

Instituto Tecnológico de Buenos Aires  
Teoría de Lenguajes y Autómatas

Trabajo Práctico 1:  
**Analizador Léxico Genómico**

**Integrantes:**

Altamiranda Enzo, 51265

Camisay Santiago, 47583

Costesich Pablo, 50109

## **Índice:**

<b>Introducción</b>	<b>2</b>
<b>Instrucciones</b>	<b>3</b>
<b>Decisiones de Diseño</b>	<b>3</b>
<b>Descripción de la Gramática</b>	<b>5</b>
<b>Dificultades Encontradas</b>	<b>6</b>
<b>Futuras Extensiones</b>	<b>6</b>

## **Introducción:**

En el siguiente informe se explicará brevemente el uso del programa para analizar secuencias genómicas, así como también las decisiones de diseño más importantes que se utilizaron para su creación, la secuencia genómica que genera cadenas permitidas por el programa, algunos problemas encontrados en la creación del mismo y qué mejoras se le podría hacer para mejorar su rendimiento.

## Instrucciones:

En la carpeta 'root' existe el directorio 'src'. En el mismo se encuentra el archivo *syntax.l* en donde está el código del programa. En la raíz del directorio existe el archivo *test\_file.sh* el cual contiene varias pruebas del programa. Al ejecutarse compila el programa y automáticamente corre varios tests para probar el correcto funcionamiento del programa. Para compilarlo, se utiliza un archivo *makefile* y los tests más importantes se encuentran en el directorio 'examples'.

Dentro del directorio 'examples', se encuentran archivos con cadenas genómicas. Entre los mismos se pueden encontrar las cadenas provistas por la cátedra, así como también otras cadenas obtenidas de una fuente de internet<sup>1</sup>.

Al ejecutar el archivo *test\_file.sh*, se obtiene en el *stdout* las cadenas de aminoácidos traducidas a partir de las cadenas genómicas.

## Decisiones de Diseño:

El programa intenta imitar de forma rudimentaria el diseño de una máquina de estados, donde se analizan la combinación de caracteres que podrían ocurrir en una secuencia genómica válida según el reglamento de la cátedra. Cada caso se podría considerar como un estado y dado ese estado sólo existe un conjunto de posibles caracteres que el programa espera encontrar. Si una secuencia de caracteres no concuerda con lo esperado, entonces como en una máquina de estados, donde se encuentra una entrada que no lleva a ningún estado válido, la cadena es rechazada.

La cadena es leída de a un carácter por vez y el mismo es almacenado en un buffer. Cuando este es completado con tres caracteres, se procede a verificar qué caracteres son y a mostrar el aminoácido correspondiente en el *stdout* dado el caso de que sea una combinación válida.

Para realizar la codificación de cadena de tres genomas a aminoácidos, se pensó en utilizar la estructura de datos llamada *trie*. El mismo tendría en cada una de sus hojas el aminoácido que corresponde a la secuencia que se siguió desde la raíz hasta la misma. Cada nodo tendría un arreglo de punteros a nodos hijo, en donde habría una relación directa entre una letra del abecedario y una posición del arreglo. Bajo estas condiciones la complejidad de la búsqueda siempre sería  $f(n) = 3$ , donde  $n$  es algún lexema válido, ya que sólo habría que recorrer tres nodos antes de llegar a la hoja y conseguir el aminoácido correspondiente al lexema.

---

<sup>1</sup> Parte de la secuencia genómica de una gallina:  
<http://www.ncbi.nlm.nih.gov/nuccore/J04598.1>

Posteriormente, se decidió utilizar una tabla de codificación de lexemas<sup>2</sup>, en donde algunos lexemas quedarían codificados en sólo uno o dos pasos. Si bien ahorarse un paso puede parecer trivial, cuando se tenga que analizar cadenas de longitudes grandes, el tiempo ahorrado por no tener que completar algunos pasos puede resultar en una diferencia significativa.

---

<sup>2</sup> La misma se puede encontrar en:  
[http://en.wikipedia.org/wiki/DNA\\_codon\\_table](http://en.wikipedia.org/wiki/DNA_codon_table)

## Descripción de la Gramática:

Se creó la siguiente gramática para representar las secuencias genómicas aceptadas por el programa:

$G = \langle VN, VT, P, \langle START \rangle \rangle$

$VN = \{ \langle START \rangle, \langle HEADER \rangle, \langle BODY \rangle, \langle END \rangle, \langle CODON \rangle, \langle POSITION \rangle, \langle DISJUNCTION \rangle, \langle BASE \rangle, \langle DIGIT \rangle, \langle D \rangle, \langle BASE\_STREAM \rangle, \langle DISJUNCTION\_PART \rangle, \langle DISJUNCTION\_OPTIONAL \rangle \}$

$VT = \{ A, U, G, T, S, O, P, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, [, ], \{, \}, * \}$

$P = \{$   
     $\langle START \rangle \rightarrow \langle HEADER \rangle \langle BODY \rangle \langle END \rangle$   
    ,  
     $\langle HEADER \rangle \rightarrow AUG$   
    ,  
     $\langle END \rangle \rightarrow STOP \mid TAA \mid TAG \mid TGA$   
    ,  
     $\langle BODY \rangle \rightarrow \langle CODON \rangle \langle BODY \rangle \mid \langle POSITION \rangle \langle BODY \rangle \mid \langle DISJUNCTION \rangle \langle BODY \rangle \mid \lambda$   
    ,  
     $\langle CODON \rangle \rightarrow \langle BASE \rangle \langle BASE \rangle \langle BASE \rangle$   
    ,  
     $\langle BASE \rangle \rightarrow C \mid G \mid A \mid T \mid -$   
    ,  
     $\langle POSITION \rangle \rightarrow [ \langle DIGIT \rangle ]$   
    ,  
     $\langle DIGIT \rangle \rightarrow 1 \langle D \rangle \mid 2 \langle D \rangle \mid 3 \langle D \rangle \mid 4 \langle D \rangle \mid 5 \langle D \rangle \mid 6 \langle D \rangle \mid 7 \langle D \rangle \mid 8 \langle D \rangle \mid 9 \langle D \rangle$   
    ,  
     $\langle D \rangle \rightarrow 0 \langle D \rangle \mid 1 \langle D \rangle \mid 2 \langle D \rangle \mid 3 \langle D \rangle \mid 4 \langle D \rangle \mid 5 \langle D \rangle \mid 6 \langle D \rangle \mid 7 \langle D \rangle \mid 8 \langle D \rangle \mid 9 \langle D \rangle \mid \lambda$   
    ,  
     $\langle DISJUNCTION \rangle \rightarrow \{ \langle BASE\_STREAM \rangle \langle DISJUNCTION\_PART \rangle \}$   
     $\langle DISJUNCTION\_OPTIONAL \rangle$   
    ,  
     $\langle DISJUNCTION\_PART \rangle \rightarrow , \langle BASE\_STREAM \rangle \langle DISJUNCTION\_PART \rangle$   
    ,  
     $\langle DISJUNCTION\_OPTIONAL \rangle \rightarrow * \mid \lambda$   
    ,  
     $\langle BASE\_STREAM \rangle \rightarrow \langle BASE \rangle \langle BASE\_STREAM \rangle \mid \lambda$   
     $\}$

### **Dificultades Encontradas:**

En la creación del programa se encontró sólo una dificultad técnica. Se quiso modularizar el código fuente en varios archivos, ya que sería conveniente que en archivo de *lex* quede sólo el código necesario para hacer el parseo. Sin embargo, no se pudo hacer que el archivo de *lex, syntax.l*, reconociera variables globales declaradas en otros módulos y ya que en la escritura del código se usaron exténsamente estas variables era impráctico rediseñar el programa. El grupo debe averiguar a qué se debe esta limitación de no poder importar variables globales de otros módulos.

### **Futuras Extensiones:**

Actualmente el programa utiliza un buffer en donde se almacenan tres caracteres, que luego se verifican para saber si hay correspondencia entre la secuencia y algún codón. Sería más conveniente que dentro de las reglas de parseo que utiliza *lex*, esté contemplado el análisis de estos codones. De tal forma que al identificar alguno, *lex* imprima directamente el aminoácido correspondiente al *stdout*. Como el código de *lex* está optimizado y no se tiene que manejar un buffer, el análisis de las secuencias sería más rápido.

Realizar este cambio no conlleva una gran complejidad ya que sólo habría que adicionar estas reglas y en los casos especiales, en donde algún codón se encuentra separado, utilizar el buffer. Como la gran mayoría de elementos que componen a la cadena son codones válidos, la rapidez del parseo puede aumentar considerablemente.