

Universidad Nacional Autónoma de México



**Instituto de Investigaciones Biomédicas
Centro de Ciencias Genómicas**



“Predicción de Residuos Funcionalmente Importantes en el dominio de unión al ligando de la familia de factores transcripcionales Crp/Fnr en bacterias”

Tesis que para obtener el grado de:

Licenciatura en Investigación Biomédica Básica

presenta:

Mónica Ivonne Peñaloza Spínola

Asesores de Tesis:

**Dr. Ernesto Pérez Rueda
Dr. Julio Collado Vides**



Universidad Nacional
Autónoma de México

Dirección General de Bibliotecas de la UNAM

Biblioteca Central



UNAM – Dirección General de Bibliotecas
Tesis Digitales
Restricciones de uso

DERECHOS RESERVADOS ©
PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis esta protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

En memoria del abuelo
Alberto Spínola Spínola (1925 – 2006)

En recuerdo del tío abuelo
Ángelo Spínola Spínola (1923 – 2006)

Abuelo, me hubiera gustado tanto compartir
esto y más en vida contigo

A todas aquellas personas que con su compañía en diferentes momentos de mi vida,
me han dado tanto y enseñado otro mucho

Agradecimientos

Definitivamente, este trabajo no lo hubiera concluido en este tiempo y no hubiera logrado concretar los resultados y discusión sin la guía, discusiones, pláticas, ayuda en programación e incansable apoyo de Irma Lozada Chávez. Mil gracias amiga y colega de Ciencia.

En cuanto al desarrollo de los algoritmos y programas principales de este trabajo, le debo un gran agradecimiento a Hely y Bruno. Sin sus conocimientos y larga experiencia en estos menesteres, no lo hubiera logrado.

Por otra parte, la paciencia y compañía de mi familia (José Luís, Lourdes y Pamela) durante el lento y tortuoso desarrollo de esta tesis ha sido un gran apoyo. Sé que aunque no entendían por qué tardaba, nunca duraron que lo lograría (al menos eso quiero pensar). Finalmente pago mi deuda con ustedes.

También la compañía, apoyo y paciencia de uno de los mejores amigos que me ha dado la vida, Christián, ayudaron mucho a que no desfalleciera en mis intentos de terminar este trabajo.

A esta tesis misma, por haberme enseñado mucho sobre lo que se puede lograr al perseverar y luchar por la realización de las ideas.

Demás compañeros de laboratorio (sin orden especial: Vero, Irma M, Albert, Soco, Martín, Sarath, Romualdo, Víctor, Conchita, Juan, Arturo) mil gracias por todos los momentos que hemos compartido, pláticas, reflexiones y desarrollo de los proyectos bioinformáticos y de la vida.

Finalmente, gracias a mis tutores Julio Collado y Ernesto Pérez Rueda; a Julio por haberme dejado realizar este proyecto (sé que no le ve mucha relación con lo demás que hacemos en el laboratorio, pero espero convencerlo que tiene conexión y aplicación directa) y a Ernesto por las recomendaciones, guía, todas las opciones que buscaba para métodos y opciones de realización de esta idea y su paciencia para que terminara este proyecto.

Índice	Página
0. Resumen	1
0.1 Summary	2
1. Introducción	3
1.1. Transcripción celular y regulación transcripcional	3
1.1.1. Elementos de la regulación transcripcional	5
1.1.2. Los factores Sigma	5
1.1.3. Los Factores transcripcionales	6
a) Dominios que conforman a los factores transcripcionales	9
1.2. Clasificación de los factores transcripcionales	10
a) Análisis de secuencia	10
b) Análisis estructural	11
1.3. Familias de reguladores transcripcionales	12
1.3.1. Familia de factores transcripcionales CRP/FNR	13
a) Especificaciones estructurales de Crp	15
b) Especificaciones estructurales de Fnr	16
1.4. Análisis de residuos funcionalmente importantes	20
a) Análisis experimental	20
b) Análisis computacional	21
c) Comparación de métodos computacionales	22
2. Hipótesis	24
3. Material y métodos	24
3.1. Colección de secuencias y estructuras tridimensionales	24
3.2. Detección de los miembros de la familia Crp/Fnr	24
3.3. Análisis evolutivo para la división en subgrupos	25
3.4. Identificación de residuos funcionalmente importantes	26
3.5. Cálculo de significancia estadística	27
a) Obtención de secuencias pseudos-aleatorias	27
b) Variabilidad intra-grupos	28
c) Medición de significancia estadística con valores Z	28
3.6. Caracterización de estructura terciaria	28

4. Resultados y Discusión	29
4.1. Identificación y agrupamiento de miembros de la familia Crp/Fnr	29
4.2. Identificación de residuos funcionalmente importantes	31
4.3. Identificación de la significancia estadística del CI	37
4.3.1. Comparando entre grupos los residuos con <i>CI</i> s significativos, que se predicen como RFIs	44
4.3.2. Residuos con evidencia experimental, ligandos y conformaciones de Crp y Fnr	45
4.3.3. De los residuos con evidencia experimental en Crp _{ECO} y los RFIs predichos para su grupo en <i>γ-Proteobacteria</i>	46
4.3.4. De los residuos con referencia en Fnr _{ECO} y los RFIs predichos para su grupo en <i>γ-Proteobacteria</i>	49
4.4. Eficiencia del método, comparando con lo ya reportado y las características de las proteínas	49
5. Conclusiones generales y perspectivas	50
5.1. Conclusiones generales	50
5.2. Perspectivas	52
6. Referencias	53
7. Anexos	57
7.1. Tabla de las secuencias de la familia Crp/Fnr	58
7.2. Alineamiento de las secuencias de la familia Crp/Fnr	60
7.3. Representación rectangular de árbol filogenético	62
7.4. Gráficas de distribución Normal	63
7.5. Tablas de los valores de <i>CI</i> , con los valores Z de cada grupo de la familia Crp/Fnr	Anexo 7.5: 1-10
7.6. Códigos de los programas en lenguaje PERL	65
7.6.1. ContInfo-aa-alignment-PRO.pl	65
7.6.2. Calc-seqs-PseudoRandom.pl	69

7.6.3. Shuffling_alingment.pl	75
7.6.4. MedsStats.pl	78
7.6.5. Calculo_rho.pl	83
7.6.6. Color_SeqRelation.pl	85

0 Resumen

Actualmente, con la obtención de cientos de miles de secuencias de genes de una gran diversidad de organismos a través de los proyectos genómicos, se hace necesario utilizar estrategias alternativas a las experimentales que permitan avanzar rápidamente en la caracterización funcional de estas secuencias. Dentro de estas estrategias se encuentran las computacionales, que se valen del entendimiento de los procesos biológicos y conocimiento de sus componentes, modelandolos lo más real posible.

Con el uso de estas estrategias computacionales se puede, inicialmente, conocer el tipo de moléculas funcionales a las que codifican. Por ejemplo, RNAs, proteínas, entre otras. Posteriormente, se puede identificar las características propias de estos productos que los definen y los hacen específicos a sus funciones; por ejemplo: las proteínas están conformadas por cierto tipo de aminoácidos que les confieren una estructura específica, la capacidad de reconocer a sus sustratos, la capacidad de realizar las reacciones metabólicas que transforman a los sustratos, etc. Uno de los resultados finales con el uso de estas estrategias es que se logra conocer la función de los productos de las secuencias, sin la necesidad de hacer los experimentos exhaustivos individualmente para cada gen.

El trabajo que se presenta en esta tesis es el avance en la caracterización de una familia de proteínas que está involucrada en la expresión genética, la familia de factores transcripcionales Crp/Fnr. El objetivo principal radica en la predicción de residuos funcionalmente importantes (RFIs) en su dominio de unión al ligando.

Los resultados obtenidos aquí se basan en una estrategia que considera el contenido informacional, y su análisis estadístico para todas las posiciones del alineamiento de estas secuencias. Los resultados muestran diferentes RFIs para los grupos analizados, unos ya caracterizados experimentalmente y otros que aún no se caracterizan experimentalmente; estos últimos se proponen que podrían ser funcionalmente indispensables para los factores transcripcionales. Estos resultados han permitido identificar, dentro de las secuencias y estructuras, varios residuos que se proponen como importantes para las interacciones con los ligandos y/o entre subunidades, mantenimiento del plegamiento, y otras funciones que hacen común la respuesta a la regulación transcripcional de estos factores en diferentes organismos.

0.1 Summary

Currently, after obtaining hundreds of thousands of gene sequences from a great diversity of organisms through genomic projects, it has become necessary to use alternative strategies to the experimental ones, in order to quickly advance in the functional characterization of such amount of sequences. Within the alternative strategies are the computational ones, which considering the understanding of the biological processes and the knowledge of their components, model them the realiest as possible.

With the use of these computational strategies it is possible, first, to get to know the type of functional molecules that they code for (e.g. RNAs, proteins, etc.). Later, it can be identified the properties that define and make specific, these products, for their functions (e.g. the type of amino acids that conform a protein confer a specific folding, a specific substrate recognition, the ability to do metabolic reactions to transform substrates, etc.). One of the final results of these strategies is to know the function of the sequence products, without doing exhaustive individual experiments.

The current work is an advance in the characterization of a protein family that is involved in the transcriptional regulation, the Crp/Fnr family. The main objective has been to predict the 'residues functionally important' (RFI) in their ligand binding domain.

The results here obtained are based on a strategy that considers the informational content of each position in the alignment of the domain sequences and its statistical analysis. There were identified different RFIs for the analyzed groups of sequences, ones of these RFIs are already characterized, but other RFIs remain to be experimentally characterized. These last RFIs are proposed to be functionally indispensable for the domain fold or structure, for different ligand and/or intersubunits interactions, or for other functions that allow the regulatory response of these transcription factors.

1 Introducción

1.1 Transcripción celular y regulación transcripcional

Todos los organismos tienen su información biológica codificada en tres tipos de moléculas básicas, DNA, RNA y proteínas. Cada tipo de molécula lleva a cabo funciones distintas y que se complementan entre sí. El DNA almacena la información genética que codifica para todas las funciones y características del organismo; el RNA es el transmisor y regulador de la información entre el DNA y las proteínas; y las proteínas son las ejecutoras de las funciones estructurales y de catálisis que la célula lleva a cabo, así como los vehículos de las reacciones metabólicas, las que dan estructura, transmiten señales, transportan, entre otras (Lewin 2000).

Durante la vida de la célula, se pasa de un tipo de molécula a otro para obtener diferentes tipos de moléculas funcionales a través de varios procesos biológicos que se representan en el Dogma Central de Biología Molecular (DCBM) (Crick 1970).

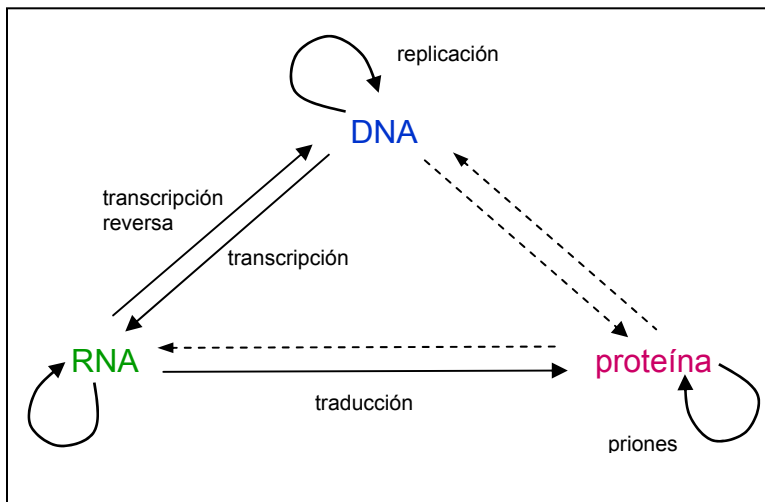


Figura 1. Dogma central de la biología molecular. Representa los nueve flujos de información biológica entre los tres diferentes tipos de moléculas básicas de la vida (DNA, RNA y proteínas), tanto los conocidos (líneas sólidas) como los que aún no se han descrito o inexistentes (líneas punteadas). Los procesos conocidos son: *Replicación DNA* ($DNA \rightarrow DNA$), *Transcripción DNA* ($DNA \rightarrow RNA$), *Transcripción reversa* ($RNA \rightarrow DNA$), *Traducción* ($RNA \rightarrow proteína$), *Replicación RNA* ($RNA \rightarrow RNA$) y los recientemente descritos *priones* ($proteína \rightarrow proteína$).

Uno de los flujos del DCBM es la transcripción (figura 2), la cual se define como el proceso mediante el cual los genes, codificados en el DNA, son leídos por la enzima RNA polimerasa (y co-factores) para transcribir su código de DNA a RNA; y así poder ser procesados para obtener de ellos los productos funcionales para la célula (diferentes tipos de RNA, proteínas estructurales, enzimas, transportadores, etc.) (Lewin 2000).

Ya que el paso de información de DNA a RNA es aquel mediante el cual se obtienen los productos necesarios para el resto de los procesos celulares, es propio de las células regular este proceso para coordinar las diferentes tareas a realizar. Así, se entiende a la regulación transcripcional como “*la coordinación del proceso de la expresión genética controlando cuánto y cuándo un gen debe ser transcrito.*” Esta regulación ocurre en diferentes etapas del proceso de transcripción y con diferentes mecanismos (etapas de transcripción y regulación, en figura 2). De ésta, de la que más se conoce es la regulación al inicio de la transcripción, comúnmente llamada ‘regulación transcripcional’. Esta regulación es muy directa y eficiente, ya que permite el ahorro de energía, tiempo y recursos antes de empezar el proceso de transcripción, así como también hace la conexión directa con diversos tipos de señales del metabolismo (Lewin 2000).

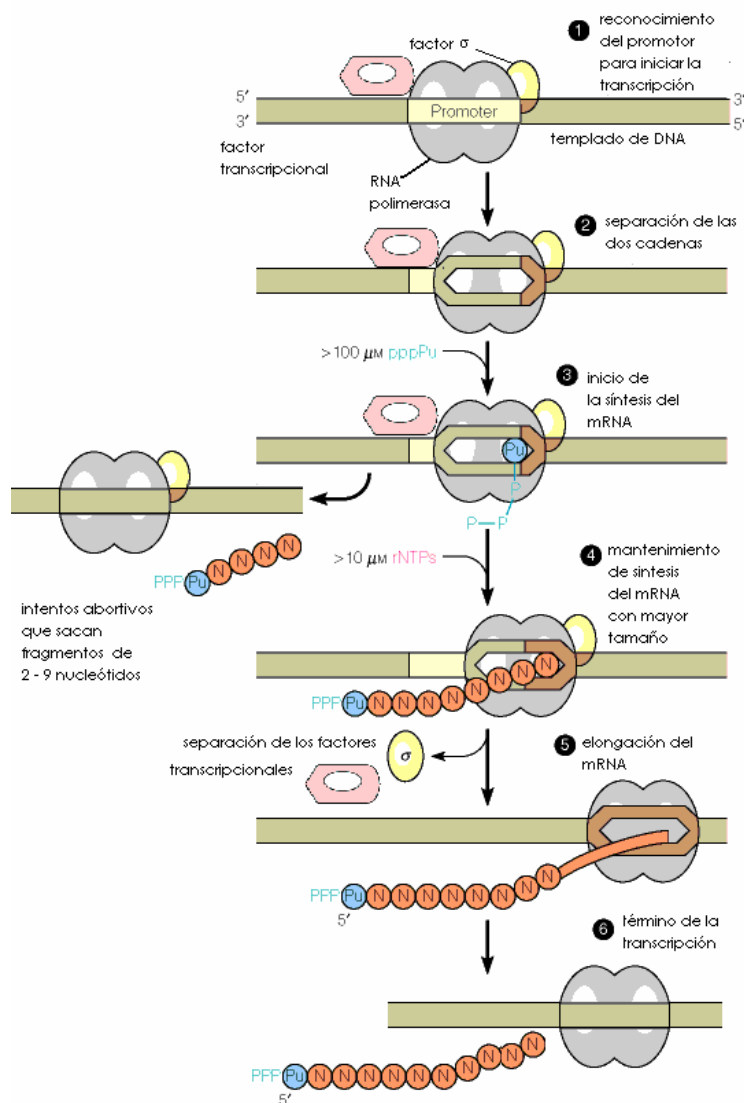


Figura 2. **Pasos del proceso de transcripción:**

1) Reconocimiento del sitio de inicio de la transcripción, con ayuda de factores sigma y factores transcripcionales, 2) separación de las dos cadenas de DNA, 3) inicios cortos de transcripción con intentos abortivos y desprendimiento de la RNA-polimerasa de los factores, 4,5) elongación del transcrito, 6) termino de la transcripción

Etapas en las que ocurre la regulación transcripcional:

a) al inicio (etapa 1), b) posterior al inicio (etapa 2), c) durante la elongación (etapa 4), d) posterior al término (posterior a etapa 6).

Las formas en que ocurre la regulación son:

- *inhibición/promoción a iniciar el proceso o para mantenerlo ocurriendo,*
- *disminución de la velocidad con que está ocurriendo,*
- *bloqueo del proceso* (Lewin 2000; Browning and Busby 2004).

La regulación transcripcional coordina la transcripción de los genes cuyos productos son necesarios para los procesos metabólicos que la célula está realizando en ese momento, mientras que también inhibe la transcripción de los genes cuyos productos son innecesarios en ese momento fisiológico (Lewin 2000; Browning and Busby 2004).

1.1.1 Elementos de la regulación transcripcional

La regulación transcripcional es llevada a cabo por proteínas que interaccionan entre ellas, con la RNA-polimerasa y/o con regiones del DNA, llamadas comúnmente reguladores transcripcionales y factores sigma (ver figura 3) (Lewin 2000; Browning and Busby 2004).

1.1.2 Los factores Sigma

Los factores sigma son una subunidad proteínica temporal de la RNA polimerasa (RNA-p) que aumentan la afinidad de la interacción RNA-p con la región del promotor del DNA en 10^7 veces (figura 3) (Perez-Rueda 1999) y promueven que la RNA-p abra las dos cadenas de DNA para que inicie la transcripción (Burgess and Anthony 2001). Para ello, los factores sigma tienen dos regiones, una donde reconocen a los promotores y otra donde interaccionan con la RNA-p y los factores transcripcionales (los cuales intervienen en la interacción de la RNA-p con el DNA). Sin los factores sigma, la RNA-p se une y despega del DNA en cualquier parte y no transcribe (Burgess, Travers et al. 1969).

A su vez, los factores sigma funcionan como reguladores modulares de la transcripción, pues hay diferentes factores sigma para diferentes condiciones globales de la célula. De los factores y sus mecanismos que más se han estudiado en procariontes, se han caracterizado diecisiete factores sigma en *Bacillus subtilis* (Haldenwang 1995; Moreno-Campuzano, Janga et al. 2006) y siete en *Escherichia coli* (Gruber and Gross 2003).

1.1.3 Los Factores transcripcionales

En la etapa de reconocimiento del DNA por la RNA-p además intervienen las proteínas que son factores transcripcionales (FT), también llamados 'reguladores transcripcionales'. Estos factores controlan la transcripción de manera más fina que los factores sigma, modulan la actividad transcripcional de los diferentes procesos celulares en respuesta a señales celulares específicas (ver más adelante en pág. 7). Activan o reprimen la expresión de los genes de acuerdo a que tanto se requieren sus productos para los diferentes procesos celulares (figura 3) (Browning and Busby 2004).

Los FT mejor caracterizados en procariontes son los de *E. coli* (Perez-Rueda and Collado-Vides 2000) y *B. subtilis* (Moreno-Campuzano, Janga et al. 2006); mientras que para el resto de los procariontes los estudios se realizan por comparaciones con FTs de estos y otros organismos (Mironov, Koonin et al. 1999; Makarova, Mironov et al. 2001; Tan, Moreno-Hagelsieb et al. 2001; Rodionov, Vitreschak et al. 2004). Para *E. coli* se ha caracterizado aproximadamente el 30% de su red de regulación (Salgado, Gama-Castro et al. 2006); con esto se sabe que del total de sus genes cerca del siete por ciento codifican para FTs¹ (314 genes son FT) (Perez-Rueda and Collado-Vides 2000), de los cuales 163 han sido determinados experimentalmente y 151 computacionalmente (RegulonDB [<http://regulondb.ccg.unam.mx/>]). A su vez, del primer conjunto de FTs, se ha determinado que el 24% son activadores, 46% son represores y 20% son duales [activan/reprimen] (el resto son de efecto desconocido) (Perez-Rueda and Collado-Vides 2000).

¹ El porcentaje de genes de un genoma que codifican para factores transcripcionales correlaciona con la observación de Pérez-Rueda et.al 2003, de que un organismo de ambientes y modos de vida versátiles tiene mayor cantidad de FTs (7-10% de sus genes, e.g. *E. coli* misma), mientras que un organismo que sólo está adecuado a uno sólo o pocos habitats tiene menos TFs (2-3%, e.g. *H. pylori*) (Perez-Rueda E and Collado-Vides J 2004).

Tabla 1. Ejemplos de factores transcripcionales de *E. coli*.

Factor transcripcional	Efector	Efecto en la transcripción	Genes – tipo de procesos regulados
OmpR	Fosfato	Activa	Principales proteínas de membrana, en respuesta a cambios de osmolaridad
MaiT	Maltotriosa, ATP	Activa	Metabolismo de maltosa
LacI	Alolactosa	Reprime	Metabolismo de lactosa
PhoP	Fosfato	Reprime	Respuesta a disponibilidad de Magnesio en el medio
TrpR	Triptofano	Reprime	Biosíntesis de triptófano
AraR	Arsenito	Reprime	Resistencia a arsénico
GalR	Galactosa	Dual	Metabolismo de galactosa
CRP	AMP-cíclico	Dual	Represión catabólica
Lrp	Leucina	Dual	Metabolismo de leucina

Fuente: RegulonDB (<http://regulondb.ccg.unam.mx/index.html>)

Los FTs realizan su función al interactuar con el DNA en sus “sitios de unión” u “operadores”, ubicados alrededor del inicio de la transcripción (entre las posiciones -200 y +40 pares de bases alrededor del inicio de la transcripción) (Madan Babu and Teichmann 2003). El reconocimiento de los operadores por parte de los FT se da sin abrir las cadenas de DNA, uniéndose a ellas por medio de puentes de hidrógeno, y suele ocurrir en el surco mayor del DNA (donde están las bases que dan la especificidad de reconocimiento) (Harrison 1991; Suzuki and Yagi 1996).

Las señales celulares a las que responden los FT pueden ser: a) sustratos o productos de las mismas vías metabólicas cuyos genes regulan ellos mismos, por ejemplo metabolitos como: vitaminas, aminoácidos, azúcares, etc (Baker, Tomlinson et al. 2001; Danot 2001; Lewis 2005); b) cambios ambientales externos como cambios de pH, de osmolaridad, de oxigenación, etc. (Gunsalus and Park 1994); y c) diferentes topologías del DNA, tales como la curvatura o su superhelicidad (Yasuzawa, Hayashi et al. 1992).

Estas señales lo que hacen es cambiar la conformación o el estado químico (óxido-reducción) de los FT, que los inducen a unirse o despegarse del DNA para hacer su efecto de regulación. Actualmente se sabe que cerca de la mitad de estas señales para modular a los FT son moléculas pequeñas (Madan Babu and Teichmann 2003). A estas señales de los TFs se les llama comúnmente ligandos o efectores.

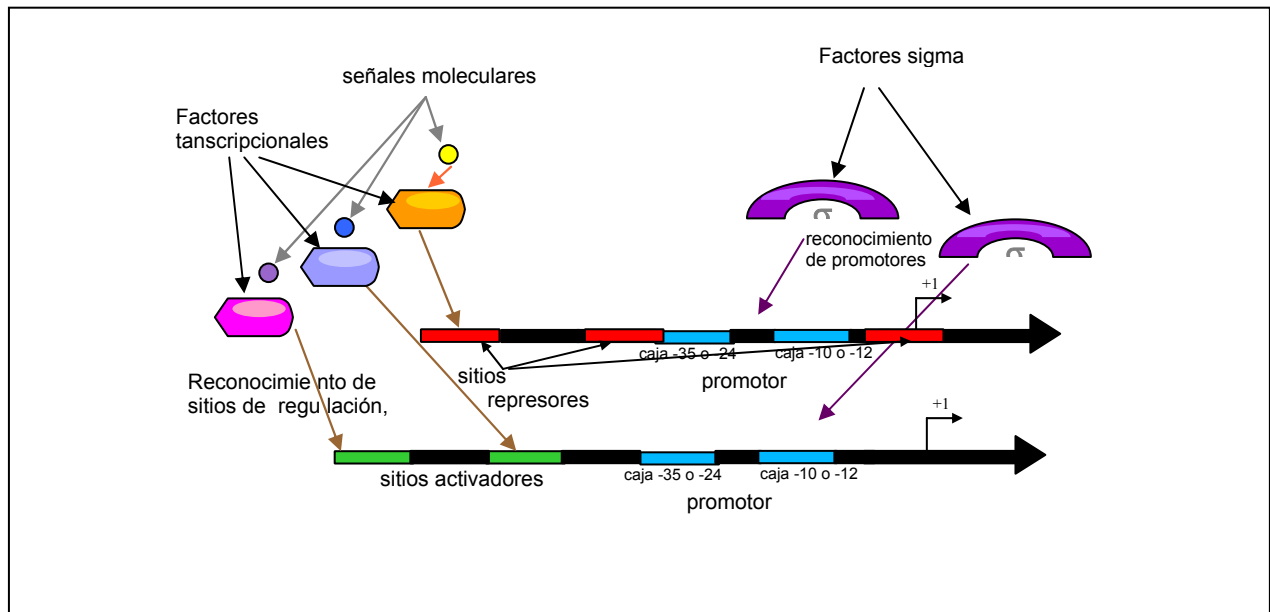


Figura 3. **Representación de ejemplo de los elementos de la regulación transcripcional junto a sus elementos en el DNA.** Los *factores sigma* interaccionan con las cajas de los *promotores* (localizadas en distancias centradas a -35 o -24 y -10 o -12 bases arriba del inicio de la transcripción +1) para promover que la RNA-polimerasa inicie la transcripción; mientras que los *factores transcripcionales*, en respuesta a diferentes *señales moleculares*, interaccionan con los *operadores* o *sitios de regulación* para facilitar o impedir que la RNA-polimerasa se una al DNA y/o inicie la transcripción.

a) Dominios que conforman a los factores transcripcionales

Los FT generalmente están conformados por dos dominios², uno de reconocimiento al DNA y el otro para unir al ligando o interactuar con otras moléculas, incluyendo subunidades de la misma proteína u otras proteínas (Madan Babu and Teichmann 2003).

De los dominios de unión al DNA (DNA-binding domain, DBD) de los FT se tienen identificados cinco tipos de motivos (ver tabla 2), de los cuales el más común en bacterias es el tipo hélice-vuelta-hélice (helix-turn-helix, HTH) (Perez-Rueda and Collado-Vides 2001; Perez-Rueda, Collado-Vides et al. 2004). El HTH consiste de aproximadamente veinte residuos, se divide en dos α -hélices unidas por una vuelta de 3 o 5 residuos que forma entre ellas un ángulo de 120° (Brennan and Matthews 1989).

Tabla 2. Estructuras de dominios de unión al DNA descritas para los factores transcripcionales.

Dominio de unión a DNA	Características	Distribución
Hélice vuelta hélice (helix-turn-helix, HTH)	Formada por dos α -hélices de aproximadamente 10 residuos cada una, conectadas por una "vuelta" de alrededor 3 a 5 residuos de aminoácidos. La hélice C-terminal es la que da la especificidad del dominio uniéndose al surco mayor del DNA.	Procariontes, Eucariontes (modificada)
Dedos de Zinc	Consta de α -hélices y β -plegadas unidas por uno o más átomos de Zinc. Las α -hélices son las que interactúan con el DNA.	Eucariontes, Arqueas
β -Plegada antiparalela	Conformada por tres o cuatro hojas β -plegadas, que interactúan con el surco mayor del DNA.	Bacterias
Zipper de leucina	Formada por dos α -hélices que interactúan entre sí a través de varios residuos de leucina.	Eucariontes
Hélice-asa-hélice (helix-loop-helix, HLH)	Consiste en una α -hélice corta conectada por un asa a una segunda α -hélice larga.	Eucariontes

Referencia: Harrison 1991.

² En este trabajo, se toma la definición de dominio como la unidad mínima independiente estructural y funcionalmente (Lewin 2000).

Tipos de dominios presentes en los FT: El 90% de los FT están conformados por un dominio de unión a DNA (DNA Binding Domain, DBD); aproximadamente el 75% de los FT tienen además un dominio extra donde interactúan con su efector, entre sus subunidades o con otras proteínas (Dominio para Otras Interacciones, DOI); cerca del 12% de los FT tienen dos DOI; y alrededor del 3% tienen tres DOI (Madan Babu and Teichmann 2003).

De los dominios para otras interacciones (DOI) de los FT, tanto su función como estructura varían; su estudio y clasificación no es tan específico con los DBD, pero se pueden agrupar en cuatro tipos conocidos: a) de unión a moléculas pequeñas, b) tipo enzima, c) interacción con otras proteínas, d) receptor, y el resto siguen siendo desconocidos (Madan Babu and Teichmann 2003).

1.2 Clasificación de los factores transcripcionales

En general, las proteínas se pueden clasificar y analizar basándose en su secuencia, estructura y función (Holm 1998). Esta clasificación genera agrupamientos, de los cuales a su vez, se generan inferencias evolutivas, funcionales, de localización celular, de conformaciones, de estructuras (Mount 2001).

El estudio de secuencias, del que se comienza cualquier análisis de genes o proteínas, permite hacer el resto de las inferencias (Wilson, Kreychman et al. 2000). El estudio de estructuras usualmente apoya el análisis de secuencia para hacer las inferencias funcionales y de localización celular (Li, Jaroszewski et al. 2002).

a) Análisis de secuencia

Este tipo de análisis se realiza con la comparación de secuencias de genes o proteínas. Para esto se han diseñado diferentes algoritmos, entre los más utilizados se encuentran FASTA (Pearson 2000), BLAST (psi-blast o blast) (Altschul, Gish et al. 1990) que hacen búsquedas de similitud de secuencias a partir del alineamiento de las mismas; y los Hidden Markov Models (HMMs) (Eddy 1998) que utilizan cadenas de Markov que toman en cuenta la memoria de las posiciones anterior y posterior de un perfil de secuencias. Las comparaciones pueden ser entre un par de secuencias (alineamientos pareados) o varias secuencias a la vez (alineamientos múltiples); y las identidades entre secuencias se consideran de dos maneras, a lo largo de toda la secuencia (global), o por regiones o dominios (local) (Mount 2001).

Las secuencias de proteínas se relacionan en grupos de acuerdo a su porcentaje de similitud. Si las secuencias son muy similares, desde el 100% hasta un 40% de

similitud, es altamente probable que las proteínas conserven la misma conformación estructural y conforme más similares sean las proteínas será más probable que realicen las mismas funciones bioquímicas (Wilson, Kreychman et al. 2000); si la similitud de secuencia disminuye hasta un 25%, estructuralmente son muy parecidas pero sus funciones ya son variantes (Chothia and Lesk 1986; Wilson, Kreychman et al. 2000).

Actualmente las secuencias de genes y proteínas se pueden consultar y obtener de bases de datos como GeneBank (Benson, Karsch-Mizrachi et al. 2003), SWISSProt (Boeckmann, Bairoch et al. 2003) y PIR (Protein Information Resource) (Barker, Garavelli et al. 1999). Por otra parte, los agrupamientos de secuencias de proteínas en familias, basados en sus alineamientos y dominios se pueden consultar y obtener de la base de datos Pfam (Sonnhammer, Eddy et al. 1998). Pfam utiliza los algoritmos de BLAST para toda la secuencia, y HMM para los dominios.

b) Análisis estructural

En este otro tipo de análisis, las proteínas se relacionan por la similitud de sus estructuras secundarias y/o terciarias de sus dominios y de la proteína completa. Esta similitud se determina por el ordenamiento de los dominios y mediciones de distancias entre estructuras, principalmente. Este análisis también, se apoya en la similitud de secuencias (Li, Jaroszewski et al. 2002).

Actualmente, las estructuras tridimensionales se pueden consultar y obtener del repositorio Protein Data Bank (PDB) (Berman, Westbrook et al. 2000). Aparte, los agrupamientos y clasificaciones estructurales se tienen en varias bases de datos: SCOP, clasificación evolutiva y estructural (Murzin, Brenner et al. 1995); CATH, clasificación jerárquica de estructuras de dominios (Pearl, Bennett et al. 2003); y Superfamily, clasificación estructural de secuencias de proteínas en familias evolutivas (Gough 2002).

La clasificación de los FT se basa principalmente en la comparación de sus secuencias, divididos en dominios. A la fecha, hay dos grandes estudios para la clasificación de los

TF: uno que se basa en análisis de los dominios de unión al DNA (DBD), considerando una similitud de secuencia mínima de 25%, del cual deriva su clasificación en 75 familias evolutivas (Perez-Rueda and Collado-Vides 2000); y el otro que además de las secuencias de los DBD considera las estructuras secundarias y terciarias de estos dominios, y los clasifica en 11 familias estructurales (Madan Babu and Teichmann 2003).

1.3 Familias de reguladores transcripcionales

Los FT que se agrupan en familias relacionándose por similitud de secuencias, con un 25% de identidad mínimo, se observa que tienen: misma ubicación del DBD (región carboxilo o amino), poseen un grupo de arquitecturas comunes, mecanismos de acción comunes, tamaños de secuencia de entre los 150 hasta los 400 aminoácidos, regulación de procesos similares, e inclusive la tendencia a comportarse como activadores y/o represores (Perez-Rueda and Collado-Vides 2000; Madan Babu and Teichmann 2003). Algunos ejemplos de familias de TFs se muestran en la tabla 4.

Tabla 4. Ejemplos de familias de FT, con las características comunes más relevantes.

Familia	tipo de procesos regulados	efecto en la regulación	ubicación de su DBD
AraC/Xyls	Metabolismo de carbono, respuesta a estrés y patogénesis.	Activación/Represión	Carboxilo
LysR	Biosíntesis de aminoácidos	Activación/Represión	Amino
MerR	Estrés oxidativo, presencia de metales pesados y antibióticos	Activación	Amino
GalR/LacI	Metabolismo de azúcares alternativos.	Represión	Amino
GntR	Metabolismo de carbonos y ácidos grasos, biosíntesis de aminoácidos y varios más.	Represión	Amino

Fuente y referencias en: RegulonDB (<http://regulondb.ccg.unam.mx/>)

1.3.1 Familia de factores transcripcionales CRP/FNR

La familia Crp/Fnr³ está representada por los dos reguladores globales que le dan nombre a la familia (Shaw, Rice et al. 1983). Crp es un FT de respuesta global que le permite a la célula contender con la falta de glucosa facilitando la activación de vías alternativas para alimentarse con otros carbohidratos (fenómeno conocido como represión catabólica) (Zubay, Schwartz et al. 1970); y Fnr es otro FT de respuesta global que le permite a la célula responder rápidamente a la ausencia de oxígeno re-dirigiendo su vía de respiración aeróbica a anaeróbica (Spiro and Guest 1990).

Debido a las importantes repercusiones que tienen los efectos de la regulación por Crp y Fnr, los integrantes de esta la familia han sido tema de diversos estudios desde su descubrimiento (Li, Wing et al. 1998; Green, Scott et al. 2001; Korner, Sofia et al. 2003). A la fecha, estos estudios se han enfocado mucho en la caracterización de las proteínas Crp y Fnr. Se han identificado nuevos homólogos de Crp, Fnr y de los nuevos miembros de la familia (Perez-Rueda and Collado-Vides 2000; Korner, Sofia et al. 2003). Se tiene la estructura tridimensional de Crp y de otros dos miembros de la familia, CooA y PrfA. Se sabe que los miembros de esta familia además de regular la represión catabólica y la respuesta a anaerobiosis, regulan otro tipo de procesos como metabolismo de aminoácidos, plegamiento de proteínas, producción de toxinas y *pilli*, respuestas a estrés oxidativo o por óxido Nítrico (NO), fijación de nitrógeno, fotosíntesis y factores de virulencia (Korner, Sofia et al. 2003). Son subdivididos en varios subgrupos: **ArcR**, **CooA**, **CprK**, **Dnr**, **FixK**, **Flp**, **FnrN**, **MalR**, **NnrR**, **NtcA**, **PrfA** y **YeiL**, propuestos por Körner *et. al* 2003 (Korner, Sofia et al. 2003). Para regular a los diferentes procesos los TFs de la familia reconocen diferentes ligandos (figura 4); sin embargo, aún no se caracterizan por completo estos ligandos a los que responde cada subgrupo, ni que determina la

³ *Antecedentes de los estudios de la familia Crp/Fnr.*

*A principios de los 70's se estudiaba el primer sistema de control positivo de la regulación, **CRP** (cAMP Receptor Protein) de *E. coli*. CRP se estudiaba como un FT activador único en su tipo (Zubay *et. al* 1970) y como el responsable de la "represión catabólica". Al mismo tiempo se estudiaba al FT **FNR** (Fumarate and Nitrate-reductase Regulator) de *E. coli* como el primer sistema de control de respuesta a cambios en la concentración de oxígeno (Lambden and Guestl 1976). Años después (1983), se descubrió la homología de CRP con FNR (como parálogos) (Shaw, Rice *et. al* 1983), lo que derivó en considerarlos como una sola familia, **CRP/FNR**.*

especificidad para su reconocimiento.

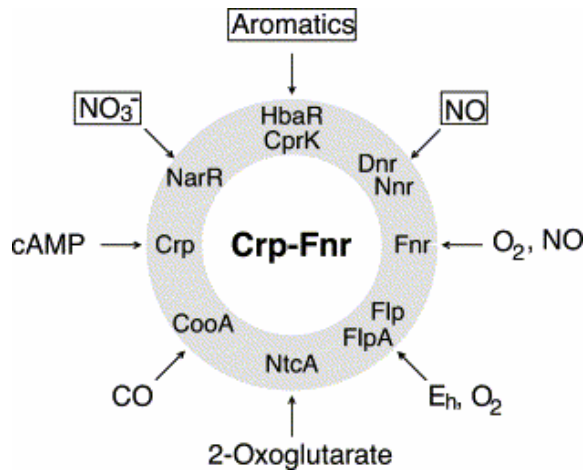


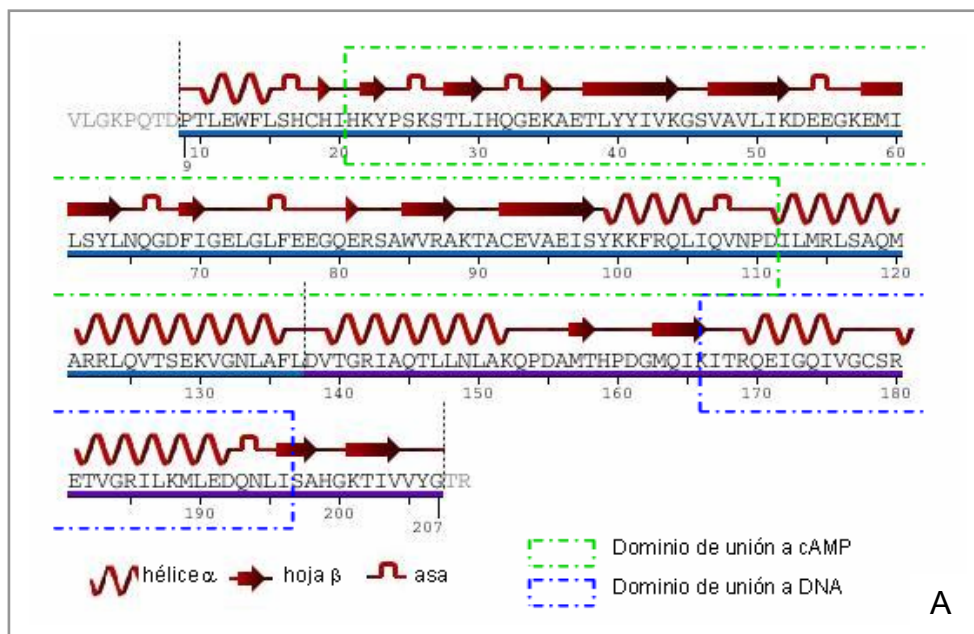
Figura 4. **Ligandos conocidos para los subgrupos de la familia Crp/Fnr**, definidos por Körner *et. al* 2003 (Figura tomada del mismo artículo).

Las características más relevantes de la familia Crp/Fnr son que:

- las proteínas son de tamaños similares (entre 230 – 250 aminoácidos) (Korner, Sofia et al. 2003);
- su DBD está ubicado en el carboxilo-terminal, es de aproximadamente 60 aminoácidos (a.a.) y su motivo es tipo hélice-vuelta-hélice alado, wHTH (wing-HTH, una variante de HTH un poco más abierto en la vuelta con lo que las hélices están un poco más separadas) (Perez-Rueda, Collado-Vides et al. 2004);
- su DOI está ubicado en la región amino, es de aproximadamente 170 a.a. y está clasificado estructuralmente como “cNMP-binding domain” (dominio de unión a un nucleótido mono-fosfatado cíclico) según Pfam (Sonnhammer, Eddy et al. 1998), aunque cada subgrupo reconoce diferentes ligandos (Korner, Sofia et al. 2003);
- tienden a reconocer operadores localizados entre las posiciones -40 o -44 arriba del inicio del gene; estos operadores son cajas palindrómicas perfectas o casi-perfectas de 5 o 6 nucleótidos separados por 4 u 8 nucleótidos; y los operadores aunque son de secuencias similares varían de acuerdo al FT que lo reconoce (Korner, Sofia et al. 2003).

a) Especificaciones estructurales de Crp

Crp de *E. coli* es una proteína de 210 aminoácidos, conformada por un DBD de tipo wHTH en la región de los residuos 166 al 197 y un dominio para unir AMPc (Adenina Mono-Fosfato cíclico) en la región de los residuos 21 al 112 (fig. 4a). La estructura secundaria de Crp combina α -hélices y hojas β -plegadas intercaladas (fig. 4.A). En la estructura terciaria se observa la región amino con una cavidad en el centro donde entra el AMPc y la región carboxilo con las dos α -hélices del HTH que contactan al DNA (fig. 4.C) (Lawson, Swigon et al. 2004). Crp normalmente se encuentra formando dimeros y sólo cuando tiene unida una molécula de AMPc en cada subunidad se modifica su conformación para unirse al DNA y llevar a cabo la regulación transcripcional (Garges and Adhya 1988) (fig. 4.B). A la fecha se han identificado varios residuos involucrados de diferentes formas en la función de la proteína, por mutagénesis dirigida (seleccionados por observaciones posteriores a la obtención y análisis de la estructura de Crp): Lys52 (K52), Asp53 (D53), Ser62 (S62), Glu72 (E72), Arg82 (R82), Tir99 (Y99), Leu124 (L124), Treo127 (T127) y Ser128 (S128) (Weber and Steitz 1987; Heyduk, Heyduk et al. 1992; Belduz, Lee et al. 1993; Lee, Glasgow et al. 1994; Baker, Tomlinson et al. 2001; Lin, Kovac et al. 2002; Gekko, Obu et al. 2004; Lawson, Swigon et al. 2004). En la figura 4.C se indican estos residuos dentro de la estructura de la proteína y su función se describe en Resultados.



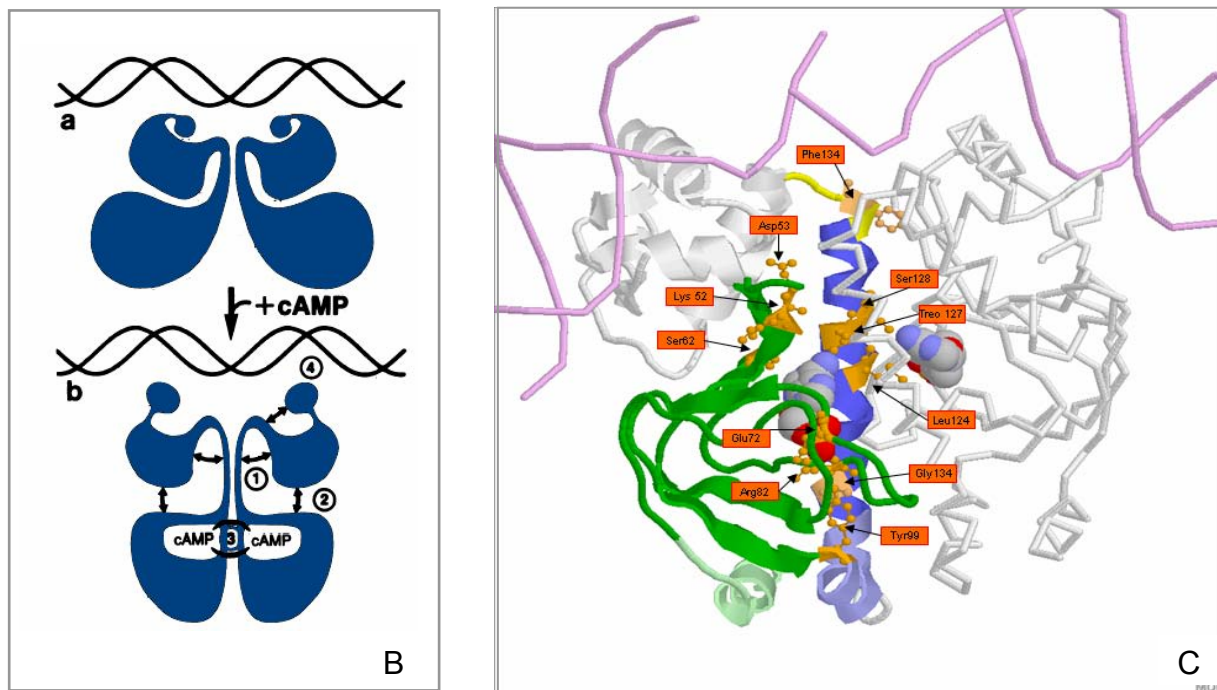
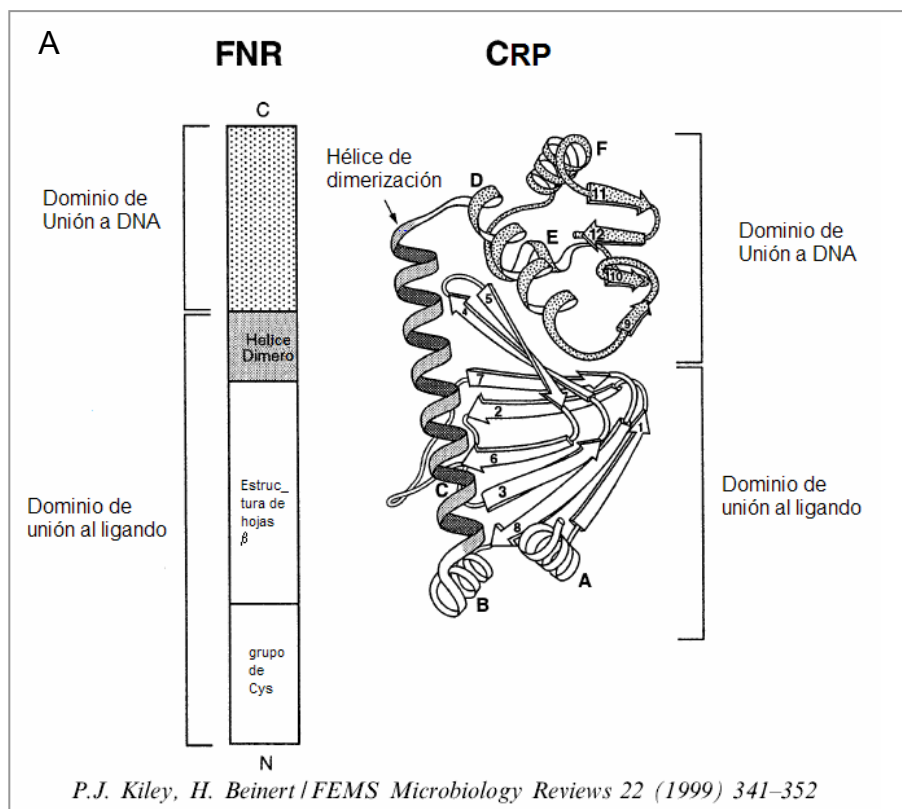


Figura 4. A) **Secuencia y estructura secundaria de CRP de *E. coli***, se indican las hélices α , hojas β y asas, junto con los dominios de unión a AMPc y a DNA (obtenida del Protein Data Bank, en la sección de “sequence details”). B) **Representación del cambio conformacional de CRP al unirse a AMPc**, resultado de esta unión se llevan a cabo 4 contactos con aminoácidos en diversas regiones de los dominios DUL y DBD, facilitando así la unión al DNA (figura tomada de Garge et. al 1988). C) **Estructura tridimensional de CRP** con los residuos probados experimentalmente con mutagénesis: Lys52 (K52), Asp53 (D53), Ser62 (S62), Glu72 (E72), Arg82 (R82), Tir99 (Y99), Leu124 (L124), Trea127 (T127) y Ser128 (S128). Cadena A monómero de Crp indicando: región N-terminal \rightarrow verde claro, dominio de unión al ligando \rightarrow verde, α -hélice de interacción inter-subunidades \rightarrow azul, región bisagra \rightarrow amarillo, dominio de unión al DNA \rightarrow blanco; cadena B Crp en blanco; cadenas de DNA en lila. (figura generada con el software *Protein Explorer 2.8* [www.proteinexplorer.org]).

b) Especificaciones estructurales de Fnr

Fnr de *E. coli* es una proteína de 250 aminoácidos, conformada por al menos dos dominios estructurales: un DBD y uno de unión al ligando tipo cNMP. Fnr en vez de unir algún cNMP, más bien interacciona con el grupo prostético Hierro-Azufre 4Fe4S por medio de cuatro cisteínas ya identificadas (Cys22, 23, 26 y 122). Para Fnr aún no se ha obtenido el cristal de su estructura tridimensional, pero se han hecho comparaciones de su secuencia con la estructura de Crp para aproximar la ubicación las diferentes

regiones de los dominios de la proteína (figura 5.A) (Kiley and Beinert 1999). En este trabajo se predijo la estructura secundaria de Fnr con *PsiPred 2.5* para poder ubicar espacialmente sus residuos funcionalmente importantes (Jones 1999; McGuffin, Bryson et al. 2000) (fig. 5.B). Se propone que Fnr sigue un arreglo estructural muy similar al de Crp, pero en la parte Amino-terminal FNR presenta una región adicional de ~35 a.a (figura 5.A y 5.B). Esta región es importante porque justo ahí se encuentran tres de las cuatro cisteínas necesarias para interaccionar con su grupo prostético 4Fe4S. La estructura de esta región se definió con estructura secundaria en forma de 'asa' con un bajo nivel de confianza (fig. 5.B), así que es muy probable que se acomode de forma distinta (pero para definir correctamente esta región hace falta información estructural de proteínas similares a su tipo).



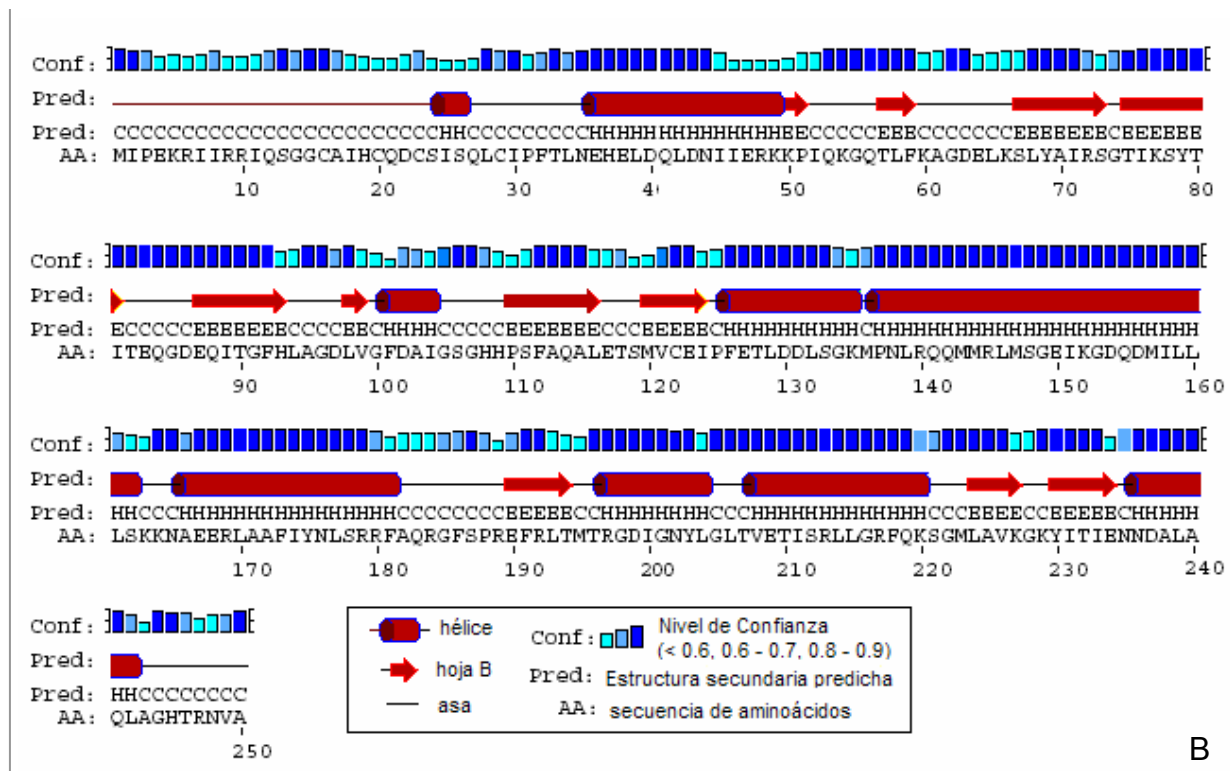


Figura 5. A) **Comparación esquemática de Fnr con la estructura de Crp.** En el esquema de Fnr se anotan las diferentes regiones con su correspondencia en la estructura de Crp en diferentes tonos de gris y blanco. Figura extraída de Kiley *et. al* 1999, sólo se tradujo el texto correspondiente. B) **Predicción de estructura secundaria de FNR de E. coli**, realizada con PsiPred. En filas, la posición en la secuencia, la predicción de estructura y la secuencia de la proteína. El nivel de confianza está en barras azules: 0.9 – 0.8 azul fuerte, 0.7 – 0.6 azul y > 0.6 azul claro; el tipo de estructura: asa (C), hoja β (E) y α hélice (H).

El grupo prostético 4Fe4S siempre está unido a cada subunidad de Fnr provocando que Fnr se acople/desacople al DNA según su oxidación-reducción en respuesta a los niveles de Oxígeno (Lazazzera, Beinert et al. 1996; Khoroshilova, Popescu et al. 1997). Cuando aumenta mucho la concentración de oxígeno, el grupo prostético se desacopla de Fnr, y en cuanto la célula detecta esto degrada a estas moléculas de Fnr-apo (Fnr sin 4Fe4S) (Dibden and Green 2005). Fnr sólo se acopla como dímero cuando se une al DNA, el resto del tiempo está como monómero (Kiley and Beinert 1999).

Con estos dos casos representativos de la familia, se expone que aunque sean de la misma familia e interaccionen con el DNA para regular la transcripción, estos factores reconocen ligandos diferentes, los residuos que intervienen en las interacciones, transmisión de señales, reconocimiento del ligando, acoplamiento con el DNA, etc. son distintos; así como también sus estados conformacionales son diferentes (uno puede estar como monómero, mientras que el otro está como dímero).

Por otra parte, se ha explicado también que cada subgrupo interacciona con ligandos diferentes y regulan procesos distintos (Korner, Sofia et al. 2003).

Como quedará explicado en la siguiente sección, los residuos que intervienen en las diferentes funciones y roles estructurales de los dominios están mucho más conservados a lo largo de la evolución respecto al resto de los residuos.

1.4 Análisis de residuos funcionalmente importantes

El proceso de caracterización de las proteínas, que comienza con los análisis de secuencia y conformación estructural, se sigue con la determinación de la función.

La función y la conformación de las proteínas son dos características que se ha visto que están directamente correlacionadas (James and Tawfik 2003). Si la proteína no se conforma de forma específica, no logra ser funcional. A lo largo de la evolución se ha mantenido la selección de ambas características conjuntamente (Saito, Sasai et al. 1997). En una aproximación teórica, una proteína funcional cumple dos características mínimas que se mantienen a lo largo de la evolución:

a) Sigue el principio de “capacidad de plegamiento” (foldability) que es aquel que sigue la proteína para conformarse espacialmente con el mínimo de energía (Saito, Sasai et al. 1997; James and Tawfik 2003), dicha propiedad es estable ante mutaciones, así como a variaciones ambientales (temperatura, solvente, etc.) (Mirny and Shakhnovich 1999).

b) Tiene un “núcleo de plegamiento” que está constituido por residuos que conforman un grupo espacialmente contiguo en la estructura (esto no implica que estén contiguos en la secuencia) (Mirny and Shakhnovich 1999). Este núcleo(s) es el primero en conformarse cuando la proteína se empieza a plegar; posterior a estos puntos espaciales, el resto de la proteína se pliega gastando el mínimo de energía (Mirny and Shakhnovich 1999).

Gran parte de los residuos que se conocen como ‘Residuos Funcionalmente Importantes’ (RFI) son los que generan y mantienen estas características (Mirny and Shakhnovich 1999). La diversidad y variación de estos RFIs derivan en funciones y acomodos conformacionales diferentes. Actualmente, hay dos estrategias que se siguen para lograr determinar a estos RFI, una experimental y otra computacional.

a) Análisis experimental. Los estudios se realizan experimental e individualmente para cada proteína que se requiere conocer a sus RFIs. Esta estrategia consiste en realizar ‘análisis mutagénicos’ (Hutchison, Phillips et al. 1978) que implican mutar los residuos que se ubican en el sitio de unión, de los que se tiene idea que pueden participar en la interacción o de cada uno de los residuos de la proteína, y después

se evalúa como funciona la proteína cambiada, y de acuerdo a ello, se determina la relevancia del residuo mutado para la función (Hutchison, Phillips et al. 1978). (Mucha de esta información sobre mutaciones hechas en proteínas y los efectos que tienen en la función y/o estructura de las proteínas se encuentra disponible en la PMD, Protein Mutant Database, (Kawabata, Ota et al. 1999)).

b) Análisis computacional. Los RFIs se pueden identificar computacionalmente en familias de proteínas que sean homólogas⁴ (identificables por su similitud de secuencia) o que tengan un mismo plegamiento (identificables por similitud estructural), debido a que la evolución de las funciones y conformaciones de las proteínas sigue ciertos comportamientos evolutivos. Por una parte, a las proteínas homólogas es que proviniendo de un ancestro común, aunque la similitud de sus secuencias sea de sólo 25% su estructura varía muy poco, dentro de 2Å (Chothia and Lesk 1986). Por otra, a las proteínas con el mismo plegamiento, pueden tener un “núcleo de plegamiento” similar (Mirny and Shakhnovich 1999). En cuanto a los RFIs, la conservación de los residuos tiende a correlacionar directamente con su importancia para la función (Zvelebil, Barton et al. 1987), y/o mantenimiento de la conformación estructural; pero pueden también conservarse por razones históricas (falta de tiempo para divergir) (Mirny and Shakhnovich 1999). Así que evaluando la conservación de los residuos y analizando su posible funcionalidad con diferentes medidas se pueden identificar a los RFIs en grupos de proteínas homólogas o con el mismo plegamiento, de las que se tiene información estructural de alguna(s) de ellas.

La estrategia computacional es la más utilizada, ya que es rápida y se puede realizar para muchas proteínas a la vez; aunque tiene algunos inconvenientes: a) depende inevitablemente de la diversidad y del número de secuencias con que se trabaje y b) del tipo de mediciones que se hagan para determinar la conservación/variabilidad de los residuos (Mirny and Shakhnovich 1999).

⁴ Homología se define como la relación de dos caracteres cualquiera que han descendido, usualmente por divergencia, de un carácter ancestral común (Fitch 2000).

Ortología se define como la relación de dos caracteres homólogos posterior a la divergencia de su ancestro común, donde estos caracteres se encuentran en diferentes organismos que divergieron del ancestro común (Fitch 2000).

Paralogía se define como la relación de dos caracteres homólogos cualquiera que surgen de una duplicación y usualmente han permanecido en un mismo organismo (Fitch 2000).

c) Comparación de métodos computacionales

Los métodos computacionales siguen varios principios comunes, pero cada uno varía en ciertas consideraciones. Los que se han considerado más relevantes se explican en el siguiente cuadro:

Autores de diferentes métodos	Hipótesis	Criterios para identificación de residuos funcionales	Medidas de funcionalidad / especificidad de los residuos
Método de Traza Evolutiva (Lichtarge, Yao et al. 2003)	<p>Las funciones comunes se retienen o evolucionan de una misma proteína ancestral, a las descendientes con cambios en o alrededor del mismo sitio estructural.</p> <p>Los cambios a nuevas funciones se observan como cambios de residuos que están implicados en la función.</p>	<p>1) Alineamiento múltiple de secuencias de una familia.</p> <p>2) Construcción de árbol por sustituciones que separa diferencia de secuencias por porcentajes</p>	<p>La conservación es la medida del tipo y/o posición del residuo en el grupo de secuencias alineadas. Esta se observa como la presencia o ausencia del residuo.</p>
(Johnson and Church 2000)	<p>Una familia de proteínas con el mismo esqueleto (backbone) estructural que reconoce ligandos distintos, tiene residuos variantes alrededor de la cavidad del ligando. Estos residuos son críticos para el reconocimiento específico de cada ligando.</p>	<p>1) Familia de proteínas con el mismo <i>backbone</i> estructural.</p> <p>2) Alineamientos múltiples de secuencias y de estructuras</p>	<p>Ubicación del residuo alrededor de la cavidad, con una distancia menor a 4.5 Å del ligando.</p>
(Mirny and Shakhnovich 1999), (Mirny and Shakhnovich 1999)	<p>Se asume que las secuencias parálogas y ortólogas⁴ realizan las mismas funciones bioquímicas de forma general, pero los ortólogos tienen la misma especificidad. Proteínas evolutivamente relacionadas tienden a conservar el mismo plegamiento.</p>	<p>1) Alineamiento múltiple de secuencias, donde el porcentaje de identidad agrupa ortólogos y separa parálogos</p> <p>2) Medición del Contenido Informacional de cada posición del alineamiento</p>	<p>Las posiciones con Contenido Informacional significativo estadísticamente corresponden a residuos funcionalmente importantes (ver en métodos)</p>

De los métodos descritos en el cuadro anterior, el de Traza Evolutiva ha sido utilizado ampliamente en otros trabajos para identificar RFIs. En cada caso, las variantes del método están dadas en la manera de evaluar la significancia de los RFIs encontrados (Landgraf, Fischer et al. 1999; Innis, Shi et al. 2000; Armon, Graur et al. 2001; Landgraf, Xenarios et al. 2001; Lichtarge, Yao et al. 2003; Berezin, Glaser et al. 2004). Este método se limita a sólo identificar a los RFIs que interaccionan con los ligandos, además que estas proteínas deben reconocer a sus ligandos en las mismas regiones o sitios activos.

Alternativamente, el método de Johnson y Church (2000), es útil sólo para casos en que la familia de proteínas sigue exactamente el mismo plegamiento, que interaccionen con su ligando en el mismo sitio, tan sólo variando el ligando que reconocen. Igualmente, este método se centra sólo en la identificación de RFIs que interaccionan con el ligando a partir de proteínas con información estructural y/o experimental.

Por último, el método de Mirny y Gelfand (2002) propone que la identificación de RFIs no se debe centrar solo en aquellos residuos que unan a un ligando; además, que los RFIs no tienen que ubicarse en los mismos sitios para todos los miembros de la familia, es decir pueden localizarse en diferentes regiones del alineamiento, no se limita a que las proteínas sigan exactamente el mismo plegamiento, y la evaluación de los RFIs encontrados no depende de observaciones sobre una estructura o su distancia al ligando, sino de su significancia estadística.

Dado que los miembros de la familia Crp/Fnr reconocen diferentes ligandos y no lo unen en una misma región, se decidió implementar el método de Mirny y Gelfand 2002 para identificar los RFIs que podrían interaccionar con un potencial ligando.

2 Hipótesis

Las diferencias funcionales dentro de una familia de factores transcripcionales parálogos pueden identificarse computacionalmente como residuos específicos en diferentes posiciones del dominio de unión al ligando.

3 Material y métodos

3.1 Colección de secuencias y estructura tridimensional

En este trabajo se decidió utilizar la clasificación de familias de factores transcripcionales (FTs) de Pérez-Rueda *et. al* 2000, ya que la información disponible de estos es principalmente de secuencias completas y sus dominios.

Todas las secuencias de las proteínas usadas en este trabajo para la familia Crp/Fnr, se extrajeron de la base de datos Genbank del National Center for Biotechnology Information (NCBI) [<http://www.ncbi.nlm.nih.gov/entrez/> (sección 'protein')] hasta el 2005. Estas secuencias pueden ser consultadas utilizando los identificadores GI (gene identifier) de esta misma base de datos (ver Anexo 1: lista de secuencias utilizadas).

La estructura tridimensional que se utilizó como referencia para la familia fue la de Crp de *E. coli* obtenida de la base de datos Protein Data Bank (PDB), cuyo identificador es 1J59 ([<http://www.rcsb.org/pdb/>] Bermann *et. al* 2000).

3.2 Detección de los miembros de la familia Crp/Fnr

Se obtuvieron 276 proteínas homólogas a las secuencias de Crp y Fnr de *E. coli* de la base de datos SwissProt ([<http://expasy.org/sprot/>] (Boeckmann, Bairoch *et al.* 2003) utilizando el programa *blastp* (E-value 10^{-6}) (Altschul, Gish *et al.* 1990). De los candidatos de la búsqueda anterior, todas aquellas secuencias que tuvieran el mismo dominio de unión al DNA de tipo Helix-Turn-Helix (SSF46785) y el dominio de unión al ligando cNMP (SSF51206), según la base de datos SuperFamily (Gough 2002), fueron considerados homólogos (ortólogos y parálogos) pertenecientes a la familia CRP/FNR. Se encontraron y eliminaron 9 secuencias que no cumplieron con estos requisitos mínimos.

Con todas estas secuencias se generó el alineamiento múltiple (multiple sequence alignment, MSA) de la familia con el programa *clustalx* (Thompson, Gibson *et al.* 1997). Con el alineamiento obtenido se generó su perfil de secuencia con una matriz de peso utilizando los programas *hmmbuild* y *hmmcalibrate* del paquete *HMMer* 2.3 (Eddy

1998). Posteriormente, se mapeo el perfil de la familia Crp/Fnr a través de 200 genomas de bacterias y arqueas utilizando el programa *hmmsearch* de *HMMer 2.3* (E-value $\leq 10^{-3}$) (Eddy 1998). De esta búsqueda se obtuvieron 171 secuencias (Anexo 7.1: Tabla de secuencias).

Como el propósito de este trabajo fue estudiar el dominio de unión al ligando (DUL), se eliminó del alineamiento el dominio de unión al DNA de la región carboxilo y se dejó sólo el DUL de la región amino, de aproximadamente 150 residuos. Posteriormente, se eliminaron las secuencias del DUL que eran 100% idénticas entre ellas para evitar redundancia en el grupo de datos estudiados, utilizando el programa *CD-HIT* (Lin, Kovac et al. 2002); se eliminaron 37 secuencias. Con las 134 secuencias del DUL homólogas y no redundantes, finalmente se volvió a hacer el MSA con *clustalx* (valores por default, ver Anexo 7.2: alineamiento de secuencias estudiadas).

3.3 Análisis evolutivo para la división en subgrupos

A las 134 secuencias del DUL se les hizo un análisis evolutivo para conocer su agrupamiento en subgrupos dentro de la familia utilizando el software para cálculos filogenéticos y graficación *MEGA* (Kumar, Tamura et al. 2004). Se seleccionó el modelo evolutivo para generar el árbol con el software *ProtTest* (REF). Se construyó su árbol filogenético siguiendo el método de Neighbor Joining utilizando el modelo JTT (parámetro $\alpha = 2.0$) (Jones 1999) y se sometió a 100 replicas de bootstrap y una prueba de consenso para verificar que los agrupamientos fueran consistentes (ver Resultados: figura 1R, Anexo 7.3 : gráfica del árbol de la familia).

3.4 Identificación de residuos funcionalmente importantes

Ya que los residuos funcionalmente importantes (RFI) se comparten entre ortólogos (misma especificidad), pero no con parálogos (diferente especificidad), estos RFI se pueden diferenciar entre estos grupos. En un alineamiento de secuencias de proteínas homólogas, las columnas donde se ubican estos residuos (RFI) tienen un *contenido informacional (CI)* diferenciable al del resto de las columnas. El *CI* es una medida que

permite la asociación de la conservación de los residuos con la especificidad de la proteína (Mirny and Shakhnovich 1999).

Por lo tanto, la identificación de RFI se realizó calculando el CI de cada columna del alineamiento de la familia Crp/Fnr con la siguiente fórmula:

$$(1) \quad CI_i = \sum_{\substack{x=1,\dots,20 \\ y=1,\dots,Y}} f_i(x,y) \log (f_i(x,y) / f_i(x) f(y))$$

Donde,

- x = cada uno de los 20 aminoácidos (A,C,D,G,...,W,Y)
- y = número de grupos en que se dividen las secuencias homólogas
- $f_i(x)$ = frecuencia del residuo tipo x en la posición i del alineamiento múltiple
- $f(y)$ = fracción de proteínas que pertenecen al grupo y
- $f_i(x,y)$ = frecuencia del residuo x en la posición i dentro del grupo y
- $f_i(x,y) / f_i(x) f(y)$ = razón de la frecuencia de un residuo en el grupo respecto al alineamiento múltiple, considerando la cantidad de proteínas en el grupo

Este cálculo se hizo con un programa en el lenguaje de programación *PERL*: ContInfo-aa-alignment-PRO.pl (Anexo 7.6: códigos de programas). El cálculo de CI de las columnas del alineamiento con *gaps* se consideró igual a cero ($CI_{gap} = 0$).

Dado que el *contenido informacional* se puede sesgar por una muestra pequeña o por la composición de aminoácidos (Mirny and Shakhnovich 1999), no se puede tomar sólo este valor para identificar a los determinantes de especificidad; sino que se debe calcular su significancia estadística (descrita abajo) y utilizar esta medida junto con el valor de CI -reales (los obtenidos para las secuencias de la familia Crp/Fnr) para predecir los RFIs.

3.5 Cálculo de significancia estadística

Para conocer el valor estadístico de los CI -reales de cada posición de los grupos del alineamiento, se utilizó la estrategia de comparación de cada valor de CI -reales contra CI -referencias de una población mayor de secuencias (que sigue una distribución Normal, Anexo 7.4), obteniendo su valor Z . Posteriormente, los CI -reales que tuvieron un

valor $Z \geq 3$ con un área ≤ 0.05 se consideraron *CI-reales* que corresponden a residuos funcionalmente importantes.

La población utilizada fue una generada con secuencias “pseudoaleatorias” que tienen la misma variabilidad intra-grupos (llamada aquí ρ) que la de los grupos del alineamiento (ver explicación en 3.5.b y discusión en resultados 4.3).

a) Obtención de secuencias pseudos-aleatorias

La generación de secuencias pseudoaleatorias, con variabilidad ρ (ver en punto 3.5.b) se hizo siguiendo los puntos 1 y 2 (explicados enseguida abajo) utilizando el programa en PERL Calc-seqs-PseudoRandom.pl (Anexo 7.6: códigos de programas). Se generaron igual número de secuencias como miembros de cada subgrupo y el proceso se repitió para la familia completa 10^3 veces.

1) Generación de secuencias parentales para cada subgrupo: Para cada grupo, se considera la distribución de aminoácidos en cada columna (que es $f_i(x)$) y con probabilidad $P(f_i(x))$ se hereda alguno de esos aminoácidos.

2) Generación de secuencias pseudo-aleatorias: Se toma la secuencia parental de cada grupo y la $f_i(x)$ de cada columna para ese grupo. De acuerdo con la probabilidad de ρ las secuencias pseudo-aleatorias heredan el mismo a.a. de la secuencia parental o cambian por algún otro a.a. de la $f_i(x)$ para esa columna.

b) Variabilidad intra-grupos (ρ)

La variabilidad intra-grupos (ρ) se utiliza para determinar la similitud dentro de los grupos de las secuencias pseudos-aleatorias y así dar una medida umbral para la herencia o cambio de aminoácido de las secuencias parentales a las nuevas secuencias pseudos-aleatorias. Si $\rho = 1$ implica que todas las secuencias generadas

son idénticas, por el contrario, si $\rho = 0$ las secuencias son tan distintas como parálogos. La ρ se obtuvo posterior a la comparación de identidades por pares cada secuencia del alineamiento, utilizando el programa en PERL `Calculo_rho.pl` (Anexo 7.6: códigos de programas), y luego de analizarlas por agrupamientos (clustering) a través del Centroit Linkage Clustering, con una métrica de distancia “uncentered correlation” del software Cluster 3.0 (Eisen 1998) y con el programa en PERL `ColorSeqRelation.pl` (Anexo 7.6: códigos de programa) (ver figura 4R en Resultados). Para la familia de secuencias de Crp/Fnr se determinó utilizar una $\rho = 0.75$ para obtener una similitud intra-grupo cercana a las proteínas conocidas.

c) Obtención de los valores Z

Los valores Z de los *CI-reales* para cada posición del alineamiento se obtuvieron con el programa en PERL `MedsStats.pl` (Anexo 7.6: códigos de programas) utilizando la fórmula:

$$(2) \quad Z = \frac{X - \mu}{\sigma}$$

Donde: x = valor *CI-reales* a evaluar
 μ = media de la población de *CI-referencias*
 σ = desviación estándar de la población de *CI-referencias*

3.6 Caracterización de estructura secundaria y terciaria

Para la familia Crp/Fnr se tienen reportadas varias estructuras de Crp de *E. coli* en la base de datos PDB, de las cuales se seleccionó una de las que tiene una mayor resolución, la 1J59 a 2.5Å. El mapeo de los residuos se hizo utilizando el software ProteinExplorer (www.proteinexplorer.org).

4 Resultados y discusión

Bajo la hipótesis de que las diferencias funcionales dentro de una familia de factores transcripcionales homólogos (ortólogos y parálogos) pueden identificarse computacionalmente como diferentes residuos específicos en diferentes posiciones del dominio de unión al ligando, la subdivisión de la familia en grupos podría hacer más sensible la identificación de los diferentes residuos funcionalmente importantes para cada uno de los grupos.

4.1. Identificación y agrupamiento de miembros de la familia Crp/Fnr

Para el estudio de residuos funcionalmente importantes (RFIs) del dominio de unión al ligando (DUL) de la familia Crp/Fnr se comenzó por la identificación de secuencias homólogas de CRP y FNR de *Escherichia coli* (CRP_{ECO} y FNR_{ECO}) a través de 200 genomas de Bacterias y Arqueas. De las 171 secuencias encontradas, se tomaron únicamente las regiones Amino donde se encuentran los Dominios de Unión al Ligando (DUL) de aproximadamente 150 a.a. (se eliminaron las regiones carboxilo de las secuencias). Posteriormente, se eliminaron 37 secuencias que eran redundantes (idénticas en los mismos organismos). Finalmente las 134 secuencias de los DUL se alinearon y se obtuvo su árbol filogenético como se explica en Métodos.

En la figura 1R se muestra el árbol filogenético de las 134 secuencias de los DUL. Lo primero a notar es que estas 134 secuencias se separan en dos grandes grupos correspondientes a los dos parálogos que definieron inicialmente a la familia: los que se diferencian con Crp (región izquierda del árbol; marcados con etiqueta naranja “CRP”) y los que se diferencian con Fnr (región derecha del árbol; marcados con etiqueta roja “FNR”). En seguida se observa que a la vez estas 134 secuencias se separan en 9 grupos (diferenciados en recuadros de diferentes colores) que representan secuencias de 4 agrupamientos filogenéticos: *Proteobacteria* (alfa- α (G0: rojo, G8: azul), beta- β (G3: verde) y gama- γ (G1: amarillo, G2: naranja), *Firmicutes* (G4: verde seco, G5: azul claro), *Actinobacteria* (G6: rosa) y *Cyanobacteria* (G7: lila) (ver también Anexo 3).

Es importante hacer notar que la subfamilia Fnr posee un mayor número de miembros subdivididos en seis grupos pertenecientes a γ (G2), β (G3) y α *Proteobacteria* (G0, G8), *Firmicutes* (G5) y *Cyanobacteria* (G7); mientras que los miembros de la subfamilia

CRP son menos y se dividen en tres grupos pertenecientes a γ -Proteobacteria (G1), Firmicutes (G4) y Actinobacteria (G6). Esto no sólo indica una mayor divergencia de los miembros de la subfamilia Fnr, sino de una mayor diversidad funcional y reconocimiento a diferentes ligandos en los organismos aquí descritos (ver Körner *et. al* 2003 que lo han reportado previamente).

Además de los ortólogos a Crp y Fnr (grupos G1 y G2 respectivamente) se identificaron algunos homólogos a FixK y YeiL, que son otros dos miembros ya caracterizados de la familia (ver Introducción) (Körner *et. al* 2003). En el grupo G0 hay 9 secuencias anotadas como FixK de α -Proteobacteria, y en el grupo G4 hay cinco secuencias anotadas como YeiL de γ -Proteobacteria. Es importante remarcar que estas secuencias de YeiL se agruparon con el grupo 4 de Firmicutes (cuadro punteado amarillo, figura 1R).

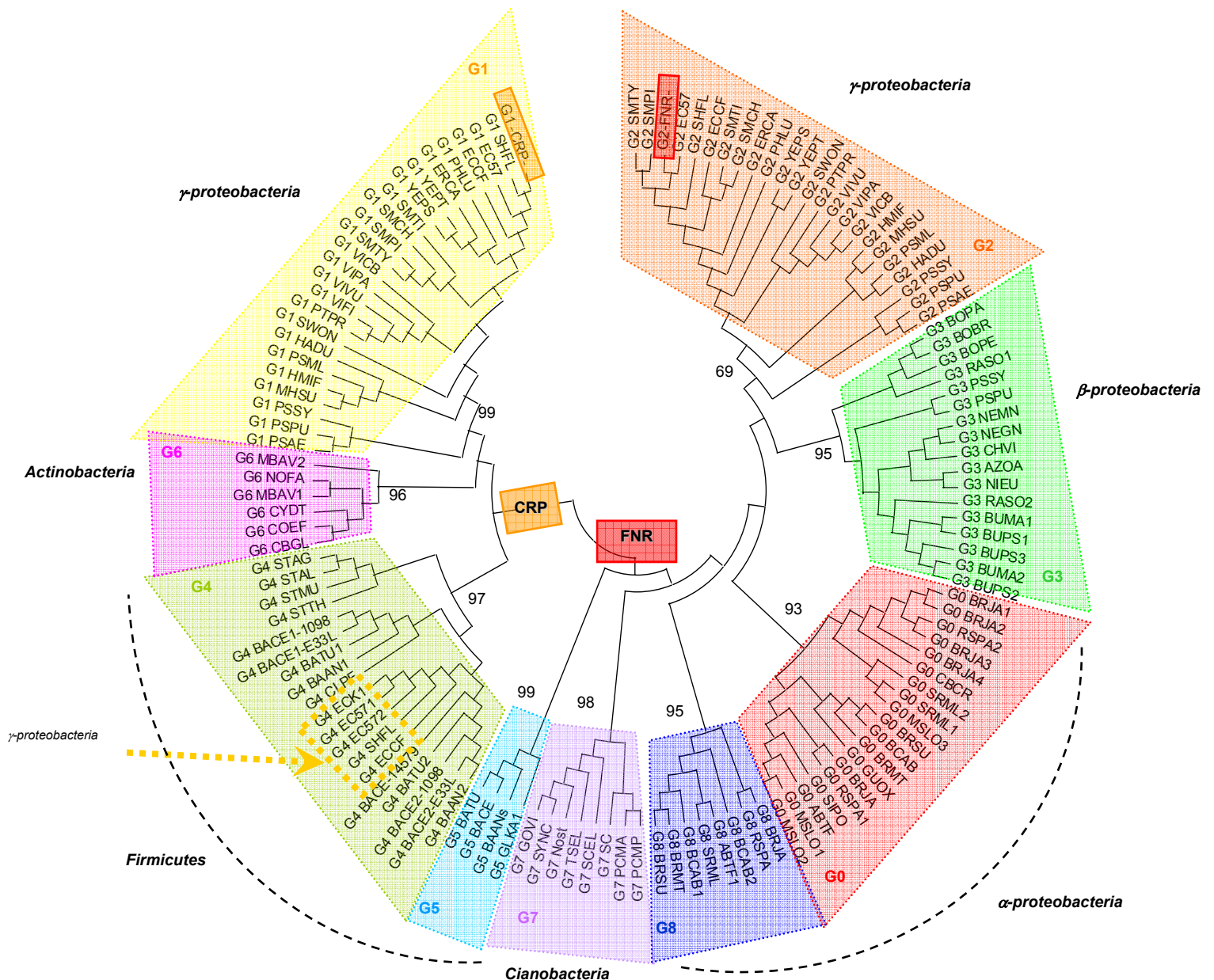


Figura 1R. **Representación circular del árbol de los 134 miembros identificados para la familia Crp/Fnr** divididos en 9 grupos, con valores bootstrap mayores a 65 (números negros en ramas que divergen a cada uno de los 9 grupos). Estos nueve grupos corresponden a cuatro grupos filogenéticos: *Proteobacteria* (α , β , γ), *Firmicutes*, *Actinobacteria* y *Cyanobacteria* (nombres indicados al lado de los grupos por la parte externa del árbol circular). Los homólogos a Crp se separan en las ramas marcadas con etiqueta naranja “CRP” y se subdividen en tres grupos marcados recuadros amarillo (*Proteobacteria* γ G1), rosa (*Actinobacteria* G6) y verde seco (*Firmicutes* G4). Los homólogos a Fnr se separan en las ramas marcadas con etiqueta roja “FNR” y se subdividen en seis grupos marcados en recuadros azul claro (*Firmicutes* G5), lila (*Cyanobacteria* G7), azul y rojo (α -*Proteobacteria* G8 y G0), verde (β -*Proteobacteria* G3).

Con este árbol se pudo observar que el DUL ha sido suficiente para diferenciar las proteínas homólogas a la familia Crp/Fnr de acuerdo a los diferentes grupos de *Bacterias* mencionados; a la vez, esta separación de los DUL en diferentes grupos es indicativa de diferentes posibles funciones.

4.2. Identificación de residuos funcionalmente importantes

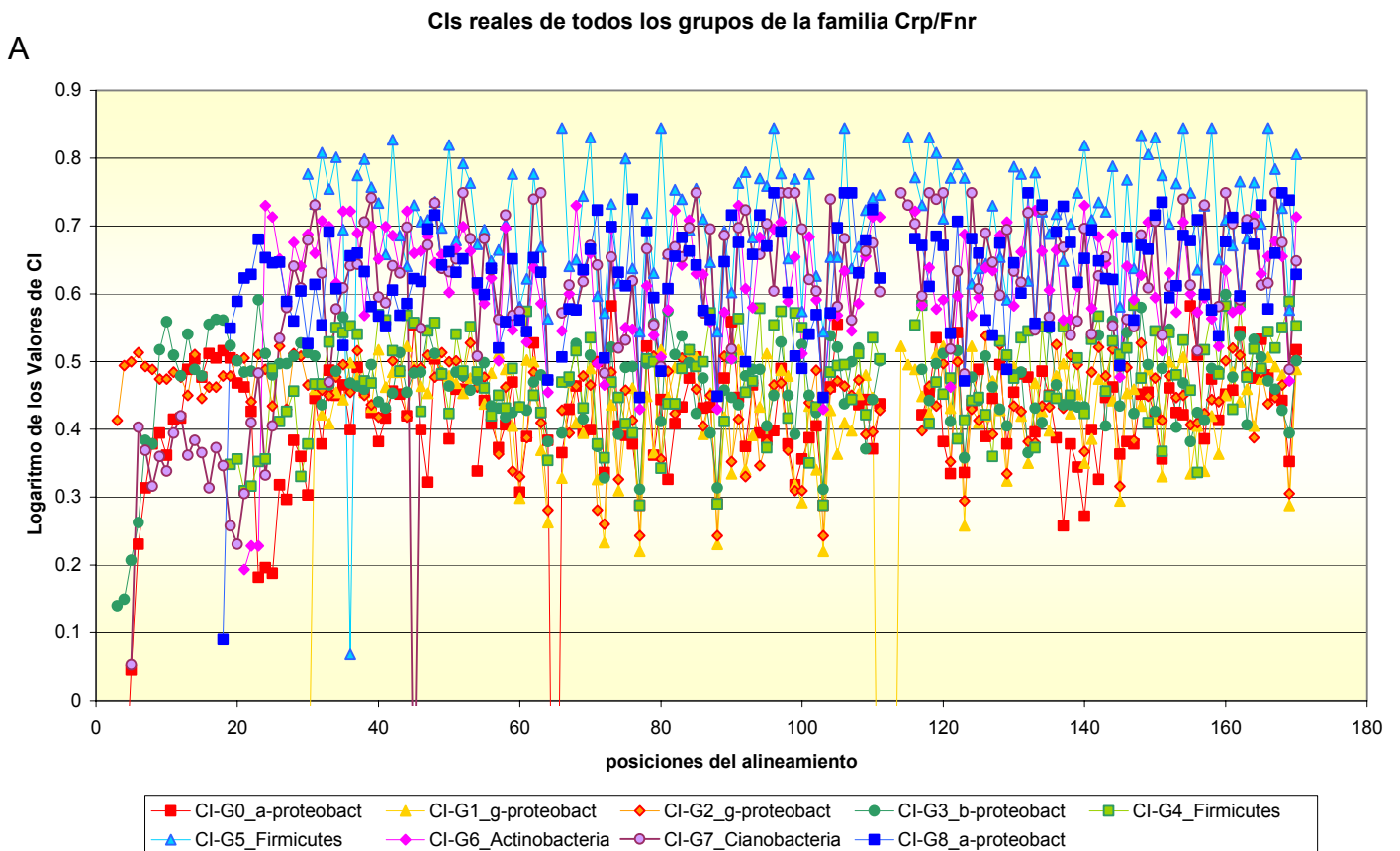
Ya que con la aproximación del ‘Contenido Informacional’ (*CI*) de los residuos de secuencias homólogas se pueden identificar posiciones en el alineamiento que tienen mayor carga de información, que implica una importancia para la función (Mirny and Shakhnovich 1999), se procedió a hacer el cálculo del *CI* para las ~180 posiciones del alineamiento de los DUL de los 9 grupos de la familia Crp/Fnr para identificar sus Residuos Funcionalmente Importantes (RFI). El cálculo de *CI* con la fórmula 1 explicada en Métodos (3.4) e implementada en el programa en el lenguaje de programación PERL ‘ContInfo-aa-alignment-PRO.pl’ (ver Anexo 7.6.1).

La distribución de los 9 grupos de secuencias DUL en 4 grupos filogenéticos distintos le dan diversidad (heterogeneidad) al universo de secuencias estudiadas. Por otra parte, el número de secuencias que constituyen a cada grupo no es homogéneo debido a la disponibilidad sesgada de secuencias en las bases de datos. En los grupos de *Actinobacteria* hay 9 secuencias [G6], de *Cyanobacteria* 8 [G7], de *Firmicutes* 24 [G4, G5] y de *Proteobacteria* 93 (α tiene 27 [G0, G8], β tiene 17 [G3] y γ tiene 49 [G1, G2]).

Es importante notar que el grupo de *Proteobacteria* es el más representado.

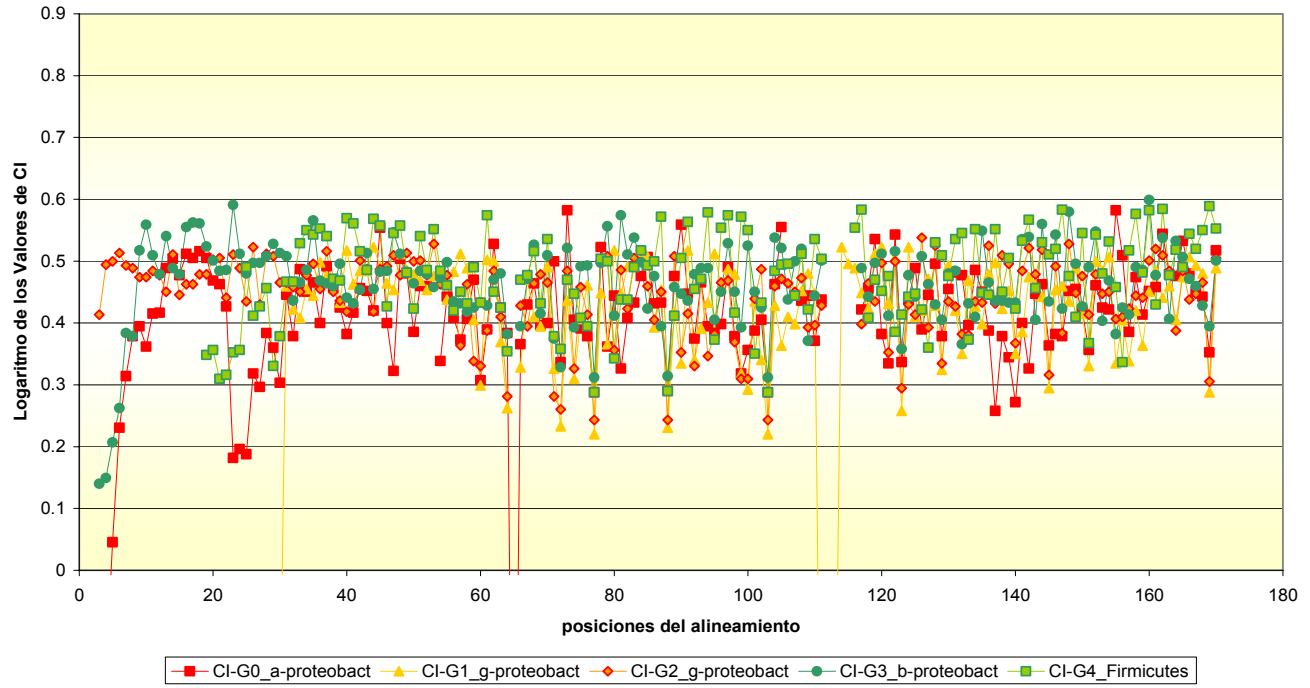
En la figura 2R se grafican los valores logarítmicos del *CI* para cada una de las ~180 posiciones que conforman el alineamiento de secuencias DUL de los 9 grupos de la familia.

Lo primero que se puede notar es que los valores de *CI* se separan en dos tendencias distintas, una de valores altos y otra de valores medios (ver 2R-B y 2R-C). En la tendencia de valores altos, gráfica 2R-B, se encuentran los grupos de *Actinobacteria* [G6], *Cyanobacteria* [G7], uno de *Firmicutes* [G5] y uno de α -*Proteobacteria* [G8]; que coinciden en ser grupos con menos de 10 secuencias y que tres de los cuatro grupos no son *Proteobacteria*. En la tendencia de valores medios, gráfica 2R-C, se encuentra el otro grupo de α -*Proteobacteria* [G0] junto con los de γ -*Proteobacteria* [G1, G2], β -*Proteobacteria* [G3] y otro de *Firmicutes* [G4]; que coinciden ser grupos con 17 o más secuencias y cuatro de los cinco grupos son *Proteobacteria*.



CIs reales de los grupos *Proteobacteria* (α , β , γ) y *Firmicutes* de la familia Crp/Fnr

B



CIs reales de los grupos de *Firmicutes*, *Actinobacteria*, *Cyanobacteria* y α -*proteobacteria* de la familia Crp/Fnr

C

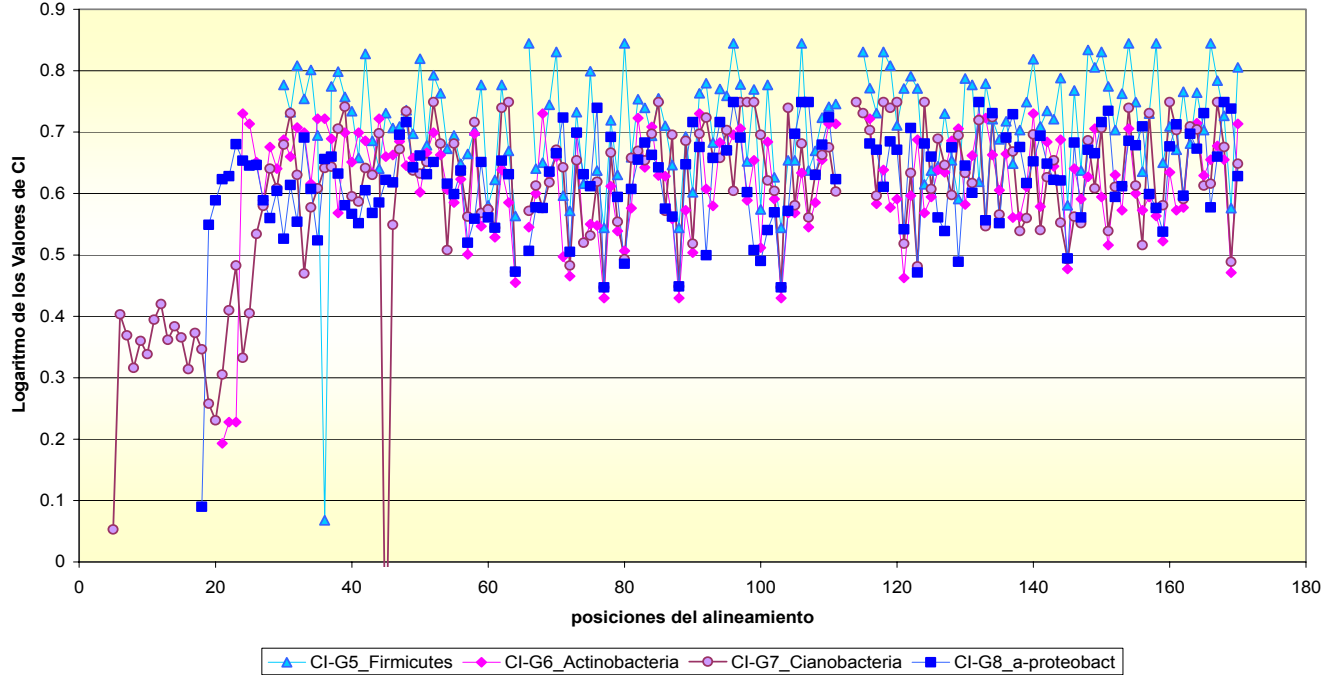


Figura 2R. A) Logaritmos del Contenido informacional de los 9 grupos de la familia Crp/Fnr. B) Logaritmos del Contenido Informacional de cinco grupos con valores altos. C) Logaritmos del Contenido Informacional de cuatro grupos con valores medios. Nota: El valor de CI de cada posición es discreto e independiente, pero se decidió trazar una línea que une a los puntos de cada grupo para facilitar su observación.

Dos de las propiedades de la métrica de CI (ver Métodos) son que CI es igual a 0 si y sólo si x y y son estadísticamente independientes, y de que un valor grande de CI indica una fuerte asociación entre x y y (donde x es uno de los 20 a.a. y y es un grupo dentro del alineamiento); así que al observar estas dos tendencias de CI de las secuencias DUI estudiadas en la figura 2R, se quiso conocer de cual de las dos variables ($x \rightarrow$ tipo de aminoácidos, $y \rightarrow$ número de secuencias en cada grupo) consideradas en el cálculo del CI están influyendo más para tener estas diferencias:

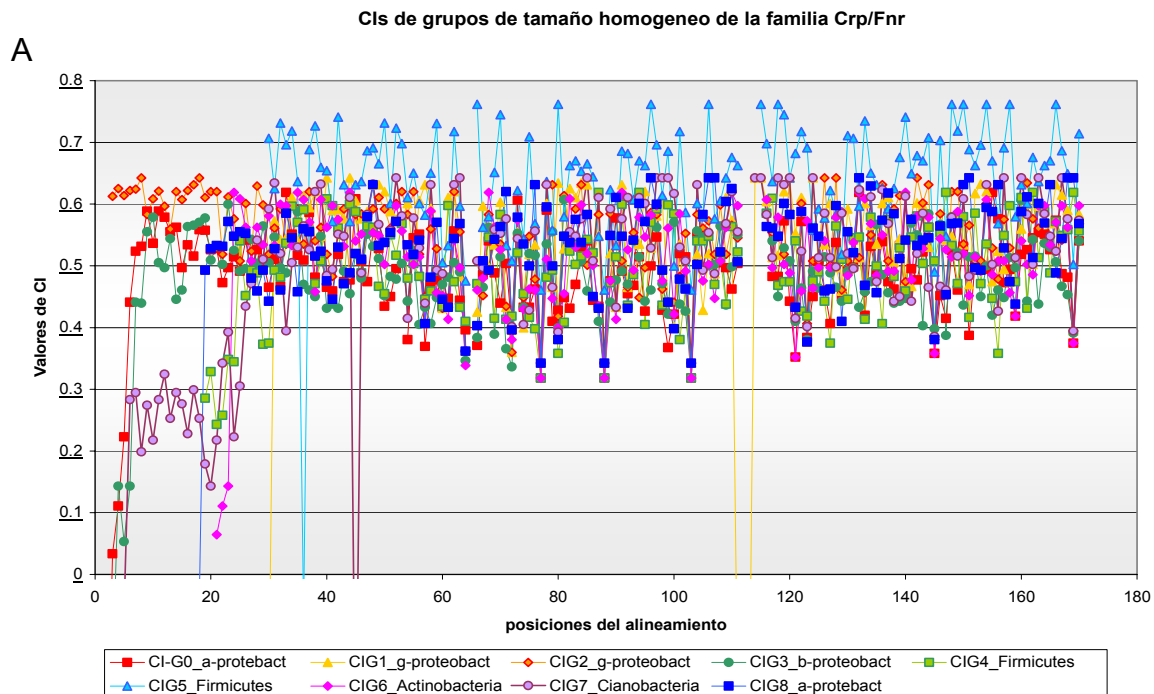
- a) Si se debe a que el número de miembros que conforman a los grupos es diferente (CI depende más de y);
- b) Si se debe a la diferente composición de aminoácidos en los grupos filogenéticos en el alineamiento (CI depende más de x).

Para la opción (a) se realizó otro análisis de CI para la misma cantidad de grupos de la familia, pero sin que su número de miembros variara (figura 3R-A). Para esto, se dejó el mismo número de miembros de los grupos que tienen 9 o menos secuencias y para los grupos con más de 17 miembros, se decidió que tuvieran al menos 9 secuencias seleccionadas al azar (ver tabla 1R). En la tabla 1R se muestra el número de secuencias que conforman a los grupos en la Familia Crp/Fnr (como están en el alineamiento) y el de los grupos con casi el mismo número de miembros seleccionados al azar.

Tabla 1R. Número de secuencias en los grupos de la familia Crp/Fnr como están en el alineamiento y de los grupos que se procuraron tuvieran el mismo número de miembros seleccionados al azar.

Grupo	Número de secuencias de los grupos de la familia (observados)	Redistribución de miembros de grupos (azar)
Actinobacteria G6	9	9
Cianobacteria G7	8	8
α -proteobacteria G8	8	8
Firmicutes G5	5	5
Firmicutes G4	19	9
α -proteobacteria G0	19	9
β -proteobacteria G3	17	9
γ -proteobacteria G1	25	9
γ -proteobacteria G2	24	9
Total	134	75

Para la opción (b) se realizó otro análisis de *CI* para la misma cantidad de grupos de la familia, pero la composición de aminoácidos de cada posición fue cambiada de forma aleatoria (figura 3R-B). El cambio de composición de aminoácidos se hizo con el programa en el lenguaje de programación PERL Shuffling-alin-PRO.pl (ver Anexo 7.6.3).



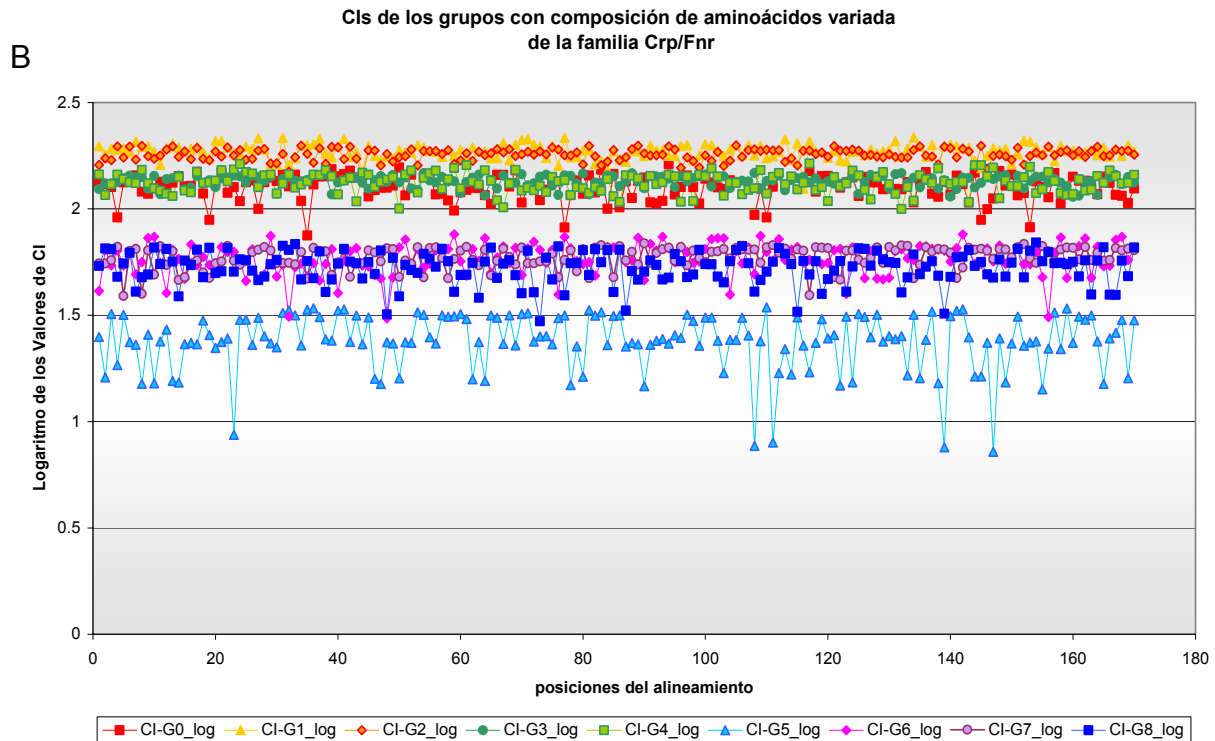


Figura 3R. A) Logaritmos de CI calculados para 9 grupos con el mismo número de secuencias. B) Logaritmos de CI calculados para 9 grupos con composición de aminoácidos variada. Nota: El valor de CI de cada posición es discreto e independiente, pero se decidió trazar una línea que une a los puntos de cada posición de los nueve grupos para facilitar su observación.

En la figura 3R-A, donde los grupos tienen aproximadamente el mismo número de secuencias, se observa que los valores de CI de todos los grupos son similares, y además que en general sus valores aumentan con una media aproximada de 0.5. Con esto se ve que los valores de CI dependen del número de secuencias en los grupos (CI depende más de y).

En la figura 3R-B, donde la composición de aminoácidos de los grupos está cambiada, se observa que aunque los valores de CI han cambiado, se siguen separando por el número de secuencias que tienen los diferentes grupos y no reflejan las relaciones filogenéticas de los diferentes grupos (CI depende más de y).

En el acomodo de las secuencias encontradas para la familia Crp/Fnr se observa claramente que el grupo 8 de α -Proteobacteria que sólo tiene 8 miembros se encuentra

con los grupos de valores altos (figura 2R-B), mientras que el resto de los grupos de *Proteobacteria* (G0, G1, G2 y G3) que tienen al menos 17 miembros presentan valores de *CI* medios (figura 2R-C); y uno de los grupos de *Firmicutes* con 19 secuencias (G4) tiene un valor también medio de *CI* (figura 2R-C) y el otro grupo de *Firmicutes* con 5 secuencias (G5) tiene un valor alto (figura 2R-B).

Con estos dos ensayos se demuestra que las diferencias en los valores de *CI* dependen del número de secuencias en los grupos (y) de las secuencias estudiadas. También se muestra que los valores de *CI* no se ponderan siendo mayores o menores a causa del número de secuencias o de la composición de aminoácidos. Por lo tanto, para poder diferenciar los valores de *CI* que se relacionan con un alta conservación (funcionalidad) de los residuos se requirió hacer una evaluación estadística.

4.3. Identificación de la significancia estadística del *CI*

Ya que el valor de *CI* no es suficiente indicador por si solo de la asociación del residuo con su funcionalidad, se procedió entonces a calcular su significancia estadística comparándolo con el *CI* de secuencias que representen un universo “azaroso” de la familia Crp/Fnr. Este universo se conforma de secuencias pseudo-aleatorias que son variables pero mantienen la misma composición de aminoácidos de la familia, lo cual permite que sigan siendo secuencias similares. Si el universo de secuencias más bien se conformara de secuencias totalmente aleatorias, se estarían incluyendo secuencias cuyas funciones y RFIs son totalmente ajenas a las de la familia y no se podrían detectar *CI*s de RFIs comunes; este universo se alejaría del modelo biológico real estudiado, la familia de Crp/Fnr.

Para determinar la variabilidad de secuencias para el estudio se graficaron las identidades entre las 134 secuencias DUL de la familia Crp/Fnr (gráfica 4R). En la gráfica 4R, se muestran sobre la diagonal en rojo las identidades que indican secuencias similares de 100 a 95%, y alrededor de la diagonal se indican identidades menores: en naranja de 85 a 75%, en naranja claro de 74 a 70%, amarillo de 69-60%, y

de verde-azul-negro las identidades menores de 60 hasta 0%. Con esta gráfica se puede ver que el universo para explorar las secuencias las pseudos-aleatorias se encuentran entre 60 y 90%. Para mayoría de los grupos el intervalo de 80 a 90% (naranja cercano a rojo) no existe. Tomar una $\rho \geq 0.9$ (que equivale a 90%) implica generar secuencias casi idénticas a las reales, y tomar una $\rho \leq 0.6$ implica considerar un intervalo ya demasiado alejado del modelo biológico real estudiado. Por esto, se utilizó una métrica de variabilidad intra-grupos de $\rho \leq 0.75$.

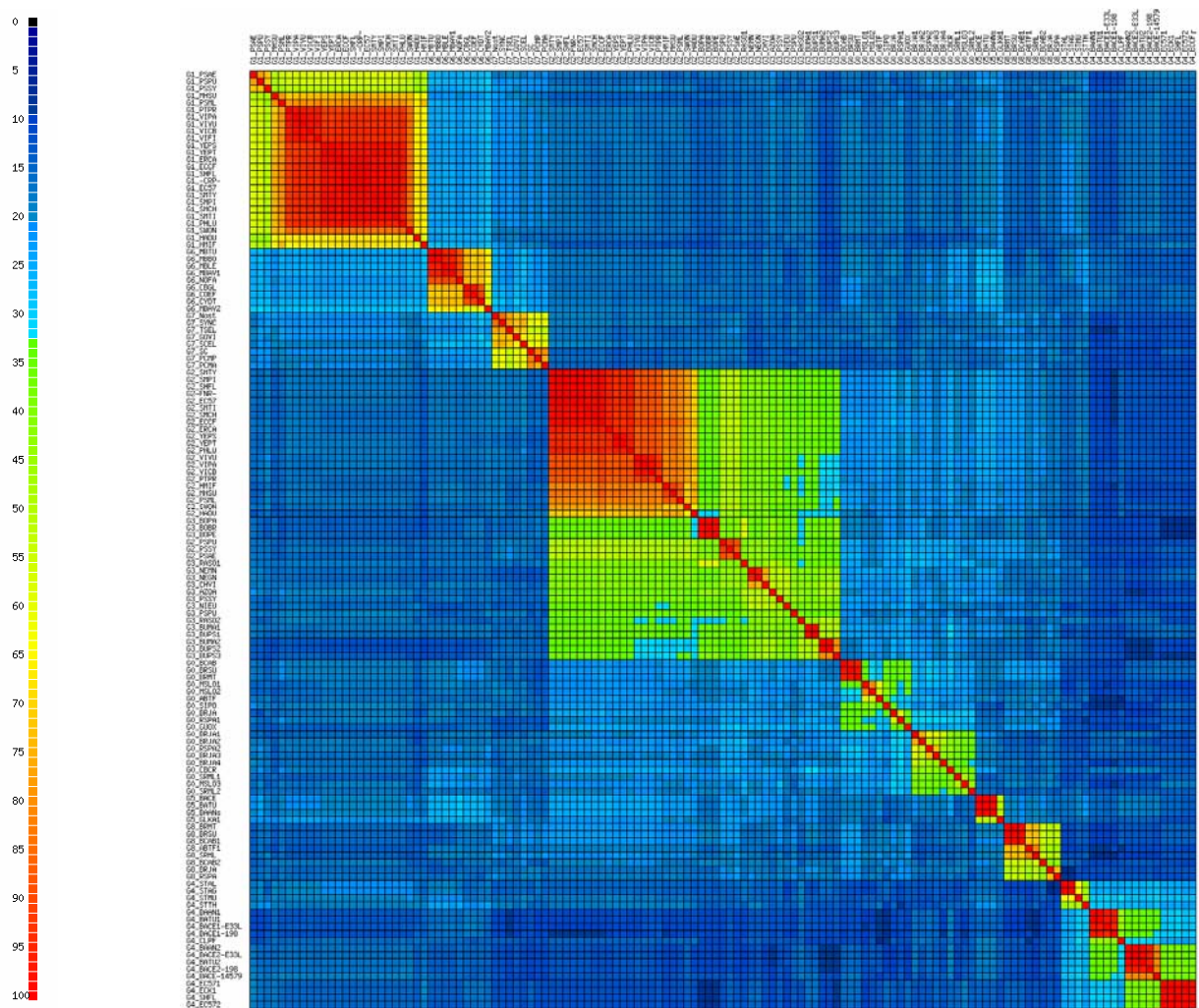


Figura 4R. Gráfica de comparación de identidades entre pares de secuencias. En la diagonal roja están las identidades máximas (1) y > 0.85 (rojo), de 0.75 a 0.65 (naranja-amarillo), y < 0.65 (tonos azules).

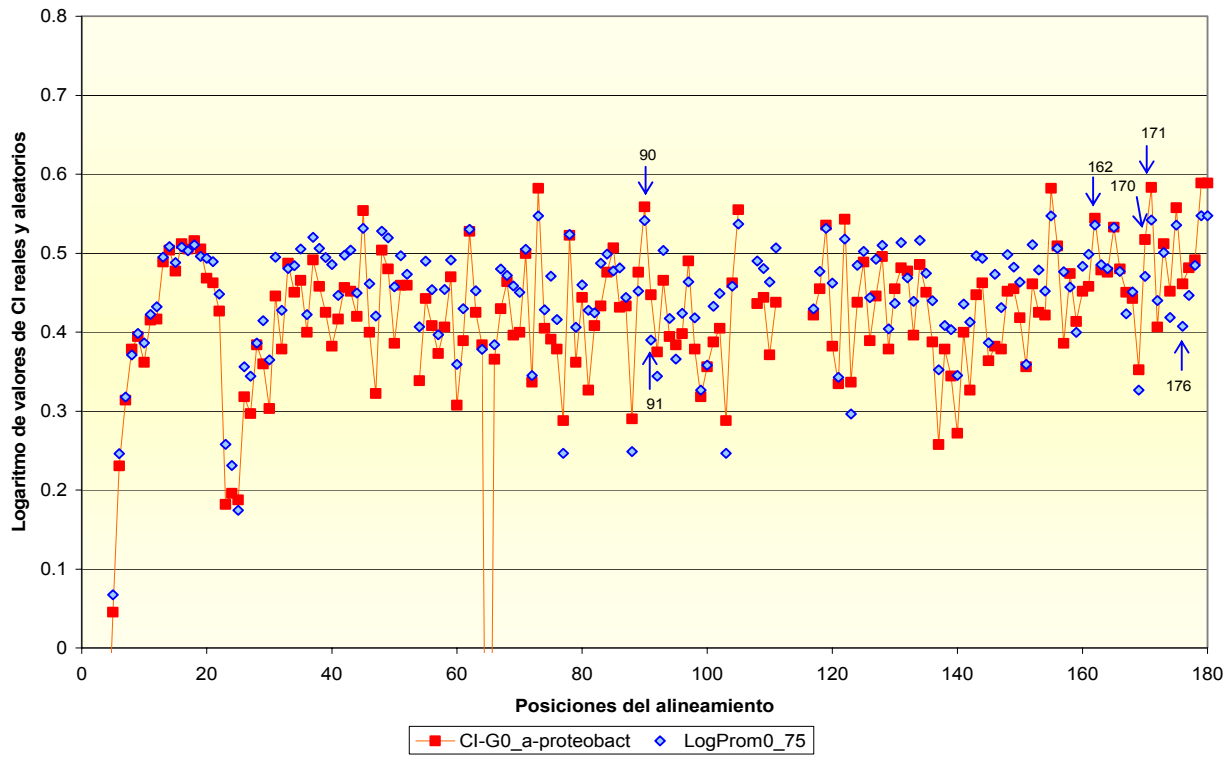
El cálculo de estas secuencias pseudos-aleatorias se realizó conforme a lo explicado en Métodos (3.5.a y 3.5.b) con la $\rho = 0.75$ con el programa en lenguaje de programación PERL Calc-PseudoRandom.pl (Anexo 7.6.2). Posteriormente se calcularon sus valores de *CI* con el mismo programa ContInfo-aa-alignment.pl (Anexo 7.6.1). A estos se les llamó *CI-referencia*, porque respecto a ellos se hizo la comparación estadística de los CIs de la familia real de Crp/Fnr, los *CIS-reales*.

Posteriormente, se procedió a la obtención de los valores Z (el valor con que el *CI-real* es tan similar o diferente al promedio de los *CIS-referencia*) para cada uno de los *CIS* de cada posición del alineamiento con la fórmula 2 explicada en Métodos (3.5.c), utilizando el programa MedsStats.pl (Anexo 7.6.4).

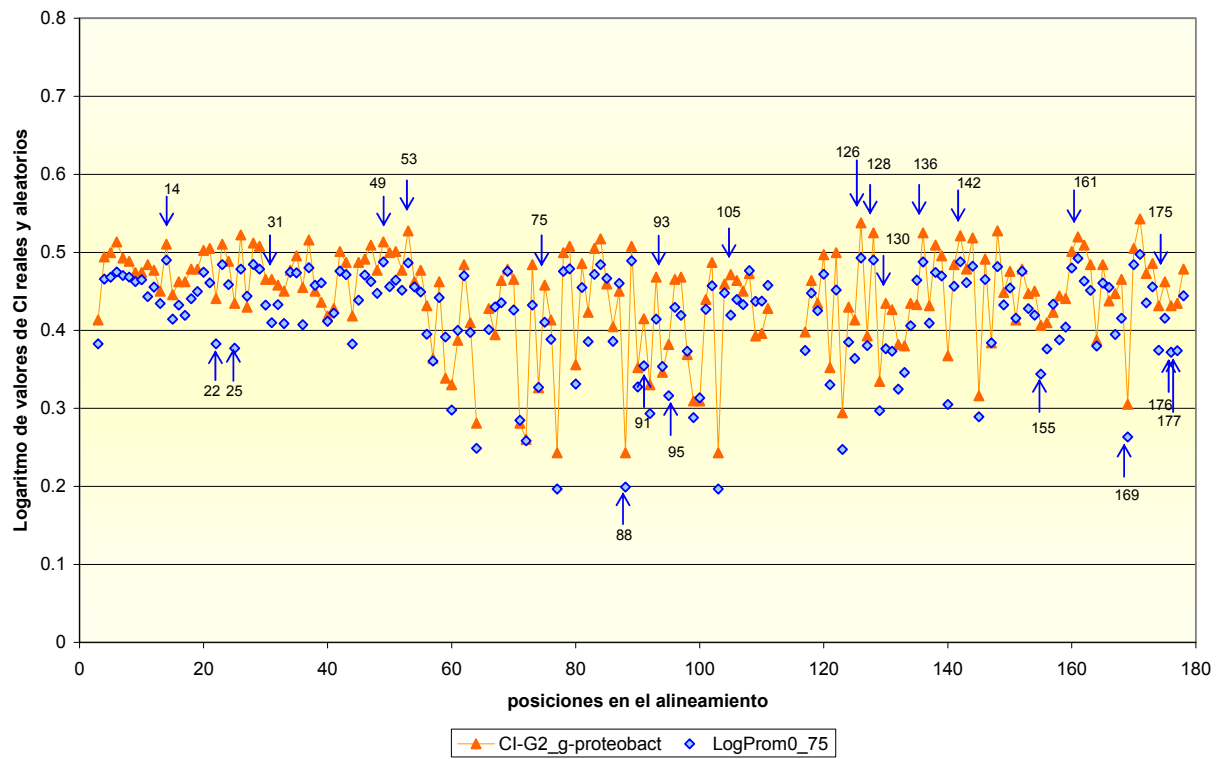
De los valores Z obtenidos para los *CI* (ver tablas en anexo 7.5), se consideraron *CIS* de residuos funcionalmente importantes aquellos que fueran ≥ 3 desviaciones estándar con un P-valor ≤ 0.05 (obtenido de tablas de valores Z, (Triola 1999)).

Para los grupos G1, G2, G3, G4 y G0, que son los grupos con más de 17 secuencias, sí se encontraron *CIS* con valores $Z \geq 3$ con un P-valor ≤ 0.05 , y se muestran en las gráficas 5R. Para los grupos G5, G6, G7 y G8, que tienen menos de 10 secuencias, se obtuvieron scores Z con valor máximo de ± 1.8 , que quedan cerca del valor promedio y no se pueden diferenciar significativamente del azar (ver Anexo 7.5, tablas con valores Z). Por lo tanto, en este estudio no se puede reportar ninguna predicción de RFIs para estos grupos. Se necesitará esperar a que haya más genomas secuenciados de estos grupos filogenéticos (*Actinobacteria*, *Cyanobacteria*, α -*proteobacteria* y *Firmicutes*), para volver a buscar ortólogos de Crp/Fnr y aumentar así el número de secuencias.

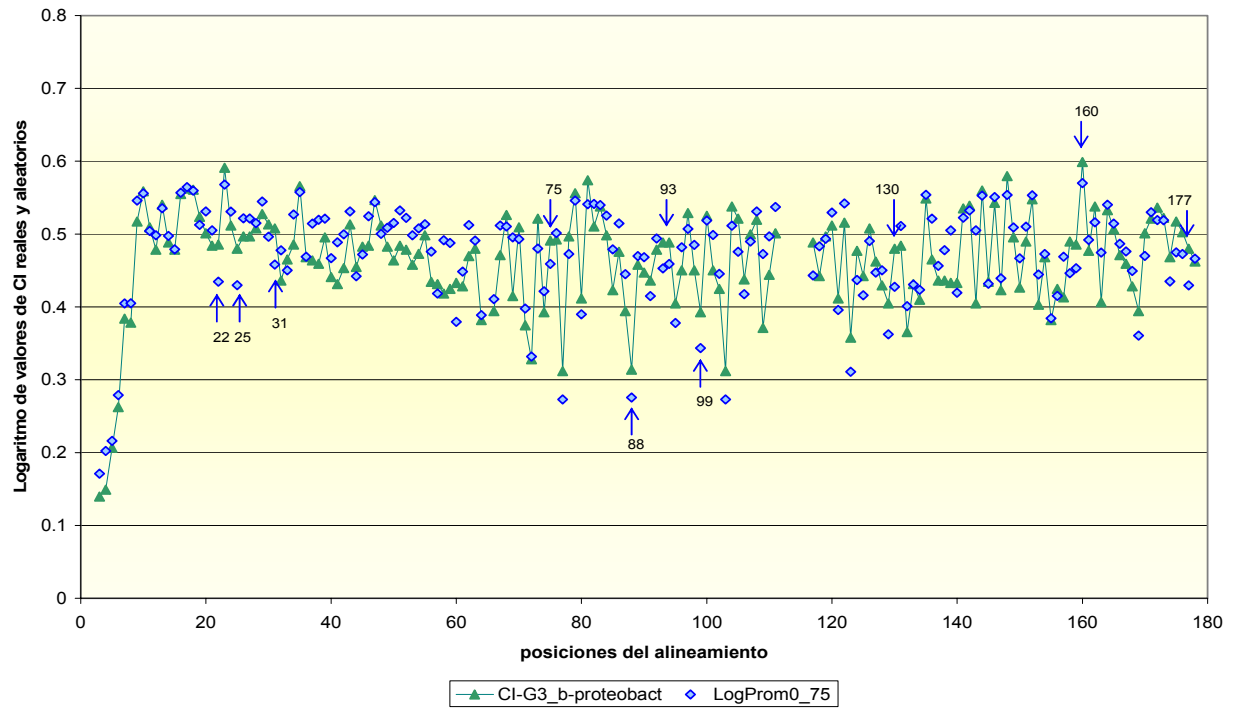
CIs reales y aleatorios del Grupo 0 - α -proteobacteria - de la familia Crp/Fnr



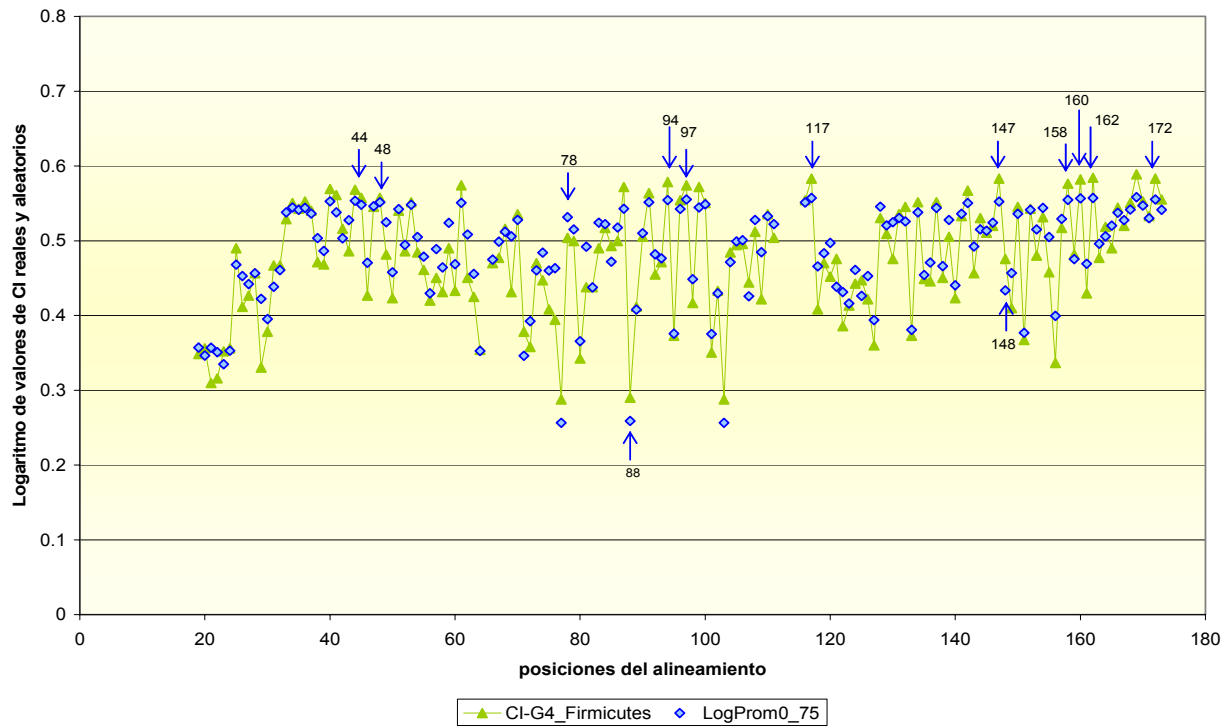
CIs reales y aleatorios del Grupo 2 - γ -proteobacteria - de la familia Crp/Fnr



Cis reales y aleatorios del Grupo 3 - β -proteobacteria - de la familia Crp/Fnr



Valores de CI reales y aleatorios del Grupo 4 -Firmicutes - de la familia Crp/Fnr



CIs reales y aleatorios del Grupo 1 -*γ-proteobacteria* - de la familia Crp/Fnr

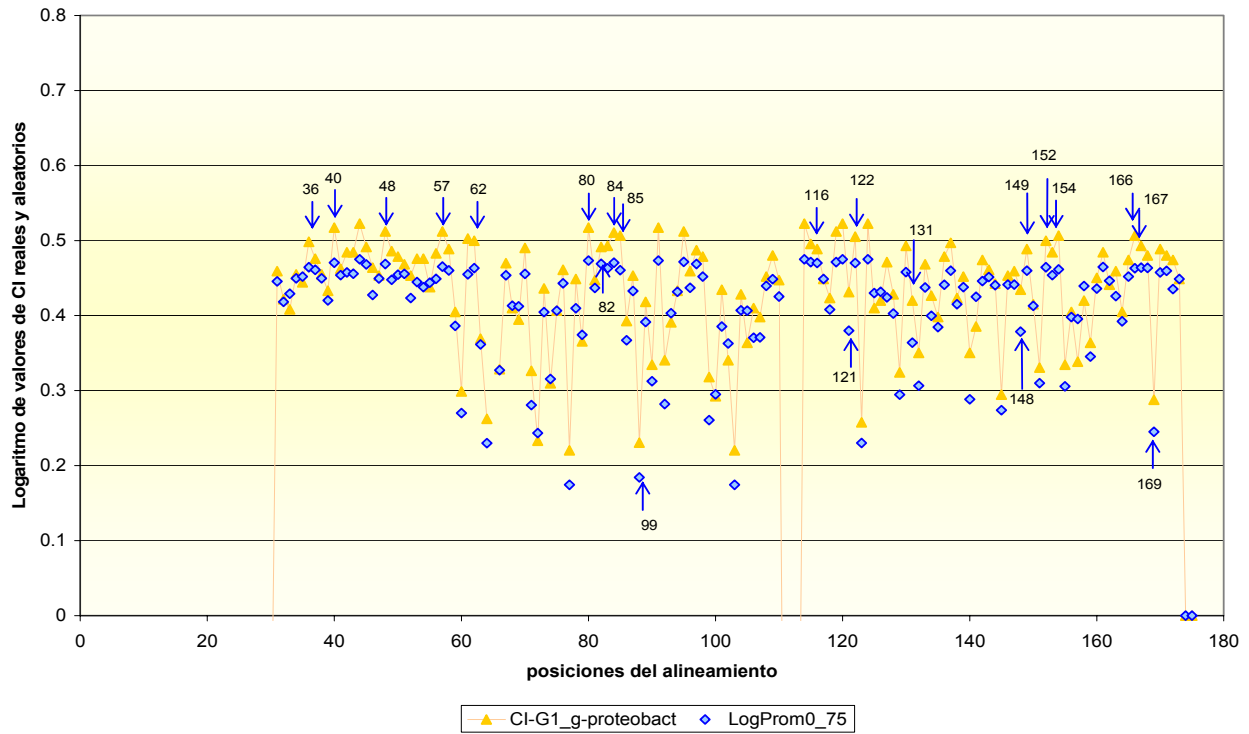


Figura 5R. Gráficas de los *CI-reales* y *CI-referencia* (con variación < 0.75) de las 180 posiciones del alineamiento para los grupos 2, 3, 0, 4 y 1. Los *CI* indicados con ↑ tienen valores $Z \geq 3$ y corresponden a las posiciones donde se encuentran residuos funcionalmente importantes.

Posterior a la identificación de los *CI-reales* con valores $Z \geq 3$ (con un valor $p \leq 0.05$), se procedió a identificar a que residuos corresponden en el alineamiento para cada grupo (figura 6R).

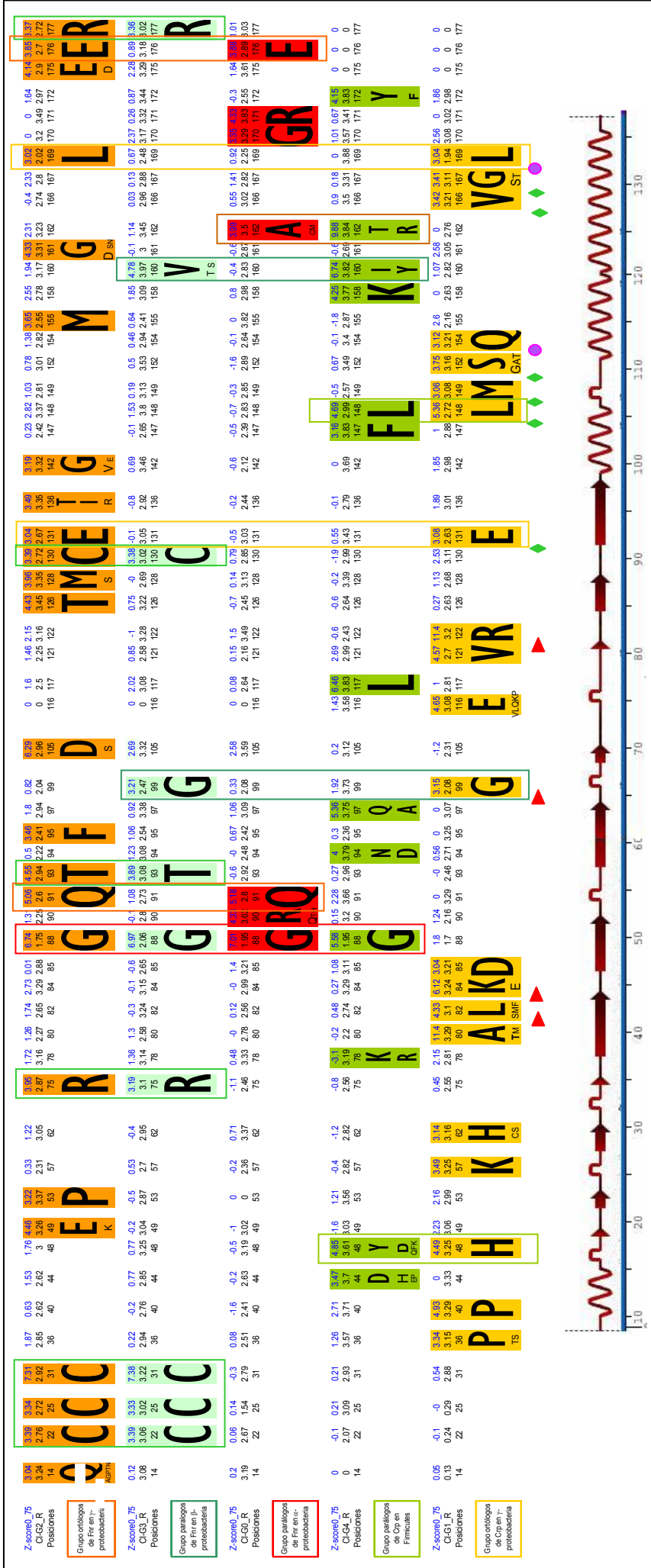


Figura 6R. Residuos Funcionalmente Importantes identificados para los grupos del alineamiento de la familia Crp/Fnr:

- G2 (naranja, tipo Fnr en γ -Proteobacteria),
- G3 (acua, tipo Fnr en β -Proteobacteria),
- G0 (rojo, tipo Fnr en α -Proteobacteria),
- G4 (verde, tipo Crp en Firmicutes) y
- G1 (amarillo, tipo Crp en γ -Proteobacteria).

Las posiciones, los valores de CI y Z se indican arriba de cada grupo.

Los residuos comunes entre grupos están enmarcados en rectángulos entre los grupos correspondientes.

En la parte inferior está la estructura secundaria de Crp, puesta como referencia de posible ubicación de los RFIs encontrados para los otros grupos.

4.3.1. Comparando entre grupos los residuos con *CI*s significativos, que se predicen como RFIs

Se encontró que entre cuatro de los cinco grupos, los grupos 2, 3, 0 y 4, comparten la Glicina 88 con *CI* significativo (figura 6R, marcada en rectángulo rojo). En particular, este aminoácido se ve conservado en todas las secuencias del alineamiento (ver Anexo 7.2). Es interesante que a pesar de ello, esta Glicina no haya resultado un residuo con *CI* significativo para el grupo 1 de Crp (ver tabla de G1 en Anexo 7.4).

En cuanto al grupo 2 de *γ-Proteobacteria* de los 24 residuos encontrados, se encontró que comparte específicamente siete RFIs iguales con el grupo 3 de *β-Proteobacteria* (figura 6R, rectángulos verde acua): Cisteinas 22, 25, 31 y 130, Argininas 75 y 177, y Treonina 93; dos específicamente con el grupo 0 de *α-Proteobacteria* (figura 6R, rectángulo rojo-naranja): Glutamina 91 y Aspartato 176; y dos específicamente con el grupo 1 de *γ-Proteobacteria* (figura 6R, rectángulos naranjas): Aspartato 131 y 176. Los restantes 12 residuos son específicos del grupo.

En cuanto a los 10 residuos encontrados para el grupo 3 de *β-Proteobacteria*, además de los ya mencionados que comparte específicamente con el grupo 2, comparte la Glicina 99 con el grupo 1 de *γ-Proteobacteria* (figura 6R, rectángulo verde agua) y la posición 160 con el grupo 4 de *Firmicutes* (figura 6R, rectángulos verde agua; Valina para el grupo 3, Isoleucina o Tirosina para el grupo 4); sólo queda un residuo único de este grupo.

Para los 7 residuos del grupo 0 de *α-Proteobacteria* comparte, además de los ya mencionados con los grupos 2 y 3, la posición 162 con el grupo 4 *Firmicutes* (figura 6R, rectángulos verde rojo-naranja; Alanina para grupo 0 y Treonina o Arginina para el grupo 4). Los otros tres residuos son únicos de este grupo.

De los 13 residuos encontrados para el grupo 4 de *Firmicutes*, además de los ya mencionados el grupo 0, comparte específicamente con el grupo 1 de *γ-Proteobacteria*

(figura 6R, rectángulos verde) la leucina 148, y la posición 48 (Tirosina o Aspartato para el grupo 4, e Histidina para el 1); los restantes 8 residuos son únicos de este grupo.

Finalmente, de los 21 residuos encontrados en el grupo 1, 16 son específicos de este grupo; el resto como ya se mencionó se comparten con los grupos 2,3 y 4 (figura 6R, rectángulos amarillos, azul y verdes).

Antes de poder hacer el análisis de estas coincidencias y diferencias mencionadas en esta sección de los RFIs encontrados en los cinco grupos de la familia, es necesario explicar lo que se conoce sobre Crp_{ECO} y Fnr_{ECO}, y mapear los RFIs encontrados para estos dos grupos en un modelo de estructura tridimensional.

4.3.2. Residuos con evidencia experimental, ligandos y conformaciones de Crp y Fnr

En general los estudios de mutagénesis de proteínas se han seguido de dos formas: con un análisis previo de la estructura para seleccionar los residuos a estudiar, o mutando residuos al azar. Con estas estrategias, los residuos que se identifican con algún efecto en la función o conformación, provienen de una selección si se parte de un estudio de estructura y dependen del tipo y condiciones del experimento realizado. Las características que se determinan a partir de estos experimentos pueden ser ambiguas funcionalmente, de función específica o sin relevancia en su función.

De los cinco grupos (2, 3, 0, 4 y 6) a los que se les pudo identificar residuos con *CI* con valores significativos, solo se cuenta con la referencia de residuos ya caracterizados para el DUL de Crp_{ECO} y algunos de los de Fnr_{ECO}; también se conocen detalles de su estado dentro de la célula y de cómo regulan.

Para Crp_{ECO}, la mayoría de los estudios parten de análisis de su estructura e inicialmente se habían enfocado en la identificación de los residuos que interaccionan específicamente con su ligando AMPc (Weber and Steitz 1987; Belduz, Lee et al. 1993;

Lee, Glasgow et al. 1994), luego para los que específicamente hacen la interacción entre las subunidades (Heyduk, Heyduk et al. 1992) y más recientemente para conocer como pasa la señal de interacción con el AMPc al dominio de interacción con el DNA y que modulan el cambio de conformación de la proteína (Baker, Tomlinson et al. 2001; Lin, Kovac et al. 2002; Gekko, Obu et al. 2004) (ver figura 4c en Introducción). Además, se sabe que Crp_{ECO} interacciona con AMPc sólo cuando la concentración de este metabolito aumenta, se mantiene como dímero unido o despegado del DNA y sólo cambia un poco su conformación para unirse al DNA (ver figura 4 en Introducción).

Para Fnr_{ECO}, se tienen descritas cuatro Cisteínas indispensables para la interacción con el grupo prostético 4Fe4S (Lazazzera, Beinert et al. 1996; Khoroshilova, Popescu et al. 1997); además Fnr_{ECO} siempre tiene unido a su grupo prostético 4Fe4S que cambia sólo de estado oxido-reducción, lo que provoca que Fnr_{ECO} cuando no está unido al DNA esté como monómero y que se conforme como dímero cuando se une al DNA (Kiley et. al. 1999).

4.3.3. De los residuos con evidencia experimental en Crp_{ECO} y los RFIs predichos para su grupo en *γ-Proteobacteria*

Para apoyar el análisis de los residuos encontrados para el grupo 1 de Crp en *γ-Proteobacterias*, estos se mapearon en una de las estructuras tridimensionales ya reportadas para Crp de *E. coli*, la 1J59 en el PDB, utilizando el software de Protein Explorer 2.8 [www.proteinexplorer.org], (tabla 2R y figura 7R).

Tabla 2R. Residuos encontrados en este trabajo con *CI*s significativos

Región en Estructura	Orientación	Posición en Tabla RFIs	Residuo	Posic Estructura	Descripción Función
	E	36	Prolina	--	
α-hélice	E	40	Prolina	9	
	E	48	Histidina	17	
hojas β y asas	E	57	Lisina	26	

	E	62	Histidina	31	
	O	80	Alanina	48	
	Cúmulo espacial	82	Leucina	50	
	Cúmulo espacial	84	Lisina	52	Se han determinado como residuos importantes en la transmisión de la señal de activación posterior a la unión de AMPc (Gekko 2004, Lin 2002)
	poco E	85	Aspartato	53	Se han determinado como residuos importantes en la transmisión de la señal de activación posterior a la unión de AMPc (Gekko 2004, Lin 2004)
	Cúmulo espacial	99	Glicina	67	
			Glutamato	72	Después del estudio de la estructura de Crp, se probó por mutagénesis que es importante para la interacción con el AMPc (Belduz et. al. 1993)
	E	116	Glutamato	81	
			Arginina	82	Después del estudio de la estructura de Crp, se probó por mutagénesis que es importante para la interacción con el AMPc (Belduz et. al. 1993)
			Serina	83	Después del estudio de la estructura de Crp, se probó por mutagénesis que es importante para la interacción con el AMPc (Lee et. al. 1994)
	orientada AMPc	121	Valina	86	
	Cúmulo espacial	122	Arginina	87	
	poco E	131	Glutamato	96	
			Tirosina	99	Contribuye en la estabilidad de la estructura dimerica, en la asociación entre las subunidades (Baker 2004)
α -hélice	contacto I-subU	148	Leucina	113	
	contacto I-subU	149	Metionina	114	
	contacto I-subU	152	Serina	117	es blanco de proteólisis cuando es mutada (Heyduk et. al. 1992)
	hacia adentro	154	Glutamina	119	
			Treonina	127	Después del estudio de la estructura de Crp, se probó por mutagénesis que es importante para la interacción con el AMPc (Lee et. al. 1994)
			Serina	128	Después del estudio de la estructura de Crp, se probó por mutagénesis que es importante para la interacción con el AMPc (Lee et. al. 1994)
	contacto I-subU	166	Valina	131	
	contacto I-subU	167	Glicina	132	residuo esencial para transmitir las señales de las interacciones alostéricas (Yu et. al 2004)
	hacia adentro	169	Leucina	134	
región visagra			Fenilalanina	136	es blanco de proteólisis cuando es mutada (Heyduk et. al. 1992)

Residuos ordenados según su ubicación (columna 1) y orientación (columna 2; E: externo, O: oculto, I-subU: inter-subunidad) en la estructura de Crp (de la figura 7R); el número de su posición en el alineamiento (columna 3) (figura 6R); su número de posición en la estructura (columna 5); tipo de función encontrada con estudios de mutagénesis (columna 6). Las líneas en gris corresponden a los residuos no identificados con el análisis de *CIS*.

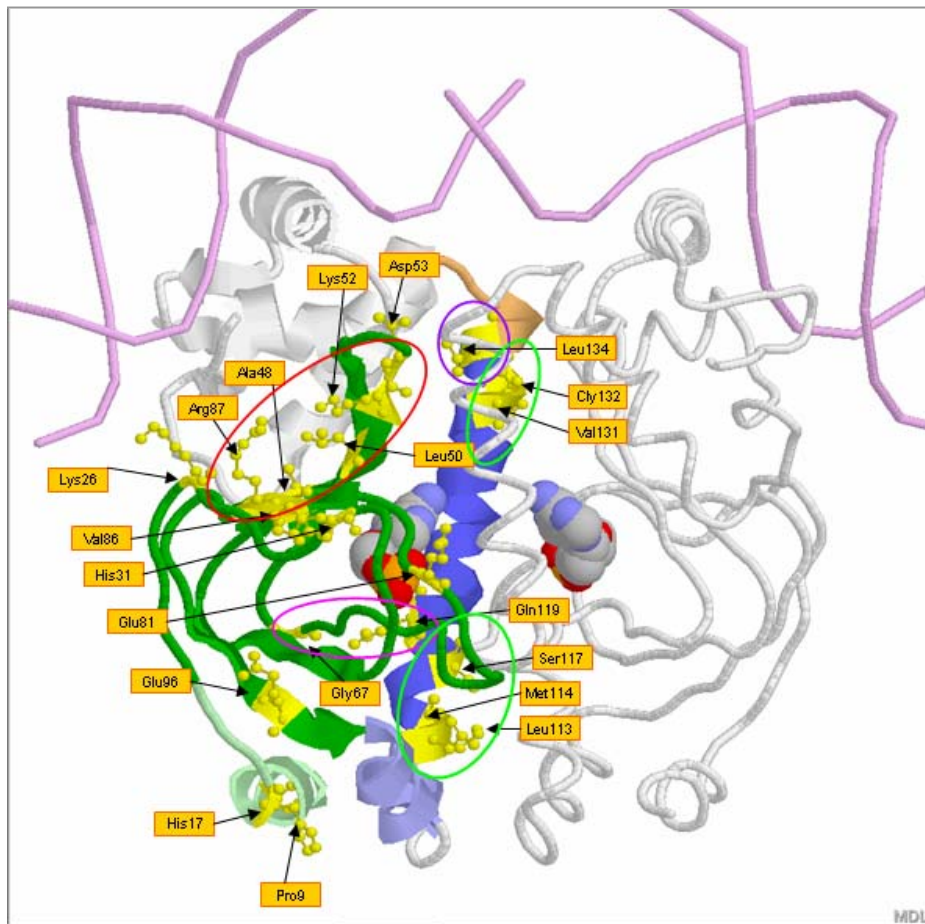


Figura 7R. Gráfica de la estructura de Crp, 1J59, con los RFIs predichos en este trabajo mapeados en color morado, en la gráfica se anotan con las posiciones en la estructura y aquí se anota la posición a la que corresponden en el alineamiento: Pro9 = P40, His17 = H48, Lys26 = K57, His31 = H62, Ala48 = A80, Leu50 = L82, Lys52 = K84, Asp53 = D85, Gly67= G99, Glu81 = E116, Val86 = V121, Arg87 = R122, Glu96 = E131, Leu113 = L148, Met114 = M149, Ser117 = S152, Gln119 = Q154, Val131 = V166, Gly132 = G167, Leu134 = L169. Ver en texto anotaciones de los círculos de colores.

Los residuos dentro de la α -hélice con que interaccionan los dominios se marcan dentro de ovalos verdes (marcados con \blacktriangleright en figura 6R) y los que están orientados hacia el propio monómero se marcan en ovalos morado y rosa (marcados con \bullet en la figura 6R), el grupo que se encuentra formando un cúmulo de residuos cercanos espacialmente se indican dentro de círculo rojo (marcados con \blacktriangle en figura 6R). Los residuos comunes con el grupo de Fnr en γ -Proteobacteria se indican en los dos círculos rojos, y los comunes

con el grupo 4 tipo Crp en *Firmicutes* en círculos verdes.

4.3.4. De los residuos con referencia en Fnr_{ECO} y los RFIs predichos para su grupo en γ -Proteobacteria

Aunque para Fnr *E. coli* se cuenta con poca información respecto a su estructura y conformación, respecto a los residuos ya reportados que interaccionan con su grupo prostético, para el grupo de Fnr en γ -Proteobacteria (G2) se encontraron las mismas Cisteínas con *CI* significativos; y no sólo en el grupo 2, sino que también se comparten con *CI* significativos con el grupo G3 en β -Proteobacteria.

Se propone que las proteínas del grupo G3 también unen un grupo prostético 4Fe4S dado que también tienen las cuatro cisteínas con *CI* significativos y en las mismas posiciones que el grupo 2 de Fnr.

4.4. Eficiencia del método, comparando con lo ya reportado y las características de las proteínas

De acuerdo a la comparación de los residuos ya reportados *versus* los residuos aquí encontrados como funcionalmente importantes (con valores de *CI* significativos), se encontraron pocas coincidencias. Por una parte, de los residuos conocidos para Crp_{ECO} hicieron falta varios que aunque en el alineamiento se conservan, no se califican como significativos. Esto puede deberse a que los RFIs detectados con este método no están involucrados con la interacción del ligando sino a que están más relacionados con la interacción inter-subunidades, el núcleo de conformación y/o la velocidad de plegamiento.

Por otra parte, para Fnr_{ECO} se encontraron las cuatro cisteínas que interaccionan con su grupo prostético y otros residuos, que por falta de caracterización estructural, no se puede decir si están relacionados con el núcleo de conformación y/o la habilidad de plegamiento rápido.

La diferencia de haber identificado a los residuos que interaccionan con el ligando para el grupo 2 de Fnr y no para el grupo 1 de Crp se puede explicar por el hecho de que Fnr en cuanto es sintetizada se une a un grupo prostético Hierro-Azufre que se oxida y reduce ($2\text{Fe}-2\text{S} \rightarrow 4\text{Fe}_4\text{S}$) en respuesta los niveles de Oxígeno, mientras que Crp sólo se une al AMPc cuando la concentración de este metabolito aumenta (ver Introducción). En el caso de Fnr, estas cuatro cisteínas permanecen ocultas por su interacción constante con el grupo prostético; además que cuando la célula detecta que Fnr no tiene a su grupo unido (Fnr-apo) lo degrada. En el caso de Crp, los residuos que interaccionan con el AMPc están expuestos al medio, si alguno de ellos se muta Crp puede perder la especificidad (reconocer otros nucleótidos monofosfatados) o quedar en la conformación activa y en algunos casos podría perder su capacidad de regulación, pero no es degradada. Esta diferencia puede marcar la diferente presión selectiva en estas proteínas para conservar ciertos residuos funcionales.

5. Conclusiones generales y perspectivas

5.1. Conclusiones generales

En el presente estudio se calculó el Contenido Informacional para identificar los residuos funcionalmente importantes en los dominios de unión al ligando para los miembros de la familia de los factores de transcripción Crp/Fnr.

Esta herramienta es dependiente del número de secuencias (tamaño de la muestra de las secuencias de la familia Crp/Fnr), por lo que su análisis estadístico es indispensable. Esta dependencia hizo que la identificación de estos residuos fuera inviable para los grupos que tuvieron menos de 10 secuencias, los grupos 5, 6, 7 y 8. El hecho de tener pocas secuencias implicó que no tuvieran una significancia estadística. Por esta dependencia, sólo se predijeron diferentes residuos potencialmente funcionales y estadísticamente significativos para los grupos 2, 3, 0, 4 y 6.

La dependencia de *CI*s al número de secuencia fue observable desde su graficación (figura 3R.B) con *CI*s que se gráficán separados a la tendencia general los otros grupos.

En general los RFIs identificados corresponden a todos los aminoácidos biológicamente conocidos, excepto triptófano.

Todavía falta completar la caracterización estructural de los diferentes grupos de la familia; sin embargo, se tiene la certeza de que los valores de *CI* que han sido calificados como significativos, han sido aquellos que cumplen con criterios estadísticos estrictos (con valores *Z* por arriba de 3 desviaciones estándar con un *p*-valor menor a 0.05), lo cual nos permite hacer predicciones sobre estos residuos en los grupos aún sin conocer su ubicación espacial dentro de la estructura.

Para varios grupos no se tiene referencia de residuos caracterizados. Falta una mayor caracterización experimental de las interacciones de estos dominios; así que tener predicciones de RFIs con valores de CI estadísticamente significativos, se pueden utilizar como guía para identificar residuos que puedan ser importantes para la conformación de estas proteínas, así como aquellos que no participen en ninguna función.

Pocos de los residuos identificados han sido reportados previamente con evidencia experimental. La gran mayoría corresponde a residuos que tentativamente se encuentran fuera del sitio físico de interacción con el ligando, pero que pueden estar involucrados en el núcleo de conformación y/o en la habilidad de plegamiento rápido de estas proteínas.

En general, los residuos que se predicen como funcionalmente importantes en este trabajo pertenecen a regiones de las proteínas que aún no se han descrito por completo bioquímicamente, en un momento en que las herramientas con las que se cuenta son sólo la observación de su conservación y posición en la estructura terciaria.

La identificación de residuos funcionalmente importantes compete al área de la genómica que nos permite completar la caracterización de los dominios de unión al ligando de proteínas como factores transcripcionales. Esto nos permite conocer los detalles bioquímicos de las interacciones con los ligandos, entre subunidades y otras funciones que hacen específica la respuesta a la regulación transcripcional de estos factores en diferentes organismos.

El método propuesto y utilizado en este trabajo permitió la identificación de RFIs sin depender de una conservación de posiciones en un alineamiento múltiple entre los miembros de los diferentes grupos. Inclusive se pudieron identificar varios residuos en posiciones comunes entre los grupos de forma independiente; además, de ser un método independiente de la disponibilidad de estructuras tridimensionales de las proteínas de la familia.

5.2. Perspectivas

El análisis computacional para identificar posibles Residuos Funcionalmente Importantes (RFIs) en la familia Crp/Fnr nos ha permitido abrir diversas preguntas que pueden continuar, tales como, la caracterización estructural de diversos miembros representativos de cada grupo (obtención de su cristal, construcción y comparación de modelos tridimensionales entre los diferentes grupos de proteínas) y la comprobación experimental de los residuos aquí determinados como importantes para la función. A su vez, estos estudios podrían arrojar luz en el conocimiento del comportamiento de estas proteínas en cuanto a sus interacciones específicas con diversos compuestos, o aquellos involucrados en los cambios de conformación como consecuencia de la unión a su ligando. Por otra parte, este enfoque se puede extender a las otras familias de factores transcripcionales que se han descrito hasta la fecha.

Finalmente, el método podría ser mejorado implementando una matriz de peso para medir los cambios de aminoácidos en lugar de sólo medir los cambios de identidades de aminoácidos.

6. Referencias

- Altschul, S. F., W. Gish, et al. (1990). "Basic local alignment search tool." J Mol Biol **215**(3): 403-10.
- Armon, A., D. Graur, et al. (2001). "ConSurf: an algorithmic tool for the identification of functional regions in proteins by surface mapping of phylogenetic information." J Mol Biol **307**(1): 447-63.
- Baker, C. H., S. R. Tomlinson, et al. (2001). "Amino acid substitution at position 99 affects the rate of CRP subunit exchange." Biochemistry **40**(41): 12329-38.
- Barker, W. C., J. S. Garavelli, et al. (1999). "The PIR-International Protein Sequence Database." Nucleic Acids Res **27**(1): 39-43.
- Belduz, A. O., E. J. Lee, et al. (1993). "Mutagenesis of the cyclic AMP receptor protein of Escherichia coli: targeting positions 72 and 82 of the cyclic nucleotide binding pocket." Nucleic Acids Res **21**(8): 1827-35.
- Benson, D. A., I. Karsch-Mizrachi, et al. (2003). "GenBank." Nucleic Acids Res **31**(1): 23-7.
- Berezin, C., F. Glaser, et al. (2004). "ConSeq: the identification of functionally and structurally important residues in protein sequences." Bioinformatics **20**(8): 1322-4.
- Berman, H. M., J. Westbrook, et al. (2000). "The Protein Data Bank." Nucleic Acids Res **28**(1): 235-42.
- Boeckmann, B., A. Bairoch, et al. (2003). "The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003." Nucleic Acids Res **31**(1): 365-70.
- Brennan, R. G. and B. W. Matthews (1989). "The helix-turn-helix DNA binding motif." J Biol Chem **264**(4): 1903-6.
- Browning, D. F. and S. J. Busby (2004). "The regulation of bacterial transcription initiation." Nat Rev Microbiol **2**(1): 57-65.
- Burgess, R. R. and L. Anthony (2001). "How sigma docks to RNA polymerase and what sigma does." Curr Opin Microbiol **4**(2): 126-31.
- Burgess, R. R., A. A. Travers, et al. (1969). "Factor stimulating transcription by RNA polymerase." Nature **221**(5175): 43-6.
- Crick, F. (1970). "Central dogma of molecular biology." Nature **227**(5258): 561-3.
- Chothia, C. and A. M. Lesk (1986). "The relation between the divergence of sequence and structure in proteins." Embo J **5**(4): 823-6.
- Danot, O. (2001). "A complex signaling module governs the activity of MalT, the prototype of an emerging transactivator family." Proc Natl Acad Sci U S A **98**(2): 435-40.
- Dibden, D. P. and J. Green (2005). "In vivo cycling of the Escherichia coli transcription factor FNR between active and inactive states." Microbiology **151**(Pt 12): 4063-70.
- Eddy, S. R. (1998). "Profile hidden Markov models." Bioinformatics **14**(9): 755-63.
- Garges, S. and S. Adhya (1988). "Cyclic AMP-induced conformational change of cyclic AMP receptor protein (CRP): intragenic suppressors of cyclic AMP-independent CRP mutations." J Bacteriol **170**(4): 1417-22.
- Gekko, K., N. Obu, et al. (2004). "A linear correlation between the energetics of allosteric communication and protein flexibility in the Escherichia coli cyclic AMP receptor protein revealed by mutation-induced changes in compressibility and amide hydrogen-deuterium exchange." Biochemistry **43**(13): 3844-52.
- Gough, J. (2002). "The SUPERFAMILY database in structural genomics." Acta

Crystallogr D Biol Crystallogr **58**(Pt 11): 1897-900.

- Green, J., C. Scott, et al. (2001). "Functional versatility in the CRP-FNR superfamily of transcription factors: FNR and FLP." Adv Microb Physiol **44**: 1-34.
- Gruber, T. M. and C. A. Gross (2003). "Multiple sigma subunits and the partitioning of bacterial transcription space." Annu Rev Microbiol **57**: 441-66.
- Gunsalus, R. P. and S. J. Park (1994). "Aerobic-anaerobic gene regulation in Escherichia coli: control by the ArcAB and Fnr regulons." Res Microbiol **145**(5-6): 437-50.
- Haldenwang, W. G. (1995). "The sigma factors of Bacillus subtilis." Microbiol Rev **59**(1): 1-30.
- Harrison, S. C. (1991). "A structural taxonomy of DNA-binding domains." Nature **353**(6346): 715-9.
- Heyduk, E., T. Heyduk, et al. (1992). "Intersubunit communications in Escherichia coli cyclic AMP receptor protein: studies of the ligand binding domain." Biochemistry **31**(14): 3682-8.
- Holm, L. (1998). "Unification of protein families." Curr Opin Struct Biol **8**(3): 372-9.
- Hutchison, C. A., 3rd, S. Phillips, et al. (1978). "Mutagenesis at a specific position in a DNA sequence." J Biol Chem **253**(18): 6551-60.
- Innis, C. A., J. Shi, et al. (2000). "Evolutionary trace analysis of TGF-beta and related growth factors: implications for site-directed mutagenesis." Protein Eng **13**(12): 839-47.
- James, L. C. and D. S. Tawfik (2003). "Conformational diversity and protein evolution--a 60-year-old hypothesis revisited." Trends Biochem Sci **28**(7): 361-8.
- Johnson, J. M. and G. M. Church (2000). "Predicting ligand-binding function in families of bacterial receptors." Proc Natl Acad Sci U S A **97**(8): 3965-70.
- Jones, D. T. (1999). "Protein secondary structure prediction based on position-specific scoring matrices." J Mol Biol **292**(2): 195-202.
- Kawabata, T., M. Ota, et al. (1999). "The Protein Mutant Database." Nucleic Acids Res **27**(1): 355-7.
- Khoroshilova, N., C. Popescu, et al. (1997). "Iron-sulfur cluster disassembly in the FNR protein of Escherichia coli by O2: [4Fe-4S] to [2Fe-2S] conversion with loss of biological activity." Proc Natl Acad Sci U S A **94**(12): 6087-92.
- Kiley, P. J. and H. Beinert (1999). "Oxygen sensing by the global regulator, FNR: the role of the iron-sulfur cluster." FEMS Microbiol Rev **22**(5): 341-52.
- Korner, H., H. J. Sofia, et al. (2003). "Phylogeny of the bacterial superfamily of Crp-Fnr transcription regulators: exploiting the metabolic spectrum by controlling alternative gene programs." FEMS Microbiol Rev **27**(5): 559-92.
- Kumar, S., K. Tamura, et al. (2004). "MEGA3: Integrated software for Molecular Evolutionary Genetics Analysis and sequence alignment." Brief Bioinform **5**(2): 150-63.
- Landgraf, R., D. Fischer, et al. (1999). "Analysis of heregulin symmetry by weighted evolutionary tracing." Protein Eng **12**(11): 943-51.
- Landgraf, R., I. Xenarios, et al. (2001). "Three-dimensional cluster analysis identifies interfaces and functional residue clusters in proteins." J Mol Biol **307**(5): 1487-502.
- Lawson, C. L., D. Swigon, et al. (2004). "Catabolite activator protein: DNA binding and transcription activation." Curr Opin Struct Biol **14**(1): 10-20.
- Lazazzera, B. A., H. Beinert, et al. (1996). "DNA binding and dimerization of the Fe-S-

containing FNR protein from *Escherichia coli* are regulated by oxygen." J Biol Chem **271**(5): 2762-8.

- Lee, E. J., J. Glasgow, et al. (1994). "Mutagenesis of the cyclic AMP receptor protein of *Escherichia coli*: targeting positions 83, 127 and 128 of the cyclic nucleotide binding pocket." Nucleic Acids Res **22**(15): 2894-901.
- Lewin, B. (2000). Genes VII. New York, Oxford University Press.
- Lewis, M. (2005). "The lac repressor." C R Biol **328**(6): 521-48.
- Li, B., H. Wing, et al. (1998). "Transcription activation by *Escherichia coli* FNR protein: similarities to, and differences from, the CRP paradigm." Nucleic Acids Res **26**(9): 2075-81.
- Li, W., L. Jaroszewski, et al. (2002). "Tolerating some redundancy significantly speeds up clustering of large protein databases." Bioinformatics **18**(1): 77-82.
- Lichtarge, O., H. Yao, et al. (2003). "Accurate and scalable identification of functional sites by evolutionary tracing." J Struct Funct Genomics **4**(2-3): 159-66.
- Lin, S. H., L. Kovac, et al. (2002). "Ability of *E. coli* cyclic AMP receptor protein to differentiate cyclic nucleotides: effects of single site mutations." Biochemistry **41**(9): 2946-55.
- Madan Babu, M. and S. A. Teichmann (2003). "Evolution of transcription factors and the gene regulatory network in *Escherichia coli*." Nucleic Acids Res **31**(4): 1234-44.
- Madan Babu, M. and S. A. Teichmann (2003). "Functional determinants of transcription factors in *Escherichia coli*: protein families and binding sites." Trends Genet **19**(2): 75-9.
- Makarova, K. S., A. A. Mironov, et al. (2001). "Conservation of the binding site for the arginine repressor in all bacterial lineages." Genome Biol **2**(4): RESEARCH0013.
- McGuffin, L. J., K. Bryson, et al. (2000). "The PSIPRED protein structure prediction server." Bioinformatics **16**(4): 404-5.
- Mirny, L. A. and E. I. Shakhnovich (1999). "Universally conserved positions in protein folds: reading evolutionary signals about stability, folding kinetics and function." J Mol Biol **291**(1): 177-96.
- Mironov, A. A., E. V. Koonin, et al. (1999). "Computer analysis of transcription regulatory patterns in completely sequenced bacterial genomes." Nucleic Acids Res **27**(14): 2981-9.
- Moreno-Campuzano, S., S. C. Janga, et al. (2006). "Identification and analysis of DNA-binding transcription factors in *Bacillus subtilis* and other Firmicutes--a genomic approach." BMC Genomics **7**: 147.
- Mount, D. W. (2001). Bioinformatics: Sequence and Genome Analysis., Cold Spring Harbor Laboratory.
- Murzin, A. G., S. E. Brenner, et al. (1995). "SCOP: a structural classification of proteins database for the investigation of sequences and structures." J Mol Biol **247**(4): 536-40.
- Pearl, F. M., C. F. Bennett, et al. (2003). "The CATH database: an extended protein family resource for structural and functional genomics." Nucleic Acids Res **31**(1): 452-5.
- Pearson, W. R. (2000). "Flexible sequence similarity searching with the FASTA3 program package." Methods Mol Biol **132**: 185-219.
- Perez-Rueda, E. (1999). Reconocimiento de motivos estructurales en reguladores transcripcionales. Centro de Investigación sobre Fijación de Nitrógeno.

Cuernavaca, Universidad Nacional Autónoma de México. **Doctorado:** 70.

- Perez-Rueda, E. and J. Collado-Vides (2000). "The repertoire of DNA-binding transcriptional regulators in *Escherichia coli* K-12." *Nucleic Acids Res* **28**(8): 1838-47.
- Perez-Rueda, E. and J. Collado-Vides (2001). "Common history at the origin of the position-function correlation in transcriptional regulators in archaea and bacteria." *J Mol Evol* **53**(3): 172-9.
- Perez-Rueda, E., J. Collado-Vides, et al. (2004). "Phylogenetic distribution of DNA-binding transcription factors in bacteria and archaea." *Comput Biol Chem* **28**(5-6): 341-50.
- Rodionov, D. A., A. G. Vitreschak, et al. (2004). "Comparative genomics of the methionine metabolism in Gram-positive bacteria: a variety of regulatory systems." *Nucleic Acids Res* **32**(11): 3340-53.
- Saito, S., M. Sasai, et al. (1997). "Evolution of the folding ability of proteins through functional selection." *Proc Natl Acad Sci U S A* **94**(21): 11324-8.
- Salgado, H., S. Gama-Castro, et al. (2006). "RegulonDB (version 5.0): *Escherichia coli* K-12 transcriptional regulatory network, operon organization, and growth conditions." *Nucleic Acids Res* **34**(Database issue): D394-7.
- Shaw, D. J., D. W. Rice, et al. (1983). "Homology between CAP and Fnr, a regulator of anaerobic respiration in *Escherichia coli*." *J Mol Biol* **166**(2): 241-7.
- Sonnhammer, E. L., S. R. Eddy, et al. (1998). "Pfam: multiple sequence alignments and HMM-profiles of protein domains." *Nucleic Acids Res* **26**(1): 320-2.
- Spiro, S. and J. R. Guest (1990). "FNR and its role in oxygen-regulated gene expression in *Escherichia coli*." *FEMS Microbiol Rev* **6**(4): 399-428.
- Suzuki, M. and N. Yagi (1996). "An in-the-groove view of DNA structures in complexes with proteins." *J Mol Biol* **255**(5): 677-87.
- Tan, K., G. Moreno-Hagelsieb, et al. (2001). "A comparative genomics approach to prediction of new members of regulons." *Genome Res* **11**(4): 566-84.
- Thompson, J. D., T. J. Gibson, et al. (1997). "The CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools." *Nucleic Acids Res* **25**(24): 4876-82.
- Triola, M. F. (1999). *Estadística*. Mexico, Pearson Educación.
- Weber, I. T. and T. A. Steitz (1987). "Structure of a complex of catabolite gene activator protein and cyclic AMP refined at 2.5 Å resolution." *J Mol Biol* **198**(2): 311-26.
- Wilson, C. A., J. Kreychman, et al. (2000). "Assessing annotation transfer for genomics: quantifying the relations between protein sequence, structure and function through traditional and probabilistic scores." *J Mol Biol* **297**(1): 233-49.
- Yasuzawa, K., N. Hayashi, et al. (1992). "Histone-like proteins are required for cell growth and constraint of supercoils in DNA." *Gene* **122**(1): 9-15.
- Zubay, G., D. Schwartz, et al. (1970). "Mechanism of activation of catabolite-sensitive genes: a positive control system." *Proc Natl Acad Sci U S A* **66**(1): 104-10.
- Zvelebil, M. J., G. J. Barton, et al. (1987). "Prediction of protein secondary structure and active sites using the alignment of homologous sequences." *J Mol Biol* **195**(4): 957-61.

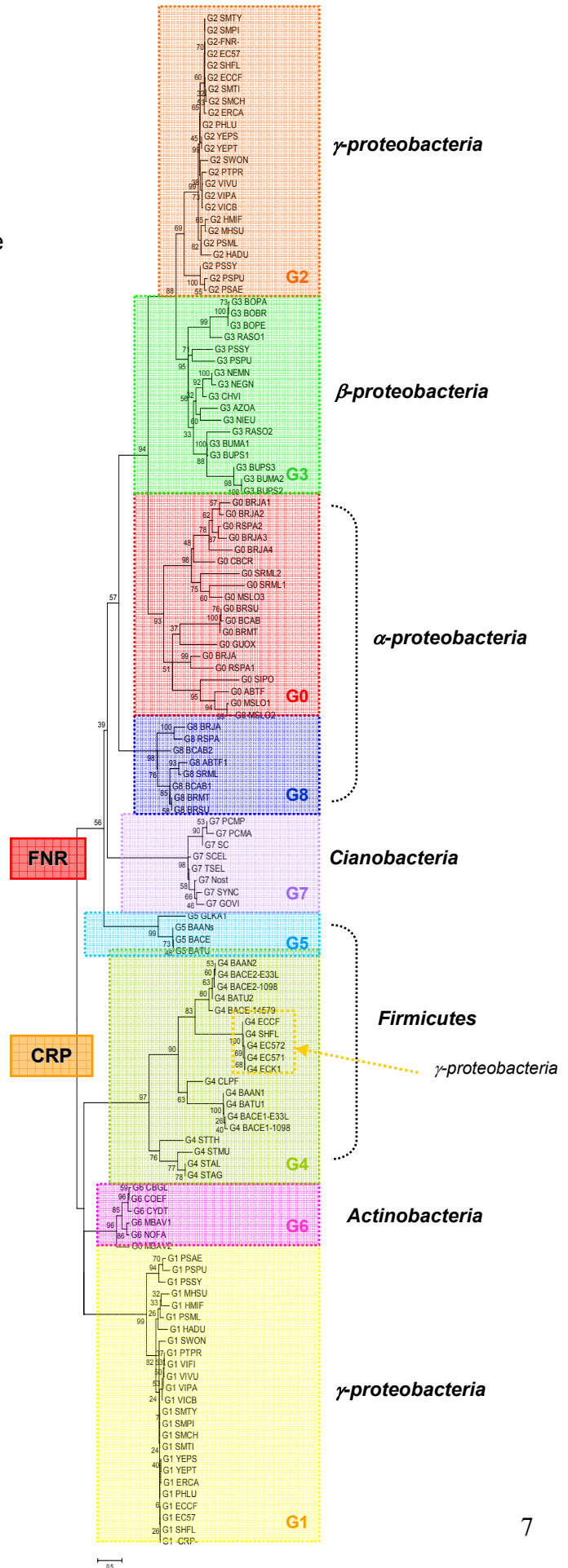
7. Anexos	57
7.1. Tabla de las secuencias de la familia Crp/Fnr	58
7.2. Alineamiento de las secuencias de la familia Crp/Fnr	60
7.3. Representación rectangular de árbol filogenético	62
7.4. Gráficas de distribución Normal	63
7.5. Tablas de los valores de CI , con los valores Z de cada grupo de la familia Crp/Fnr	Anexo 7.5: 1-10
7.6. Códigos de los programas en lenguaje PERL	65
7.6.1. ContInfo-aa-alignment-PRO.pl	65
7.6.2. Calc-seqs-PseudoRandom.pl	69
7.6.3. Shuffling_alingment.pl	75
7.6.4. MedsStats.pl	78
7.6.5. Calculo_rho.pl	83
7.6.6. Color_SeqRelation.pl	85

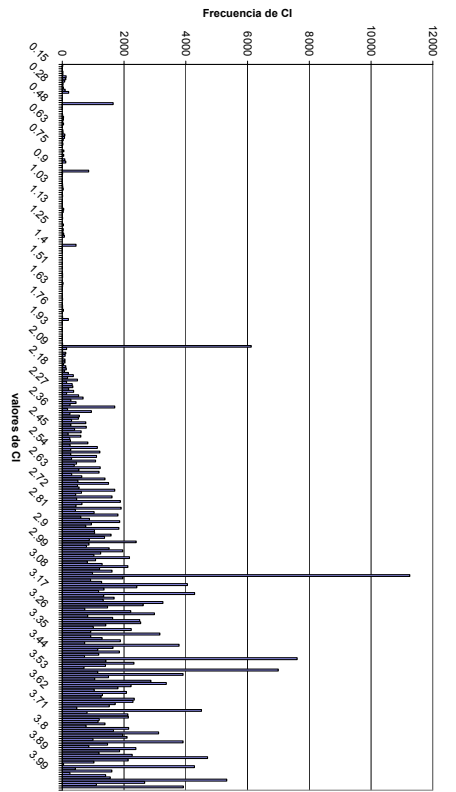
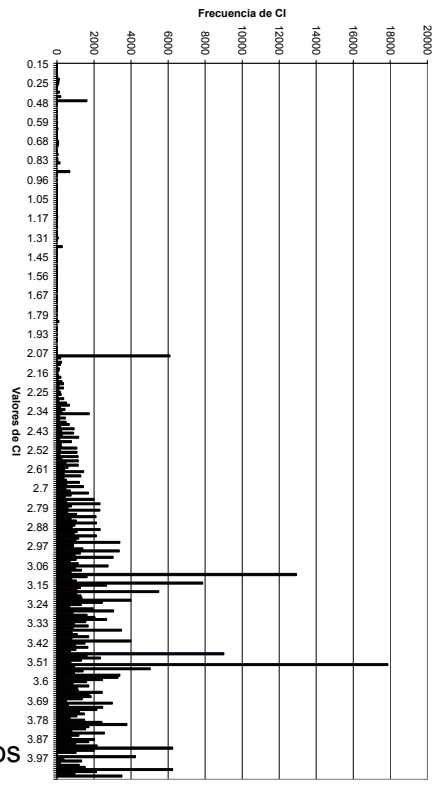
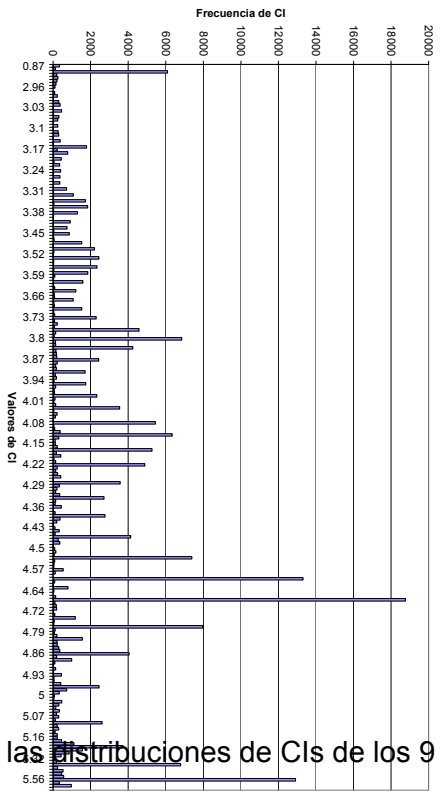
Anexo 7.1 Tabla de secuencias de la familia Crp/Fnr

MSLO3	SRML2	19	(Mesorhizobium loti MAFF530309)	transcriptional regulator FixK	[13475493]
G0	BRMT	19	(Sporichobium medii 1021)	transcriptional regulator FixK	[16262801]
G8	BCBJ	19	(Bacillus subtilis 104)	transcriptional regulator, cap box	[17892931]
G8	BCAB1	19	(Bacillus subtilis 1030)	transcriptional regulator, CapT	[23500019]
G8	ABT1	19	(Bacillus subtilis 1030)	transcriptional regulator, CapT	[62317811]
G8	SRML	19	(Agrobacterium tumefaciens str. C58)	hypothetical protein AGR_L_9	[16283132]
G8	BCA22	19	(Sporichobium medii 1021)	putative crp/fnr-like transcript	[62317794]
G8	BRJA	19	(Bacillus subtilis 1030)	NrrR, transcriptional regulator	[27382195]
G8	RSPA	19	(Rhodospirillum rubrum)	transcriptional regulatory prot	[399837203]
G5	BACE	19	(Bacillus cereus ATCC 14579)	Transcription regulator, Crp Ia	[30020257]
G5	BACE	19	(Bacillus thuringiensis serovar konkukian str. 97-27)	transcriptional regulator, Crp Ia	[52143301]
G5	BAAN	19	(Bacillus anthracis str. Sterne)	transcriptional regulator, Crp Ia	[49184969]
G5	GLKA	19	(Geobacillus kaustophilus HTA428)	transcriptional regulator of arc	[586419303]
G4	STAL	19	(Streptococcus agalactiae 2603VR)	cyclic nucleotide-binding dom.	[22836032]
G4	STAG	19	(Streptococcus agalactiae NEM316)	hypothetical protein gbs1882	[25011821]
G4	STMU	19	(Streptococcus mitis U4169)	putative transcriptional regulator	[24378661]
G4	STTH	19	(Streptococcus thermophilus LMG 18311)	transcriptional regulator, putat	[95820250]
G4	BAANI	19	(Bacillus anthracis str. Ames)	cyclic nucleotide-binding dom.	[30263413]
G4	BACE-E33	19	(Bacillus thuringiensis serovar konkukian str. 97-27)	putative transcriptional regulator	[30263413]
G4	BACE-E33	19	(Bacillus thuringiensis serovar konkukian str. 97-27)	cyclic nucleotide-binding dom.	[52142514]
G4	CLPF	19	(Clostridium perfringens str. 13)	probable transcriptional regul.	[18309529]
G4	BAAN2	19	(Bacillus anthracis str. Ames)	cyclic nucleotide-binding dom.	[30265430]
G4	BACE-E33	19	(Bacillus cereus E331)	cyclic nucleotide-binding dom.	[52145218]
G4	BATU2	19	(Bacillus thuringiensis serovar konkukian str. 97-27)	cyclic nucleotide-binding dom.	[49480446]
G4	BACE-109	19	(Bacillus cereus ATCC 10987)	cyclic nucleotide-binding dom.	[42784633]
G4	BACE-1457	19	(Bacillus cereus ATCC 14579)	Catabolite gene activator	[300234339]
G4	EC271	19	(Escherichia coli O157:H7 EDL83)	putative transcriptional regulator	[15802719]
G4	EC272	19	(Escherichia coli K12)	DNA-binding transcriptional ar	[38704654]
G4	EC273	19	(Escherichia coli O157:H7 str. Sakai)	putative transcriptional regulator	[38704654]
G4	SHFL	19	(Shigella flexneri 2a str. 2467T)	putative transcriptional regulator	[300639604]
G4	ECCF	19	(Escherichia coli CF1073)	Regulatory protein nsr	[26248545]

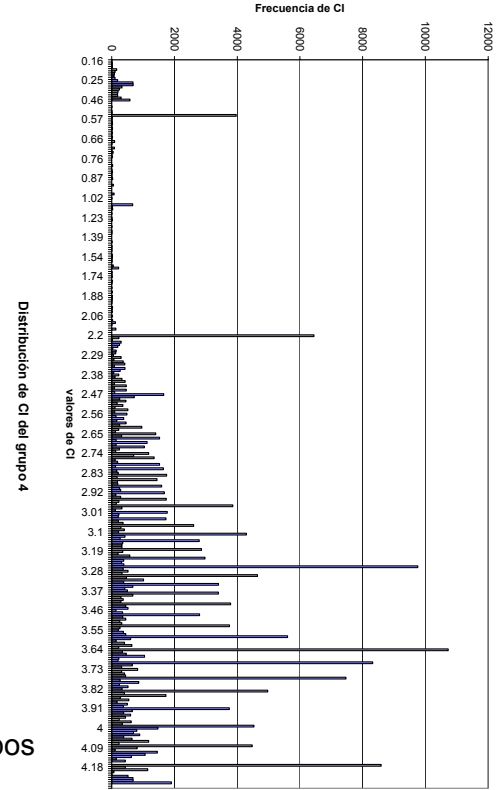
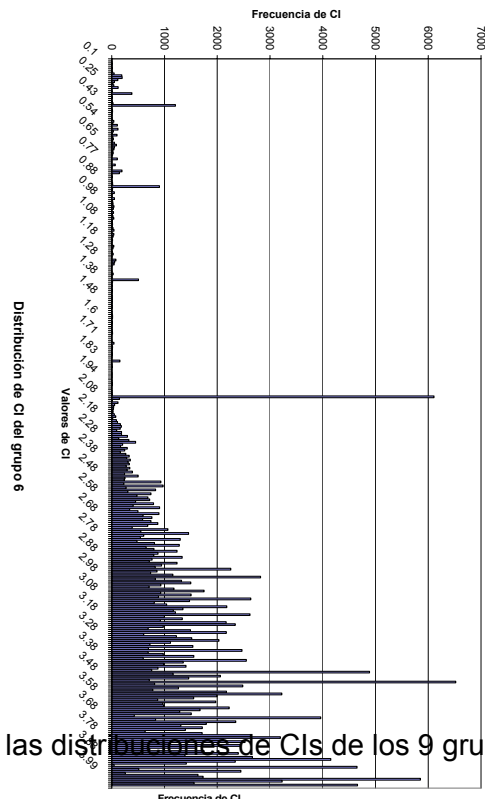
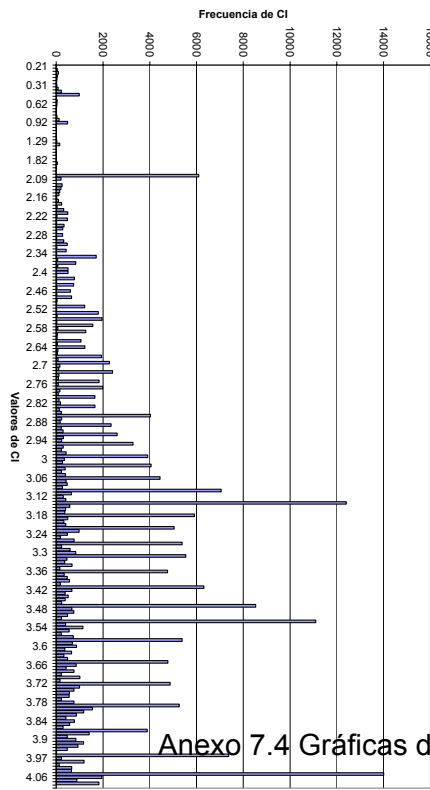
Anexo 7.3 Representación del árbol filogenético de la familia Crp/Fnr

Árbol de 134 secuencias de la familia Crp/Fnr, con los valores de 'bootstrap' que ayudaron a determinar los nueve grupos. Generado por Neighbor Joining siguiendo el modelo JTT (gama = 2.0) con el software Mega3.1

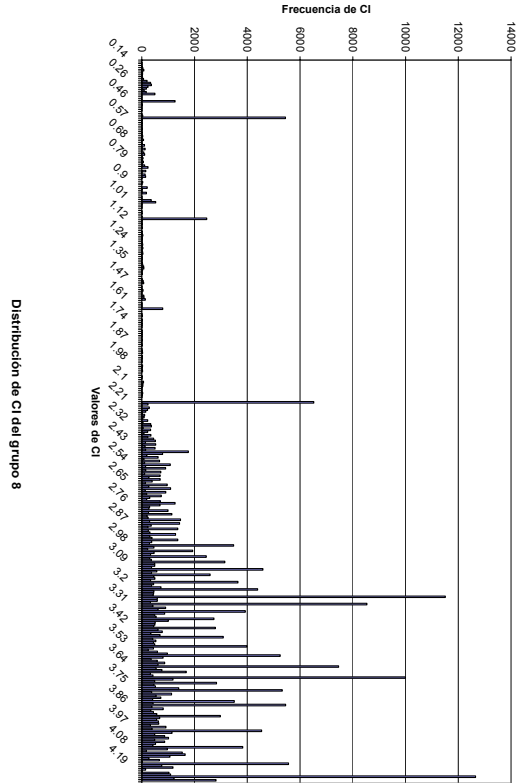
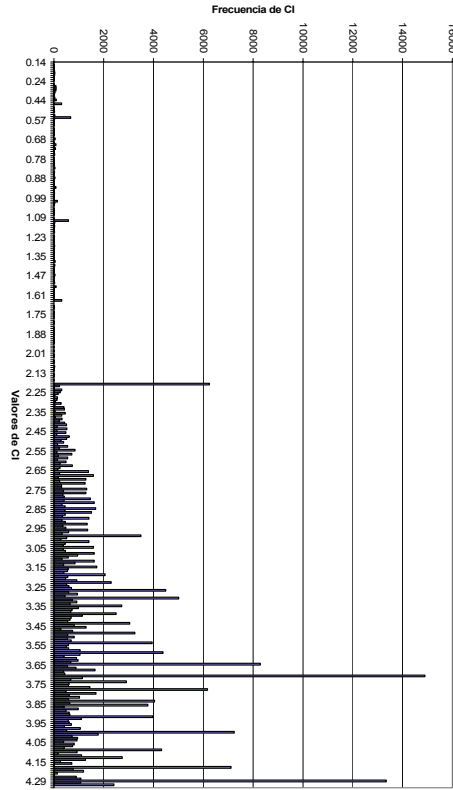




as de las distribuciones de CIs de los 9 grupos



Anexo 7.4 Gráficas de las distribuciones de CIs de los 9 grupos



Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.1 ContInfo-aa-alignment-PRO.pl

```
=====
=====
NAME
    ContInfo-aa-alignment-PRO.pl

DESCRIPTION
    Calculates and assigns a score for the amino acids content of each
    column in an alignment

AUTHORS
    heladia@ccg.unam.mx; mipspin@ccg.unam.mx

CATEGORY
    Multiple Sequence Alignment analysis; Informational Content

USAGE
    CalculateLogAligment [-i AlignmentFileName] [-help] [-h] [-f
    InputFormat] [-o OutputFile]

    perl ContInfo-aa-alignment-PRO.pl -i alignment.txt -f
    pipe/id-tab/tab3/tab2 -o OupptputFilename.txt

Formula Used
    ScoreCol = PropAAGroupCol * (log (PropAAGroupCol/ (PropAACol *
    totalElementByGroup)))

Where

    PropAAGroupCol = ratio of the amino acid in each column by the group it
    belongs to

    PropAACol = ratio of the amino acid in each column

    ScoreCol = score of amino acid

    SumScoreCol = sumatory of the scores of the aminoacid by column
```

```

-help
Same as -h

-i AlignmentFileName
name of the file [don't use spaces in the name]

-f InputFormat
input Formats: "pipe" group|identifier|name|organism sequence "id-tab"
group_identifier sequence "tab3" group identifier sequence "tab2" group
sequence

-o output
The output file for the CI results

=====
#!/usr/local/bin/perl

=pod

=head1 NAME

ContInfo-aa-alignment-PRO.pl

=head1 DESCRIPTION

Calculates and assigns a score for the amino acids content of each column
in an alignment

=head1 AUTHORS

heladia@ccg.unam.mx; mipspin@ccg.unam.mx

=head1 CATEGORY

Multiple Sequence Alignment analysis; Informational Content

=head1 USAGE

CalculateLogAlignment [-i AlignmentFileName] [-help] [-h] [-f InputFormat]
[-o OutputFileName]
perl ContInfo-aa-alignment-PRO.pl AlignmentFile pipe/id-tab/tab3/tab2 -o
OutputFileName.txt

=cut

#####
###          MAIN
#####
&ReadArguments;
&OpenOutPut;
&ReadData;
&CalculateLog;

#####
### Calcula el score para cada columna

=pod

=item B<Formula Used>

ScoreCol = PropAAGroupCol * (log (PropAAGroupCol/
                                (PropAACol * totalElementByGroup)))

Where

PropAAGroupCol = ratio of the amino acid in each column by the group it
                belongs to

PropAACol = ratio of the amino acid in each column

ScoreCol = score of amino acid

SumScoreCol = sumatory of the scores of the aminoacid by column

=cut

sub CalculateLog {

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.1 ContInfo-aa-alignment-PRO.pl

```

$totalElementByGroup = 0;
@SumScore_tot = ();

print OUT "Posiciones\t";
foreach $scol ( sort { $a <=> $b } keys %{ $FreqAAByGroupCol{$group} } ){
    $scol_imprimir = $scol + 1;
    print OUT "$scol_imprimir\t";
}
print OUT "\n";

# recorremos los grupos
foreach $group ( sort keys %GroupGICol ) {

    # total de elementos por grupo
    $totalElementByGroup = keys %{ $GroupGICol{$group} };

    print OUT "$group\t";

    ## recorrer las columnas de manera numerica 0,1,2,3..
    foreach $scol ( sort { $a <=> $b } keys %{ $FreqAAByGroupCol{$group} } ){

        $SumlogAACol = 0;
        $PropAACol = 0;
        $PropAAGroupCol = 0;
        $SumScoreCol = 0;

        foreach $aa ( sort keys %{ $FreqAAByGroupCol{$group}{$scol} } ){

            $FreAA_Gn_ColN = $FreqAAByGroupCol{$group}{$scol}{$aa};

            # proporcion del aminoacido por grupo y en columna
            $PropAA_GroupCol = $FreAA_Gn_ColN / $totalElementByGroup;

            # proporcion del aminoacido por columna
            $PropAA_AllCol = $FreAA_All_ColN / $TotalElementsAllGroup;

            #proporcion de secuencias de cada grupo respecto al total de secs
            $propElemEnGrup = $totalElementByGroup / $TotalElementsAllGroup;

            $FE = ($PropAA_AllCol * $propElemEnGrup);

            $FO_FE = ($PropAA_GroupCol / $FE );

            $log = log ($FO_FE);

            #score del aminoacido
            $ScoreCol = ($PropAA_GroupCol * $log);

            if ($aa eq '-') { $ScoreCol = 0; }

            #sumatoria de todos los scores de los aminoacidos en la columna n
            $SumScoreCol += $ScoreCol;

        }

        $SumScore_tot[$scol] += $SumScoreCol;

        $SumScoreCol_dec = sprintf("%3.2f", $SumScoreCol);
        print OUT "$SumScoreCol_dec\t";

    }

    print OUT "\n";

}

print OUT "GT\t";

$num_cols = scalar (@SumScore_tot);
$scol = 0;

while ($scol < $num_cols) {

    $SumScore_tot_dec = sprintf ("%3.2f", $SumScore_tot[$scol]);

    print OUT "$SumScore_tot_dec\t";
    $scol++;
}

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.1 ContInfo-aa-alignment-PRO.pl

```

print OUT "\n";
}

#####
### Read the info form file name
sub ReadData {

    $TotalElementsAllGroup = 0;

    open(IN,"sort $filename |") || die "I can't open the file\n";

    while(<IN>){
        chomp;

        if ($informat =~ /pipe/i) {
            #ARVHIVO EN FORMATO +|+|+|s+*

            next if $_ !~ /^(^|+)|(^|+)|(^|+)|(^|s+)|s+(.*)$/;
            #print "$1 .. $2 .. $3 .. $4 ..$5\n";      #CHECKPOINT
            $group = $1;
            $gi = $2;
            $alignment = $5;
        }

        if ($informat =~ /id-tab/i) {
            #ARCHIVO EN FORMATO + +s+*
            next if $_ !~ /^(^|+)|(^|+)|(^|+)|(^|s+)|s+(.*)$/;
            $group = $1;
            $gi = $2;
            $alignment = $3;
            #print "$group .. $gi .. $alignment\n";    #CHECKPOINT
        }

        if ($informat =~ /tab3/i) {
            #OCUPAR ESTAS LINEAS EN CASO DE ARCHIVO EN FORMATO 'identif de grupo \t secuencia'

            ($group, $gi, $alignment) = split (/\\t/, $_);
            #print "$group .. $gi .. $alignment\n";    #CHECKPOINT
        }

        if ($informat =~ /tab2/i) {
            #OCUPAR ESTAS LINEAS EN CASO DE ARCHIVO EN FORMATO 'identif de grupo \t secuencia'

            ($group, $alignment) = split (/\\t/, $_);
            #print "$group .. $alignment\n"; #CHECKPOINT
        }
    }

    $TotalElementsAllGroup++;

    # separamos los aminoacidos en columnas
    $i = 0;
    foreach $aa (split(//,$alignment)) {

        #print "$group $i = $aa\n";

        # lleva los aminoacidos por columna
        $GroupGICol{$group}{$gi}{$i} = $aa;

        # Frecuencia de AA en cada grupo x columna
        $FreqAAByGroupCol{$group}{$i}{$aa}++;

        # Frecuencia de AA en todo el alineamiento x columna
        $FreqAAByColAllGroup{$i}{$aa}++;

        $i++;
    }
}

}

#####
### Open Output File

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.1 ContInfo-aa-alignment-PRO.pl


```

sub OpenOutPut {
    open (OUT, ">$outFile") || die "No pude abrir el archivo de salida";
}

#####
### display full help message
sub PrintHelp {
    system "pod2text -c $0";
    exit()
}

#####
### Read arguments
sub ReadArguments {

    my $arg = "";
    while ($arg = shift (@ARGV)) {
        if (($arg eq "-h") || ($arg eq "-help")) {
            &PrintHelp();

            ## List of options

=pod
=item B<-help>

Same as -h
=cut
7.6.1 ContInfo-aa-alignment-PRO.pl
7.6.2 Calc-seqs-PseudoRandom.pl

=pod
=item B<-i AlignmentFileName>

name of the file [don't use spaces in the name]

=cut
    } elsif ($arg eq "-f") {
        $informat = lc (shift (@ARGV));

=pod
=item B<-f InputFormat>

input Formats: "pipe" group|identifier|name|organism      sequence "id-tab"
group_Identifier sequence "tab3" group      identifier      sequence "tab2" group      sequence

=cut
    } elsif ($arg eq "-o") {
        $outFile = shift (@ARGV);

=item B<-o output>

The output file for the CI results

=cut

    } else {
        &PrintHelp();
    }
}

}

=====
NAME
    Calc-seqs-PseudoRandom.pl

DESCRIPTION
    Generates pseudo-aleatory sequences starting from a multiple sequence
    alignment, in accordance to a probability given

```

AUTHORS
 mipspin@ccg.unam.mx;heladia@ccg.unam.mx

CATEGORY
 Multiple Sequence Alignment analysis; Generate Pseudo-random sequences

USAGE
 Calc-seqs-PseudoRandom.pl [-help] [-h] [-i AlignmentFileName] [-f
 InputFormat] [-p probability] [-n NumberRepetitions] [-o OutputFile]

perl Calc-seqs-PseudoRandom.pl [-help] [-h] [-i AlignmentFileName] [-f
 InputFormat] [-p probability] [-n NumberRepetitions] [-o OutputFile]

-help
 Same as -h

-i AlignmentFileName
 name of the file [don't use spaces in the name]

-f inputFormat
 inputFormat:"pipe" group|identifier|name|organism sequence "id-tab"
 group_identifier sequence "tab3" group identifier sequence "tab2" group
 sequence

-p probability
 probability: real number between 0 and 1

-n NumberRepetitions
 the times that the generation of sequences will be done

-o OutputFilename
 name of the output file [don't use spaces in the name]

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.2 Calc-seqs-PseudoRandom.pl

```
=====
#!/usr/local/bin/perl

=pod

=head1 NAME

Calc-seqs-PseudoRandom.pl

=head1 DESCRIPTION

Generates pseudo-aleatory sequences starting from a multiple sequence
alignment, in accordance to a probability given

=head1 AUTHORS

mipspin@ccg.unam.mx;heladia@ccg.unam.mx

=head1 CATEGORY

Multiple Sequence Alignment analysis; Generate Pseudo-random sequences

=head1 USAGE

Calc-seqs-PseudoRandom.pl [-help] [-h] [-i AlignmentFileName]
[-f InputFormat] [-p probability] [-n NumberRepetitions] [-o OutputFile]

perl Calc-seqs-PseudoRandom.pl [-help] [-h] [-i AlignmentFileName]
[-f InputFormat] [-p probability] [-n NumberRepetitions] [-o OutputFile]

=cut

#####
###      MAIN
#####
&ReadArguments;
&ReadData;
&CalculateProb_AA;

for ($r = 0; $r <= $numberRep; $r++) {

    $out_file = "$r"."out";
    open (OUT_S, ">$out_file") || die "No fue posible abrir archivo salida num $r\n";

    &Genera_Seq_Parental;

```

```

        &Genera_Seqs_hijas;

        close (OUT_S);
    }
#####
###

#####
###
sub Genera_Seqs_hijas {

    @seq_hija = ();

    foreach $group ( sort keys %GroupGICol ) {

        $totalSecsByGroup{$group} = keys %{$GroupGICol{$group}};

        $g = 0;
        while ($g < ($totalSecsByGroup{$group} - 1)) {

            $scol = 0;
            $prob_a_busc = 0;

            $largo_sec = (keys %{ $AcumAAGroupCol{$group} });

            while ($scol < $largo_sec) {

                $prob_a_busc = rand ();

                $scol = $AcumAAGroupCol{$group}[$scol];
                $seq_hija{$group}[$scol] = $aa_ahora;

                print OUT_S "$seq_hija{$group}[$scol]";

            }

            else {

                $prob_a_busc_2nd = rand ();

                $i = 0;

                $aa_ahora = $aa_array{$group}[$scol][$i];
                $prob_a_probar = $AcumAAGroupCol{$group}[$scol][$i];

                while ($prob_a_busc_2nd >= $prob_a_probar && $prob_a_busc_2nd != 1) {

                    $i++;
                    $aa_ahora = $aa_array{$group}[$scol][$i];
                    $prob_a_probar = $AcumAAGroupCol{$group}[$scol][$i];

                    if ($prob_a_probar == "") {
                        #Este ciclo lo pongo
                        #por si llega a pasar
                        #que la probab de algun
                        #aa se pierda por ahí y
                        #asi saber al menos
                        #donde buscar y reparar

                        print STDERR "Error de P_aa en G:$group C:$scol i:$i\n";
                    }

                    if ($prob_a_probar == 1 || $prob_a_busc_2nd == 1 ) { last }

                }

                $seq_hija{$group}[$scol] = $aa_ahora;

                print OUT_S "$seq_hija{$group}[$scol]";
            }

            $scol++;
        }

        print OUT_S "\n";
        $g++;
    }
}
#####

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.2 Calc-seqs PseudoRandom.pl

```

####
sub Genera_Seq_Parental {
    @seq_parental = ();
    foreach $group ( sort keys %GroupGICol ) {
        print OUT_S "$group\t00\t";
        $col = 0;
        $prob_a_busc = 0;
        $largo_sec = (keys %{ $AcumAAGroupCol{$group} });
        while ($col < $largo_sec) {
            $i = 0;
            $prob_a_busc = rand ();
            $aa_ahora = $aa_array{$group}{$col}{$i};
            $prob_a_probar = $AcumAAGroupCol{$group}{$col}{$i};
            while ($prob_a_busc >= $prob_a_probar && $prob_a_busc != 1) {
                $i++;
                $aa_ahora = $aa_array{$group}{$col}{$i};
                $prob_a_probar = $AcumAAGroupCol{$group}{$col}{$i};
            }
            if ($prob_a_probar == "") {
                #Este ciclo lo pongo
                #por si llega a pasar
                #que la probab de algun
                #aa se pierda por ahí y
                #asi saber al menos
                #donde buscar y reparar
                print STDERR "Error de P_aa en G:$group C:$col i:$i\n";
            }
            if ($prob_a_probar == 1 || $prob_a_busc == 1 ) { last }
        }
        $seq_parental{$group}{$col} = $aa_ahora;
        print OUT_S "$seq_parental{$group}{$col}";
        $col++;
    }
    print OUT_S "\n";
}

#####
### Calcula probabilidad de aa en cada columna

sub CalculateProb_AA {
    @totalSecsByGroup = ();
    @AcumAAGroupCol = ();
    @aa_array = ();
    $g = 0;
    # recorremos los grupos
    foreach $group ( sort keys %GroupGICol ) {
        # total de elementos por grupo
        $totalElementByGroup{$group} = keys %{ $GroupGICol{$group} };
        $totSecs_x_group{$group}[$g] = $totalElementByGroup{$group};
        ## recorrer las columnas de manera numerica 0,1,2,3..
        foreach $col ( sort { $a <=> $b } keys %{ $FreqAAByGroupCol{$group} } ){
            $PropAAGroupCol_lx1 = 0;
            $acumulando_propAA = 0;

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.2 Calc-seqs-PseudoRandom.pl

```

        $i = 0;
        foreach $aa ( sort keys %{ $FreqAABByGroupCol{$group}{$col} } ) {

            $aa_array{$group}{$col}[$i] = $aa;

            $FreAA_Gn_ColN = $FreqAABByGroupCol{$group}{$col}{$aa};

            # proporcion del aminoacido por grupo y en columna
            $PropAAGroupCol_1x1 = $FreAA_Gn_ColN / $totalElementByGroup{$group};

            $acumulando_propAA += $PropAAGroupCol_1x1;

            $GroupAAGroupCol{$group}{$col}[$i] = $acumulando_propAA;
        }

        $g++;
    }

}

#####
### Read the info form file name
sub ReadData {

    $TotalElementsAllGroup = 0;

    open(IN,"$filename") || die "I can't open the file\n";
    while(<IN>){
        chomp;

        if ($informat =~ /pipe/i) {
            #ARVHIVO EN FORMATO +|+|+|s+*

            next if $_ !~ /^([\^|]+)\|([\^|]+)\|([\^|]+)\|([\^s]+)\s+(.*)$/;
            #print "$1 .. $2 .. $3 .. $4 .. $5\n";      #CHECKPOINT
            $group = $1;
            $gi = $2;
            $alignment = $5;
        }

        if ($informat =~ /id-tab/i) {
            #ARCHIVO EN FORMATO +_+s+*

            next if $_ !~ /^([\^_|]+)\_([\^s]+)\s+(.*)$/;
            $group = $1;
            $gi = $2;
            $alignment = $3;
            #print "$group .. $gi .. $alignment\n";      #CHECKPOINT
        }

        if ($informat =~ /tab3/i) {
            #ARCHIVO EN FORMATO 'grupo \t identif de grupo \t secuencia'

            ($group, $gi, $alignment) = split (/\\t/, $_);
            #print "$group .. $gi .. $alignment\n";      #CHECKPOINT
        }

        if ($informat =~ /tab2/i) {
            #ARCHIVO EN FORMATO 'identif de grupo \t secuencia'

            ($group, $alignment) = split (/\\t/, $_);
            #print "$group .. $alignment\n"; #CHECKPOINT
        }

        $TotalElementsAllGroup++;

        # separamos los aminoacidos en columnas
        $i = 0;
        foreach $aa (split(//,$alignment)) {

            #print "$group $gi $i = $aa\n";

            # lleva los aminoacidos por columna
            $GroupGICol{$group}{$gi}[$i] = $aa;
        }
    }
}

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.2 Calc-seqs-PseudoRandom.pl

```

        # lleva la frecuencia por columna
        $FreqAAByGroupCol{$group}{$i}{$aa}++;

        # Frecuencia de AA en todos los grupos
        $FreqAAByColAllGroup{$i}{$aa}++;

        $i++;
    }
}

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.2 Calc-seqs-PseudoRandom.pl

```

#####
### display full help message
sub PrintHelp {
    system "pod2text -c $0";
    exit()
}

#####
### Read arguments
sub ReadArguments {

    my $arg = "";
    while ($arg = shift (@ARGV)) {
        if (($arg eq "-h") || ($arg eq "--help")) {
            &PrintHelp();

            ## List of options

            =pod

            =item B<-help>

            Same as -h

            =cut

        } elsif ($arg eq "-i") {
            $filename = shift(@ARGV);

            =pod

            =item B<-i AlignmentFileName>

            name of the file [don't use spaces in the name]

            =cut

        } elsif ($arg eq "-f") {
            $informat = shift(@ARGV);

            =pod

            =item B<-f inputFormat>

            inputFormat:"pipe" group|identifier|name|organism sequence
                        "id-tab" group_identifier sequence
                        "tab3" group identifier sequence
                        "tab2" group sequence

            =cut

        } elsif ($arg eq "-p") {
            $prob = shift(@ARGV);
            $probname = ($prob * 100);

            =pod

            =item B<-p probability>

            probability: real number between 0 and 1

            =cut

        } elsif ($arg eq "-n") {
            $numberRep = shift(@ARGV);

```

```

=pod
=item B<-n NumberRepetitions>
the times that the generation of sequences will be done
=cut

    } elsif ($arg eq "-o") {
        $outFile = shift(@ARGV);
    }

=pod
Anexo 7.6 Códigos de programas en lenguaje PERL
7.6.2 Calc-seqs-PseudoRandom.pl
7.6.3 Shuffling_alin_PRO.pl
name of the output file [don't use spaces in the name]
=cut

    } else {
        &PrintHelp();
    }
}

=====
NAME
    Shuffling_alin_PRO.pl

DESCRIPTION
    TOMAR UN ALINEAMIENTO Y MEZCLAR SUS RESIDUOS ALEATORIAMENTE PARA SER
    COLOCADOS ALEATORIAMENTE; Y ASI GENERAR SECUENCIAS 'SHUFFLED'

AUTHORS
    mipspin@ccg.unam.mx

CATEGORY
    Multiple Sequence Alignment analysis; To shuffle an alignment

USAGE
    perl Shuffling_alin_PRO.pl -h -i [alignment-file] -f [inputFormat]

    -help
    Same as -help

    -i AlignmentFileName
    name of the file [don't use spaces in the name]

    -f inputFormat
    inputFormat:"pipe" group|identifier|name|organism sequence "id-tab"
    group_identifier sequence "tab3" group identifier sequence "tab2" group
    sequence
=====

#!/usr/bin/perl
=pod
=head1 NAME
Shuffling_alin_PRO.pl

=head1 DESCRIPTION
TOMAR UN ALINEAMIENTO Y MEZCLAR SUS RESIDUOS ALEATORIAMENTE
PARA SER COLOCADOS ALEATORIAMENTE; Y ASI GENERAR SECUENCIAS
'SHUFFLED'

=head1 AUTHORS
mipspin@ccg.unam.mx

=head1 CATEGORY
Multiple Sequence Alignment analysis; To shuffle an alignment

=head1 USAGE

```

```
perl Shuffling_alin_PRO.pl -i [alignment-file] -f [inputFormat]
```

```
=cut
```

```
#####  
###          MAIN  
#####  
&ReadArguments;  
&ReadData;  
&RandomizeData;  
&DoRandAlin;
```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.3 Shuffling_alin_PRO.pl

```
#####  
### Read the info form file name and split alignment  
sub ReadData {  
  
    @aa_x_posic = ();  
    @ids = ();  
  
    open (ALIN, "$filename") || die "No se pudo abrir el alineamiento:( \n";  
  
    $j = 0;  
    while (<ALIN>) {  
        chomp;  
        if ($informat =~ /pipe/i) {  
            #ARVHIVO EN FORMATO +|+|+|s+*  
  
            next if $_ !~ /^(^[^|]+)\|([^\|]+)\|([^\|]+)\|([^\s]+)\s+(.*)$/;  
            #print "$1 .. $2 .. $3 .. $4 ..$5\n";      #CHECKPOINT  
            $group = $1;  
            $gi = $2;  
            $alignment = $5;  
        }  
  
        if ($informat =~ /id-tab/i) {  
            #ARCHIVO EN FORMATO +_+s+*  
  
            next if $_ !~ /^(^[^_]+)\_([^\s]+)\s+(.*)$/;  
            $group = $1;  
            $gi = $2;  
            $alignment = $3;  
            #print "$group .. $gi .. $alignment\n";      #CHECKPOINT  
        }  
  
        if ($informat =~ /tab3/i) {  
            #ARCHIVO EN FORMATO 'grupo \t identif de grupo \t secuencia'  
  
            ($group, $gi, $alignment) = split (/t/, $_);  
            #print "$group .. $gi .. $alignment\n";      #CHECKPOINT  
        }  
  
        if ($informat =~ /tab2/i) {  
            #ARCHIVO EN FORMATO 'identif de grupo \t secuencia'  
  
            ($group, $alignment) = split (/t/, $_);  
            #print "$group .. $alignment\n"; #CHECKPOINT  
        }  
  
        $ids[$j] = $id;  
  
        # se separan las columnas de la secuencia  
        $i = 0;  
        foreach $aa (split(//,$alignment)) {  
            $aa_x_posic{$id}{$i} = $aa;  
            #print "aas: $aa_x_posic{$id}{$i}\t";      #CHECK_POINT  
            $i++;  
        }  
        #print "\n";      #parte de CHECKPOINT ultimo-anterior  
    }  
    #print "$ids[$j]\n"; #CHECK_POINT
```



```

        $j++;
    }

close ALIN;

} #####FIN_ReadData

#Randomizing an array as Perl says it can be done
sub RandomizeData {

    $l = 0;
    %bolsa_nums movida = ();
    while ($l < @bolsa_nums) {
        @bolsa_nums = (1..180);
        @moviendo_nums = ();

        while (@bolsa_nums) {

            push (@moviendo_nums, splice (@bolsa_nums, rand (@bolsa_nums), 1))

        }
        #print "nums-mov: @moviendo_nums \n";      #CHECKPOINT

        $h = 0;
        foreach $moviendo_nums (@moviendo_nums) {
            $bolsa_nums_movida{$l}{$h} = $moviendo_nums[$h];

            #print "$bolsa_nums_movida{$l}{$h} \t";  #CHECKPOINT

            $h++;
        }
        #print "\n"; #salto de linea del CHECKPOINT de la linea 90

        $l++;
    }

}

#
} #####FIN_RandomizeData

####SUB-MODULO 3####
#Assigning the random order to the secuencia of each @aa_x_posic{$identificador}
#
sub DoRandAlin {
    $n = 0;
    open (OUT, ">alin_shuf.aln" ) || die "no se pudo hacer un archivo de salida \n";

    while ($n < $j){

        $id = $ids[$n];
        print OUT "$id\t";      #CHECKPOINT

        $col = 0;
        while ($col < 180) {

            $num = $bolsa_nums_movida{$n}{$col};
            #print "numero: $num\t"; #CHECKPOINT

            $aa_x_mov{$id}{$col} = $aa_x_posic{$id}{$num};
            #print "aa: $aa_x_posic{$id}{$num} de:$aa_x_mov{$id}{$col}\t";      #CHECKPOINT

            print OUT "$aa_x_posic{$id}{$num}";
            $col++;
        }
        print OUT "\n";
        $n++;
    }

}

close (OUT);

} #####Fin DoRandAlin

#####
### display full help message
sub PrintHelp {
    system "pod2text -c $0";
}

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.3 Shuffling_alin_PRO.pl

Se muestra el código de los programas según el número de columnas del alineamiento

```

    exit()
}

#####
### Read arguments
sub ReadArguments {

    my $arg = "";
    while ($arg = shift (@ARGV)) {
        if (($arg eq "-h") || ($arg eq "-help")) {
            &PrintHelp();
        }
    }

}

#####

=====
NAME
    MedsStats.pl

DESCRIPTION
    It obtains the average, standard deviation and Z-score of the
    Informational Content of the real alignment compared with the
    Informational Content of the pseudo-random sequences generated with
    Calc-PseudoRandom.pl

AUTHORS
    mipspin@ccg.unam.mx

CATEGORY
    Multiple Sequence Alignment analysis; Statistic Measures

USAGE
    MedsStats.pl[-help] [-h] [-ir InfContFileReal] [-ip
    InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

    perl Calc-seqs-PseudoRandom.pl [-help] [-h] [-ir InfContFileReal] [-ip
    InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

    Output is automatically saved in a file named Meds-Stats0_[percentage of
    variability of the pseudorandom seqs]_G[group identifier].out

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.3 Shuffling_alin_PRO.pl

7.6.4 Meds-Stats.pl

Same as -h

=cut

```

    } elsif ($arg eq "-i") {
        $filename = shift(@ARGV);
    }

```

=pod

=item B<-i AlignmentFileName>

name of the file [don't use spaces in the name]

=cut

```

    } elsif ($arg eq "-f") {
        $inputFormat = shift(@ARGV);
    }

```

=pod

=item B<-f inputFormat>

```

inputFormat:"pipe" group|identifier|name|organism sequence
            "id-tab" group_identifier sequence
            "tab3" group identifier sequence
            "tab2" group sequence

```

=cut

```

    } else {
        &PrintHelp();
    }

```

}

NAME

MedsStats.pl

DESCRIPTION

It obtains the average, standard deviation and Z-score of the Informational Content of the real alignment compared with the Informational Content of the pseudo-random sequences generated with Calc-PseudoRandom.pl

AUTHORS

mipspin@ccg.unam.mx

CATEGORY

Multiple Sequence Alignment analysis; Statistic Measures

USAGE

```

MedsStats.pl[-help] [-h] [-ir InfContFileReal] [-ip
InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

```

```

perl Calc-seqs-PseudoRandom.pl [-help] [-h] [-ir InfContFileReal] [-ip
InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

```

Output is automatically saved in a file named Meds-Stats0_[percentage of variability of the pseudorandom seqs]_G[group identifier].out

```

-help
Same as -h

-ir InfContFileReal
name of the file with CI of real data set [don't use spaces in the name]

-ip InfContFilePseudoRandSeqs
name of the file with CI of pseudo random seqs [don't use spaces in the
name]

-g GroupIdentifier
The group identifier, usually a single number or letter
-p probab ility
probability: real number between 0 and 1 The probability at which the
pseudo-random sequences were done, just to have the reference of the
corresponding population and calculus

-n NumberRepetitions
the times that the calculation must be done
=====

#!/usr/bin/perl

=pod

=head1 NAME

MedsStats.pl

=head1 DESCRIPTION

It obtains the average, standard deviation and Z-score of the Informational
Content of the real alignment compared with the Informational Content of the
speudo-random sequences generated with Calc-PseudoRandom.pl

=head1 AUTHORS

mipsp@ccg.unam.mx

=head1 CATEGORY

Multiple Sequence Alignment analysis; Statistic Measures

=head1 USAGE

MedsStats.pl[-help] [-h] [-ir InfContFileReal]
[-ip InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

perl Calc-seqs-PseudoRandom.pl [-help] [-h] [-ir InfContFileReal]
[-ip InfContFilePseudoRandSeqs] [-g GroupIdentifier] [-n NumberRepetitions]

Output is automatically saved in a file named Meds-Stats0_[percentage of
variability of the pseudorandom seqs]_G[group identifier].out

=cut

#####
##
&ReadArguments;
&OpenOutPut;
&ArrayContInfReal;
&OpenCIPseudoRand;
&MedsStats;

#####
#
sub OpenOutPut {
    open (SALIDA, ">Meds-Stats0_${probname}_G${group}.out") || die "No pude abrir el archivo";

    print SALIDA "CI-G${group}_R\t";
}

#####
### M O D U L O 1 ##### Abre y pone en arreglo los Cont.Inf. real del grupo
#####

sub ArrayContInfReal {

```

Anexo 7.6 Códigos de programas en lenguaje PERL
7.6.4 Meds-Stats.pl

```

@ci_Gs_lines =();
open (GS, "$filenameCI_real") || die "Can't open the file :( \n"; # prueba_sta

while(<GS>){
    chomp($_);
    push(@ci_Gs_lines,$_);
    #print "$_\n"; #CHECKPOINT
}

close GS;

@part_elem_real_col = ();
Anexo 7.6 Códigos de programas en lenguaje PERL
7.6.4 Meds-Stats.pl
foreach $curr_line (@ci_Gs_lines){

    next if $curr_line !~ /^G$(group)\t(.*)$/; #busca que se cumpla el
                                                # patron 'GX' donde X= numero o #letra,
                                                #seguido de un tabulador, #seguido de cualquier
                                                #caracter #0 o mas veces y el resto de la
                                                #linea

    $valores = $1; #asigna como $valores lo siguiente a la
                  # primer columna, en este caso
                  #los aminoacidos de la secuencia

    print SALIDA "$valores\n"; #CHECKPOINT

    @part_elem_real_col = split (/t/, $valores);

}

}
#####
##FiN M O D U L O 1 #####
#####

###MINI-M O D U L O ### Abre y pone en arreglo el archivo con los "Random"
#de ContInf

sub OpenCIPseudoRand {

@out_lines =();
open (IN, "$filenameCI_pseudoR") || die "No pude abrirlo :( \n";

while(<IN>){
    chomp($_);
    push(@out_lines,$_);
}

close IN;
}
### FIN MINI-M O D U L O ###

#####
### M O D U L O 2 ##### Pone en arreglo y hace todas las operaciones con
#####'Random' de ContInf

sub MedsStats {
    chomp @out_lines;
    $tot_secs_alin = scalar (@out_lines);#asigno variable con numero de Cont.
                                                #Infs de secuencias pseudo-aleatorias
                                                # que tengo

    $veces_repetir = $numberRep;

    @Random_en_arreglo = ();
    @part_elem_col = ();

    $i = 0;
    while ($i < $veces_repetir) {

        next if $out_lines[$i] !~ /^G$(group)\s+(.*)$/; #busca el patron 'G
                                                #seguido de un número y
                                                # uno o mas espacios,
                                                #seguido de cualquier
                                                #caracter 0 o mas veces

```

```

#y el resto de la linea.
$valores = $1;          #asigna como $valores lo
                        #siguiente a la primer
                        #columna, en este
                        #caso los numeritos.

@part_elem_col = split (/\/t/, $valores);  #parte $valores por
#columnas y los pone
#en arreglo

$sh = 0;
foreach $part_elem_col (@part_elem_col) { #aquí se pondra cada
#del arreglo
#part_elem_col dentro
# de un arreglo
#bidimensional
$Random_en_arreglo[$i][$sh] = $part_elem_col[$sh];

$sh++;
}

$si++;
}

$tot_cols = scalar (@part_elem_col);  #CALCULO DE SUMA por COLUMNA
$tot_lines = $veces_repetir;

print SALIDA "Posiciones\t";
$il = 0;
while ($il < $tot_cols) {
    $col_imprimir = $il + 1;
    print SALIDA "$col_imprimir\t";
    $il++;
}

print SALIDA "\nSum_Rand0_{probname}\t";

@suma_rand_arreglo = ();

$il = 0;
while ($il < $tot_cols) {

    $Suma_random_col = 0;
    $si = 0;
    while ($si < $tot_lines) {

        $Suma_random_col += $Random_en_arreglo[$i][$il];

        $si++;
    }

    $suma_rand_arreglo[$il] = $Suma_random_col;

    printf SALIDA "%.3f\t", "$suma_rand_arreglo[$il]\t";  #CHECKPOINT

    $il ++;
}

print SALIDA "\n";

print SALIDA "Prom_Rand0_{probname}\t";

@promedio_arreglo = ();

$il = 0;
while ($il < $tot_cols) {          #CALCULO DE PROMEDIO por COLUMNA

    $promedio = ($suma_rand_arreglo[$il] / $veces_repetir);
    $promedio_arreglo[$il] = $promedio;

    printf SALIDA "%.3f\t", "$promedio_arreglo[$il]\t";  #CHECKPOINT
    $il ++;
}

print SALIDA "\n";

print SALIDA "\Desv_Est_0_{probname}\t";

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.4 Meds-Stats.pl

```

@desv_est_arreglo = ();

$l = 0;
while ($l < $tot_cols) {

    $cuadrado_resta = 0;
    $Resta_x_promedio = 0;
    $sumatoria_cuadrados = 0;

    $i = 0;
    while ($i < $tot_lines) {
        $Resta_x_promedio += ($arreglo[$l][$i] - $promedio_arreglo[$l]);
        $Resta_x_promedio **= 2;
        $sumatoria_cuadrados += $cuadrado_resta;

        $i++;
    }

    $div_suma_cuad_N = $sumatoria_cuadrados / $veces_repetir;
    $desv_est[$l] = sqrt($div_suma_cuad_N);

    #printf "%.6f\t", "$desv_est[$l]\n"; #CHECKPOINT
    $l ++;
}

print SALIDA "\n";

print SALIDA "Z-score0_{probname}\t"; #imprime en archivo la linea deseada

@z_score = ();

$l = 0;
while ($l < $tot_cols) { #CALCULO DE Z-SCORE por COLUMNA

    if ($desv_est[$l] > 0) {

        $z_score[$l] = (($part_elem_real_col[$l] - $promedio_arreglo[$l]) / $desv_est[$l]);

    } else {$Z_score[$l] = 0;}

    printf SALIDA "%.3f\t", "$z_score[$l]";

    $l++;
}

print SALIDA "\n\n";

}
#####
##FIN M O D U L O 2 #####
#####

#####
### display full help message

sub PrintHelp {
    system "pod2text -c $0";
    exit()
}

#####
### Read arguments
sub ReadArguments {

    my $arg = "";
    while ($arg = shift (@ARGV)) {
        if (($arg eq "-h") || ($arg eq "--help")) {
            &PrintHelp();
        }

        ## List of options
    }

}

=pod

=item B<-help>

Same as -h

=cut

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.4 Meds-Stats.pl

```

    } elsif ($arg eq "-ir") {
        $filenameCI_real = shift(@ARGV);
=pod
=item B<-ir InfContFileReal>
name of the file with CI of real data set [don't use spaces in the name]
=cut
    } elsif ($arg eq "-ip") {
=pod
7.6.4 Meds-Stats.pl
7.6.5 Calculo_rho.pl
=item B<-ip InfContFilePseudoRandSeqs>
name of the file with CI of pseudo random seqs [don't use spaces in the name]
=cut
    } elsif ($arg eq "-g") {
        $group = shift(@ARGV);
=pod
=item B<-g GroupIdentifier>
The group identifier, usually a single number or letter
=cut
    } elsif ($arg eq "-p") {
        $prob = shift(@ARGV);
        $probname = ($prob * 100);
=pod
=item B<-p probability>
probability: real number between 0 and 1
The probability at which the pseudo-random sequences were done, just
to have the reference of the corresponding population and calculus
=cut
    } elsif ($arg eq "-n") {
        $numberRep = shift(@ARGV);
=pod
=item B<-n NumberRepetitions>
the times that the calculation must be done
=cut
    } else {
        &PrintHelp();
    }
}

=====
=====
Calculo_rho.pl
-----

#! /usr/bin/perl

### PROGRAMA QUE CALCULA LAS RELACIONES DE IDENTIDAD DE UN ALINEAMIENTO DE SECUENCIAS DE FORMA PAREADA.

### Programa escrito en lenguaje PERL
### Irma Lozada Chavez: ilozada@ccg.unam.mx

### Ejecucion:
### .\Calculo_rho.pl

```

```
# /home/ilozada/GENETIC_REGULATION/TRN_ORTHOLOGY/TFs_alignments/TFs_ECO_alignments/
# crp      son_S00624      DVAGRIAQTLHLAKQPDAMTHPDGMQIKITRQEIGQIVGCSRETIVGRILKMLEEQNLIQ
# crp      hin_HI0957      DVAGRIAQTLMLNLAQPEAMTHPDGMQIKITRQEIGQIVGCSRETIVGRILKMLEDQNLIH
# lexA     sty_STY4433      MKALTARQQEVFDLIRDHISQTGMPPTRAEIAQRLGFRSPNAAEEHLKALARKGVLEIVS
# lexA     stm_STM4237      MKALTARQQEVFDLIRDHISQTGMPPTRAEIAQRLGFRSPNAAEEHLKALARKGVLEIVS
# lexA     eco_b4043|1-72| MKALTARQQEVFDLIRDHISQTGMPPTRAEIAQRLGFRSPNAAEEHLKALARKGVLEIVS
```

```
system "rm Gs-G_ORG-1a_SL.rho.txt";
Anexo 7.6 Códigos de programas en lenguaje PERL
7.6.5 Cálculo_rho.pl
open ENTRA, "Gs-G_ORG-1a_SL.txt" || die "No pude abrir Gs-G_ORG-1a_SL.txt\n";
open SALE, ">>Gs-G_ORG-1a_SL.rho.txt" || die "No pude abrir Gs-G_ORG-1a_SL.rho.txt\n";

@aln= ();
@aln= <ENTRA>;
chomp @aln;

print SALE "\t";
print "\t";

@horizontal= ();
$z= 0;
foreach $all (@aln) {
    chomp $all;
    $all=~ s/\s+/\t/gis;
    ($idh0, $sec0)= split "\t", $all;
    chomp $idh0, $sec0;

    #print "\t$idh0";
    $horizontal[$z++]= $idh0;
}

foreach $h (@horizontal){
    print SALE "$h\t";
    print "$h\t";
}
print SALE "\n";
print "\n";

@aln2= ();
@aln2= @aln;

$j= 0;
$k= 0;
@vertical= ();
foreach $query (@aln) {
    chomp $query;
    $query=~ s/\s+/\t/gis;
    ($idh, $sec)= split "\t", $query;
    chomp $idh, $sec;
    #print "$sec\n";
    @aas= ();
    @aas= split "", $sec;          # A V K G K Y I T I E N S D A L A A L A G H T R
    chomp @aas;
    $aas= @aas;

    print SALE "$idh";
    print "$idh";
    #@aln2= ();
    #@aln2= @aln;

    %hashline= ();
    foreach $subject (@aln2) {
        chomp $subject;
        $subject=~ s/\s+/\t/gis;
        ($idv, $secs)= split "\t", $subject;
        chomp $idv, $secs;
        #print "$secs\n";
        @aas2= ();
        @aas2= split "", $secs;    # A V K G K Y I T I E N S D A L A A L A G H T R
        chomp @aas2;

        $aln= 0;
        $total= 0;
        $iden= 0;
        $i= 0;
        AA:

```



```

for ($i>= 0; $i<= $aas; $i++) {
    #print "$i-$aas\t";
    #print "$aas[$i]-$aas2[$i]\t";

    if ($aas[$i] eq "-" || $aas2[$i] eq "-" || $aas[$i] eq "-" && $aas2[$i] eq "-") {
        #print "$aas[$i]-$aas2[$i]\t";
        next AA;
    }
    if ($aas[$i] eq $aas2[$i]) {
        #print "$aas[$i]-$aas2[$i]\t";
        $iden+= 1;
    }
}
#print "$idh=$iden/$total\n";
$identity= $iden/$total;
$identity= sprintf "%.2f",$identity);
print SALE "\t$identity";
print "\t$identity";
#<STDIN>;
}
print SALE "\n";
print "\n";
#shift @aln;
}

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.5 Calculo_rho.pl

7.6.6 color_SeqRelation.pl

```

close ENTRA;
close SALE;

```

```

=====
NAME
    color_SeqRelation.pl

DESCRIPTION
    It colors a pairwise-identity matrix table within the format: . line1:
    OrganismsIdentifiers[tab separated] Further lines:
    SequenceIdentifier[tab ]identity1[tab]identity2[tab]identity...
    Separation Line between groups:-----[tab]0[tab]0[tab] Further lines:
    SequenceIdentifier[tab ]identity1[tab]identity2[tab]identity...

```

```

AUTHORS
    Irma Lozada Chavez: ilozada@ccg.unam.mx

```

```

CATEGORY
    To color sequence relationships

```

```

USAGE
    color_SeqRelation.pl -h -i [InputMatrix]

    -help
    Same as -h

    -i InputMatrix
    name of the file [don't use spaces in the name]

```

```

#!/usr/bin/perl

=pod

=head1 NAME

color_SeqRelation.pl

=head1 DESCRIPTION

It colors a pairwise-identity matrix table within the format:
line1: OrganismsIdentifiers[tab separated]
Further lines: SequenceIdentifier[tab ]identity1[tab]identity2[tab]identity...
Separation Line between groups:-----[tab]0[tab]0[tab]
Further lines: SequenceIdentifier[tab ]identity1[tab]identity2[tab]identity...

=head1 AUTHORS

```

Irma Lozada Chavez:
ilozada@ccg.unam.mx

=head1 CATEGORY

To color sequence relationships

=head1 USAGE

color_SeqRelation.pl -h -i [InputMatrix]

=cut **Anexo 7.6 Códigos de programas en lenguaje PERL**

7.6.6 color_SeqRelation.pl

```
#####
&ReadArguments;

#####
# GD library
use GD;

open (IN, "$filename") or die "Can't open ENTRADA!!\n";
@arr1 = (<IN>);
$arr2 = join ("", @arr1);
$arr2 =~ s/^\t+/\t/;
@arr = split ("\n", $arr2);

#print "Mandar salida a:\nPantalla\t(1)\nArchivo\t(2)\n";
#$opt = <STDIN>;
#chomp $opt;

#### Sacando dimensiones
$filaMasGrande = 1;
$columnaMasGrande = 1;
foreach (@arr) {
    $ac = 0;
    if ($_ =~ /^(\S+)/) {
        @a = split ("", $1);
        foreach (@a) {
            $ac++;
        }

        if ($ac > $filaMasGrande) {
            $filaMasGrande = $ac;          ### Tamano del encabezado mas grande de las filas
        }

        $filas++;                          ### Numero de filas
    }
}

$an = join ("", @arr[0]);
@anc = split ("\t", $an);
foreach (@anc) {
    $bc = 0;
    if ($_ =~ /(\S+)/) {
        @b = split ("", $1);
        foreach (@b) {
            $bc++;
        }

        if ($bc > $columnaMasGrande) {
            $columnaMasGrande = $bc;      ### Tamano del encabezado mas grande de las columnas
        }

        $columnas++;                       ### Numero de columnas
    }
}

$alto = ((($filas * 10) + ($columnaMasGrande * 7)) - 10);  ### ANCHO DE LA FIGURA
$ancho = (20 + ($columnas * 10) + ($filaMasGrande * 7));  ### ALTO DE LA FIGURA

#print "Ancho: $ancho\nAlto: $alto\nFilas: $filas\nColumnas: $columnas\nfilaMasGrande =
$filaMasGrande\ncolumnaMasGrande = $columnaMasGrande\n";

##### Definiendo las variables
$box = 10;                                     #TAMANO DE CUADRO

#if ($opt == 2) {
open(OUT,"$filename.png") || die "No display... arghh!\n";
}
```

```

#}else{
#open(OUT,"| display") || die "No display... arghh!\n";
#}

```

```
binmode OUT;
```

```
$image = new GD::Image($ancho,$alto);
```

```

$whit = $image->colorAllocate(255,255,255);
DEFINIDO ES EL DEL FONDO
$gray = $image->colorAllocate(120,120,120);
$graydark = $image->colorAllocate(120,120,120);
$black = $image->colorAllocate(0,0,0);

```

```
$text_color = $black;
```

```
#EL PRIMER COLOR
```

```
#COLOR DEL TEXTO
```

```
#POSITIVE
```

```

$c[0] = $image->colorAllocate(0,0,0);
$c[1] = $image->colorAllocate(0,0,150);
$c[2] = $image->colorAllocate(0,10,150);
$c[3] = $image->colorAllocate(0,17,150);
$c[4] = $image->colorAllocate(0,24,150);
$c[5] = $image->colorAllocate(0,31,150);
$c[6] = $image->colorAllocate(0,38,150);
$c[7] = $image->colorAllocate(0,45,150);
$c[8] = $image->colorAllocate(0,52,150);
$c[9] = $image->colorAllocate(0,59,150);
$c[10] = $image->colorAllocate(0,66,200);
$c[11] = $image->colorAllocate(0,73,200);
$c[12] = $image->colorAllocate(0,80,200);
$c[13] = $image->colorAllocate(0,87,200);
$c[14] = $image->colorAllocate(0,94,200);
$c[15] = $image->colorAllocate(0,101,200);
$c[16] = $image->colorAllocate(0,108,200);
$c[17] = $image->colorAllocate(0,115,200);
$c[18] = $image->colorAllocate(0,122,200);
$c[19] = $image->colorAllocate(0,129,200);
$c[20] = $image->colorAllocate(0,136,200);
$c[21] = $image->colorAllocate(0,143,255);
$c[22] = $image->colorAllocate(0,150,255);
$c[23] = $image->colorAllocate(0,157,255);
$c[24] = $image->colorAllocate(0,164,255);
$c[25] = $image->colorAllocate(0,171,255);
$c[26] = $image->colorAllocate(0,178,255);
$c[27] = $image->colorAllocate(0,185,255);
$c[28] = $image->colorAllocate(0,192,255);
$c[29] = $image->colorAllocate(0,199,255);
$c[30] = $image->colorAllocate(0,206,255);
$c[31] = $image->colorAllocate(0,213,255);
$c[32] = $image->colorAllocate(0,220,255);
$c[33] = $image->colorAllocate(95,255,0);
$c[34] = $image->colorAllocate(100,255,0);
$c[35] = $image->colorAllocate(105,255,0);
$c[36] = $image->colorAllocate(110,255,0);
$c[37] = $image->colorAllocate(115,255,0);
$c[38] = $image->colorAllocate(120,255,0);
$c[39] = $image->colorAllocate(125,255,0);
$c[40] = $image->colorAllocate(130,255,0);
$c[41] = $image->colorAllocate(135,255,0);
$c[42] = $image->colorAllocate(140,255,0);
$c[43] = $image->colorAllocate(145,255,0);
$c[44] = $image->colorAllocate(150,255,0);
$c[45] = $image->colorAllocate(155,255,0);
$c[46] = $image->colorAllocate(160,255,0);
$c[47] = $image->colorAllocate(165,255,0);
$c[48] = $image->colorAllocate(170,255,0);
$c[49] = $image->colorAllocate(175,255,0);
$c[50] = $image->colorAllocate(180,255,0);
$c[51] = $image->colorAllocate(185,255,0);
$c[52] = $image->colorAllocate(190,255,0);
$c[53] = $image->colorAllocate(195,255,0);
$c[54] = $image->colorAllocate(200,255,0);
$c[55] = $image->colorAllocate(205,255,0);
$c[56] = $image->colorAllocate(210,255,0);
$c[57] = $image->colorAllocate(215,255,0);
$c[58] = $image->colorAllocate(220,255,0);
$c[59] = $image->colorAllocate(225,255,0);
$c[60] = $image->colorAllocate(230,255,0);
$c[61] = $image->colorAllocate(235,255,0);
$c[62] = $image->colorAllocate(240,255,0);
$c[63] = $image->colorAllocate(245,255,0);

```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.6 color_SeqRelation.pl

Anexo 7.6 Códigos de programas en lenguaje PERL
 7.6.6 color_SeqRelation.pl

```

$c[64] = $image -> colorAllocate(250,255,0);
$c[65] = $image -> colorAllocate(255,248,0);
$c[66] = $image -> colorAllocate(255,241,0);
$c[67] = $image -> colorAllocate(255,234,0);
$c[68] = $image -> colorAllocate(255,227,0);
$c[69] = $image -> colorAllocate(255,220,0);
$c[70] = $image -> colorAllocate(255,213,0);
$c[71] = $image -> colorAllocate(255,206,0);
$c[72] = $image -> colorAllocate(255,199,0);
$c[73] = $image -> colorAllocate(255,192,0);
$c[74] = $image -> colorAllocate(255,185,0);
$c[75] = $image -> colorAllocate(255,178,0);
$c[76] = $image -> colorAllocate(255,171,0);
$c[77] = $image -> colorAllocate(255,164,0);
$c[78] = $image -> colorAllocate(255,157,0);
$c[79] = $image -> colorAllocate(255,150,0);
$c[80] = $image -> colorAllocate(255,143,0);
$c[81] = $image -> colorAllocate(255,136,0);
$c[82] = $image -> colorAllocate(255,129,0);
$c[83] = $image -> colorAllocate(255,122,0);
$c[84] = $image -> colorAllocate(255,115,0);
$c[85] = $image -> colorAllocate(255,108,0);
$c[86] = $image -> colorAllocate(255,101,0);
$c[87] = $image -> colorAllocate(255,94,0);
$c[88] = $image -> colorAllocate(255,87,0);
$c[89] = $image -> colorAllocate(255,80,0);
$c[90] = $image -> colorAllocate(255,73,0);
$c[91] = $image -> colorAllocate(255,66,0);
$c[92] = $image -> colorAllocate(255,59,0);
$c[93] = $image -> colorAllocate(255,52,0);
$c[94] = $image -> colorAllocate(255,45,0);
$c[95] = $image -> colorAllocate(255,38,0);
$c[96] = $image -> colorAllocate(255,31,0);
$c[97] = $image -> colorAllocate(255,24,0);
$c[98] = $image -> colorAllocate(255,17,0);
$c[99] = $image -> colorAllocate(255,10,0);
$c[100] = $image -> colorAllocate(255,3,0);

foreach (@arr) {
    chomp;
    unless (/^\S/) {
        my ($id,@data) = (split/\t/);
        foreach $i (1..$#data+1) {
            if ($i =~ /\S/) {
#COORDENADAS DEL TEXTO EN X (ESPECIES)
                $inicio_texto_X_horizontal = (8 + $filaMasGrande * 6);
                $inicio_texto_X_vertical = ($columnaMasGrande * 6); #igual que
                $inicio_vertical_cuadro_1
                $image ->
stringUp(gdSmallFont,$i*$box+$inicio_texto_X_horizontal,$inicio_texto_X_vertical,$data[$i-1],$text_color);
            }
        }
        next;
    }
    $line++;
    my ($id,@data) = (split/\t/);

#COORDENADAS DEL TEXTO EN Y (DOMINIOS)
    $inicio_texto_Y_vertical = (($columnaMasGrande * 6) - 2);
    $inicio_texto_Y_horizontal = ($filaMasGrande * 6);

    $image -> string(gdSmallFont,$box,$line*$box+$inicio_texto_Y_vertical,$id,$text_color);

    foreach $i (1..$#data+1) {

#COORDENADAS DE LOS CUADROS DEL PERFIL

        $inicio_horizontal_cuadro_1 = (10 + $filaMasGrande * 6);
        $inicio_vertical_cuadro_1 = ($columnaMasGrande * 6);

        $image ->
filledRectangle($i*$box+$inicio_horizontal_cuadro_1,$line*$box+$inicio_vertical_cuadro_1,$i*$box+$box+$inicio_h
orizontal_cuadro_1,$line*$box+$box+$inicio_vertical_cuadro_1,$c[$data[$i-1]]);
        $image ->
rectangle($i*$box+$inicio_horizontal_cuadro_1,$line*$box+$inicio_vertical_cuadro_1,$i*$box+$box+$inicio_horizon
tal_cuadro_1,$line*$box+$box+$inicio_vertical_cuadro_1,$black);
    }
}

#print "\n";

```

```
##### ESCALA DE POSITIVOS #####
```

```
foreach $i (0..100) {  
#foreach $i (0..35) {
```

```
#COORDENADAS DE LOS CUADROS DE LA ESCALA  
if ($i =~ /\d+$/) {
```

Anexo 7.6 Códigos de programas en lenguaje PERL

7.6.6 color_SeqRelation.pl

```
$inicio_horizontal_cuadro_escal_1 = 10;  
$inicio_vertical_cuadro_escal_1 = 20;  
$fin_vertical_cuadro_escal_1 = 30;  
$image ->  
filledRectangle($inicio_horizontal_cuadro_escal_1+$i*$box,$inicio_vertical_cuadro_escal_1,$inicio_horizontal_  
cuadro_escal_1+$i*$box+$box,$fin_vertical_cuadro_escal_1,$c{$i});  
$image ->  
rectangle($inicio_horizontal_cuadro_escal_1+$i*$box,$inicio_vertical_cuadro_escal_1,$inicio_horizontal_cuadro_  
_escal_1+$i*$box+$box,$fin_vertical_cuadro_escal_1,$white);  
  
$text = $i;  
  
}
```

```
#COORDENADAS DEL TEXTO DE LA ESCALA
```

```
if ($i =~ /\d+$/) {  
if ($i =~ /^0$/ or $i =~ /^5$/ or $i =~ /^10$/ or $i =~ /^15$/ or $i =~ /^20$/ or $i =~ /^25$/  
or $i =~ /^30$/ or $i =~ /^35$/ or $i =~ /^40$/ or $i =~ /^45$/ or $i =~ /^50$/ or $i =~ /^55$/ or $i =~ /^60$/  
or $i =~ /^65$/ or $i =~ /^70$/ or $i =~ /^75$/ or $i =~ /^80$/ or $i =~ /^85$/ or $i =~ /^90$/ or $i =~ /^95$/  
or $i =~ /^100$/) { #Inicio y final de los valore de la escala  
$text = "$i";  
$inicio_texto_escal_horizontal = 8;  
$inicio_texto_escal_vertical = 45;  
  
$image ->  
stringUp(gdSmallFont,$inicio_texto_escal_horizontal+$i*$box,$inicio_texto_escal_vertical,$text,$text_color);  
}  
}
```

```
# Send to output  
print OUT $image->png;
```

```
#####  
### display full help message  
sub PrintHelp {  
system "pod2text -c $0";  
exit()  
}
```

```
#####  
### Read arguments  
sub ReadArguments {
```

```
my $arg = "";  
while ($arg = shift (@ARGV)) {  
if (($arg eq "-h") || ($arg eq "--help")) {  
&PrintHelp();  
  
## List of options  
=pod  
=item B<-help>  
  
Same as -h  
  
=cut  
} elsif ($arg eq "-i") {  
$filename = shift(@ARGV);
```

```
=pod  
=item B<-i InputMatrix>  
  
name of the file [don't use spaces in the name]
```

```
=cut  
    } else {  
        &PrintHelp();  
    }  
}  
}
```

=====

Anexo 7.4 Tablas de CIs con valores Z de los 9 grupos de la familia Crp/Fnr

Posics	CI-r	G8	Sum-0.75	Prom0.75	Des.E0.75	valor-Z0.75	Area acum
1	0	0	0	0	0	0	
2	0	0	0	0	0	0	
3	0	0	0	0	0	0	
4	0	0	0	0	0	0	
5	0	0	0	0	0	0	
6	0	0	0	0	0	0	
7	0	0	0	0	0	0	
8	0	0	0	0	0	0	
9	0	0	0	0	0	0	
10	0	0	0	0	0	0	
11	0	0	0	0	0	0	
12	0	0	0	0	0	0	
13	0	0	0	0	0	0	
14	0	0	0	0	0	0	
15	0	0	0	0	0	0	
16	0	0	0	0	0	0	
17	0	0	0	0	0	0	
18	1.23	1377.09	1.38	1.93	-0.08		
19	3.54	4062.23	4.06	1.37	-0.38		
20	3.88	4236.3	4.24	1.42	-0.25		
21	4.2	4472.62	4.47	1.66	-0.16		
22	4.25	4560.61	4.56	1.47	-0.21		
23	4.79	4822.84	4.82	1.51	-0.02		
24	4.5	4737.85	4.74	1.48	-0.16		
25	4.42	4626.46	4.63	1.49	-0.14		
26	4.43	4705.07	4.71	1.54	-0.18		
27	3.88	4261.74	4.26	1.45	-0.26		
28	3.63	3951.44	3.95	1.41	-0.23		
29	4.02	4468.6	4.47	1.48	-0.30		
30	3.36	3439.84	3.44	1.12	-0.07		
31	4.11	4526.42	4.53	1.45	-0.29		
32	3.58	3946.37	3.95	1.35	-0.27		
33	4.91	4903.83	4.90	1.50	0.00		
34	4.05	4384.66	4.39	1.45	-0.23		
35	3.34	3702.54	3.70	1.31	-0.28		
36	4.53	4884.69	4.89	0.61	-0.58		
37	4.57	5201.6	5.20	0.41	-1.54	0.06	
38	4.29	4736.37	4.74	0.49	-0.91		
39	3.81	4376.74	4.38	0.58	-0.99		
40	3.69	4181.03	4.18	0.56	-0.88		
41	3.56	3800.93	3.80	0.29	-0.84		
42	4.03	4232.9	4.23	0.27	-0.76		
43	3.7	4560.59	4.56	0.55	-1.57	0.06	
44	3.85	4165.84	4.17	0.69	-0.46		
45	4.19	4678.18	4.68	0.44	-1.11		
46	4.15	4942.13	4.94	0.47	-1.67	0.05	
47	4.96	5325.87	5.33	0.36	-1.01		
48	5.2	5304.94	5.31	0.40	-0.26		
49	4.39	4949.32	4.95	0.47	-1.20		
50	4.59	4967.98	4.97	0.47	-0.80		
51	4.28	4797.53	4.80	0.53	-0.97		
52	4.48	4712.81	4.71	0.52	-0.45		
53	0	0	0	0	0		
54	4.13	4399.06	4.40	0.46	-0.59		
55	3.97	3972.48	3.97	0.16	-0.02		
56	4.34	4662.42	4.66	0.77	-0.42		
57	3.31	3658.16	3.66	0.55	-0.63		
58	3.62	4242.79	4.24	0.56	-1.11		
59	4.48	4981.16	4.98	0.48	-1.04		
60	3.64	3997.16	4.00	0.54	-0.66		
61	3.5	3510.45	3.51	0.13	-0.08		
62	4.5	4828.24	4.83	0.52	-0.63		
63	4.28	4410.97	4.41	0.51	-0.26		
64	2.97	2965.78	2.97	0.11	0.04		
65	0	0	0.00	0.00	0.00		
66	3.21	4013.84	4.01	0.76	-1.06		
67	3.78	4511.65	4.51	0.54	-1.35		
68	3.77	3790.36	3.79	0.21	-0.10		
69	4.32	5079.41	5.08	0.45	-1.69	0.05	
70	4.63	4892.95	4.89	0.67	-0.39		
71	5.29	5362.18	5.36	0.38	-0.19		
72	3.2	3598.52	3.60	0.78	-0.51		
73	5	5168.08	5.17	0.30	-0.55		
74	4.28	4386.92	4.39	0.42	-0.25		
75	4.09	4641.5	4.64	0.56	-0.99		
76	5.49	5544.24	5.54	0.20	-0.27		
77	2.8	2800	2.8	0	0		
78	4.92	5086	5.09	0.45	-0.37		
79	3.93	4316.81	4.32	0.70	-0.56		
80	3.06	3388.24	3.39	0.38	-0.87		
81	4.05	4329.1	4.33	0.45	-0.62		
82	4.52	5142.52	5.14	0.44	-1.42		
83	4.82	5055.11	5.06	0.47	-0.50		
84	4.6	4800.97	4.80	0.24	-0.84		
85	4.39	4449.92	4.45	0.40	-0.15		
86	3.76	4367.99	4.37	0.56	-1.09		
87	3.65	3961.85	3.96	0.52	-0.61		
88	2.81	2808.46	2.81	0.03	0.05		
89	4.44	4754.74	4.76	0.42	-0.75		
90	5.2	5309.81	5.31	0.40	-0.27		
91	4.74	4887.11	4.89	0.12	-1.26		
92	3.16	3469.94	3.47	0.45	-0.69		
93	4.55	4603.83	4.60	0.30	-0.18		
94	5.2	5200	5.20	0.00	0.00		
95	4.68	5036.73	5.04	0.53	-0.67		
96	5.61	5610	5.61	0.00	0.00		
97	4.91	5256.2	5.26	0.46	-0.76		
98	4	4046.35	4.05	0.34	-0.14		
99	3.22	3224.4	3.22	0.11	-0.04		

Anexo 7.4 Tablas de CIs con valores Z de los 9 grupos de la familia Crp/Fnr

100	3.09	3456.89	3.46	0.41	-0.89	
101	3.47	3992.62	3.99	0.55	-0.95	
102	3.71	3968.73	3.97	0.30	-0.86	
103	2.8	2800	2.80	0.00	0.00	
104	3.73	4180.72	4.18	0.37	-1.24	
105	4.98	5121.49	5.12	0.47	-0.30	
106	5.61	5610	5.61	0.00	0.00	
107	5.61	5610	5.61	0.00	0.00	
108	4.27	4615.62	4.62	0.45	-0.76	
109	4.78	5066.62	5.07	0.51	-0.56	
110	5.3	5536.91	5.54	0.18	-1.29	
111	4.2	4862.02	4.86	0.42	-1.56	0.06
112	0	0	0.00	0.00	0.00	
113	0	0	0.00	0.00	0.00	
114	0	0	0.00	0.00	0.00	
115	0	0	0.00	0.00	0.00	
116	4.8	4948.23	4.95	0.19	-0.78	
117	4.69	4771.23	4.77	0.40	-0.20	
118	4.08	4104	4.10	0.29	-0.08	
119	4.84	5092.53	5.09	0.56	-0.45	
120	4.69	4798.6	4.80	0.48	-0.23	
121	3.48	3517.57	3.52	0.55	-0.07	
122	5.09	5171.16	5.17	0.32	-0.26	
123	2.96	3180.11	3.18	0.47	-0.47	
124	4.8	5086.74	5.09	0.39	-0.73	
125	4.57	5115.93	5.12	0.57	-0.96	
126	3.64	4017.7	4.02	0.42	-0.91	
127	3.46	3707.31	3.71	0.30	-0.84	
128	4.73	5255.31	5.26	0.39	-1.35	
129	3.08	3568.07	3.57	0.59	-0.82	
130	4.42	4693.03	4.69	0.42	-0.65	
131	3.99	4519.83	4.52	0.45	-1.18	
132	5.61	5610	5.61	0.00	1.00	
133	3.6	3647.28	3.65	0.49	-0.10	
134	5.38	5535.62	5.54	0.21	-0.75	
135	3.56	4296.53	4.30	0.49	-1.50	0.07
136	4.91	5252.37	5.25	0.40	-0.86	
137	5.36	5515.54	5.52	0.22	-0.70	
138	4.74	4934.26	4.93	0.59	-0.33	
139	4.14	4742.66	4.74	0.62	-0.98	
140	4.49	4800.27	4.80	0.85	-0.37	
141	4.94	5224.68	5.23	0.40	-0.72	
142	4.45	4750.8	4.75	0.46	-0.66	
143	4.19	4621.46	4.62	0.49	-0.88	
144	4.18	4472.34	4.47	0.47	-0.62	
145	3.12	3129.41	3.13	0.13	-0.07	
146	4.82	5193.84	5.19	0.32	-1.18	
147	3.64	3670.15	3.67	0.20	-0.15	
148	4.69	4767.28	4.77	0.39	-0.20	
149	4.63	4945.62	4.95	0.58	-0.55	
150	5.2	5392.15	5.39	0.34	-0.57	
151	5.43	5573.89	5.57	0.12	-1.18	
152	3.93	4468.32	4.47	0.57	-0.94	
153	4.09	4373.54	4.37	0.49	-0.59	
154	4.85	5141.5	5.14	0.45	-0.64	
155	4.77	4975.95	4.98	0.37	-0.55	
156	5.12	5240.18	5.24	0.42	-0.29	
157	3.97	4803.87	4.80	0.61	-1.36	
158	3.77	3770	3.77	0.00	0.00	
159	3.45	3457.59	3.46	0.13	-0.06	
160	4.75	5090.63	5.09	0.55	-0.62	
161	5.16	5383.88	5.38	0.34	-0.66	
162	3.95	4357.88	4.36	0.47	-0.87	
163	4.98	5084.53	5.09	0.43	-0.24	
164	4.71	4958.17	4.96	0.28	-0.88	
165	5.38	5466.05	5.47	0.30	-0.29	
166	3.78	4015.96	4.02	0.26	-0.91	
167	4.57	4943.28	4.94	0.50	-0.75	
168	5.61	5610	5.61	0.00	1.00	
169	5.47	5580.81	5.58	0.10	-1.07	
170	4.25	4250	4.25	0.00	0.00	
171	5.2	5319.13	5.32	0.39	-0.31	
172	5.31	5443.77	5.44	0.30	-0.45	
173	5.61	5610	5.61	0.00	0.00	
174	0	0	0	0	0	
175	0	0	0	0	0	
176	0	0	0	0	0	
177	0	0	0	0	0	
178	0	0	0	0	0	
179	0	0	0	0	0	
180	0	0	0	0	0	