



T.C. Fırat Üniversitesi
Bilgisayar Mühendisliği Bölümü

BMÜ450 - Biyoinformatik Dersi
BLAST İle Dizi Hizalama Ödevi Raporu

Ders Sorumlusu: Prof. Dr. Mehmet KAYA

Öğrenci İsim: Mert İNCİDELEN
Öğrenci No.: 170260101

ÖDEV TANIMI

Ödevde, BLAST Algoritmasının araştırılması ve çalışma prensibinin anlaşılması; BLAST kullanan bir araç üzerinde örnek dizilerin hizalanıp, nasıl yapıldığının gösterilmesi beklenmektedir.

1. BLAST ALGORİTMASI

BLAST (Basic Local Alignment Search Tool), biyoinformatik çalışmalarında yaygın kullanılan veri tabanı arama ve hizalama aracıdır. Dizi veri tabanlarını aramak için birçok algoritma mevcuttur, ancak BLAST Algoritmaları hızı itibarıyla diğerlerine kıyasla en popüler olanlarıdır. BLAST Algoritmasının temel stratejisi amino asitlerden oluşan sorgu (BLASTP) üzerinden örneklenmiştir.

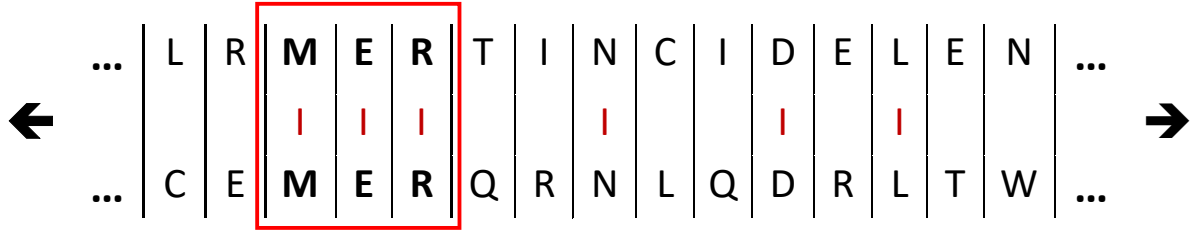
“... **M E R T I N C I D E L E N** ...” sorgusu, boyutu 3 olan (BLASTP için varsayılan boyut 3 amino asittir) kelimelere ayrılır. Böylece MER, ERT, RTI, TIN, INC, NCI, CID, IDE, DEL, ELE, LEN kelimeleri elde edilir. Bir sorgu için elde edilen kelime $n-2$ adet olmaktadır. Proteinler 20 farklı amino asit içerdiğinden, bir kelimenin her 20^3 tripeptitte bir kez (8000 tripeptit, 24000 amino asit içinde) rastgele oluşması beklenir, bu da herhangi bir proteinden daha uzundur ve dolayısıyla idealdir.

Sorgu: ... M E R T I N C I D E L E N ...
Kelimeler: ... M E R
 E R T
 R T I
 T I N
 I N C
 N C I
 C I D
 I D E
 D E L
 E L E
 L E N ...

Elde edilen boyutu 3 olan bu kelimelerin komşu kelimeleri bulunur ve bir eşik değeri belirlenir. Bir skor tablosu kullanarak eşik değerinin üzerinde skora sahip komşu kelimeler de sorgulamaya dahil edilir. Sonraki aşamada “MER” kelimesi için komşu kelimeler, eşik değeri 11 belirlenerek BLOSUM 62 tablosuna göre bulunmuştur. Eşik değeri düşürülerek arama kapsamı genişletilebilir.

KELİME	SKOR (BLOSUM62)
MER	15
MEK	12
LER	12
MDR	12
MQR	12
MEQ	11
IER	11
VER	11
MKR	11
MEN	10
QER	10
...	...
Eşik değeri = 11	

Daha sonra hedef dizilerde aranan kelimeler bulunarak, isabetin olduğu bölgede her iki yönde eşleştirmeler genişletilir.



Genişletilen bölgenin HSP (high-scoring segment pair) puanı, skor tablosu üzerinden hesaplanır ve en iyi eşleşmeler bu şekilde saptanabilir. Yukarıdaki bölge için HSP puanı BLOSUM62 matrisine göre hesaplanmıştır.

EŞ.	SKOR	EŞ.	SKOR	EŞ.	SKOR	EŞ.	SKOR	EŞ.	SKOR
L-C	-1	E-E	5	I-R	-3	I-Q	-3	L-L	4
R-E	0	R-R	5	N-N	6	D-D	6	E-T	-1
M-M	5	T-Q	-1	C-L	-1	E-R	0	N-W	-4

Pozitif skorlu eşleşmelerin değerlerini toplarsak ($5+5+5+6+6+4 = 31$), HSP puanı 31 olarak bulunur. Negatif skorlu eşleşmeler de dahil edilerek toplam puan 17 olarak bulunabilir.

2. BLAST ile Dizi Hizalama Uygulamaları

2.1. BLASTn ile Uygulama

Dizi hizalama işlemi için blast.ncbi.nlm.nih.gov/Blast.cgi adresinde bulunan araçtan faydalanılmıştır. Eşleştirme için ikinci örnek uygulamada dizi olarak **ACTA1** kimliğine sahip (ID No. 58) 2852 nükleotit uzunluğunda bir dizi kullanılmıştır. BLASTN (nükleotit-nükleotit karşılaştırma ve eşleştirme) ile ACTA1 dizisi standart veri tabanı üzerinde benzerlik derecesi yüksek olacak şekilde sorgulanmış ve aşağıdaki veriler elde edilmiştir.

Sorgulama sonrasında karşılaştırılan dizilerden benzerlik oranı yüksek olandan düşüğe doğru sonuç ekranında aşağıda görüldüğü şekilde sıralanmıştır.

Job Title

Biyoinformatik Dersi Hizalama Odevi - 170260101

RID

[A3KFVCY1016](#) Search expires on 04-25 01:18 am [Download All](#) ▼

Program

BLASTN [Citation](#) ▼

Database

nt [See details](#) ▼

Query ID

lcl|Query_12683

Description

None

Molecule type

dna

Query Length

2852

Other reports

[Distance tree of results](#) [MSA viewer](#) [?](#)

Filter Results

Organism

only top 20 will appear ☐ exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

Percent Identity

to

E value

to

Query Coverage

to

Filter

Reset

Descriptions

Graphic Summary

Alignments

Taxonomy

Sequences producing significant alignments

Download ▼

Manage Columns ▼

Show 100 ▼

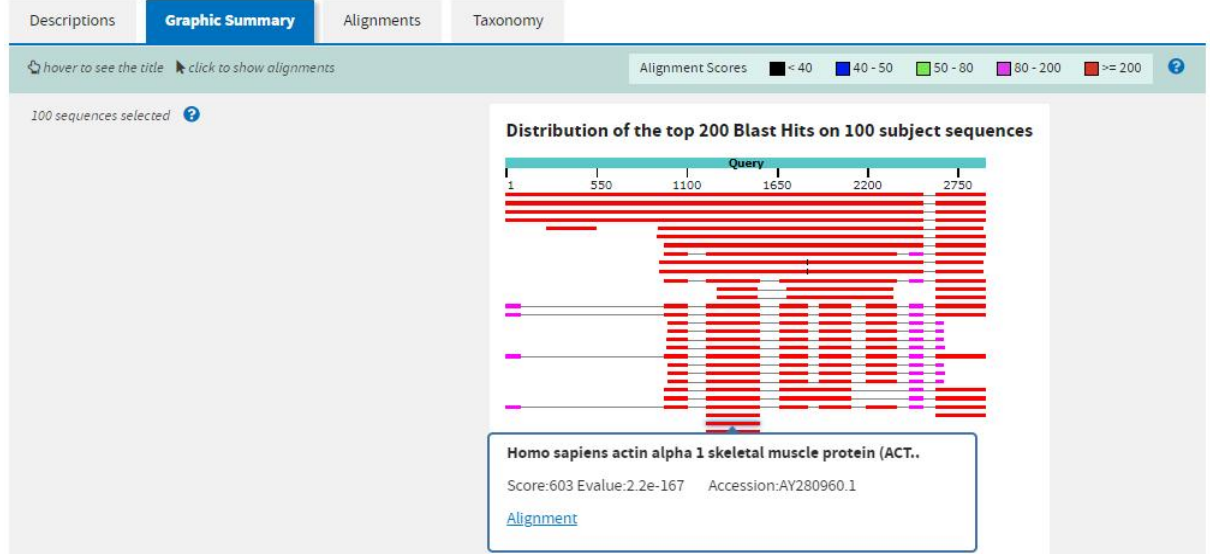
[?](#)

☐ select all 0 sequences selected

[GenBank](#) [Graphics](#) [Distance tree of results](#)

	Description	Max Score	Total Score	Query Cover	E value	Per Ident	Accession
<input type="checkbox"/>	Homo sapiens actin alpha 1 skeletal muscle (ACTA1), RefSeqGene (LRG_429) on chromosome 1	5267	5267	100%	0.0	100.00%	NG_008672.1
<input type="checkbox"/>	Human DNA sequence from clone RP5-1068B5 on chromosome 1q42.11-43, complete sequence	5267	5267	100%	0.0	100.00%	AL160004.18
<input type="checkbox"/>	Homo sapiens skeletal muscle alpha-actin gene (ACTA1), complete cds	5188	5188	100%	0.0	99.54%	AF182035.1
<input type="checkbox"/>	Human skeletal alpha-actin gene, complete cds	5116	5116	100%	0.0	99.12%	M20543.1
<input type="checkbox"/>	Lutra lutra genome assembly, chromosome_14	2113	2540	66%	0.0	86.83%	LR738416.1
<input type="checkbox"/>	Sus scrofa skeletal alpha actin gene, complete cds	1905	1905	67%	0.0	84.81%	U16368.1
<input type="checkbox"/>	Bos taurus alpha skeletal actin precursor gene, complete cds	1851	1851	66%	0.0	84.63%	U02285.1
<input type="checkbox"/>	PREDICTED: Ovis aries actin alpha 1, skeletal muscle (ACTA1), transcript variant X1, mRNA	1310	1903	59%	0.0	87.29%	XM_027962294.1
<input type="checkbox"/>	Canis lupus familiaris breed Labrador retriever chromosome 04b	1173	2199	66%	0.0	86.91%	CP050625.1
<input type="checkbox"/>	Canis lupus familiaris breed Labrador retriever chromosome 04a	1151	2177	66%	0.0	86.52%	CP050572.1

Bu diziler bir grafik üzerinde hizalanma konumları ile kabaca gösterilmiştir. Kırmızı renk ile gösterilenler eşleşme skoru 200 ve üzerinde olan hizalamalardır.



Eşleşen dizilerden bir tanesi için eşleşen kısımların bir bölümünün sonuç dökümü ise aşağıdaki gibi gösterilmiştir.

Lutra lutra genome assembly, chromosome: 14

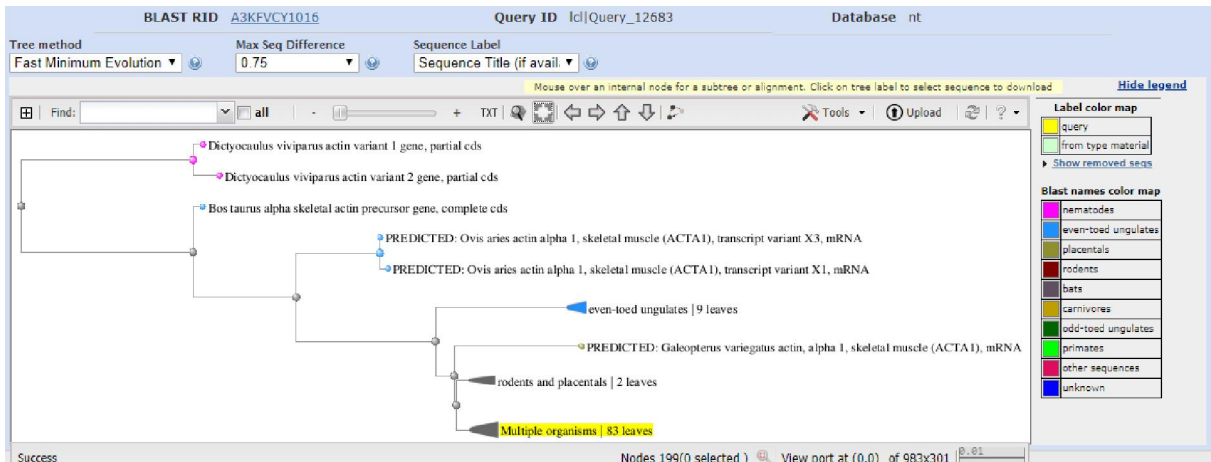
Sequence ID: [LR738416.1](#) Length: 89080780 Number of Matches: 4

Range 1: 1388448 to 1390078 [GenBank](#) [Graphics](#)

[Next Match](#) [Previous Match](#)

Score	Expect	Identities	Gaps	Strand
1879 bits(1017)	0.0	1452/1655(88%)	57/1655(3%)	Plus/Minus
Query 924	CTGCGCTGATGCACCGCGCCTCTTCGCGGTCTCCCTGTCCTT-GCAGAACTAGACACAA	982		
Sbjct 1390078	CTGTGCTGACGCGCCGCGTGTCTCCACGGCCTTCCTATCCTTTGCAGAAACCAGACACCA	1390019		
Query 983	TGTGCGACGAAGACGAGACCACCGCCCTCGTGTGCGACAATGGCTCCGGCCTGGTGAAAG	1042		
Sbjct 1390018	TGTGTGACGAAGACGAGACCACCGCCCTCGTGTGCGACAATGGCTCCGGCCTGGTGAAAG	1389959		
Query 1043	CCGGCTTCGCCGGGGATGACGCCCTAGGGCGTGTTCCTGTCATCGTGGGCGGCCCTC	1102		
Sbjct 1389958	CCGGCTTCGCCGGCGACGACGCCCTAGGGCGTGTTCCTTCCATCGTGGGCGGCCCTC	1389899		
Query 1103	GACACCAGGTACGGCTGCCCTCCGAGAGGGAGCCGGCTCGGGGTCC-CCGCGTAAGCC	1161		
Sbjct 1389898	GCCACCAGGTACGGCTGCCCTCGCGGAGGGAGCCGGGCGGGGACCTCCGTG-GAGCC	1389840		
Query 1162	-AGCCTGGTGCCACCCGAGCGCGTTAACGGGTGCGTGGTGTCTCGGCTCTGCAGGGCG	1220		
Sbjct 1389839	GGGCCCTGTGCACTCTGAGCCGCGTAACGTGGGCGTGGTGTCTCCGCTCCGAGGGCG	1389780		
Query 1221	TCATGGTTCGGTATGGGTGAGAAAGATTCTACGTGGGCGACGAGGCTCAGAGCAAGAGAG	1280		
Sbjct 1389779	TCATGGTGGGTATGGGTGAGAAAGATTCTATGTGGGCGACGAGGCTCAGAGCAAGAGAG	1389720		
Query 1281	GTATCCTGACCCTGAAGTACCCTATCGAGCAGGCATCATCACCAACTGGGATGACATGG	1340		
Sbjct 1389719	GCATCCTGACCCTGAAGTACCCATCGAGCAGGCATCATCACCAACTGGGACGACATGG	1389660		
Query 1341	AGAAGATCTGGCACCACACCTTCTACAACGAGCTTCGCGTGGCTCCCGAGGAGCACCCCA	1400		
Sbjct 1389659	AGAAGATCTGGCACCACACCTTCTACAACGAGCTCCGCGTGGCCCTGAGGAGCACCCCA	1389600		
Query 1401	CCCTGCTCACCAGGCCCCCTCAATCCCAAGGCCAACCAGGAGAGATGACCCAGATCA	1460		
Sbjct 1389599	CCCTGCTCAGGAGGCCCCCTCAACCCCAAGCCAACCAGGAGAGATGACCCAGATCA	1389540		
Query 1461	TGTTTGAGACCTTCAACGTGCCCCGCATGTACGTGGCCATCCAGGCCGTGCTGTCCCTCT	1520		
Sbjct 1389539	TGTTTGAGACCTTCAACGTGCCCCGCATGTACGTGGCCATCCAGGCCGTGCTGTCCCTCT	1389480		

Benzerlikleri esas alınarak sorgu sonucunda elde edilen verilere göre dizilerin filogenetik haritası ise aşağıdaki gibidir.



2.2. BLASTp ile Uygulama

Eşleştirme için ikinci örnek uygulamada dizi olarak **Thrombospondin Type-1** kimliğine sahip 1657 aminoasit içeren bir dizi kullanılmıştır. BLASTp (aminoasit-aminoasit karşılaştırma ve eşleştirme) ile Thrombospondin Type-1 dizisi veri tabanı üzerinde sorgulanmış ve aşağıdaki veriler elde edilmiştir.

Sorgulama sonrasında karşılaştırılan dizilerden benzerlik oranı yüksek olandan düşüğe doğru sonuç ekranında aşağıda görüldüğü şekilde sıralanmıştır.

Job Title

Biyoinformatik Dersi Hizalama Odevi - 170260101

RID

A3R56SRW016 Search expires on 04-25 02:21 am Download All

Program

BLASTP Citation

Database

nr See details

Query ID

lcl|Query_54413

Description

None

Molecule type

amino acid

Query Length

1657

Other reports

Distance tree of results Multiple alignment MSA viewer

Filter Results

Organism

only top 20 will appear

☐ exclude

Type common name, binomial, taxid or group name

+ Add organism

Percent Identity

to

E value

to

Query Coverage

to

Filter

Reset

Descriptions

Graphic Summary

Alignments

Taxonomy

Sequences producing significant alignments

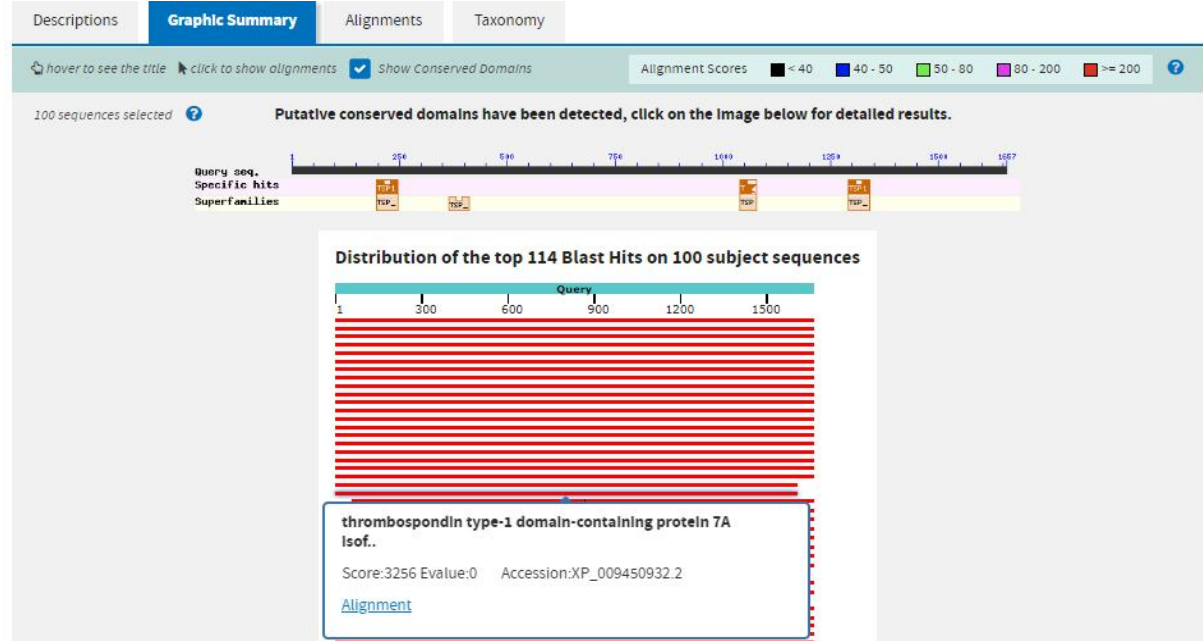
Download Manage Columns Show 100

☒ select all 100 sequences selected

GenPept Graphics Distance tree of results Multiple alignment

	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A precursor (Homo sapiens)	3395	3395	100%	0.0	100.00%	NP_056019.1
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A isoform X1 (Pan troglodytes)	3376	3376	100%	0.0	99.46%	XP_527689.3
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Gorilla gorilla gorilla)	3371	3371	100%	0.0	99.26%	XP_004045160.2
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Pongo abelii)	3343	3343	100%	0.0	98.31%	XP_002818243.2
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Nomascus leucogenys)	3343	3343	100%	0.0	98.31%	XP_003252614.2
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Hylobates moloch)	3343	3343	100%	0.0	98.31%	XP_032616349.1
<input checked="" type="checkbox"/>	THSD7A isoform 1 (Pongo abelii)	3340	3340	100%	0.0	98.25%	PNJ88887.1
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Theropithecus gelada)	3326	3326	100%	0.0	97.95%	XP_025235320.1
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A isoform X1 (Papio anubis)	3323	3323	100%	0.0	97.89%	XP_017812404.2
<input checked="" type="checkbox"/>	thrombospondin type-1 domain-containing protein 7A (Ptilocobus tephrosceles)	3322	3322	100%	0.0	97.77%	XP_026305173.1

Bu diziler için hizalama grafiği kabaca aşağıdaki gibi elde edilmiştir. Kırmızı renkli eşleşmeler, eşleşme puanı 200 ve üzerinde olanları ifade etmektedir.



Eşleşen dizilerden bir tanesi için eşleşen kısımların bir bölümünün sonuç dökümü ise aşağıdaki gibi elde edilmiştir.

PREDICTED: thrombospondin type-1 domain-containing protein 7A isoform X1 [Chlorocebus sabaeus]

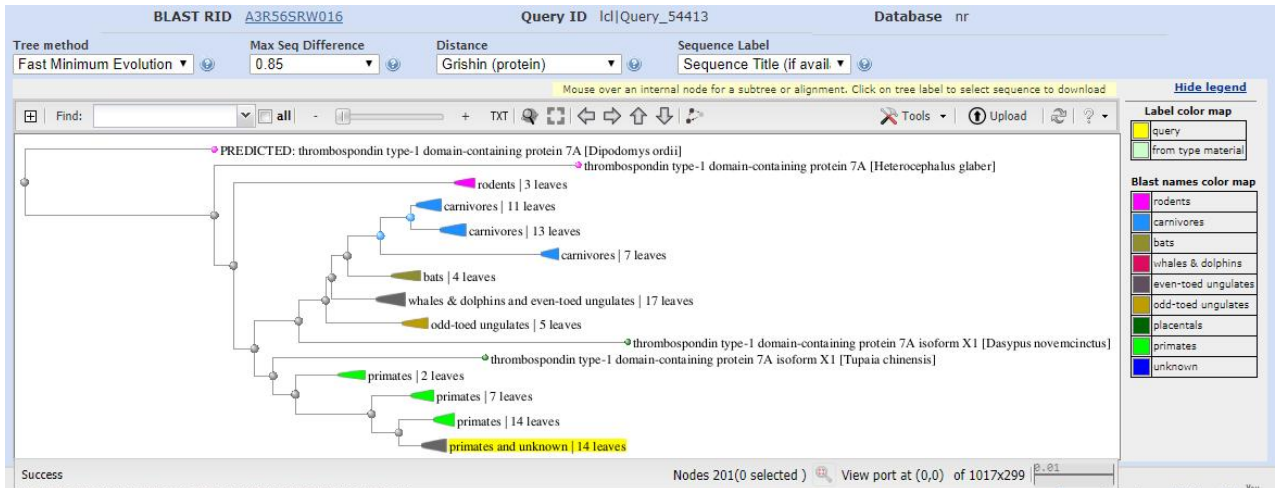
Sequence ID: [XP_007980223.1](#) Length: 1654 Number of Matches: 1

Range 1: 1 to 1654 [GenPept](#) [Graphics](#)

[▼ Next Match](#) [▲ Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
3314 bits(8592)	0.0	Compositional matrix adjust.	1620/1658(98%)	1637/1658(98%)	5/1658(0%)
Query 1	MGLQARRWASGSRGAAGPRRGVLQLLPLPLPLLLLLLLLRPGA-GRAAAQGEAEPTLY				59
Sbjct 1	MGLQAR WASGSRGAAGPRRGVLQLLPL L LPLLLLL PGA GRAAAQGEAE PTLY				56
Query 60	LWKTGPWGRCMGDECGPGGIQTRAVMCAHVEGWTTLHTNCKQAEPPNNQQNCFKVCDDWHK				119
Sbjct 57	LWKTGPWGRCMGDECGPGGIQTRAVMCAHVEGWTTLHTNCKQAEPPNNQQNCFKVCDDWHK				116
Query 120	ELYDWRGLGPNWQCQPVISKSLEKPLECIKGEEGIQVREIACIQKDKDIPAEDIICEYFEP				179
Sbjct 117	ELYDWRGLGPNWQCQPVISKSLEKPLECIKGEEGIQVREI CIQKDKDIPAEDIICEYFEP				176
Query 180	KPLLEQACLIPCCQDCIVSEFSAWSECSKTCGSGLQHRTRHVVPQFGGSGCPNLTEFQ				239
Sbjct 177	KPLLEQACLIPCCQDCIVSEFSAWSECSKTCGSGLQHRTRHVVPQFGGSGCPNLTEFQ				236
Query 240	VCQSSPCEAEELRYSLVHVPWSTCSMPHSRQVRQARRRGKNKEREKDRSGVKDPEAREL				299
Sbjct 237	VCQSSPCEAEEL YSLVHVPWSTCSMPHSRQVRQARRRGKNKEREKDRSGVKDPEAREL				296
Query 300	IKKKRRNRNRQNRQENKYWDIQIGYQTRVEMCINKTGKAADLSFCQQEKLPMTFQSCVITK				359
Sbjct 297	IKKKRRNRNRQNRQENKYWDIQIGYQTRVEMCINKTGKAADLSFCQQEKLPMTFQSCVITK				356
Query 360	ECQVSEWSEWSPCKTCHDMVSPAGTRVTRTRIRQFPFISGEKECEPEEKEPCLSQGDGV				419
Sbjct 357	ECQVSEWSEWSPCKTCHDMVSPAGTRVTRTRIRQFPFISGEKECEPEEKEPCLSQGDGV				416
Query 420	VPCATYGWRTTEWTECRVDPLLSQQDKRRGNQALCGGGIQTREYVCVQANENLLSQLST				479
Sbjct 417	APCATYGWRTTEWTECRVDPLLSQQDKRRGNQALCGGGIQTREYVCVQANENLLSQLNT				476
Query 480	HKNKEASKPMDLKLCTGPIPNNTTQLCHIPCPTCEVSPWSAWGPCTYENCNDQQGKKGFK				539
Sbjct 477	HKNKEASKP++LKLCTGPIPNNTTQLCHIPCPTCEVSPWSAWGPCTYENCNDQQGKKGFK				536
Query 540	LRKRRITNEPTGGSGVTGNCPHLLAIPCEEPACYDMKAVRLGNCEPDNGKECGPGTQVQ				599
Sbjct 537	LRKRRITNEPTGGSG TGNCPHLLAIPCEEPACYDMKAVRLG+CEPDNGKECGPGTQVQ				596
Query 600	EVVCINSDGEEVDRLCRDAIFPIPVACDAPCPKDCVLSTWSTWSSSHTCSGKTTEGKQ				659
Sbjct 597	EVVCINSDGEEVDRLCRDAIFPIPVACDAPCPKDCVLSTWSTWSSSHTCSGKTTEGKQ				656

Benzerliklerine göre sorgu sonucunda elde edilen veriler için dizilerin benzerlik haritası ise aşağıdaki gibidir.



3. YARARLANILAN KAYNAKLAR

[1][http://resources.qiagenbioinformatics.com/manuals/clcmainworkbench/800/index.php?manual=How does BLAST work.html](http://resources.qiagenbioinformatics.com/manuals/clcmainworkbench/800/index.php?manual=How%20does%20BLAST%20work.html)

[2][https://bio.libretexts.org/Bookshelves/Cell and Molecular Biology/Book%3A Investigations in Molecular Cell Biology \(O'Connor\)/9%3A Protein Conservation/9.3%3A BLAST algorithms are used to search databases](https://bio.libretexts.org/Bookshelves/Cell_and_Molecular_Biology/Book%3A_Investigations_in_Molecular_Cell_Biology_(O'Connor)/9%3A_Protein_Conservation/9.3%3A_BLAST_algorithms_are_used_to_search_databases)

[3][https://en.wikipedia.org/wiki/BLAST \(biotechnology\)](https://en.wikipedia.org/wiki/BLAST_(biotechnology))

[4]<https://biotechgo.org/tr/?view=article&id=247:lo2&catid=136>

[5][https://www.ccg.unam.mx/~vinuesa/tlem/pdfs/Bioinformatics explained BLAST.pdf](https://www.ccg.unam.mx/~vinuesa/tlem/pdfs/Bioinformatics_explained_BLAST.pdf)

[6]<http://www.yasinhoca.com/2016/11/dna-dizilerini-karşılaştırma-blast.html>

[7][https://www.comp.nus.edu.sg/~ksung/algo in bioinfo/slides/Ch5 database.pdf](https://www.comp.nus.edu.sg/~ksung/algo_in_bioinfo/slides/Ch5_database.pdf)

[8]<https://bioinf.comav.upv.es/courses/biotech3/theory/blast.html>

[9]<https://www.youtube.com/watch?v=4AcnZnZRss8>

[10]<https://www.youtube.com/watch?v=g0nSH17psDc>

[11]<https://www.youtube.com/watch?v=NnY2f5111FU>