

¿Qué es un gen?

'Gene' is not a typical four-letter word. It is not offensive. It is never bleeped out of TV shows. And where the meaning of most four-letter words is all too clear, that of gene is not. The more expert scientists become in molecular genetics, the less easy it is to be sure about what, if anything, a gene actually is.

Helen Pearson: Genetics: What is a gene? (*Nature* 441, 398-401)

¿Qué es un gen?

Concepto mendeliano:

Unidad discreta de herencia (factor) que expresa un fenotipo dado
En 1909, W. Johanson acuñó el término 'gen'

Morgan y CIA

Los genes se encuentran en loci específicos en los cromosomas

G.W. Beadle and E.L. Tatum

Confirmaron en 1941 la hipótesis de un gen una enzima postulado en 1908 por E.R. Garrod.
Los genes se heredan y pueden sufrir mutaciones.

S. Benzer (1957)

Acuñó diferentes términos que finalmente daban al dogma central de la biología molecular:

gen (DNA) --> transcrito (RNA) --> proteína

La definición de un gen

A locatable region of genomic sequence, corresponding to a unit of inheritance, which is associated with regulatory regions, transcribed regions and/or other functional sequence regions.

Helen Pearson: Genetics: What is a gene? (*Nature* 441, 398-401)

Genes RNA

- tRNA, snRNA, rRNA, snoRNA, lncRNA, vault RNA, piRNA, microRNA ...

Funciones: Ribonucleoproteína, biogénesis de otros RNAs, traducción, regulación génica

Genes codificantes de proteínas

- **Procariotas:**
 - Sin intrones
 - Regulación más simple
- **Eucariotas**
 - Los genes contienen intrones,
 - Regulación compleja

La tarea de buscar los genes codificantes

El genoma humano tiene 3,200,000,000 pares de bases pero solo entre 20,000 y 25,000 genes codificantes que ocupan una proporción muy pequeña del genoma (aprox. 1,5%)

La tarea de buscar los genes codificantes

El genoma humano tiene 3,200,000,000 pares de bases pero solo entre 20,000 y 25,000 genes codificantes que ocupan una proporción muy pequeña del genoma (aprox. 1,5%)

¿Como podemos encontrar las islitas pequeñas que son los genes dentro del océano genómico?

La tarea de buscar los genes codificantes

El punto de partida: una secuencia de DNA

GGCCTACCATTTTCTACAGATGGCATTATTCAGCAAGGGAATGTTGTAGGTATAGTAATTTATCTTTTTTCCATGAAACATGTGTT
ATCTTCTTGCTCTATTTTGAAGAGATATTATGATAAATGATGCCACTATTTAAACTTGTAAAGCTTCTTGAGGTTCTCAGAAAAA
GTACAAGGTATAGGTCAAAAAGAGTAATTTTGTCTCTCACAAAAGAATTTTCTGTTCAAAAATGCTGTAACTTGAAAAAGTCATA
TCTCCAAAAGAACGTTCTAATAGGTGGGATAATCATCATCATATTTCCAGCTATAATCCATCATATTTCCAGGTGTATCCCTTCCCT
TAAAAGTAAAAAGGAGAAAATCACTTTGCTGTGTGGAAGTGCTTTGTGTGTTTTATAAAATCCACATAAAATGTTTCGTGGTTCAAG
GAGTGGCTGCCGCATTTCTAAGGCAAATAAGCTTATCCCTGAAGAAGCTTGAAGATCTGACTCTCAAGAATCTTTGTGTTCTGGATG
AAAAAGGAATCATTATTTTTTAACTCCCTGCTCAAGTCAGTACTTCCACTGGCACCCAATCCATTGGTTATAAGTACCCATGAAT
TGCCTTGCCAATTATAGGCAGGATAACCTAGGAGGACCAACCAAAACCAACTCCTGGTAAGCTTTATTAGCCTTTTCTACAATGTG
TCTCAGTTTTCTTTTACACCTTGGAGTGAAGACAATCTCTACTCACACTGTTGTGGTATGTGTCTGTCTCAGATGTGTGTA
TAAGGTGTATTTCTCAGCTCGGTATCTGAAGTAATTTCTTTTTCTTTTTTTTTTAAAGACTGCTGTCTTAAGGCTTAAAAGTAA
TTTAAGGTTTTGTTTTTGTATTTTTAAGAGTTGCAGTTTTCTTTGGAGACAACAAAGTACTTTTTTCCAAAGATGAAAAAGGCATA
TTCAGACTTTCTATGTTATGTTCCATAATTGTTTTGCTTATCACTGTATAACATTCTCATTATCATGGGCTACTTTAAGCTAATCAG
TTGCATTTTGGGGTCCCTCATTTAGAAGTCAATCTGGAAAAAGATGGCTGTAAACTTAGATAAACAGTATGCTATTTCTCCTTCC
TTGTCTGCCTTTACAGATGTTTCTCTTTCTCTCTCTTCCCTGACTATCTCTCTCAATGATCATTGACAGCTCCCGCATCCCCTC
TTTCACTCAGCTTCACTCCACCATGACCAGAGCACCTCTCCTGCTACTATGTGTTGCCCTGGTGCTGCTTGGGCATGTGAATGGAGC
CACAGTAAGAAATGAGGACAAATGGAAGCCACTCAACAACCCAGAAAACAGAGATCTGGTAAGAACACACTTGGGGCCTGGGTTTCA
GGATAGAATAATGCCTCTTTGGATTTACGCACACAGCATTCTGTTTTCCAGAACGTAGTGGAAGGGGGCTACCAAGGGTTGAGAA
AATTCTCCAACAGGATTGGACCAGAGTTTTACTACCCTAGTGGTTGGCAACCTATTGTTTCTCTGGCCAGTTTACATTTGTTTCTGT
AAAGAGAAACAAATTTATATTTAGAAACAATCTTATACACACATAGTACAATTCATTTGTAAGTATGCTTATATAACACATTGACTTT
GTAGCCTTGATTTCTAATACATATAAAAAATAAAGAGTGTAACAGGGAATTGCTCGTGCATAGCTTAGATTGAACTGATGATAAAAC
AAGTTATCATATAACTAGAGATTGCCAAGCTCAGCTGCAAAAATTGAAAGTTGACCCTGCATGGGTATCAATTTACTGGACATATTT
TCAACATATAGTTATCTATAAAATAGAATTCAAAGATCAGTAAACATTTCTTTTTAAAAAAAATGCCAAATACAATAAACATGTTGA
CATGACTATTTTGAATATGGTTTTTCTAAAGTGAATGATGGGAAATATAACAACATAGGATTAATAAACTGAGATAAGAGACCA
GCAGGATTTGTTTTCTGGTCCAGCTGCTGACGAAAACAGGATCTGGTC

La tarea de buscar los genes codificantes

Diferentes aproximaciones

Búsqueda *ab-initio* (a partir de primeros principios)

Buscar signos o señales indicadores de genes codificantes de proteínas (promotores, composición de secuencia, existencia de una ORF (**O**pen **R**eadin**G** **F**rame, Marco de lectura abierto))

Aproximación extrínseca

- Se basa en evidencia externa, es decir se busca en el genoma objetivo secuencias similares a la de una secuencia externa
- Estas secuencias pueden provenir de experimentos de secuenciación del RNA total de la misma especie o de genes conocidos en otras especies

Las predicciones deben de ser confirmado y validado por experimentos en el laboratorio

Confirmar que realmente se transcriben y traducen

Ab-Initio

¿Que propiedades de un gen codificante podemos buscar en una secuencia de DNA?

Marco de lectura abierto (ORF)

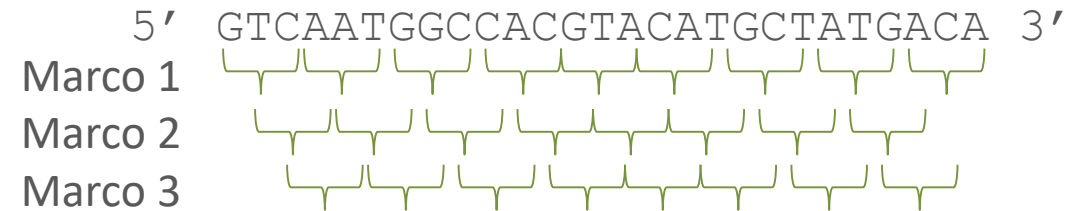
La lectura de la información genética se lleva a cabo de triplete (codón) en triplete



- Un marco abierto de lectura empieza con un codón de inicio y termina con un codón de parada que se encuentran en la misma pauta de lectura
- No puede haber codones de parada en la pauta de lectura

Búsqueda de ORFs

Buscamos primero en la **hebra directa** en los 3 marcos posibles

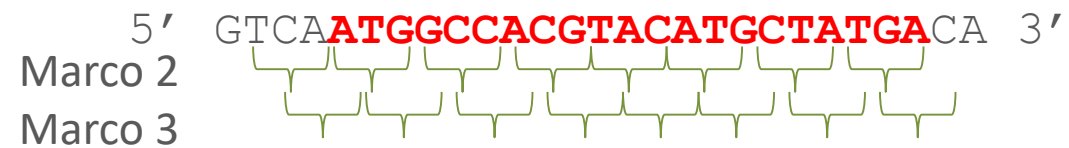
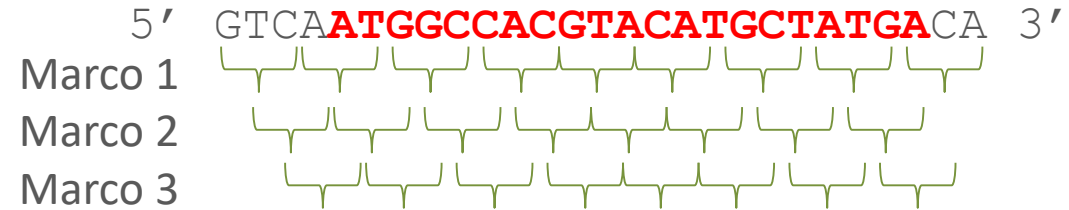


Pero también tenemos que buscar en la secuencia **inversa complementaria**



Búsqueda de ORFs

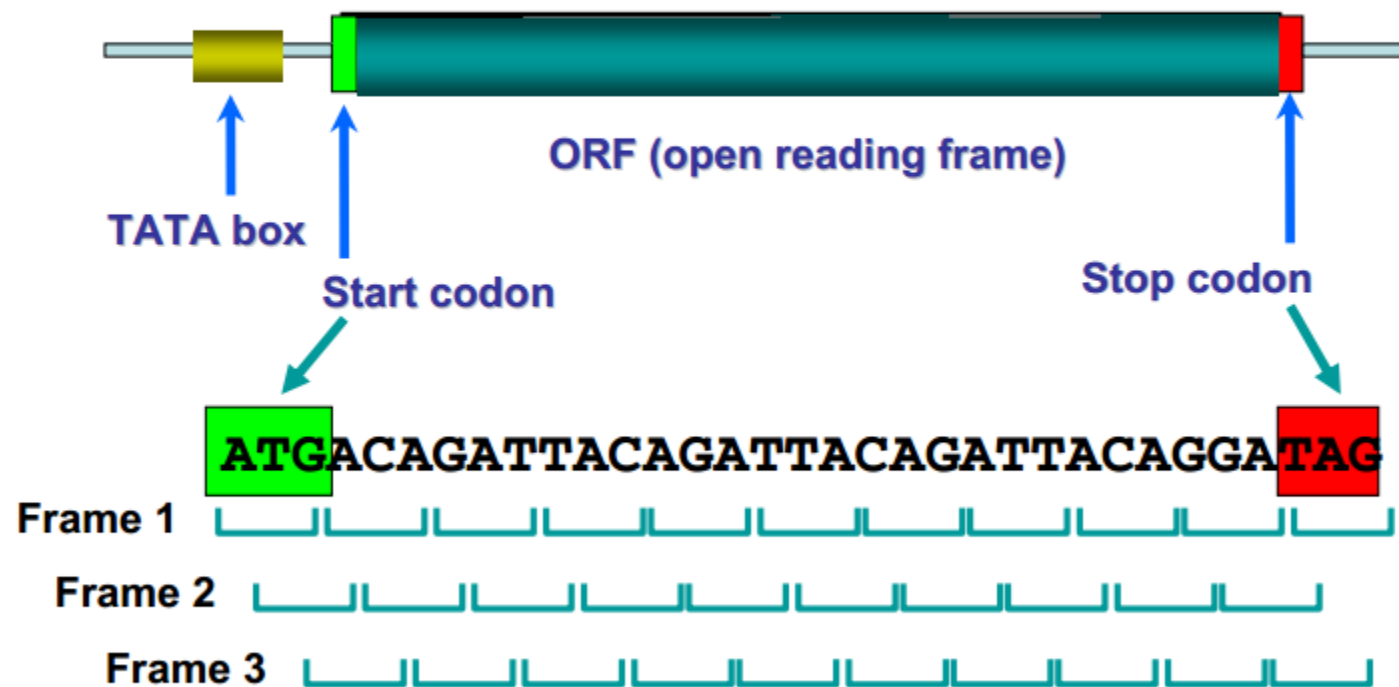
Buscamos primero en la hebra directa en los 3 marcos posibles



Hay un marco abierto de lectura en el segundo marco

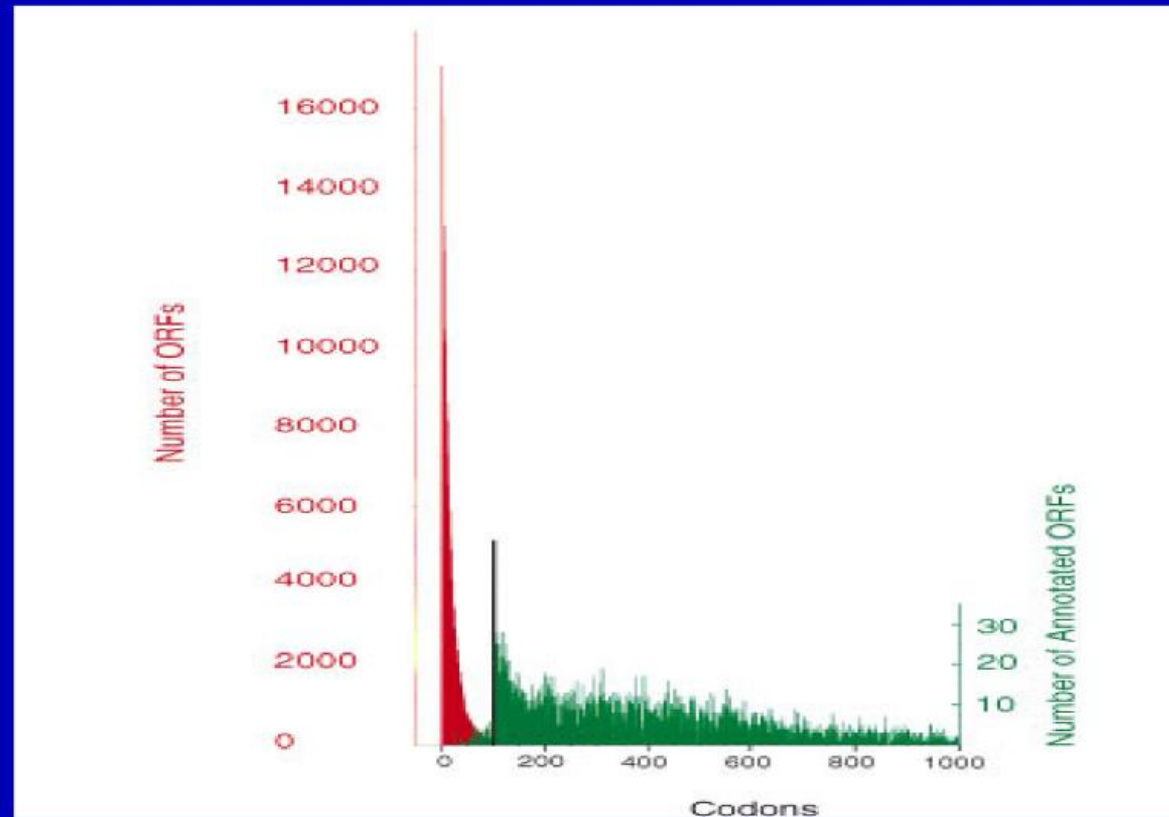
En procariontas el ORF está presente en la secuencia de DNA por la ausencia de intrones

Prokaryotic Gene Structure

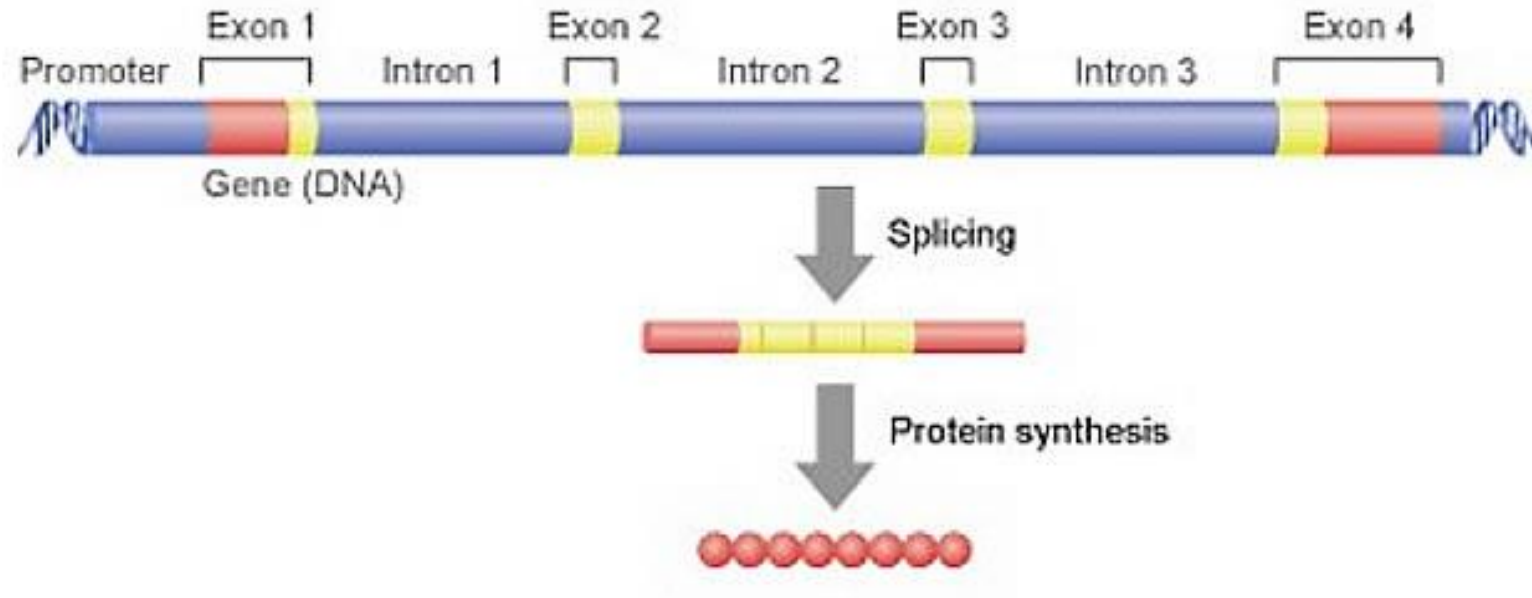


Que es la diferencia entre ORF y CDS?

Una medida simple: comparación de las longitudes de ORFs anotadas y espúreas en *S. cerevisiae*

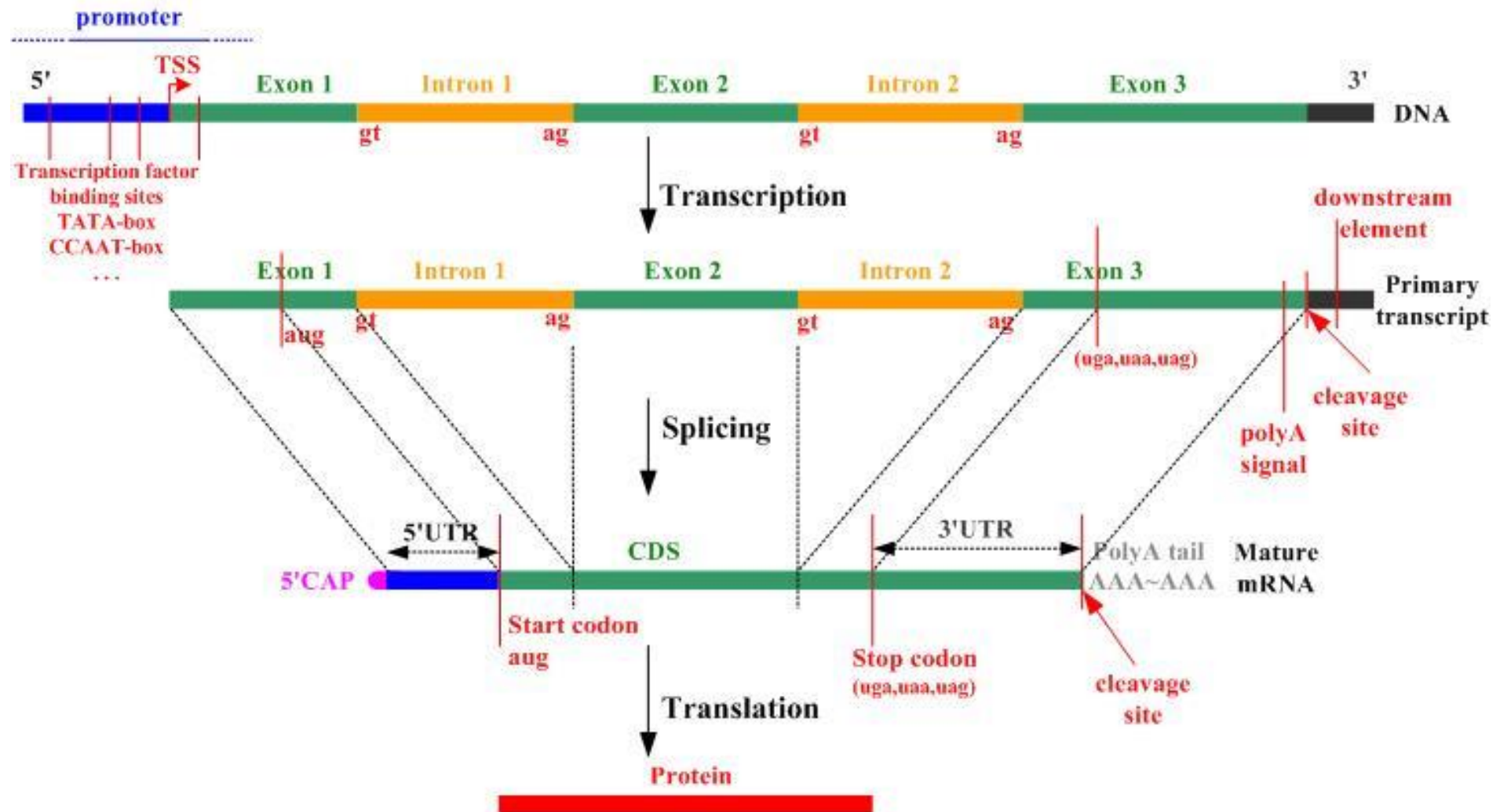


Predicción de genes en eucariotas



¡En genes con intrones, el ORF no se presenta en la secuencia de DNA!

Detectar los genes



Predicción de genes en eucariotas

¿Cómo podemos buscar/predecir genes que contienen intrones?

1) Potencial codificador de los cdsExons

- Los exones tienen otra composición de secuencia que las regiones no-exónicas (frecuencias de hexámeros).
- Esta información puede ser usada por métodos de aprendizaje automatizado

2) Señales como donadores y aceptores

- Los 2 primeros nucleótidos de un intron son en la inmensa mayoría de los casos 'GT' (Donador)
- Los 2 últimos nucleótidos de un intron son en la inmensa mayoría de los casos 'AG' (Aceptor)

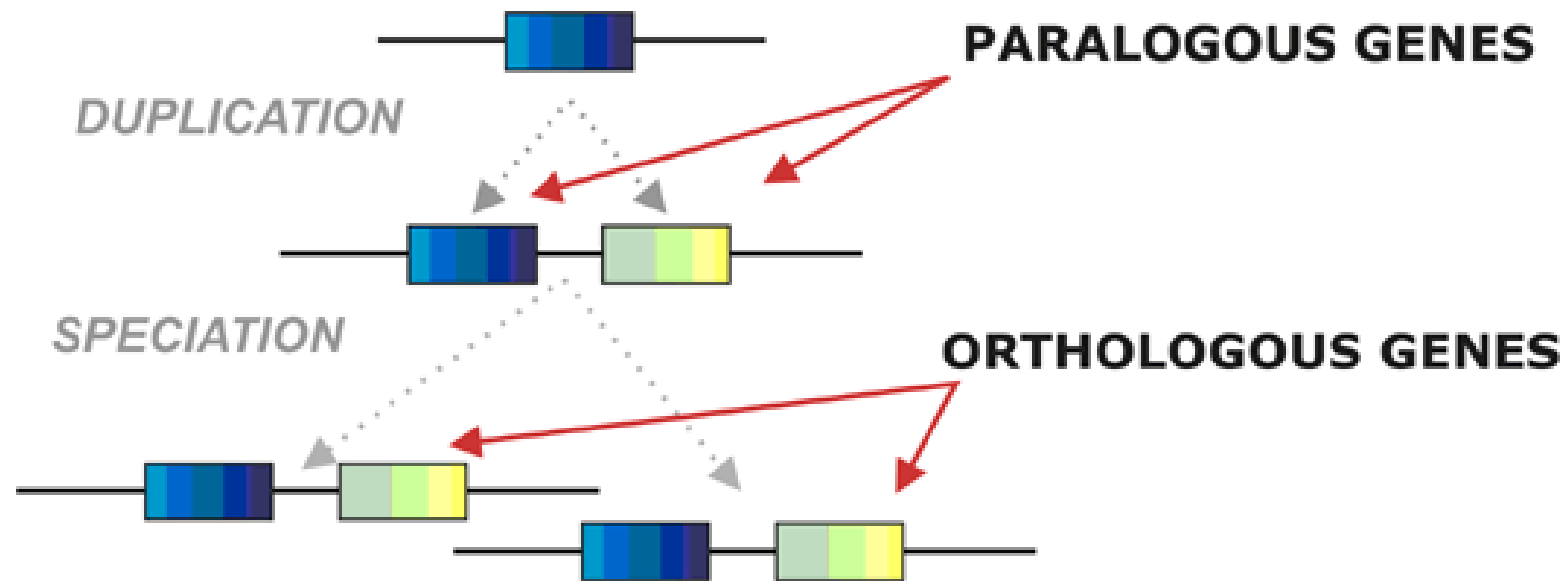
3) Promotores (islas CpG, composición de secuencia)

- La búsqueda de genes se puede refinar mediante la detección de regiones o motivos asociados al promotor (TFBS, islas CpG)

Usar información extrínseca

La idea:

Alinear secuencias de mRNAs conocidos (de otras especies o obtenidas en estudios de expresión génica) frente al genoma objetivo



Usar información extrínseca

La idea:

Alinear secuencias de mRNAs conocidos (de otras especies o obtenidas en estudios de expresión génica) frente al genoma objetivo

Por ejemplo:

- Tenemos una secuencia mRNA de humano y buscamos si esta secuencia también se encuentra en el chimpancé
- Si encontramos la secuencia y si además observamos que no hay mutaciones que cambien el ORF, podemos concluir que es probable que este gen es funcional también en el chimpancé

Pero:

¡La prueba definitiva que el posible gen se transcribe y traduce tiene que venir por parte de experimentos de laboratorio (Northern Blot, Western Blot)!

Alineamiento: similitud de secuencia

$$\text{Similitud de secuencia} = \frac{\text{Número de bases emparejadas}}{\text{Número de bases alineadas}}$$

```
ACGTATAGCG
|||||  |||
ACGTA--GCG
```

Ungapped: 100%

Gapped: 80%

```
ACGTATAGCG
|||||  |
ACGTAGCG
```

75%

Usar información extrínseca

- El BLAT de ADN está diseñado para encontrar de forma rápida una secuencia en un genoma dado
- Optimizado para las secuencias con mayor similitud de 95% y mayor longitud que 25 bp → probablemente falla detectar alineamientos mas divergentes o cortos

Human BLAT Search

BLAT Search Genome

Genome: Assembly: Query type: Sort output: Output type:

Paste in a query sequence to find its location in the the genome. Multiple sequences may be searched if separated by lines starting with '>' followed by the sequence name.

File Upload: Rather than pasting a sequence, you can choose to upload a text file containing the sequence.
Upload sequence:

Only DNA sequences of 25,000 or fewer bases and protein or translated sequence of 10000 or fewer letters will be processed. Up to 25 sequences can be submitted at the same time. The total limit for multiple sequence submissions is 50,000 bases or 25,000 letters.

For locating PCR primers, use [In-Silico PCR](#) for best results instead of BLAT.

Kent WJ.
BLAT--the BLAST-like alignment
tool.
Genome Res. 2002

BLAT

```
>panTro4_refGene_NM_001252546  
GGGTGGGAGCGCGTGCTGTTGGGAGTTGCTTGGAGGTTGGCGGCGCGGGG  
CTGAAGGCTAGCAAACCGAGCGATCATGTGCGCACAAACAAATTTACTATT  
CGGACAAATACGACGACGAGGAGTTTGAGTATCGACATGTCATGCTGCCC  
AAGGACATAGCCAAGCTGGTCCCTAAAACCCATCTGATGTCTGAATCTGA  
ATGGAGGAATCTTGGCGTTCAGCAGAGTCAGGGATGGGTCCATTATATGA  
TCCATGAACCAGAACCTCACATCTTGCTGTTCCGGCGCCCCTACCCAAG  
AAACCAAAGAAATGAAGCTGGCAAGCTACTTTTCAGCCTCAAGCTTTACA  
CAGCTGTCCTTACTTCCTAACATCTTTCTGATAACATTATTATGTTGCCT  
TCTTGTTTCTCACTTTGATATTTAAAAGATGTTCAATACACTGTTTGAAT  
GTGCTGGTAACTGCTTTGCTTCTTGAGTAGAGCCACCACCACCATAGCCC  
AGCCAGATGAGTGCTCTGTGGACCCACAGCCTCAGCTGAGTGTGACCCCA  
G
```



Alinear frente al genoma mediante BLAT

BLAT

Longitud y coordenadas de alineamiento de la secuencia de entrada

Región en el genoma
Cromosoma, hebra, inicio, final

Score/puntuación

Similitud de secuencia

Longitud de la región en el cromosoma

BLAT Search Results

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
browser details	panTro4_refGene_NM_001252546	547	1	551	551	99.9%	1	+	154947147	154951489	4343
browser details	panTro4_refGene_NM_001252546	513	13	551	551	97.6%	14	+	96689382	96689920	539
browser details	panTro4_refGene_NM_001252546	503	1	551	551	95.3%	11	+	82920366	82920915	550
browser details	panTro4_refGene_NM_001252546	502	32	551	551	98.3%	8	-	81556932	81557451	520
browser details	panTro4_refGene_NM_001252546	502	14	551	551	96.3%	X	+	30635509	30636045	537
browser details	panTro4_refGene_NM_001252546	470	76	551	551	99.4%	5	-	61807834	61808309	476
browser details	panTro4_refGene_NM_001252546	468	76	551	551	99.2%	10	+	29986864	29987339	476
browser details	panTro4_refGene_NM_001252546	388	13	551	551	91.3%	5	+	143818706	143819276	571
browser details	panTro4_refGene_NM_001252546	93	14	130	551	89.8%	19	-	9115818	9115934	117
browser details	panTro4_refGene_NM_001252546	26	492	534	551	71.5%	2	+	7182785	7182816	32
browser details	panTro4_refGene_NM_001252546	24	452	480	551	76.0%	4	-	24220330	24220354	25
browser details	panTro4_refGene_NM_001252546	24	292	316	551	100.0%	1	-	62669892	62669936	45
browser details	panTro4_refGene_NM_001252546	23	485	507	551	100.0%	11	+	907550	907572	23
browser details	panTro4_refGene_NM_001252546	20	396	415	551	100.0%	3	-	155686371	155686390	20

BLAT: Resuelve la estructura génica

Dos exones de
NM_001252546
(chimpancé) alineados
en el genoma humano
(bases en azul)

```
gtgctccgtg gggctcccat tatgttgtag ataatcaact agaagactat 154950196
aaggtccatg aaggcggaga aaatgtgcct tgaacttaag tgctcattca 154950246
ctggaattat tcaactaattt gtatttttaa tatactctcg gaatttatat 154950296
aacatagcct gtattttctag ctgctgctac cccacaataa aattatataa 154950346
actctgacat caaaaaccag tttttctttc ctttcttggc cttgtaaaca 154950396
ggcatttgtc taagagactg tatctggtac taacataaat tccccacttc 154950446
cccgtttctg ttacagACAT GTCATGCTGC CCAAGGACAT AGCCAAGCTG 154950496
GTCCTAAAA CCCATCTGAT GTCTGAATCT GAATGGAGGA ATCTTGCGT 154950546
TCAGCAGAGT CAGGGATGGG TCCATTATAT GATCCATGAA CCAGgtcagt 154950596
gcactggcta aaaacaacca tatagaactg ctacactgag agaatgaaag 154950646
aataagattg tataacccaa atagggagat aggaaatggt ttactggttc 154950696
cttccccctc cagtcgtggg ggattttttt aaaaaaaaaac tagtgaccaa 154950746
aaataagact aaaatatctg ggaagttcag agacaacctg tcaactgaaa 154950796
acctcctgta atctttcatt caatcagagg gtattctttt taaggccaca 154950846
tatagcctga tcatagcccc tgccctcattc tccatcgaaa acattcttgg 154950896
atgtgttcta aataagcaaaa ggaaagtata tttattgata agacaccaga 154950946
caccagctg ccaggcaaaa ctaataaggg acaccctggg gctgtataaa 154950996
catagcaaaa gaactgatat taacaattct gtacttggca gacagtccag 154951046
acttctgggt ctgcttctaa ggccatagc ttaagtcttt atttagttat 154951096
aaagatctga gtgggtgata gctggggagg tggtagtgga atacactata 154951146
ggttgtaacat agaatggtgt gagcttgtct cttagatttc cctcactctt 154951196
tcagAACCTC ACATCTTGCT GTCCGGCGC CCACTACCCA AGAAACCAA 154951246
GAAATGAAGC TGGCAAGCTA CTTTTCAGCC TCAAGCTTTA CACAGCTGTC 154951296
CTTACTTCCT AACATCTTTC TGATAACATT ATTATGTTGC CTTCTTGTTT 154951346
CTCACTTTGA TATTTAAAAG ATGTTCAATA CACTGTTTGA ATGTGCTGGT 154951396
AACTGCTTTG CTTCTTGAGT AGAGCCACCA CCACCATAGC CCAGCCAGAT 154951446
GAGTGCTCTG TGGACCCACA GCCTaAGCTG AGTGTGACCC CAGaagccac 154951496
gatgtgctct gtatccagaa cacacttggc agatggagga agcatctgag 154951546
tttgagacca tggctgttac agggatcatg taaacttget gtt
```