

# Find a Gene

Serina Huang, [seh004@ucsd.edu](mailto:seh004@ucsd.edu), A12245564

December 3, 2018

## Q1.

Protein: Hepatocyte Nuclear Factor 6 (HNF6)

Species: Homo sapiens

Accession number: NP\_004489

Function known: A transcription factor in the Cut homeobox family. Expression of HNF6 is enriched in the liver, which it stimulates transcription of liver-specific genes and antagonizes glucocorticoid-stimulated gene expression. May influence glucose metabolism, cell cycle regulation, and may be associated with cancer.

## Q2.

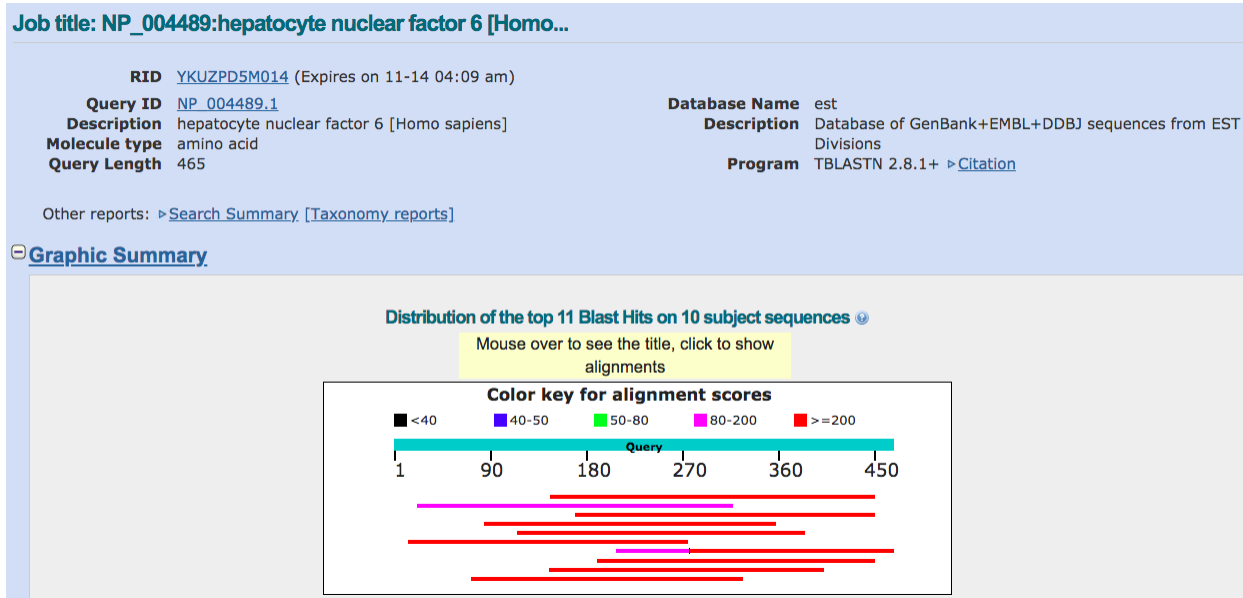
Method: TBLASTN (2.8.1)

Database: Expressed sequence tags (est)

Limit organism: All species

Alignment of choice: BLOSUM62

Chosen match: Assession CA476912.1, a 942 bp clone from *Danio rerio*.



Descriptions

Sequences producing significant alignments:

Select: AllNoneSelected:0

Alignments

Download

GenBank

Graphics

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	<a href="#">FY471421 full-length enriched tammar gravid uterus cDNA library Notamacropus eugenii cDNA clone MEGC-003B08 3', mRNA sequence</a>	407	407	64%	5e-138	71%	<a href="#">FY471421.1</a>
<input type="checkbox"/>	<a href="#">AGENCOURT_10702568 NCI_CGAP_ZEmb3 Danio rerio cDNA clone IMAGE:6800416 5', mRNA sequence</a>	198	198	63%	1e-56	46%	<a href="#">CA476912.1</a>
<input type="checkbox"/>	<a href="#">FY475673 full-length enriched tammar gravid uterus cDNA library Notamacropus eugenii cDNA clone MEGC-008N01 3', mRNA sequence</a>	508	508	60%	3e-178	90%	<a href="#">FY475673.1</a>
<input type="checkbox"/>	<a href="#">602284768F1 NIH_MGC_86 Homo sapiens cDNA clone IMAGE:4372264 5', mRNA sequence</a>	303	303	58%	1e-96	63%	<a href="#">BG111077.1</a>
<input type="checkbox"/>	<a href="#">AUF_lpPit_33_p18 Pituitary cDNA library Ictalurus punctatus cDNA 5' similar to one cut domain, family member 2; oncut 2, mRNA sequence</a>	251	251	57%	3e-77	54%	<a href="#">CK415893.1</a>
<input type="checkbox"/>	<a href="#">AGENCOURT_16879333 NCI_CGAP_ZEmb3 Danio rerio cDNA clone IMAGE:7057865 5', mRNA sequence</a>	226	226	55%	5e-68	54%	<a href="#">CK144141.1</a>
<input type="checkbox"/>	<a href="#">FDR103-P00001-DEPE-R_E20 FDR103 Danio rerio cDNA clone FDR103-P00001-BR_E20 3', mRNA sequence</a>	340	429	55%	1e-112	87%	<a href="#">EH457866.1</a>
<input type="checkbox"/>	<a href="#">FY549871 full-length enriched tammar gravid uterus cDNA library Notamacropus eugenii cDNA clone MEGC-109C14 3', mRNA sequence</a>	369	369	55%	2e-123	75%	<a href="#">FY549871.1</a>
<input type="checkbox"/>	<a href="#">FS645687 full-length enriched swine cDNA library, adult brain (frontal lobe) Sus scrofa cDNA clone BFLT10071H11 5', mRNA sequence</a>	304	304	55%	1e-98	67%	<a href="#">FS645687.1</a>
<input type="checkbox"/>	<a href="#">AGENCOURT_16916377 NCI_CGAP_ZEmb2 Danio rerio cDNA clone IMAGE:7062891 5', mRNA sequence</a>	260	260	54%	3e-81	59%	<a href="#">CK147211.1</a>

Download

GenBank

Graphics

Next

Previous

Descriptions

AGENCOURT\_10702568 NCI\_CGAP\_ZEmb3 Danio rerio cDNA clone IMAGE:6800416 5', mRNA sequence.  
Sequence ID: [CA476912.1](#) Length: 942 Number of Matches: 1

Range 1: 21 to 878

GenBank

Graphics

Next Match

Previous Match

Score	Expect	Method	Identities	Positives	Gaps	Frame
198 bits(503)	1e-56	Compositional matrix adjust.	184/312(59%)	200/312(64%)	44/312(14%)	+3
Query 22	PAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMA	all	dgsggggdyHHHHR	RAPEH-SL		79
Sbjct 21	P+ ADL+ G S H RS HR S AH RSMGMAS+LD G		+HHHR PEH L			170
Query 80	PSAADLMTGDSAHHRS---HRSSLS--AHARSMGMAS		ILDSGD-----YHHHRPPEHPGL			347
Sbjct 171	AGPLHPTMTMACETPPGMSM	ptt	tytltltp	lqplppISTVSDKF	phhhhhhhhhhhphhhQ	139
Query 140	A LHP M+MACE PPGMSM +TYTTLTPLQPLPPISTVSDKF		PHHHHHHHHHHH H Q			347
Sbjct 348	ATHLHPAMSMACEAPPGMSM	SST	YTTLTPLQPLPPISTVSDKF	PHHHHHHHHHHHHPH-Q		527
Query 200	RLAGNVSGSFTLMRDERGLASMNNLYTPYHKDVAG	ggs	slsplsssglsIHNSQQGLPH			199
Sbjct 528	R+ GNVSGSFTLMRD+RGLA MNNLY+PYHKDVA		MGQSLSPLS SGL IHNSQQGLP			527
Query 260	YAHFGAAMPTDKMLTPNGFEAHPAMLGRHGEQHLTPTSAGMVP		INGLpphhphAHLNAQ			259
Sbjct 666	YAHFGA MP +KMLTPNGFEAHPAML RHG		H Q			665
Query 304	YAHFGATMPAEKMLTPNGFEAHPAMLARHG-----GAAHERVFGEHGADQ					842
Query 304	GHGQLLG-TAREPNPSVTGAQVNSGNSGQMEETINT-----KEVAQRITT					303
Sbjct 843	H ++ P P TGA + +G+ I		KE + I T			842
Query 304	ELKRYSIQAIIF					315
Sbjct 843	ELKRYSIQPIF					878

Q3.

Chosen sequence:

>D. rerio protein (sequence taken from BLAST result)  
PSAADLMTGDSAHHRSRSHSSLSAHARSMGMASILDSGDYHHHRPPEHPGLATHLHPAMSMACEAPPGMSM  
SSTYTTLTPLQPLPPISTVSDKFPHHHHHHHHHHHHPHQRIPGNVSGSFTLMRDDRGLAPMNNLYSPYHK  
DVASMGQSLSPLSGSGLSGIHNSQQGLPPYAHFGATMPAEKMLTPNGFEAHPAMLARHGGAHERVFGE  
HGADQRHPPSPARPSQRPGPRTGAGLHSGAEPLLRSRIAAQQRVAVRVRWKRSSIPKEWPKGIPTELKRY  
IPQPIF

Name: Danio transcription factor

Species: *Danio rerio*

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Actinopterygii; Neopterygii; Teleostei;  
Ostariophysi; Cypriniformes; Cyprinidae; Danio.

Q4.

Method: BLASTP 2.8.1

Database: Non-redundant protein sequences (nr)

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	<a href="#">hepatocyte nuclear factor 6 isoform X1 [Danio rerio]</a>	414	414	99%	4e-141	72%	<a href="#">XP_005173695.1</a>
<input type="checkbox"/>	<a href="#">hepatocyte nuclear factor 6 [Danio rerio]</a>	414	414	99%	5e-141	72%	<a href="#">NP_956867.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: hepatocyte nuclear factor 6-like isoform X1 [Sinocyclocheilus rhinoceros]</a>	383	383	100%	1e-128	72%	<a href="#">XP_016416804.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: hepatocyte nuclear factor 6-like isoform X2 [Sinocyclocheilus rhinoceros]</a>	382	382	100%	1e-128	72%	<a href="#">XP_016416805.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: hepatocyte nuclear factor 6-like isoform X1 [Sinocyclocheilus anshuiensis]</a>	380	380	100%	7e-128	72%	<a href="#">XP_016326525.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: hepatocyte nuclear factor 6-like isoform X2 [Sinocyclocheilus anshuiensis]</a>	380	380	100%	9e-128	72%	<a href="#">XP_016326526.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: LOW QUALITY PROTEIN: hepatocyte nuclear factor 6-like [Cyprinus carpio]</a>	374	374	100%	1e-125	74%	<a href="#">XP_018949348.1</a>
<input type="checkbox"/>	<a href="#">onecut-1 isoform a [Danio rerio]</a>	374	374	99%	2e-125	72%	<a href="#">AAW02845.1</a>
<input type="checkbox"/>	<a href="#">onecut-1 isoform b [Danio rerio]</a>	374	374	99%	4e-125	72%	<a href="#">AAW02846.1</a>
<input type="checkbox"/>	<a href="#">hepatocyte nuclear factor 6-like isoform X2 [Carassius auratus]</a>	371	371	100%	3e-124	72%	<a href="#">XP_026143615.1</a>
<input type="checkbox"/>	<a href="#">hepatocyte nuclear factor 6-like isoform X1 [Carassius auratus]</a>	371	371	100%	4e-124	72%	<a href="#">XP_026143613.1</a>
<input type="checkbox"/>	<a href="#">PREDICTED: hepatocyte nuclear factor 6-like [Cyprinus carpio]</a>	370	370	99%	1e-123	76%	<a href="#">XP_018949346.1</a>

Download

GenPept

Graphics

Distance tree of results

Multiple alignment

NextPreviousDescriptions

hepatocyte nuclear factor 6 isoform X1 [Danio rerio]  
Sequence ID: [XP\\_005173695.1](#) Length: 456 Number of Matches: 1  
[See 1 more title\(s\)](#)

Range 1: 23 to 306

GenPept

Graphics

Next Match

Previous Match

Score	Expect	Method	Identities	Positives	Gaps
414 bits(1064)	4e-141	Compositional matrix adjust.	225/311(72%)	232/311(74%)	53/311(17%)
Query 2	SAADLMTGDSAHHRSHRSSL	SAHARSMGMASILD	SGDYHHHRPPEHPGLATHLHPAMSMA	61	
Sbjct 23	SAADLMTGDSAHHRSHRSSL	SAHARSMGMASILD	SGDYHHHRPPEHPGLATHLHPAMSMA	82	
Query 62	CEAPPGMSMSSTYTTLT	PLQLPPISTVSDKFP	HHHHHHHHHHHPHQRI	PGNVSGSFTL	121
Sbjct 83	CEAPPGMSMSSTYTTLT	PLQLPPISTVSDKFP	HHHHHHHHHHHHHPHQRI	PGNVSGSFTL	142
Query 122	MRDDRGLAPMNNLYSPYHKDV	ASMGQSLSP	SGSGLSGIHNSQQGLPPYAHPGATMPAEK	181	
Sbjct 143	MRDDRGLAPMNNLYSPYHKDV	ASMGQSLSP	SGSGLSGIHNSQQGLPPYAHPGATMPAEK	202	
Query 182	MLTPNGFEAHHPAMLARHC	-----	GAHERVFG	EHGADQ	215
Sbjct 203	MLTPNGFEAHHPAMLARHC		H +V G	+Q	261
Query 216	RHPPSPARPSQRPGRPTG	AGLHSGAEPLLR	SRIAAQQRVAVRVWRKRSIPKEWPKGIPTE	275	
Sbjct 262	NH-----	SSVPGSQLNNGSSSQ	MEVNTKEVAQR-----	ITTE	295
Query 276	LKRYISIPQPIF	286			
Sbjct 296	LKRYISIPQAIF	306			

Related Information

[Gene](#) - associated gene details  
[New Genome Data Viewer](#) - aligned genomic context  
[Identical Proteins](#) - Identical proteins to XP\_005173695.1

Q5.

Relabeled sequences:

>Novel\_zebrafish

PSAADLMTGDSAHHRSHRSSL  
SAHARSMGMASILD  
SGDYHHHRPPEHPGLATHLHPAMSMA  
CEAPPGMSMSSTYTTLT  
PLQLPPISTVSDKFP  
HHHHHHHHHHHPHQRI  
PGNVSGSFTLMRDDRGLAPMNNLYSPYHK  
DVASMGQSLSP  
SGSGLSGIHNSQQGLPPYAHPGATMPAEK  
MLTPNGFEAHHPAMLARHG  
GAHERVFG

HGADQRHPPSPARPSQRPGPRTGAGLHSGAEPLLSRIAQQRVAVRVRWKR SIPKEWPKGIPTELKRY S  
IPQPIF

>Original\_human

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGSGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACRKEQE H GKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Sumatran\_orangutan

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGSGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACRKEQE H GKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Gibbon

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGSGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACRKEQE H GKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Cheetah

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGGGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACRKEQE H GKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Chimp

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGSGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGTMTPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACRKEQE H GKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Cat

MNAQLTMEAIGELHGV SHEPVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGGGGGDYHH  
HHRAP EHS LAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPPISTVSDKFP HHHHHHHHHHPHHHQR  
LAGNVSGSF TLMRDERGLASMNNLYTPYHKDVAGMGQSLSP LSSSGLGSIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLP PPHPHAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRY SIPQAIFAQRVLCRSQGTLSDLLRNP KPWSKLKSGRETFRM

WKWLQEPEFQRMSALRLAACKRKEQEHGKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Grizzly\_bear

MNAQLTMEAIGELHGVSHPEVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGGGGGDYHH  
HHRAPESLAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPLPPISTVSDKFPHHHHHHHHHHPHHHQR  
LAGNVSGSFTLMRDERGLASMNNLYTPYHKDVAGMGQSLSPSTSGLGGIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLPPHHPAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRYISIPQAIFAQRVLCRSQGTLSDLLRNPKPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACKRKEQEHGKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Dog

MNAQLTMEAIGELHGVSHPEVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGGGGGDYHH  
HHRAPESLAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPLPPISTVSDKFPHHHHHHHHHHPHHHQR  
LAGNVSGSFTLMRDERGLASMNNLYTPYHKDVAGMGQSLSPSSSGLSGIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLPPHHPAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRYISIPQAIFAQRVLCRSQGTLSDLLRNPKPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACKRKEQEHGKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Killer\_whale

MNAQLTMEAIGELHGVSHPEVPAPADLLGGSPHARSSVAHRGSHLPPAHPRSMGMASLLDGGSGSGDYHH  
HHRAPESLAGPLHPTMTMACETPPGMSMPTTYTTLTPLQLPLPPISTVSDKFPHHHHHHHHHHPHHHQR  
LAGNVSGSFTLMRDERGLASMNNLYTPYHKDVAGMGQSLSPSGSGLSGIHNSQQGLPHYAHPGAAMPTD  
KMLTPNGFEAHHPAMLGRHGEQHLTPTSAGMVPINGLPPHHPAHLNAQGHGQLLGTAREPNPSVTGAQV  
SNGSNSGQMEEINTKEVAQRITTELKRYISIPQAIFAQRVLCRSQGTLSDLLRNPKPWSKLKSGRETFRM  
WKWLQEPEFQRMSALRLAACKRKEQEHGKDRGNTPKKPRLVFTDVQRRTLHAIFKENKRPSKELQITISQ  
QLGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTKA

>Tropical\_clawed\_frog

MNAQLTMDAIGDLHGISHESVPGTADLMGSSPHHRGVSVTHRSNHLAHPRSMGMASILDGGDYHHHHHHH  
HRPPDHALTGPLHPTMTMACDTPPGMSMSSTYTTLTPLQLPLPPISTVSDKFPHHHHHHHHPHQRIPGNVSG  
SFTLMRDDRGLASMNNLYSPYHKVETGMGQSLSPSGSGLSGIHGAQQGP PHYAHP SAAMPTEKMLTPNG  
FEAHHPAMLTRHGEQHLTPPSAGMVPINGIPHHPAHLNAQSHGQILASTRDQNPPSVTGSQINNGSNSG  
QMEEINTKEVAQRITTELKRYISIPQAIFAQRVLCRSQGTLSDLLRNPKPWSKLKSGRETFRMMKWLP  
EFQRMSALRLAALVPADPVFHSGQLPADSLVKIGYPSQSTQSNHMSCKRKEQEHGKDRGNTPKKPRLV  
TDVQRRTLHAIFKENKRPSKELQITISQQLGLELSTVSNFFMNARRRSLDKWQDEGNSGSGNTSSSSSTC  
TKA

>Sea\_squirt

MPVSLGQSVTTPSAKSASILNQHPGIATDFITMATAASGELNGFHHLHHHHHPSEQYYRHEHYHHHFH  
HPNFDGYPNYNDRDTP IAGDMQKNNLHNFASKSMSLEGEKLDENCNKSPNYLPPIGDALLRRDNRS DASK  
NNAKEEDES GCSKFVMQETDNLTELQKSSAVSEHEKKEEVQLKTNDAPEDFSVKTEQSELYQFHARNFS  
IFT PSSQRGTPDEGMNLIPVETTDHTSIDSYFRSDATNANPNSNPIDSVPSVDGPSYATLTPLQLPSI  
SSVSDKYMPTNETSYATLTNQELTDCSSYSKMGMGHSPLPLSNRMILNGLAAQTRGGMQSQA AIDAVNQ  
AAAAAVGLSHYNKPVLSNIIPPPPVSNPYDPHVFGRI DQCNDMGAGFPGGHMFPHRSTGFVSQYGLQD  
LSSSLQVSAPSERRRPTHE DIPADNGKRHSGSDRLGGSLQPHSSNSASSSRTQQIEEVNTKEVASKITQ  
ELKRYISIPQAIFAQRVLCRSQGTLSDLLRNPKPWSKLKSGRETFRMMKWLPQEPEFQRMSALRLAACKRK  
EDEKSYENS VN SPKKPRLVFTDLQRRTLHAIFKESKRPSKEMQIQISQQLGLEVTTVSNFFMNARRRSLD  
KWQDDESGYNSKENSRSNPNSSDHL SASPNHQQQQQQQQQQQQQQYQQHTQDSRLSYAPGESLLS  
PLCGSPSGHLHFP PP HHLHHHNLHQQQNTMLSASHLTSSGLVHPYQSQHQLLGSDVTGLVNPR

Alignment using EBI's MUSCLE:

CLUSTAL multiple sequence alignment by MUSCLE (3.8)

```

Sea_squirt      -----MPVSLGQSVTTSPSAKSASILNQQH-PGIATDFITMATAA
Novel_zebrafish -----PSAADLMTGDSAHHRS---HRSSL--SAHARSMGMASIL
Tropical_clawed_frog MNAQLTMDAIGDLHGISHESVPGTADLM-GSSPHHRGSVTHRSNHL-SAHPRSMGMASIL
Killer_whale    MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Chimp           MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Dog             MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Grizzly_bear    MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Cat            MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Original_human  MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Sumatran_orangutan MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Gibbon         MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL
Cheetah        MNAQLTMEAIGELHGVSEHPVPAPADLL-GGSPHARSSVAHRGSHLPPAHPRSMGMASLL

```

\* . . \* . . . . : \*\*:

```

Sea_squirt      SGELNGFHHLHH---HHHPSEQYRHEHYHHHFHHPNFDGYPNYNDRDTPDIAGDMQKNN
Novel_zebrafish DSG-----DYH-----HHRPPE-----HPGLATH-----
Tropical_clawed_frog DGG-----DYHHHHHHHHHRPPD-----H-ALTGP-----
Killer_whale    DGG-SGSGDYHH---HHRAPE-----H-SLAGP-----
Chimp           DGG-SGGGDYHH---HHRAPE-----H-SLAGP-----
Dog             DGG-GGGGDYHH---HHRAPE-----H-SLAGP-----
Grizzly_bear    DGG-GGGGDYHH---HHRAPE-----H-SLAGP-----
Cat            DGG-GGGGDYHH---HHRAPE-----H-SLAGP-----
Original_human  DGG-SGGGDYHH---HHRAPE-----H-SLAGP-----
Sumatran_orangutan DGG-SGGGDYHH---HHRAPE-----H-SLAGP-----
Gibbon         DGG-SGGGDYHH---HHRAPE-----H-SLAGP-----
Cheetah        DGG-GGGGDYHH---HHRAPE-----H-SLAGP-----

```

.. \* \*\*...: \* :

```

Sea_squirt      LHNFAKSM SLEGEKLDENCNKSPNYLPPIGDALLRRDNRSDASKNNAKEEDES GC SKFV
Novel_zebrafish LH---PAMSMA-----CEAPPG-----MS
Tropical_clawed_frog LH---PTMTMA-----CDTPPG-----MS
Killer_whale    LH---PTMTMA-----CETPPG-----MS
Chimp           LH---PTMTMA-----CETPPG-----MS
Dog             LH---PTMTMA-----CETPPG-----MS
Grizzly_bear    LH---PTMTMA-----CETPPG-----MS
Cat            LH---PTMTMA-----CETPPG-----MS
Original_human  LH---PTMTMA-----CETPPG-----MS
Sumatran_orangutan LH---PTMTMA-----CETPPG-----MS
Gibbon         LH---PTMTMA-----CETPPG-----MS
Cheetah        LH---PTMTMA-----CETPPG-----MS

```

\*\* :\*: : \*: .\*: :

```

Sea_squirt      MQETDNSLTELQKSSAVSEHEKKEEVQLKTNDAPEDFSVKTEQSELYQFHARNFSIFTPS
Novel_zebrafish MSSTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHH-
Tropical_clawed_frog MSSTYTTTLTPLQLPPISTVSDK-----FPHHHHHHH-
Killer_whale    MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP
Chimp           MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP
Dog             MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP
Grizzly_bear    MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP
Cat            MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP
Original_human  MPTYTTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHHHP

```



Sumatran_orangutan	MPTTYTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHP
Gibbon	MPTTYTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHP
Cheetah	MPTTYTTLTPLQLPPISTVSDK-----FPHHHHHHHHHHP
	* * .: ** ** .: * ..* : . * ..
Sea_squirt	SQRGTPDEGMNLIPVETTDHTSIDSYFRSDATNANPNSNPIDSVSSVDGPSYATLTPLQ
Novel_zebrafish	-----HPHQRIPGNVSG----SFTLMR
Tropical_clawed_frog	-----HPHQRIPGNVSG----SFTLMR
Killer_whale	-----HHHQRLAGNVSG----SFTLMR
Chimp	-----HHHQRLAGNVSG----SFTLMR
Dog	-----HHHQRLAGNVSG----SFTLMR
Grizzly_bear	-----HHHQRLAGNVSG----SFTLMR
Cat	-----HHHQRLAGNVSG----SFTLMR
Original_human	-----HHHQRLAGNVSG----SFTLMR
Sumatran_orangutan	-----HHHQRLAGNVSG----SFTLMR
Gibbon	-----HHHQRLAGNVSG----SFTLMR
Cheetah	-----HHHQRLAGNVSG----SFTLMR
	: : :...*. * : * ..
Sea_squirt	PLPSISSVSDKYMPTNETSYATLTNQELTDCSSYSKMGGMGHSLPPLSNRMILNLGAAQT
Novel_zebrafish	DDRGLAPMNNLYSP-----YHKDVASMGQSLSPLSG----SGL----
Tropical_clawed_frog	DDRGLASMNNLYSP-----YHKEVTGMGQSLSPLSG----SGL----
Killer_whale	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLSG----SGL----
Chimp	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Dog	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Grizzly_bear	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Cat	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Original_human	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Sumatran_orangutan	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Gibbon	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
Cheetah	DERGLASMNNLYTP-----YHKDVAGMGQSLSPLS--SGL----
	.: : : : * * : : : . *: ** . ** . **
Sea_squirt	RGGMQSQAIDA V NQAAAAVGLSHYNKP--VLSSNIIPPPPPVSNPYDPHVFGRIDQCN
Novel_zebrafish	-SGIHN-----SQQGLPPYAHPGATMPAEKMLTPNGF-EAHPA MLAR-----
Tropical_clawed_frog	-GSIHG-----AQQGPPHYAHP SAAMPTEKMLTPNGF-EAHPA MLTR-----
Killer_whale	-GGIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Chimp	-GSIHN-----SQQGLPHYAHPGT TMTPTDKMLTPNGF-EAHPA MLGR-----
Dog	-SGIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Grizzly_bear	-GGIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Cat	-GGIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Original_human	-GSIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Sumatran_orangutan	-GSIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Gibbon	-GSIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
Cheetah	-GSIHN-----SQQGLPHYAHPGAAMP TDKMLTPNGF-EAHPA MLGR-----
	.: : : : : : * . * : * .: : : : . * . : : : * : : *
Sea_squirt	DMGAGFPGGHMFPHRSTGFVSQYGLQDLSSSLQVSAPSERRRPT HEDIPADNGKRHSGSD
Novel_zebrafish	-----HGG-----AAHERVFG EHG A-----DQRHPPS-----
Tropical_clawed_frog	-----HGEQHLTPSAGMVP INGI-----PPHHPAHLNA-----QSHG
Killer_whale	-----HGEQHLTPSAGMVP INGL-----PPHHPAHLNA-----QGHG
Chimp	-----HGEQHLTPSAGMVP INGL-----PPHHPAHLNA-----QGHG
Dog	-----HGEQHLTPSAGMVP INGL-----PPHHPAHLNA-----QGHG
Grizzly_bear	-----HGEQHLTPSAGMVP INGL-----PPHHPAHLNA-----QGHG

Cat	-----HGEQHLTPTSAGMVPINGL-----PPHHPHAHLNA-----QGHG
Original_human	-----HGEQHLTPTSAGMVPINGL-----PPHHPHAHLNA-----QGHG
Sumatran_orangutan	-----HGEQHLTPTSAGMVPINGL-----PPHHPHAHLNA-----QGHG
Gibbon	-----HGEQHLTPTSAGMVPINGL-----PPHHPHAHLNA-----QGHG
Cheetah	-----HGEQHLTPTSAGMVPINGL-----PPHHPHAHLNA-----QGHG

\*                    ..                    \*

Sea_squirt	RLGGSGLQPH-----SSNSASSSRTQQIEEVNTKEVASKITQELKRYISIPQAIFAQRVLC
Novel_zebrafish	----PARPSQRP---GPRTGAGLHSGAEPLLSRIAQ-----QRVAV
Tropical_clawed_frog	QILASTRDQNPPSVTGSQINNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Killer_whale	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Chimp	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Dog	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Grizzly_bear	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Cat	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Original_human	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Sumatran_orangutan	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Gibbon	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC
Cheetah	QLLGTAREPN-PSVTGAQVSNGSNSGQMEINTKEVAQRITTELKRYISIPQAIFAQRVLC

.                    :                    ...                    . . :                    : . . .                    . \* .                    \*\*\*

Sea_squirt	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSLRLAA-----
Novel_zebrafish	RVR-----WKR-----SIPKEW-----
Tropical_clawed_frog	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAALVPADPVFQHS
Killer_whale	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Chimp	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Dog	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Grizzly_bear	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Cat	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Original_human	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Sumatran_orangutan	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Gibbon	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----
Cheetah	RSQGTLSDLLRNPKPWSKLKSGRETFRMWKWLQEPEFQRMSSALRLAA-----

\* .                    \* . .                    : : . \*

Sea_squirt	-----CKRKEDEKSYENS VN SPKKPRLVFTDLQRRTL-HAIF
Novel_zebrafish	-----PKGIPTTELKRYISIPQIF
Tropical_clawed_frog	QLPADSLVKIGYPSQSTQSNHMSCKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Killer_whale	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Chimp	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Dog	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Grizzly_bear	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Cat	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Original_human	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Sumatran_orangutan	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Gibbon	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF
Cheetah	-----CKRKEQEHGKDRG-NTPKKPRLVFTDVQRRTL-HAIF

\* . : \* : : \* : : : \* : : \* : \*

Sea_squirt	KESKRPSKEMQIQISQQLGLEVTTVSNFFMNARRRSLDKWQDESGYNSKENSRSNNPSS
Novel_zebrafish	-----
Tropical_clawed_frog	KENKRPSKELQITISQQLGLELSTVSNFFMNARRRSLDKWQDEGSSSGSGNTSSSSSTCTK
Killer_whale	KENKRPSKELQITISQQLGLELSTVSNFFMNARRRSLDKWQDEGSSSGSGNTSSSSSTCTK
Chimp	KENKRPSKELQITISQQLGLELSTVSNFFMNARRRSLDKWQDEGSSSGSGNTSSSSSTCTK



Dog	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Grizzly_bear	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Cat	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Original_human	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Sumatran_orangutan	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Gibbon	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK
Cheetah	KENKRPSKELQITISQQIGLELSTVSNFFMNARRRSLDKWQDEGSSNSGNSSSSSSTCTK

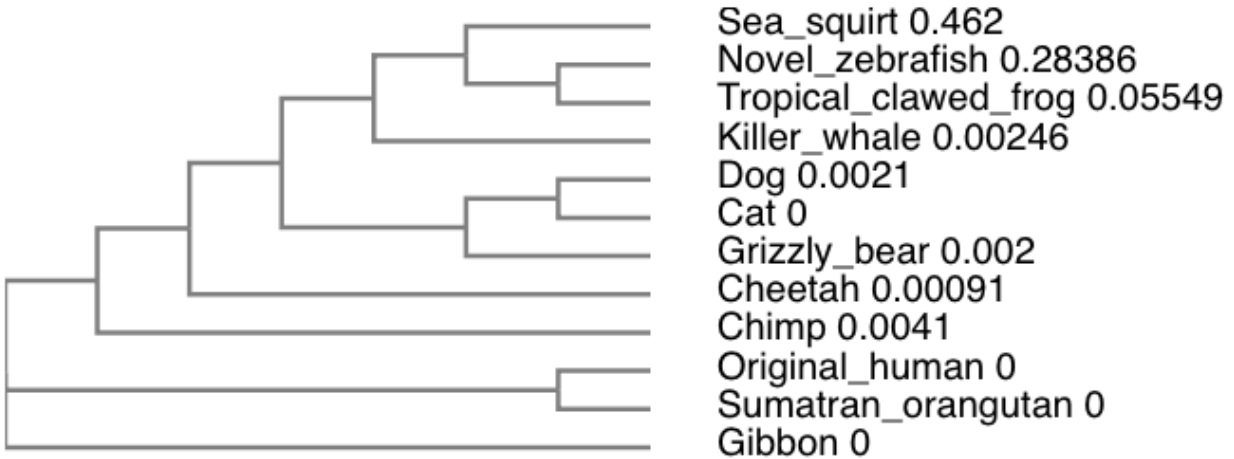
Sea_squirt	DHHLSPNHHQQQQQQQQQQQQQQQAYQQHTQDSRLSYAPGESLLSPLCGSPSGHLHF
Novel_zebrafish	-----
Tropical_clawed_frog	A-----
Killer_whale	A-----
Chimp	A-----
Dog	A-----
Grizzly_bear	A-----
Cat	A-----
Original_human	A-----
Sumatran_orangutan	A-----
Gibbon	A-----
Cheetah	A-----

Sea_squirt	PPPHLHHHNLHQQQNTMLSASHLTSSGLVHPYQSQHQLLGSDVTGLVNPR
Novel_zebrafish	-----
Tropical_clawed_frog	-----
Killer_whale	-----
Chimp	-----
Dog	-----
Grizzly_bear	-----
Cat	-----
Original_human	-----
Sumatran_orangutan	-----
Gibbon	-----
Cheetah	-----

Note: The gap overhang in this alignment file was necessary because the sea\_squirt's sequence has multiple insertions compared to other species' sequences.

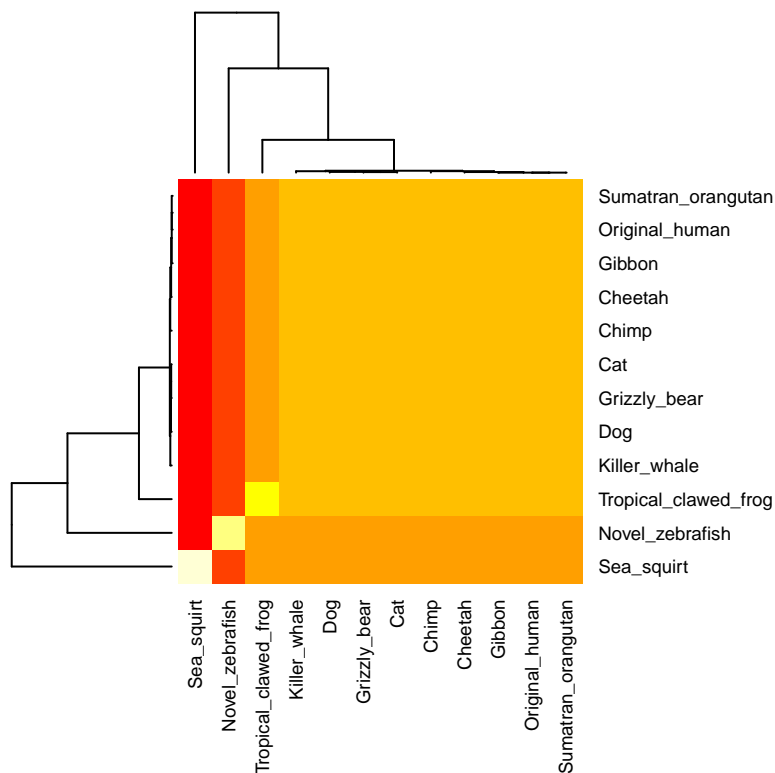
## Q6.

Phylogenetic tree using EBI's "Simple Phylogeny" feature.



Q7.

Read MUSCLE alignment file into Seaview and saved file as FASTA.



Q8.

Because my consensus sequence contains a lot of gaps, I will choose the sequence with the highest identity in the alignment. To do so, I will first calculate the sum of each row in the identity matrix. Then, I will find the first sequence with the maximum sum.

```
sums <- rowSums(identity)
which.max(sums)
```

```
## Original_human
##          9
```

Thus, in the main protein structure database, I will use the Original\_human sequence to search for the most similar atomic resolution structure to my aligned sequences.

```
human <- read.fasta("human.fasta")
hits <- blast.pdb(human, database = "pdb")
```

```
## Searching ... please wait (updates every 5 seconds) RID = OCP700WD01R
## ...
## Reporting 67 hits
```

```
# head(hits, 3)
```

We see that the top 3 unique hits are **2D5V\_A**, **1S7E\_A**, and **1WH6\_A**. I will save these to a new dataframe and add my own annotations: structure ID, method used to solve the structure, resolution, and source organism.

```
top <- hits$hit.tbl[1:3,]
anno <- pdb.annotate(top$pdb.id)
```

```
## Warning in pdb.annotate(top$pdb.id): ids should be standard 4 character
## PDB-IDs: trying first 4 characters...
```

```
top_anno <- merge(top, anno, by.x = "pdb.id", by.y = "row.names")
```

```
# Only take relevant columns
```

```
relevant_colnames <- c("pdb.id", "experimentalTechnique", "resolution", "source", "evaluate", "identity")
x <- match(relevant_colnames, colnames(top_anno))
top_anno_relevant <- top_anno[,x]
```

```
# Last thing, split PDB identifier on the underscore
```

```
ids <- top_anno_relevant$pdb.id
```

```
# Split PDB identifier on the underscore
```

```
ids_split <- strsplit(ids, "_")
```

```
# [[ is synonymous to $. They both select an element from a list
```

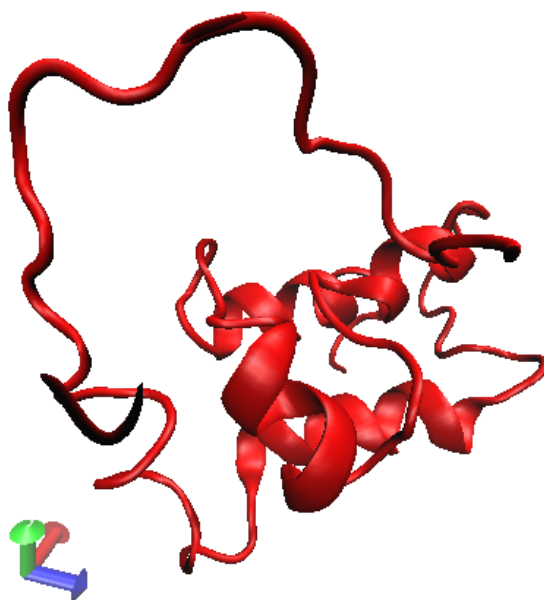
```
top_anno_relevant$pdb.id <- lapply(ids_split, "[", 1)
```

```
top_anno_relevant
```

```
##   pdb.id experimentalTechnique resolution      source  evaluate
## 1  1S7E      SOLUTION NMR      <NA> Mus musculus 5.70e-106
## 2  1WH6      SOLUTION NMR      <NA> Homo sapiens 3.02e-13
## 3  2D5V      X-RAY DIFFRACTION    2.0 Rattus norvegicus 1.78e-112
##   identity
## 1  100.000
## 2   46.341
## 3  100.000
```

## Q9.

Here, I will generate a molecular figure of the first identified PDB structure, 1S7E, using VMD. After downloading the pdb file from Protein Data Bank, load the molecule into VMD.



The sequence similarity is 100% for 1S7E. Thus, this structure from *Mus musculus* is very likely to be similar to my “novel” *Danio rerio* protein.

## Q10.

Searched ChEMBL with my novel sequence. Found 1 Binding Assay and 0 Functional Assay for ChEMBL2176818 (*Mus musculus*, which was the closest because no *Danio rerio* data was available) but 0 Ligand Efficiency Data.

Binding Assay: <https://www.ebi.ac.uk/chembl/assay/inspect/ChEMBL2186578>

Pioglitazone and rosiglitazone are two diabetes drugs that share a common functional core: glitazone. Two variants of the glitazone scaffold, pioglitazone and rosiglitazone, are tested to identify off-target binding events in the rat heart. The purpose of these tests is to explain recently reported cardiovascular risk associated with these drugs.

Results suggest that glitazone has affinity for dehydrogenases. Both drugs bind ion channels and modulators, with implications in congestive heart failure, arrhythmia, and peripheral edema. Additional proteins involved in glucose homeostasis, synaptic transduction, and mitochondrial energy production were detected and potentially contribute to drug efficacy and cardiotoxicity.

Hoffmann BR, El-Mansy MF, Sem DS, Greene AS. Chemical proteomics-based analysis of off-target binding profiles for rosiglitazone and pioglitazone: clues for assessing potential for cardiotoxicity. *J. Med. Chem.* (2012)55:8260-8271. doi: 10.1021/jm301204r.