

ANÁLISIS DE UNA SECUENCIA USANDO HERRAMIENTAS DISPONIBLES EN LA RED

Samuel Pintos González

1. SELECCIÓN DE LA SECUENCIA E IDENTIFICACIÓN DE ORFs

He utilizado “ORF finder” para el primer ejercicio, copiando la secuencia proporcionada en esa página web con los parámetros por defecto, salvo en la sección “ORF start codon to use”, ya que como no sé a qué organismo pertenece la secuencia, he seleccionado “ATG and alternative initiation codons”

Elijo ORF1 [Figura 1] ya que encapsula al resto de ORFs en la secuencia y es el más grande.

2. ANÁLISIS FUNCIONAL DE LA SECUENCIA

La secuencia de aminoácidos que viene dada en la propia herramienta de “ORF finder” al seleccionar el ORF1 es:

MNYSFTEKKRIRKSFAKRENVLEVPFLLATQIDSYAKFLQLENAFDKRTDDGLQAAFNSIFPIVSHNGYARLE
FVYYTLGEPLFDIPECQLRGITYAAPLRARIRLVILDKEASKPTVKEVRENEVYMGEIPLMTPSGSFVINGTER
VIVSQLHRSPGVFFEHDKGKTHSSGKLLFSARIIPYRGSWLDFEFDPKDLLYFRIDRRRKMPVTILLKALGYN
NEQILDIFYDKETFYLSSNGVQTDLVAGRLKGETAKVDILDKEGNVLVAKGKRITAKNIRDITNAGLTRLDVE
QESLLGKALAADLIDSETGEVLASANDEITEELLAKFDINGVKEITTLINELDQGAYISNTLRTDETAGRQAA
RVAIYRMMRPGEPPTEEAVEQLFNRLFFSEDSYDLRVGRMKFNTRTYEQLSEAQQNSWYGRLLNETF
AGAADKGGYVLSVEDIVASIATLVELRNHGHEVDDIDHLGNRRVRSVGELTENQFRSGLARVERAVKERLN
QAESENLMPHDLINAKPVSAAIKEFFGSSQLSQFMDQTNPLSEVTHKRRVSALGPGGLTRERAGFEVRDV
HPHYGRVCPIETPEGPNIGLINSLSVYARTNDYGLETYPYRRVIDGKVTEEIDYLSAIEEGRYVIAQANADL
DSDGNLIGDLVTCREKGETIMATPDRVQYMDVATGQVVSVAASLIPFLEHDDANRALMGANMQRQAVP
CLRPEKPMVGTGIERSVAVDSATAIVARRGGVVEYVDANRVVIRVHDDEATAGEVGVDIYNLVKFTRSNQS
TNINQRPVAVKAGDVLQRGDLVADGASTDLGELALGQNMTIAFMPWNGYNYEDSILISEKVAADDRYTSI
HIEELNVVARDTKLGAEDITRDIPNLSERMQNRLDESGIVYIGAEVEAGDVLVGKVTPKGETQLTPEEKLLR

AIFGEKASDVKDTSLRMPTGMSGTVIDVQVFTREGIQRDKRAQSIIDSELKRYRLDLNDQLRIFDNDADFDR
IERMIVGQKANGGPMKLAGSEITTEYLAGLPSRHDWFDIRLTDEDLAKQLELIKLSLQQKREEADELYEIK
KKKLTQGDELQPGVQKMKVVFIAIKRRLQAGDKMAGRHHGNGKVVSRLPVEDMPYMADGRPVDIVLNP
LGVPSRMNIGQILEVHLGWAAKGIGERIDRMLKERRKAGELREFLNKLYNGSGKKEDLDSLTDEEIIELASN
LRKGASFASPVFDGAKESEIREMLNLAYPSEDPEVEKLGFNDSKTQITLYDGRSGEAFDRKVTVGVMHYLK
LHHLVDEKMHARSTGPYSLVTQQPLGGKAQFGGQRFGEMEVWALEAYGAAYTLQEMLTVKSDDVNGR
TKMYENIVKGEHKIDAGMPESFNVLVKEIRSLGLDIDLERY

Copio la secuencia de aminoácidos en protein BLAST seleccionando la base de datos "Reference proteins (refseq_protein)". A continuación me aparece que su cobertura es el 100% de la secuencia con un porcentaje de identidad del 100% y un E-valor de 0.0 [Figura 2]. Pertenecería a la especie *Neisseria gonorrhoeae*, donde la descripción de esta proteína sería "DNA-directed RNA polymerase subunit beta".

Con core nucleotide database (core nt) para el nucleotide blast, me aparece la misma cobertura, E-valor y porcentaje de identidad pero con descripciones distintas una de otra, a pesar de que el organismo coincide también [Figura 3].

3. COMPARACIÓN EVOLUTIVA

[WP_003690105.1](#) *Neisseria gonorrhoeae* (el organismo al que pertenece la secuencia)

[WP_119518745.1](#) *Pseudomonas aeruginosa*

[WP_353524646.1](#) *Escherichia coli*

[Figura 4].

Haciendo el alineamiento a 3Bits para ser más estrictos que a 2Bits, las 2 regiones regiones que tienen una conservación más grande (en cuanto a longitud de la cadena de aminoácidos sin variaciones) entre estos 3 organismos, siendo estas para *Neisseria gonorrhoeae* son:

aminoácido 527 al 603 y aminoácido 1260 al 1361. Hay varias regiones donde faltan datos (en gris, siendo la más grande la del aminoácido 420 al 444) y hay diferentes regiones bastante variables [Figura 5].

4. ANÁLISIS ESTRUCTURAL

Utilizando la herramienta PROSITE, introduzco la secuencia de aminoácidos que me dio ORFfinder.

Obtengo un dominio proteico entre los aminoácidos 1104 y 1116 (correspondiendo al código "PS01166" y con la serie de aminoácidos "GDKMAGRHHGNKGV"), el cual coincide con una zona muy conservada entre el organismo de la secuencia y el resto de organismos del alineamiento (aunque no está dentro de las zonas conservadas más grandes que se han destacado previamente).

Según los datos aportados por la propia herramienta seleccionando la secuencia conservada, contiene dos lisinas muy bien conservadas que son parte del sitio activo del extremo C-terminal de todas las cadenas beta y este sería el motivo por el que varía tan poco entre especies.

ANEXO: IMÁGENES DE APOYO AL TRABAJO Y SU DESCRIPCIÓN



Figura 1. Representación de ORFfinder de todos los ORFs de la secuencia seleccionada en barras horizontales rojas, siendo ORF1 el primero y el más grande, abarcando al resto de ORFs

Query Cover	E value	Per. Ident	Acc. Len	Accession
100%	0.0	99.93%	1392	WP_404518383.1
100%	0.0	99.93%	1392	WP_218422894.1
100%	0.0	100.00%	1392	WP_003690105.1

Figura 2. Resultados del análisis por protein BLAST de la secuencia de aminoácidos que se ha obtenido previamente con ORFfinder

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
✓	Neisseria gonorrhoeae strain 9035 chromosome .complete genome	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2223133	CP104546.2
✓	Neisseria gonorrhoeae isolate G97687 genome assembly .chromosome .1	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2174841	LS999565.1
✓	Neisseria gonorrhoeae strain 10538 chromosome .complete genome	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2223795	CP104548.2
✓	Neisseria gonorrhoeae strain AT159 chromosome .complete genome	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2232771	CP097846.1
✓	Neisseria gonorrhoeae strain TH2288 chromosome .complete genome	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2227834	CP138339.1
✓	Neisseria gonorrhoeae strain CT530 chromosome .complete genome	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2169323	CP048254.1
✓	Neisseria gonorrhoeae strain JS-23-24 chromosome .1 .complete sequence	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2231116	CP195874.1
✓	Neisseria gonorrhoeae strain WHO K genome assembly .chromosome .1	Neisseria gonorrhoeae	7718	7718	100%	0.0	100.00%	2169846	LT591908.1

Figura 3. Nucleotide BLAST donde se analiza la secuencia de nucleótidos obtenida del enunciado del ejercicio donde aparece la especie a la que pertenece (*Neisseria gonorrhoeae*) y distintos valores como su cobertura y porcentaje de identidad.

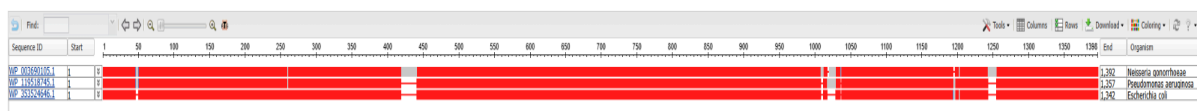


Figura 4. Selección de 3 organismos distintos para hacer un alineamiento. Las barras rojas horizontales indican zonas conservadas entre estos organismos, mientras que los espacios se deben a zonas variables o desconocidas

WP_003690105.1	1	MNYSFTEKKRIRKSF	AKRENVLEVPFL	ATQIDSYAKFLQ	LENAFDKRTDDGLQAAFN	SIFPIVSHNGYAR	LEFVYYTLG	80																																																																								
WP_119518745.1	1	MAYSYTEKKRIRKDF	SKLPDVMVPYLLA	TQLDSYREFLQ	AGATKEQFRDIGLHAAFK	SVFPIISYSGNALE	VVGYRLG	80																																																																								
WP_353524646.1	1	MVYSYTEKKRIRKDF	KRPQVLDPVYLLS	IQLDSFQKFI	EQDPEGQ----	YGLEAAFRSVFPI	IQSYSGNSELQVVS	YRLG 76																																																																								
WP_003690105.1	81	EPLFDIP	ECQLRGITYA	APLARIRL	VILDK	EASKPTV	KEVRENEVYMG	EIPLMTPSGSFVINGTERVIV	SQHLRSPGVF 160																																																																							
WP_119518745.1	81	EPAFDVK	CEVLRGVTF	AVPLRVK	VLRIIFDRESS	KAIDKEQ	EVYMG	EIPLMTE	NGTFIINGTERVIV	SQHLRSPGVF 160																																																																						
WP_353524646.1	77	EPVFDV	QECQIRGV	TVSAPLRV	KLRLVI	YERAE	EGTVKDI	KEQEVYMG	EIPLMTDNGTFVINGTERVIV	SQHLRSPGVF 156																																																																						
WP_003690105.1	161	FEHDKGK	THSSGKLL	SARIIPYRGS	WLDFFD	PKDILLYFRIDRRR	KMPVTILLK	ALGYNN	EQILDFYDKET	FYLLSSNG 240																																																																						
WP_119518745.1	161	FDHDKG	THSSGKLL	SARIIPYRGS	WLDFFD	PKDILLYFRIDRRR	KLPASVLL	RALGYST	EEILNAFYATN	VNFHIKGET 240																																																																						
WP_353524646.1	157	FDSKKG	THSSGKVL	YNARIIPYRGS	WLDFFD	PKDILLYFRIDRRR	KLPATIIL	RALNYTTE	QILDLFF	EKVIFEIRDNK 236																																																																						
WP_003690105.1	241	VQTDLV	AGRLKGET	AKVDIL	KEGNVL	AKKRITAK	NIRDI	TNAGL	TRL	DVEQESLLG	KALAADLIDSET	GEVLASAND 320																																																																				
WP_119518745.1	241	LNLELVP	QRLRGE	VASIDIKD	GSGKVI	EQRRITAR	HINQLEK	AGVSQ	LEV	FDYLI	GR	TI	AKAIVHPATGE	IIAECNT 320																																																																		
WP_353524646.1	237	LQMLVPE	RLRGET	ASF	DI--	FANGKV	VVEQRRITAR	HIRQLEK	DDVKL	E	VP	VEYI	AGKVVAKD	YIDEST	GELICAANM 315																																																																	
WP_003690105.1	321	EITEELL	AKFDING	VKEIT	TYLINE	LDQ	GAYISNT	LRD	ETAGR	QAARVAI	YRM	MRP	GEPTTE	E	EA	VEQLFN	RLFFS	EDSY 400																																																														
WP_119518745.1	321	ELTDL	LAKVAKQ	VRIET	LYTND	IDCGP	FISDTL	IKDNT	SNQLE	ALVEI	YRM	MRP	GEPTTE	K	AA	ETL	GNL	FFSAERY 400																																																														
WP_353524646.1	316	ELS	LDLAKLSQ	SGHKRI	ETLFTND	LHGP	YIS	ETLR	VDPTNDR	Q	SALVEI	YRM	MRP	GEPTTE	R	AA	ESL	FN	FFS	EDRY 395																																																												
WP_003690105.1	401	DLSRV	GRMKFN	TRTYE	QKL[22]	KGGV	YLSVED	IVAS	IATL	VELRN	GHGEV	DDI	DHLGN	RRV	SV	GELT	ENQFR	SGL	ARV 498																																																													
WP_119518745.1	401	DLS	AVGRMK	FNRRIG	RT	E	EGP	VL	SKED	I	IDL	KTLD	VRNG	KGI	VDDI	DHLGN	RRV	CVG	MAENQ	FRVGL	VRV 476																																																											
WP_353524646.1	396	DLS	AVGRMK	FNRRS	LLREEI		EGS	GIL	SKDD	I	IDV	MKLI	DIR	NGK	GEV	DDI	DHLGN	RRV	SVG	MAENQ	FRVGL	VRV 471																																																										
WP_003690105.1	499	ERAVK	ERL	NAESEN	LMPH	DLINAK	PVSA	AAIK	EFFG	SSQLS	QFMDQ	TNPLSE	VTHK	RRV	SAL	GP	GG	L	TRERAG	F	VRD	VH 578																																																										
WP_119518745.1	477	ERAVK	ERL	MAESEN	LMPQ	DLINAK	PVSA	AAIK	EFFG	SSQLS	QFMDQ	TNPLSE	ITHK	RRV	SAL	GP	GG	L	TRERAG	F	VRD	VH 556																																																										
WP_353524646.1	472	ERAVK	ERL	SLG	DL	LMPQ	DMINAK	PTSA	AAV	K	EFFG	SSQLS	QFMDQ	TNPLSE	ITHK	RRV	SAL	GP	GG	L	TRERAG	F	VRD	VH 551																																																								
WP_003690105.1	579	P	THYGRV	CP	IET	PE	GN	I	GL	IN	SL	SV	YART	ND	Y	G	F	L	E	T	P	Y	R	R	V	I	D	G	K	V	T	E	E	I	D	Y	L	S	A	I	E	E	G	R	V	I	A	Q	A	N	A	D	L	S	D	G	N	L	I	G	658																			
WP_119518745.1	557	P	THYGRV	CP	IET	PE	GN	I	GL	IN	SL	SV	YART	ND	Y	G	F	L	E	T	P	Y	R	R	V	I	D	G	K	V	T	E	E	I	D	Y	L	S	A	I	E	E	A	D	H	V	I	A	Q	A	S	A	T	L	N	E	K	G	Q	L	V	D	636																	
WP_353524646.1	552	P	THYGRV	CP	IET	PE	GN	I	GL	IN	SL	SV	YART	ND	Y	G	F	L	E	T	P	Y	R	R	V	I	D	G	K	V	T	E	E	I	D	Y	L	S	A	I	E	E	G	N	V	I	A	Q	A	N	S	N	L	D	E	E	G	H	F	V	E	631																		
WP_003690105.1	659	D	L	V	T	C	R	E	K	G	E	T	I	M	A	T	P	D	R	V	Q	Y	M	D	V	A	T	G	Q	V	V	S	A	A	S	L	I	P	F	L	E	H	D	D	A	N	R	A	L	M	G	A	N	M	Q	R	A	V	P	T	L	R	A	D	K	P	L	V	G	T	G	M	E	R	N	V	A	R	D	738
WP_119518745.1	637	E	L	V	A	R	H	L	N	E	F	T	V	K	A	P	E	D	V	T	L	M	D	V	S	P	K	Q	V	V	S	A	A	S	L	I	P	F	L	E	H	D	D	A	N	R	A	L	M	G	S	N	M	Q	R	A	V	P	T	L	R	A	D	K	P	L	V	G	T	G	M	E	R	N	V	A	R	D	716	
WP_353524646.1	632	D	L	V	T	C	R	S	K	G	E	S	S	L	F	S	R	D	Q	V	D	M	D	V	S	T	G	Q	V	V	S	G	A	S	L	I	P	F	L	E	H	D	D	A	N	R	A	L	M	G	A	N	M	Q	R	A	V	P	T	L	R	A	D	K	P	L	V	G	T	G	M	E	R	N	V	A	R	D	711	
WP_003690105.1	739	S	A	T	A	I	V	A	R	R	G	G	V	V	E	Y	D	A	N	R	V	I	R	V	H	D	E	A	T	A	G	E	V	G	D	I	N	L	V	K	F	T	R	S	N	Q	S	T	N	I	N	Q	R	P	A	V	K	A	G	D	V	L	Q	R	G	D	L	V	A	D	G	A	S	T	D	L	818			
WP_119518745.1	717	S	G	V	C	V	A	R	R	G	G	V	I	D	S	V	A	S	R	V	V	R	V	A	D	E	V	E	T	G	A	G	D	I	N	L	T	K	Y	T	R	S	N	Q	T	C	I	N	Q	R	P	L	V	S	K	G	D	V	A	R	G	D	L	A	D	G	P	S	T	D	L	796								
WP_353524646.1	712	S	G	V	T	A	V	A	K	R	G	G	V	Y	D	A	S	R	I	V	I	K	V	N	E	D	E	M	P	Y	P	E	A	G	D	I	N	L	T	K	Y	T	R	S	N	Q	T	C	I	N	Q	M	P	C	V	S	L	G	E	P	V	E	R	G	D	L	A	D	G	P	S	T	D	L	791					
WP_003690105.1	819	G	E	L	A	L	G	Q	N	M	T	I	A	F	M	P	W	N	G	N	Y	N	E	D	S	I	L	S	E	K	V	A	A	D	D	R	Y	T	S	I	H	I	E	L	N	V	A	R	D	T	K	L	G	A	E	I	T	R	D	I	P	N	L	S	E	R	M	Q	N	R	L	D	E	S	G	I	V	898		
WP_119518745.1	797	G	E	L	A	L	G	Q	N	M	R	V	A	F	M	P	W	N	G	N	F	N	E	D	S	I	L	S	E	K	V	A	A	D	D	R	Y	T	S	I	H	I	E	L	N	V	A	R	D	T	K	L	G	P	E	I	T	A	D	I	P	N	V	G	E	A	A	L	N	K	L	D	E	A	G	I	V	876		
WP_353524646.1	792	G	E	L	A	L	G	Q	N	M	R	V	A	F	M	P	W	N	G	N	F	N	E	D	S	I	L	S	E	K	V	A	A	D	D	R	Y	T	S	I	H	I	E	L	N	V	A	R	D	T	K	L	G	P	E	I	T	A	D	I	P	N	V	G	E	A	A	L	S	K	L	D	E	S	G	I	V	871		
WP_003690105.1	899	Y	I	G	A	E	V	E	A	G	D	V	L	V	G	K	T	P	K	G	E	T	Q	L	T	P	E	E	K	L	R	A	I	F	G	E	K	A	S	D	V	K	D	T	S	L	R	M	P	T	G	M	S	G	T	V	I	D	V	Q	F	T	R	E	G	I	Q	R	D	K	A	Q	S	I	D	S	978			
WP_119518745.1	877	Y	V	G	A	E	V	Q	A	G	D	I	L	V	G	K	T	P	K	G	E	T	Q	L	T	P	E	E	K	L	R	A	I	F	G	E	K	A	S	D	V	K	D	T	S	L	R	V	P	T	G	T	G	T	V	I	D	V	Q	F	T	R	D	G	V	E	R	D	S	R	A	L	S	I	E	K	M	956		
WP_353524646.1	872	Y	I	G	A	E	V	T	G	D	I	L	V	G	K	T	P	K	G	E	T	Q	L	T	P	E	E	K	L	R	A	I	F	G	E	K	A	S	D	V	K	D	S	L	R	V	P	N	G	S	G	T	V	I	D	V	Q	F	T	R	D	G	V	E	K	D	K	R	A	L	E	I	E	M	951					
WP_003690105.1	979	E	L	K	R	Y	R	L	D	L	N	Q	L	R	I	F	D	N	A	F	D	R	I	E	R	M	I	V	G[4]	G	G	P	M	K	---	L	A	K	G	S	E	I	T	E	Y	L	A	G	L	P	S	R	H	D	F	D	I	R	L	D	E	D	A	K	L	Q	L	E	L	I	K	1056								
WP_119518745.1	957	Q	L	D	Q	I	R	K	D	L	N	E	E	F	R	I	V	E	G	A	T	F	E	R	L	A	A	L	V		A	K	A	E	G	g	p	a	L	K	K	G	E	I	T	D	Y	L	D	G	L	E	-	R	G	Q	W	F	L	R	M	A	D	D	A	L	N	E	Q	L	E	K	A	Q	1032					
WP_353524646.1	952	Q	L	Q	A	K	D	L	S	E	E	L	Q	I	L	E	A	G	L	F	S	R	I	R	A	V	L	V		G	G	V	E	A	-----	E	K	L	D	K	L	P	-	R	D	R	W	L	E	L	G	L	T	D	E	K	Q	N	L	Q	L	E	Q	L	A	1015														
WP_003690105.1	1057	L	S	L	Q	Q	R	E	A	D	E	L	Y	E	I	K	K	K	L	T	Q	G	D	E	L	P	Q	G	V	K	M	K	V	F	I	A	I	K	R	R	L	Q	A	G	D	K	M	A	G	R	H	G	N	K	G	V	S	R	I	L	P	V	E	D	M	P	M	Y	A	D	G	R	P	V	D	1136				
WP_119518745.1	1033	A	Y	I	S	D	R	R	Q	L	L	D	K	F	E	D	K	R	K	L	Q	G	D	L	A	P	G	V	L	I	K	V	Y	L	A	I	K	R	R	I	Q	P	G	D	K	M	A	G	R	H	G	N	K	G	V	S	V	I	M	P	V	E	D	M	P	H	A	N	G	T	P	V	D	1112						
WP_353524646.1	1016	E	Q	Y	D	E	L	K	H	E	F	E	K	L	E	A	K	R	R	I	T	Q	G	D	L	A	P	G	V	L	I	K	V	Y	L	A	V	K	R	R	I	Q	P	G	D	K	M	A	G	R	H	G	N	K	G	V	I	S	K	I	N	P	E	D	M	P	Y	D	E	N	G	T	P	V	D	1095				
WP_003690105.1	1137	I	V	L	N	P	L	G	V	P	S	R	M	N	I	G	Q	I	L	E	V	L	H	G	A	A	K	G	I	G	E	R	I	D	R	M	L	K	E	R	R	K	A	G	E	L	R	E	F	L	N	K	L	Y	N	---	G	S	G	K	E	D	L	S	L	T	D	E	E	I	I	E	L	A	S	N	1213			
WP_119518745.1	1113	I	V	L	N	P	L	G	V	P	S	R	M	N	V	G	Q	I	L	E	T	H	L	G	A	A																																																						