Implementacion de comandos	10
4.1	bio start m	10
4.2	bio read archivo.txt	10
4.3	bio search XX	11
4.4	bio max, bio min, bio all	11
4.4.1	bio max	12
4.4.2	bio min	12
4.4.3	bio all	13
4.5	bio exit	14
4.6	help	14
5	Diseño del programa (arquitectura general)	15
6	Manejo de memoria	16
7	Conclusion	17
	Referencias	18

1. Introduccion

Este proyecto tiene como objetivo crear un programa que lea una secuencia de ADN (formada por letras A, C, G y T) y busque genes cortos dentro de ella. Para ello se utilizaran diversas estructuras de datos. Se implementara un arbol 4-ario (tipo trie) para guardar las combinaciones posibles de genes y listas enlazadas para guardar las posiciones donde aparecen. Se implementara la deteccion de comandos, entre ellos estaran:

- **bio start m:** crea un arbol 4-ario con genes de tamaño m
- **bio read archivo.txt:** lee la secuencia de ADN desde el archivo archivo.txt
- **bio search XX:** busca el gen XX en la secuencia
- **bio max:** muestra el gen mas repetido en la secuencia
- **bio min:** muestra el gen menos repetido en la secuencia
- **bio all:** muestra todos los gen que tienen frecuencia ≥ 0
- **bio exit:** libera memoria y sale del programa

El proyecto se desarrolla en el lenguaje de programacion C.

2. Objetivos

- Implementar un sistema de búsqueda de patrones utilizando árboles tries.
- Implementar y manipular estructuras de datos abstractas como árboles y listas enlazadas para almacenar y gestionar los datos manipulados.
- Desarrollar habilidades en programación en lenguaje C, centrándose en el manejo de memoria, punteros y eficiencia algorítmica.

3. Estructuras de datos

Para la implementación del proyecto se utilizaron las siguientes estructuras de datos:

- Lista enlazada (Linked List)
- Arbol 4-ario (Trie)

3.1. Lista enlazada

Una **lista enlazada** es una estructura de datos dinámica que consiste en una serie de elementos llamados nodos, donde cada nodo contiene un dato y una referencia al siguiente nodo en la secuencia. A diferencia de un array, las listas enlazadas permiten un crecimiento dinámico sin un tamaño fijo y facilitan la inserción y eliminación de elementos, ya que solo requieren actualizar referencias.

Sin embargo, las listas enlazadas presentan tanto ventajas como desventajas frente a los arreglos:

- **Ventajas:** Inserciones y eliminaciones en cualquier posición sin necesidad de mover todos los elementos. Crecimiento dinámico sin tamaño predefinido.
- **Desventajas:** La búsqueda de un dato específico requiere recorrer la lista nodo por nodo ($O(N)$), mientras que en un arreglo, si se conoce el índice, se puede acceder directamente en $O(1)$.

Para comprender mejor una lista enlazada es importante saber que es un nodo.

3.1.1. Nodo

En informática, un nodo es un elemento individual de una estructura de datos que almacena información y, en muchos casos, tiene referencias o enlaces a otros nodos. Los nodos son la base de muchas estructuras de datos dinámicas, como listas enlazadas, árboles y grafos.

Un nodo tiene dos componentes principales:

- **Dato:** La información que el nodo almacena, que puede ser de cualquier tipo de dato, como un número, una cadena de texto o un objeto más complejo.
- **Referencia:** Un puntero o referencia al siguiente nodo en la lista

En la siguiente imagen se puede apreciar de manera grafica que es un nodo y tambien un extracto de codigo en C que representa un nodo de una lista enlazada.

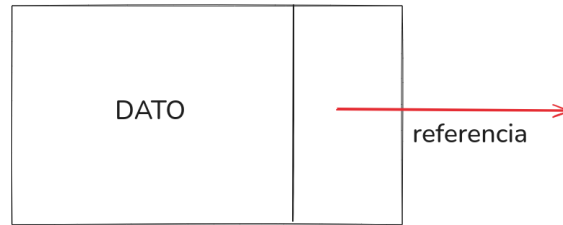


Figura 1: Nodo

```

1  /*Definición del nodo*/
2  typedef struct Nodo {
3      int dato;           /*Valor almacenado*/
4      struct Nodo* siguiente; /*Puntero al siguiente nodo*/
5  } Nodo;
6
7  /*Creación de un nuevo nodo*/
8  Nodo* crearNodo(int valor) {
9      Nodo* nuevo = (Nodo*) malloc(sizeof(Nodo));
10     nuevo->dato = valor;
11     nuevo->siguiente = NULL;
12     return nuevo;
13 }
14

```

Listing 1: Creacion de Nodo

Las listas enlazadas tienen un nodo inicial llamado **head** el cual apunta al primer elemento de la lista. Si se llega a un nodo que no apunta a ningún otro, se dice que dicho nodo apunta a **null** o final de la lista.

En la siguiente imagen se puede apreciar de manera grafica que es una lista enlazada y su implementacion en C

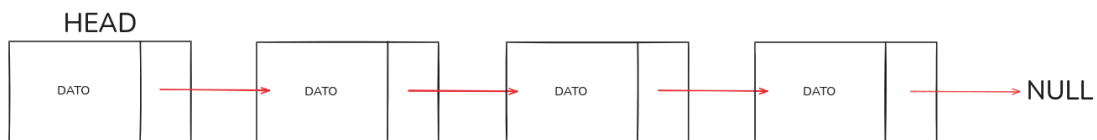


Figura 2: Lista enlazada

```

1  /*Inserta un nodo al final de la lista*/
2  void insertarFinal(Nodo** cabeza, int valor) {
3      Nodo* nuevo = crearNodo(valor);
4      if (*cabeza == NULL) {

```

```

5      *cabeza = nuevo;
6      return;
7  }
8  Nodo* actual = *cabeza;
9  while (actual->siguiente != NULL) {
10     actual = actual->siguiente;
11  }
12  actual->siguiente = nuevo;
13  }
14

```

Listing 2: Estructura lista enlazada

La lista enlazada, se utiliza dentro del arbol 4-ario para almacenar las posiciones donde aparecen los genes. En el nodo hoja del Trie (aquel que corresponde al final de un gen de longitud m), se almacena un puntero a una lista (`ListaInt* lista_posiciones`). Cada nodo de esta lista (`ListaInt`) guarda un entero: la posición (índice 0-basado) donde ese gen comienza en la secuencia S . Es decir cada vez que el programa detecta un gen (una secuencia de longitud m) dentro de la secuencia principal de ADN, se almacena su posición inicial (el índice donde empieza).

Como un mismo gen puede aparecer muchas veces en distintas partes de la secuencia, se necesita una estructura que permita agregar fácilmente nuevas posiciones sin un tamaño fijo ni un límite predefinido. Aquí es donde se utiliza la **lista enlazada**. También permite que el programa no solo diga si un gen existe, sino también cuántas veces (bio search, bio max, bio min) y dónde comienza cada aparición.

3.2. Árboles

Los Árboles son las estructuras de datos mas utilizadas, pero también una de las mas complejas, Los Árboles se caracterizan por almacenar sus nodos en forma jerárquica y no en forma lineal como las Listas enlazadas, Colas, Pilas, etc.

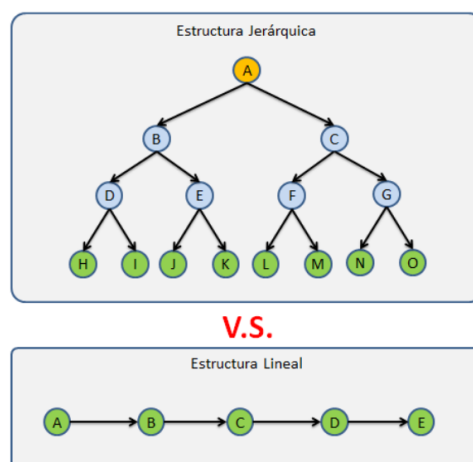


Figura 3: Visualización de árbol vs lista

Para comprender que es un árbol, es importante entender como esta estructurado y cuales son sus componentes. Algunos de ellos son:

- **Nodo:** Se le llama nodo a cada elemento que contiene un árbol.
- **Raíz:** Es el nodo principal del árbol, es el primer nodo del árbol, solo puede haber un nodo como raíz.
- **Padre:** Este termino representa a cualquier nodo que tiene al menos un hijo.
- **Hijo:** Los hijos son los nodos que tienen un padre.
- **Nodo hermano:** Los nodos hermanos son aquellos que comparten el mismo padre.
- **Hoja:** Son los nodos que no tienen hijos, por lo tanto siempre se encuentran en los extremos de un árbol.

Puede ser algo difícil de comprender, por lo que a continuacion se muestra una imagen que ejemplifica los conceptos antes mencionados.

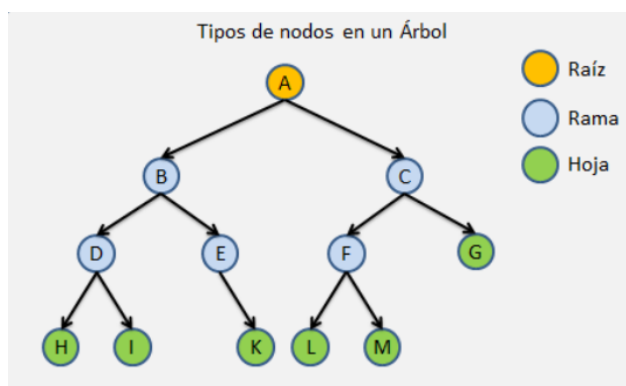


Figura 4: Componentes de un árbol

Los árboles a demas de nodos tienen otras propiedades importantes como:

- **Nivel:** El nivel es cada generacion dentro del árbol, por ejemplo: si se tiene un nodo hoja y se le agrega un hijo, entonces el nodo hoja se convierte en un nodo rama, pero a demas el árbol aumenta un nivel mas.
- **Altura:** La altura de un árbol es la cantidad de niveles que tiene el árbol.

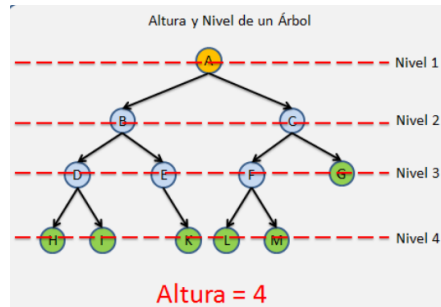


Figura 5: Nivel y altura de un árbol

- **Peso:** El peso corresponde a la cantidad de nodos que tiene un árbol, es un concepto muy importante ya que permite tener una idea de que tan grande es el árbol y el tamaño de memoria que usará en ejecución.

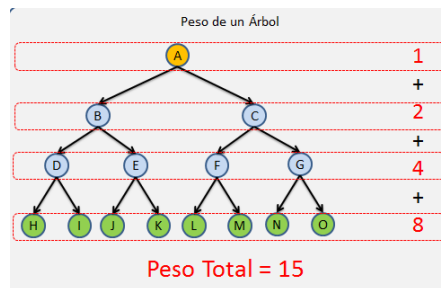


Figura 6: Peso de un árbol

A continuación un extracto de código en C de un árbol binario:

```

1  /*Estructura para un Nodo de árbol Binario (2 hijos)*/
2  typedef struct NodoArbol {
3      int dato; /*El valor que almacena el nodo*/
4      /*Punteros a los hijos (recursivos)*/
5      struct NodoArbol* izquierda; /*Puntero al sub árbol izquierdo*/
6      struct NodoArbol* derecha; /*Puntero al sub árbol derecho*/
7  } NodoArbol;
8

```

Listing 3: Árbol binario

3.2.1. Árbol 4-ario (Trie)

Existen muchas formas de implementar un árbol, por ejemplo un árbol binario, pero en este proyecto se utilizará un árbol 4-ario (tipo trie). Este tipo de árbol tiene 4 hijos posibles, uno por cada letra del ADN (A, C, G y T). Los tries se usan comúnmente en sistemas de búsqueda de patrones porque permiten recorrer y almacenar grandes cantidades de datos de forma ordenada y eficiente. Gracias a su estructura, es posible encontrar patrones dentro de una secuencia de

manera rápida sin necesidad de comparar cada combinación manualmente.

El árbol 4-ario se puede representar mediante la siguiente imagen:

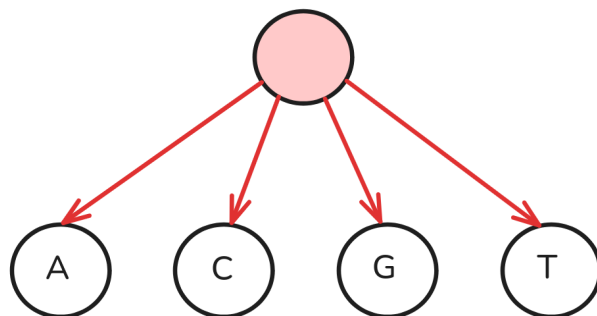


Figura 7: Arbol 4-ario

4. Implementacion de comandos

En este proyecto se implementaron principalmente siete comandos para interactuar con el programa. Estos comandos permiten al usuario iniciar el arbol, leer la secuencia de ADN, buscar genes,etc.

4.1. bio start m

Este comando crea un arbol de orden m, donde m es el tamaño de la secuencia de ADN.

```
1      Nodo* bio_start(int m) {
2          if (m <= 0) {
3              printf("Error: El tamano del gen (m) debe ser mayor a 0.\n");
4              return NULL;
5          }
6          if (TRIE_ROOT != NULL) {
7              bio_exit();
8          }
9          GEN_SIZE = m;
10         TRIE_ROOT = generar_arbol_recursivo(0, m);
11
12         if (TRIE_ROOT != NULL) {
13             printf("Arbol creado con tamano %d\n", m);
14         } else {
15             printf("Error al crear el arbol.\n");
16         }
17         return TRIE_ROOT;
18     }
19
```

Listing 4: Implementacion bio start m

4.2. bio read archivo.txt

Este comando lee la secuencia de ADN desde el archivo archivo.txt.

```
1      /*Logica para leer toda la secuencia S*/
2      fseek(file, 0, SEEK_END);
3      long n_size = ftell(file);
4      fseek(file, 0, SEEK_SET);
5
6      char* secuencia_S = (char*)malloc(n_size + 1);
7      /*Verificaci n de malloc*/
8      size_t read_len = fread(secuencia_S, 1, n_size, file);
9      fclose(file);
```

Listing 5: Fragmento clave read archivo.txt

Aginacion dinamica: Se utiliza fseek y ftell para determinar el tamaño exacto del archivo (n) y luego se asigna dinámicamente el buffer `secuencia_S` con malloc. Esto garantiza que la secuencia genética completa se cargue en la memoria (heap) de un solo golpe con fread, mejorando la eficiencia de lectura del archivo.

4.3. bio search XX

La función `bio_search` implementa el núcleo del sistema de búsqueda. Su objetivo es encontrar la frecuencia y todas las posiciones de inicio de un gen específico (de longitud m) dentro de la secuencia genética S previamente cargada en el Árbol Trie

```

1      /** Verificar que el arbol esta inicializado **/
2      if (TRIE_ROOT == NULL || GEN_SIZE == 0) return -1;
3      /** Validar que el gen tenga el largo correcto **/
4      if ((int)strlen(gen_G) != GEN_SIZE) return -1;
5      /** Recorrer el trie segun cada caracter del gen **/
6      for (int i = 0; i < GEN_SIZE; i++) {
7          int index = get_base_index(gen_G[i]);
8          if (index == -1) return -1;      /** caracter invalido **/
9          actual = actual->hijos[index];  /** avanzar en el arbol **/
10     }
11     /** Contar cuantas veces aparece el gen **/
12     int frecuencia = contar_posiciones(actual->lista_posiciones);
13     if (frecuencia == 0) return -1;
14     /** Imprimir el gen, su frecuencia y posiciones **/
15     printf("%s %d", gen_G, frecuencia);
16     imprimir_posiciones(actual->lista_posiciones);
17

```

Listing 6: Fragmento clave bio search XX

4.4. bio max, bio min, bio all

Estos comandos permiten obtener el gen mas repetido (bio max), el gen menos repetido (bio min) y todos los gen que tienen frecuencia ≥ 0 (bio all). A continuacion extractos de cada uno de estos comandos

4.4.1. bio max

Busca dentro del árbol el gen con mayor frecuencia. Se realiza un recorrido completo (DFS) y en cada hoja se compara la frecuencia actual con la maxima encontrada hasta el momento.

```
1      void bio_max() {
2          if (TRIE_ROOT == NULL || GEN_SIZE == 0) { printf("Error: Arbol no
cargado.\n"); return; }
3          char* gen_buffer = (char*)malloc(GEN_SIZE + 1);
4          if (gen_buffer == NULL) { perror("Error de asignacion de memoria");
return; }
5          int max_freq = 0;
6          recorrido_extremos(TRIE_ROOT, gen_buffer, 0, 1, &max_freq);
7          if (max_freq == 0) { printf("No hay genes presentes.\n"); free(
gen_buffer); return; }
8          imprimir_recorrido(TRIE_ROOT, gen_buffer, 0, max_freq);
9          free(gen_buffer);
10         }
11     }
```

Listing 7: Funcion bio max

4.4.2. bio min

Similar a bio max, pero busca el gene con frecuencia minima ¿0. Ignora hojas sin apariciones. Comparte la misma estrategia de recorrido pero con la condicion inversa.

```
1      void bio_min() {
2          if (TRIE_ROOT == NULL || GEN_SIZE == 0) { printf("Error: Arbol no cargado.\n
"); return;
3          char* gen_buffer = (char*)malloc(GEN_SIZE + 1);
4          if (gen_buffer == NULL) { perror("Error de asignacion de memoria"); return;
}
5          int min_freq = INT_MAX;
6          recorrido_extremos(TRIE_ROOT, gen_buffer, 0, 0, &min_freq);
7          if (min_freq == INT_MAX) { printf("No hay genes presentes.\n"); free(
gen_buffer); return; }
8          imprimir_recorrido(TRIE_ROOT, gen_buffer, 0, min_freq);
9          free(gen_buffer);
10         }
11     }
```

Listing 8: Funcion bio min

4.4.3. bio all

Recorre todas las hojas del arbol y muestra solo aquellas cuyo conteo de apariciones es mayor a cero. Es una enumeración completa de todos los genes presentes en los datos cargados.

```
1      void bio_all() {
2          if (TRIE_ROOT == NULL || GEN_SIZE == 0) { printf("Error: Arbol no cargado.\n"); return; }
3          char* gen_buffer = (char*)malloc(GEN_SIZE + 1);
4          if (gen_buffer == NULL) { perror("Error de asignacion de memoria"); return; }
5          /*Funcion recursiva interna para ALL (imprime todos los que tengan frecuencia > 0)*/
6          void bio_all_recursivo(Nodo* nodo, char* buffer, int nivel) {
7              if (nodo == NULL) return;
8              if (nivel == GEN_SIZE) {
9                  int frecuencia = contar_posiciones(nodo->lista_posiciones);
10                 if (frecuencia > 0) {
11                     buffer[GEN_SIZE] = '\0';
12                     printf("%s %d", buffer, frecuencia);
13                     imprimir_posiciones(nodo->lista_posiciones);
14                     printf("\n");
15                 }
16                 return;
17             }
18             const char bases[] = {'A', 'C', 'G', 'T'};
19             for (int i = 0; i < 4; i++) {
20                 buffer[nivel] = bases[i];
21                 bio_all_recursivo(nodo->hijos[i], buffer, nivel + 1);
22             }
23         }
24         bio_all_recursivo(TRIE_ROOT, gen_buffer, 0);
25         free(gen_buffer);
26     }
27 }
```

Listing 9: Funcion bio all

El arreglo `gen_buffer` se utiliza para reconstruir el gen letra por letra durante el recorrido del árbol. Como cada nivel del Trie representa una base del ADN, el buffer permite almacenar la 'ruta' tomada para llegar a cada hoja, reconstruyendo así el gen completo sin necesidad de almacenarlo textualmente en cada nodo (lo cual ahorra memoria).

Complejidad computacional

- El recorrido completo del Trie tiene complejidad $O(4^m)$ donde $m = GEN_SIZE$ (longitud del gen).

-
- Esto es eficiente porque el tamaño del alfabeto es fijo (A,C,G,T) y m es pequeño y constante en el contexto biologico.
 - Las funciones bio max, bio min y bio all tienen la misma complejidad porque requieren recorrer todas las hojas.

4.5. bio exit

```
void bio_exit()liberar_arbol(TRIE_ROOT);TRIE_ROOT = NULL;GEN_SIZE = 0;printf("Limpiando cacheysalendodel programa...");
```

4.6. help

Este comando es un adicional, ya que no estaba en los requerimientos iniciales, pero se considero importante para el usuario. Ya que al iniciar el programa no hay nada que indique como funciona si el usuario no tiene conocimientos previos. Para todo lo relacionado con la salida por terminal se utilizo definicion de macros con el fin de tener una interfaz mas atractiva para el usuario.

A continuacion la salida por terminal de la funcion help:

```
COMANDOS DISPONIBLES:
bio start m: crea un arbol 4-ario con genes de tamaño m
bio read archivo.txt: lee la secuencia de ADN desde el archivo archivo.txt
bio search XX: busca el gen XX en la secuencia
bio max: muestra el gen mas repetido en la secuencia
bio min: muestra el gen menos repetido en la secuencia
bio all: muestra todos los gen que tienen frecuencia > 0
bio exit: libera memoria y sale del programa
] >bio
```

Figura 8: Salida por terminal del comando help

5. Diseño del programa (arquitectura general)

Al ejecutar el programa pasa por los siguientes pasos:

- bio start 4: Crea un arbol 4-ario con genes de tamaño 4.
- bio read adn.txt: Lee la secuencia de ADN desde el archivo adn.txt.
- bio search XX: Busca el gen XX en la secuencia.
- bio maz/min/all: recorren todo el arbol y muestran los genes que cumplen con la condición.
- bio exit: libera memoria y sale del programa.

Al ejecutar se tiene la siguiente salida por terminal:

```
>bio start 4
Arbol creado con tamano 4
>bio read adn.txt
Secuencia S leida desde el archivo
>bio search GATA
GATA 2 0 9
>bio search ATAC
ATAC 2 1 10
>bio search TOCA
-1
>bio min
AATT 1 5
ACAA 1 3
ACAG 1 12
ATTG 1 6
CAAT 1 4
TGAT 1 8
TTGA 1 7
>bio max
ATAC 2 1 10
GATA 2 0 9
TACA 2 2 11
>bio all
AATT 1 5
ACAA 1 3
ACAG 1 12
ATAC 2 1 10
ATTG 1 6
CAAT 1 4
GATA 2 0 9
TACA 2 2 11
TGAT 1 8
TTGA 1 7
>bio exit
Limpiando cache y saliendo del programa...
```

Figura 9: Salida por terminal del programa

6. Manejo de memoria

El programa hace un uso intensivo de memoria dinamica, ya que tanto el arbol (trie) como las listas enlazadas se construyen en tiempo de ejecucion a partir de la secuencia de ADN cargada. Para ello se utilizan los siguientes mecanismos:

- **malloc:** Se emplea para reservar memoria de forma dinamica para cada nodo del arbol y para cada nodo de las listas enlazadas donde se almacenan las posiciones de aparicion de cada gen.
- **free:** Se utiliza para liberar toda la memoria asignada una vez que el programa finaliza o cuando se ejecuta el comando `bio exit`. Esto evita fugas de memoria (memory leaks).
- **Estructuras dinamicas:** El arbol se implementa como un trie 4-ario, donde cada nodo contiene un arreglo de cuatro punteros (A, C, G, T). En las hojas del arbol, cada gen almacena una lista enlazada con todas las posiciones donde aparece en la secuencia.

Cada nodo del arbol se reserva mediante `malloc`, incluyendo su arreglo de punteros a hijos. De forma similar, cada posicion encontrada en la secuencia se inserta en una lista enlazada, cuyos nodos tambien son reservados dinamicamente. Para asegurar una correcta administracion de memoria, el programa realiza un recorrido postorden del arbol al momento de liberar los recursos. En este proceso:

1. Se liberan primero todas las listas enlazadas asociadas a cada hoja.
2. Luego se libera cada nodo del arbol recursivamente.
3. Finalmente, se libera la memoria utilizada por la secuencia principal y demas buffers auxiliares.

Esto garantiza que ninguna referencia quede perdida y que toda la memoria solicitada al sistema sea correctamente devuelta.

7. Conclusion

El desarrollo de este proyecto permitio implementar un sistema eficiente para el analisis de secuencias de ADN mediante el uso de estructuras de datos dinamicas. El arbol 4-ario (trie) demostro ser especialmente adecuado para almacenar y buscar genes de longitud fija, permitiendo realizar consultas en tiempo $O(m)$ sin depender del tamano total de la secuencia. Complementariamente, las listas enlazadas utilizadas en cada nodo hoja permitieron registrar todas las posiciones de aparicion de cada gen, proporcionando informacion completa y facilmente accesible. Ademias, la gestion dinamica de memoria mediante `malloc` y `free` aseguro un uso responsable de los recursos, evitando fugas y permitiendo la liberacion ordenada de todas las estructuras al finalizar la ejecucion del programa. En conjunto, el proyecto integra de manera efectiva conceptos fundamentales como arboles, listas enlazadas, manejo de memoria y diseño modular. El resultado es una herramienta funcional que permite procesar y consultar secuencias geneticas de forma rapida, clara y eficiente.

Referencias

- esantos. (2024, noviembre). *¿Qué es una Lista Enlazada?* Consultado el 13 de noviembre de 2025, desde <https://www.ervinsantos.com/post/50-%C2%BFqu%C3%A9-es-una-lista-enlazada>
- oblancarte. (2021, agosto). *Estructura de datos Árboles*. Consultado el 13 de noviembre de 2025, desde <https://www.oscarblancarteblog.com/2014/08/22/estructura-de-datos-arboles/>
- wikipedia. (2024, diciembre). *Nodo (informática)*. Consultado el 13 de noviembre de 2025, desde [https://es.wikipedia.org/wiki/Nodo_\(inform%C3%A1tica\)](https://es.wikipedia.org/wiki/Nodo_(inform%C3%A1tica))