

PRIMERJAVE ZAPOREDIJ

- Primerjave zaporedij; homologije, podobnosti, identičnost; poravnava, prileganje;
- Homologi, ortologi paralogi;
- Ocenjevanje poravnav; ocenjevalne matrike, rezultat poravnave, teža vrzeli;
- Globalne, lokalne poravnave;
- Točkovni diagrami.

PRIMERJAVE ZAPOREDIJ

Primerjamo med seboj dve zaporedji: proteini, DNA, RNA

Osnoven proces. Iskanje po bazah.

query sequence (*probe*) vs **podatkovna zbirka** (*database, subject*)

Kako podobne so si zaporedja

malo počakati na analizo bioloških podatkov ali filogenetske analize

zelo lahko govoriš o družinah proteinov

npr. v neki družini proteinov je znanih 10 predstavnikov pri podgani, a le 7 pri človeku. Verjetno bodo odkriti še trije:

Farmakološko stališče: poznamo biološki efekt, a ne poznamo proteina zanj

Molekularno-biološko stališče: priložnost kloniranja novih človeških proteinov (*in silico* kloniranje).

PRIMERJAVE ZAPOREDIJ

Homologija (*homology*) Kvalitativna oznaka

Poudarja evolucijsko povezanost dveh zaporedij. Homologni proteini so se razvili iz skupnega prednika. Podobnost je med različnimi družinami homolognih proteinov različno ohranjena. Homologni proteini imajo vedno podobno 3D zgradbo. Pri sklepanju na homologije je potrebna previdnost in statistična analiza rezultatov (eksperimentiranje, ki je podobno eksperimentiranju v laboratorijih, “*wet labs*”)

Podobnost (*similarity*) Kvantitativna oznaka

Zaporedji sta si lahko zelo podobni, ni pa nujno tudi evolucijske povezanosti.

Identičnost (*identity*) Kvantitativna oznaka

Pove koliko je identičnih amino kislin ali nukleotidov med dvema zaporedjima.

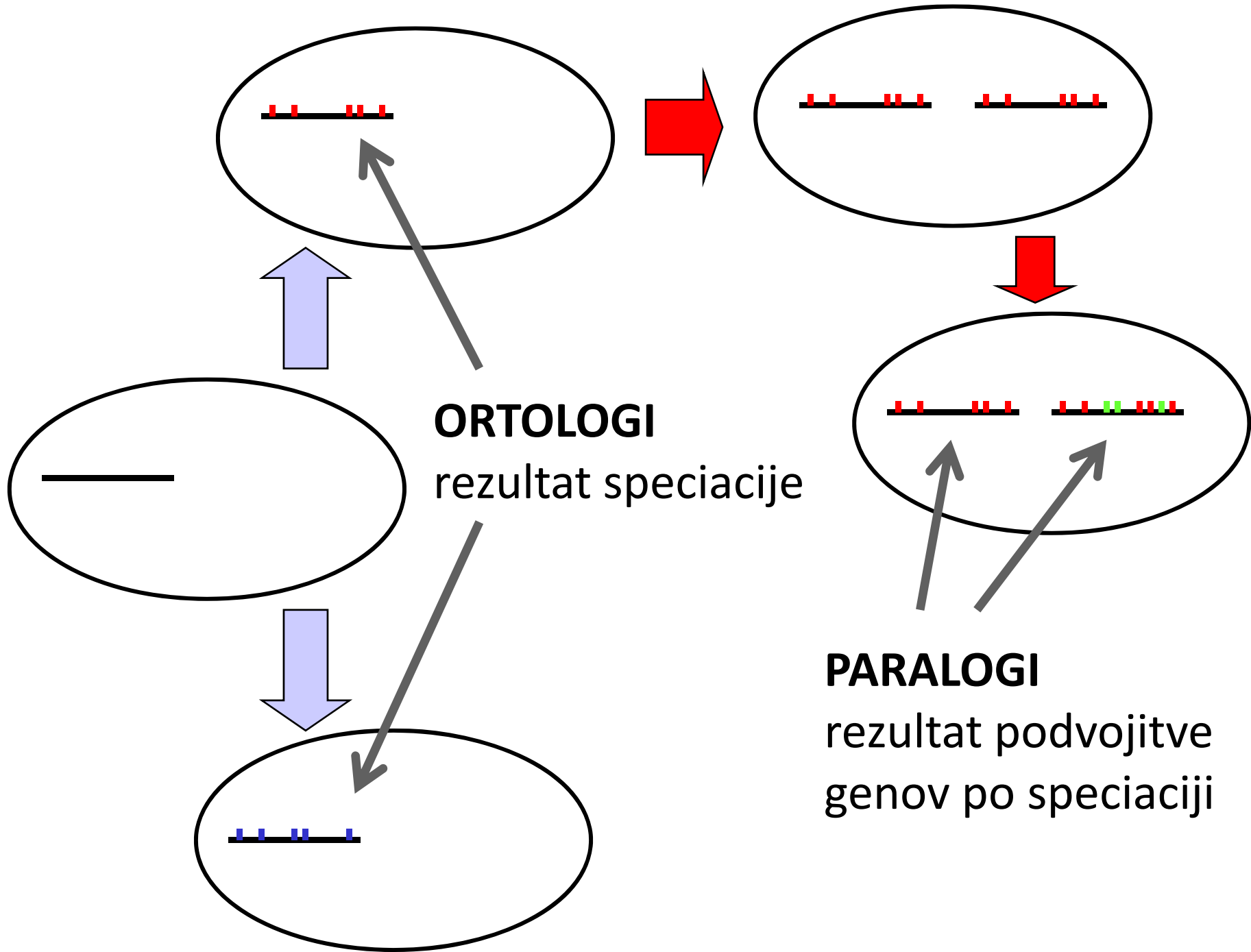
PORAVNAVANJE, PRILEGANJE

PORAVNAVA

Poseben zapis dveh ali več zaporedij. Znake izpišemo v stolpcih.

A	G	G	V	L	I	I	Q	V	G	N	M	R	T	P
:	:	:	:		:	:				:	:			
A	G	G	D	L	V	I	Q	-	-	N	M	K	S	N

Bistvo poravnave zaporedij je poiskati takšno, ki pove nekaj o evlucijskih dogodkih, ki so privedli do današnjega stanja (zamenjave, insercije, delecije).



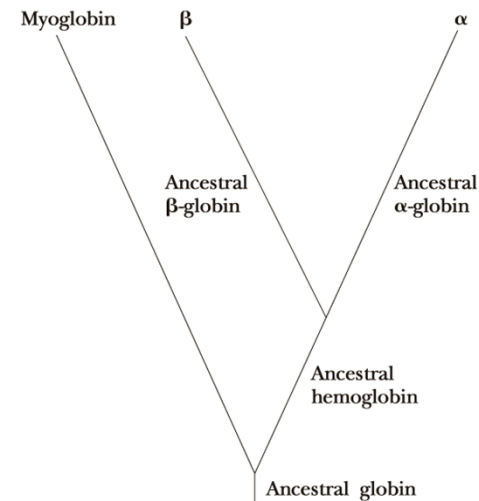
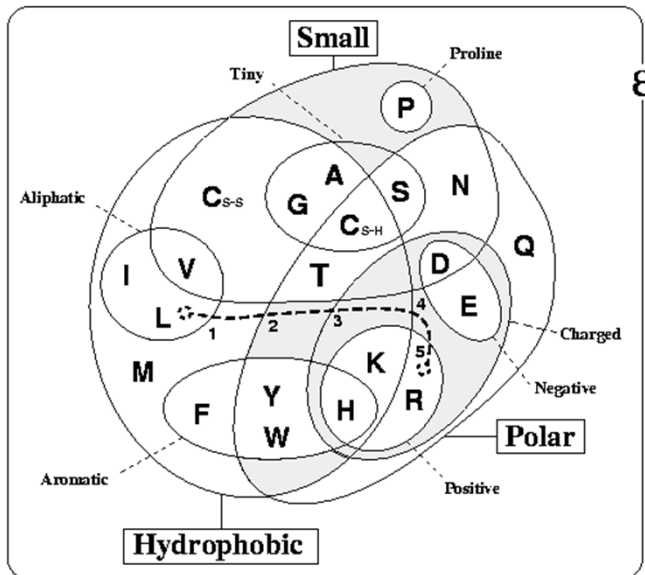
ORTOLOGI
rezultat specijacije

PARALOGI
rezultat podvojitve
genov po specijaciji




PORAVNAVA, PRILEGANJE

	10	20	30	40	50	60	70	80	
ALPHA1	-MVLSPADKTNVKA	AWGKVGAHAGEYGA	EALERMFLSFPTTK	TYFPHFDLSH----	G-SAQVKGHGKKV	ADALTNVAHVDD	MPN	80	
ALPHA2	-MVLSPADKTNVKA	AWGKVGAHAGEYGA	EALERMFLSFPTTK	TYFPHFDLSH----	G-SAQVKGHGKKV	ADALTNVAHVDD	MPN	80	
THETA	-MALSAEDRALV	RALWKKLGSNVGV	TTEALERTFLAF	PATKTYFSHLDL	SP----	G-SSQVRAHGQ	KVADALSLAVER	80	
GAMMA	MGHFTEEDKATIT	SLWGVN--VEDAG	GETLGRLLVVYP	PWTQRFFDSFG	NLSSASAIMGN	PVKVKAHGKKV	LTSLGDAIKH	85	
BETA	MVHLTPEEKSA	VTALWGVN--VDE	VGGGALGRLLVV	YPWTQRFFESF	GDLSSTPDAVM	GNPKVKAHGKK	VLGAFSDGLA	85	
EPSILON	MVHFTAEEKAA	VTSLWSKMN--	VEEAGGALGRLL	VVYPWTQRFFD	SFGNLSSPSAIL	GNPKVKAHGKK	VLTSFGDAIK	85	
DELTA	MVHLTPEEKTA	VNALWGVN--VDA	VGGGALGRLLVV	YPWTQRFFESF	GDLSSTPDAVM	GNPKVKAHGKK	VLGAFSDGLA	85	
MYOGLOBIN	-MGLSDGEWQL	VLVNVWCKVEAD	IPGHGQEVLRIR	LFKGHPETLEK	EDKEKHLKSEDE	MKASEDLKKHG	ATVLTALGGIL	86	

	90	100	110	120	130	140	150	
ALPHA1	LSALS	DLHAHKL	RVDPVNF	KLLSHCLL	VTLAAHL	PAEFTPAV	HASLDKFL	142
ALPHA2	LSALS	DLHAHKL	RVDPVNF	KLLSHCLL	VTLAAHL	PAEFTPAV	HASLDKFL	142
THETA	LSALS	SHLHACQ	LRVDPASF	OLLGHCLL	VTLARHYP	GDSPALQ	ASLDKFLSH	142
GAMMA	FAQL	SELHCD	KLHVDPE	NFKLLGN	VLVTLA	IHF	GKEFTPEV	147
BETA	FATL	SELHCD	KLHVDPE	NFRLLGN	VLVLCV	LAHHF	GKEFTPPV	147
EPSILON	FAKL	SELHCD	KLHVDPE	NFKLLGN	VMVIL	AHFG	GKEFTPEV	147
DELTA	FSQL	SELHCD	KLHVDPE	NFRLLGN	VLVLCV	LAHNF	GKEFTPQ	147
MYOGLOBIN	IKPLA	QSHATK	HKIPV	KYLEFISE	CEIIQV	LQSKHP	GDFGADAQ	154



	A	G	G	V	L	I	I	Q	V	G	N	M	R	T	P
	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:
	A	G	G	D	L	V	I	Q	-	-	N	M	K	S	N
Rezultat	1	1	1	0	1	0	1	1	-1	-1	1	1	0	0	0

		
UJEMANJE	VRZEL	NEUJEMANJE

Poravnavi lahko pripišemo **rezultat** poravnave, ki je odraz:

- ujemanj (števila identičnih elementov med obema zaporedjima) 1
- neujemanj (števila različnih elementov med obema zaporedjima) 0
- vrzeli (insercij ali delecij)
- -1

$$1+1+1+0+1+0+1+1-1-1+1+1+0+0+0=6$$

Uteži za ujemanja in neujemanja, ki jih uporabljamo za določitev rezultata poravnave imenujemo **ocenjevalna matrika**, uteži za vrzeli imenujemo **teža vrzeli**.

S kvantifikacijo ujemanj in neujemanj poskušamo poiskati poravnavo, ki da najvišji rezultat.

Kako poravnati?

“Na oko” vs računalnik z algoritmi

Primer:

AGGVLI IQVG 6

AGGVLI IQVG 9

***** ***

AGGVLIQVG

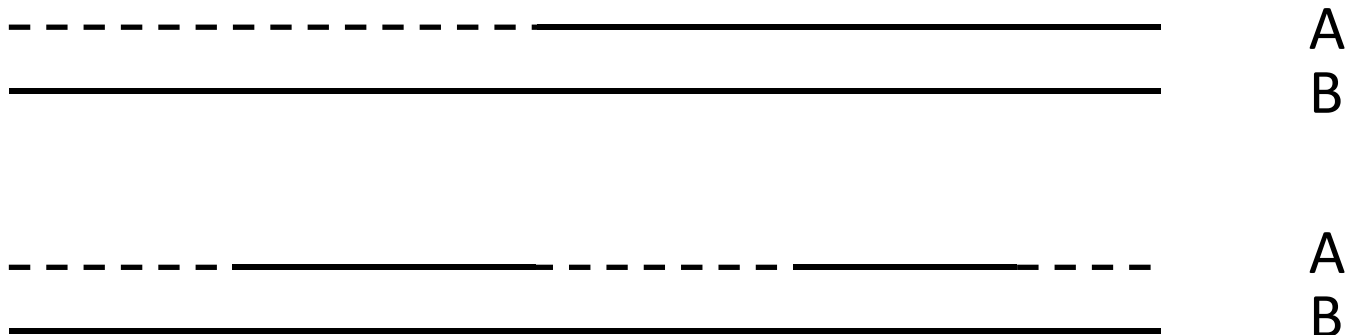
AGGVLI-QVG

Zaporedja bioloških makromolekul daleč od idealnih primerov:

dolga zaporedja (večina proteinov ima 200-500 amino kislin)

različne dolžine

možne različne poravnave

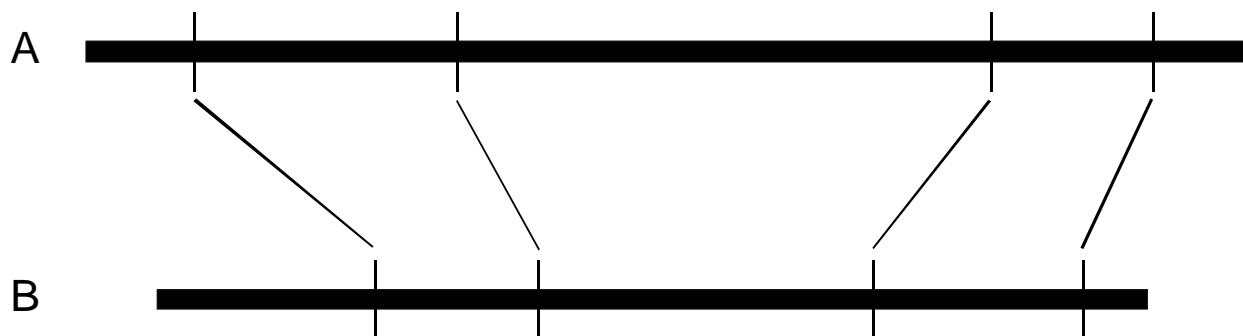




Kadar pričakujemo podobnost celih ali večine zaporedij



Globalna poravnava



Kadar so zaporedja podobna v nekem krajšem segmentu



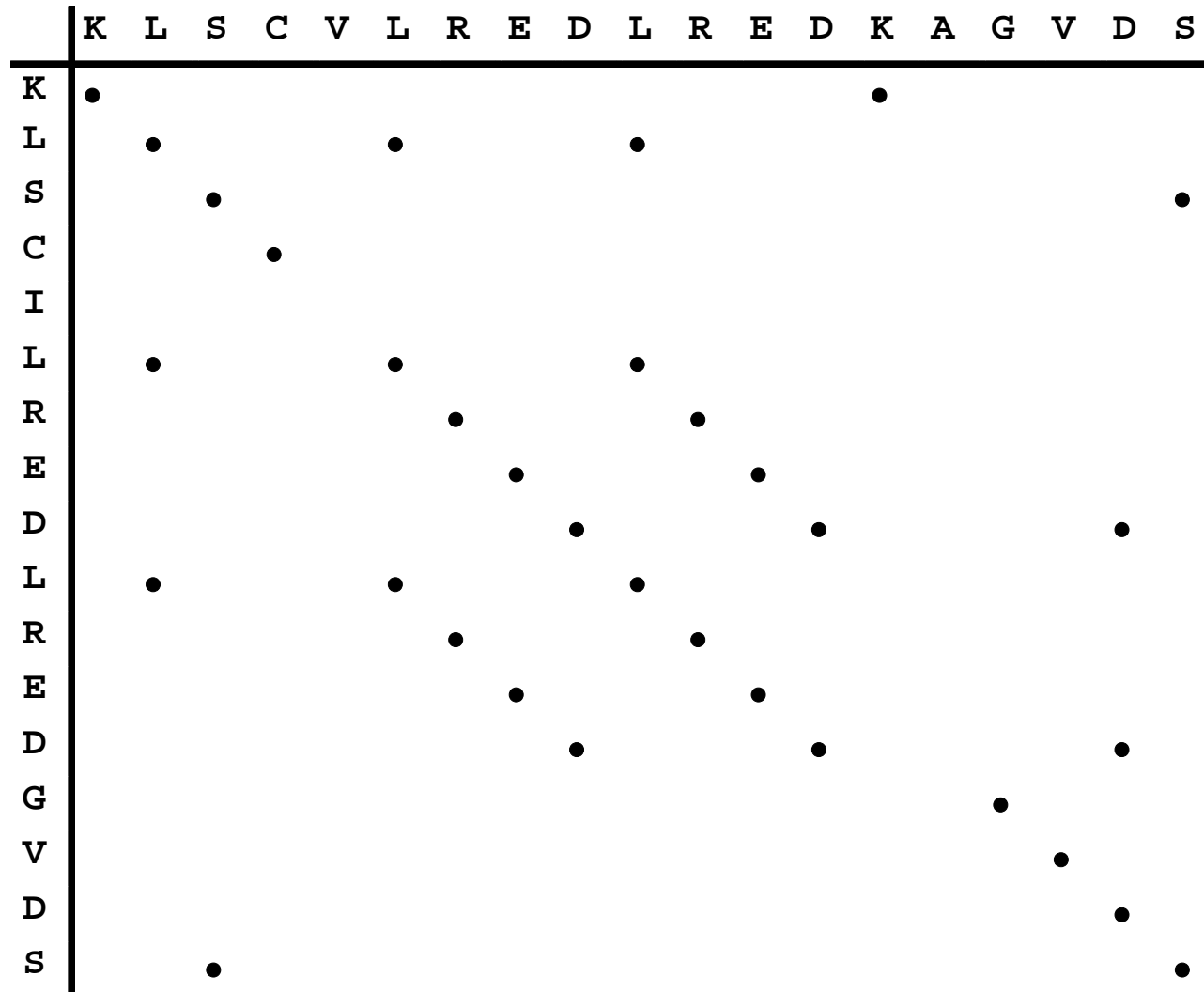
Lokalna poravnava



TOČKOVNI DIAGRAMI

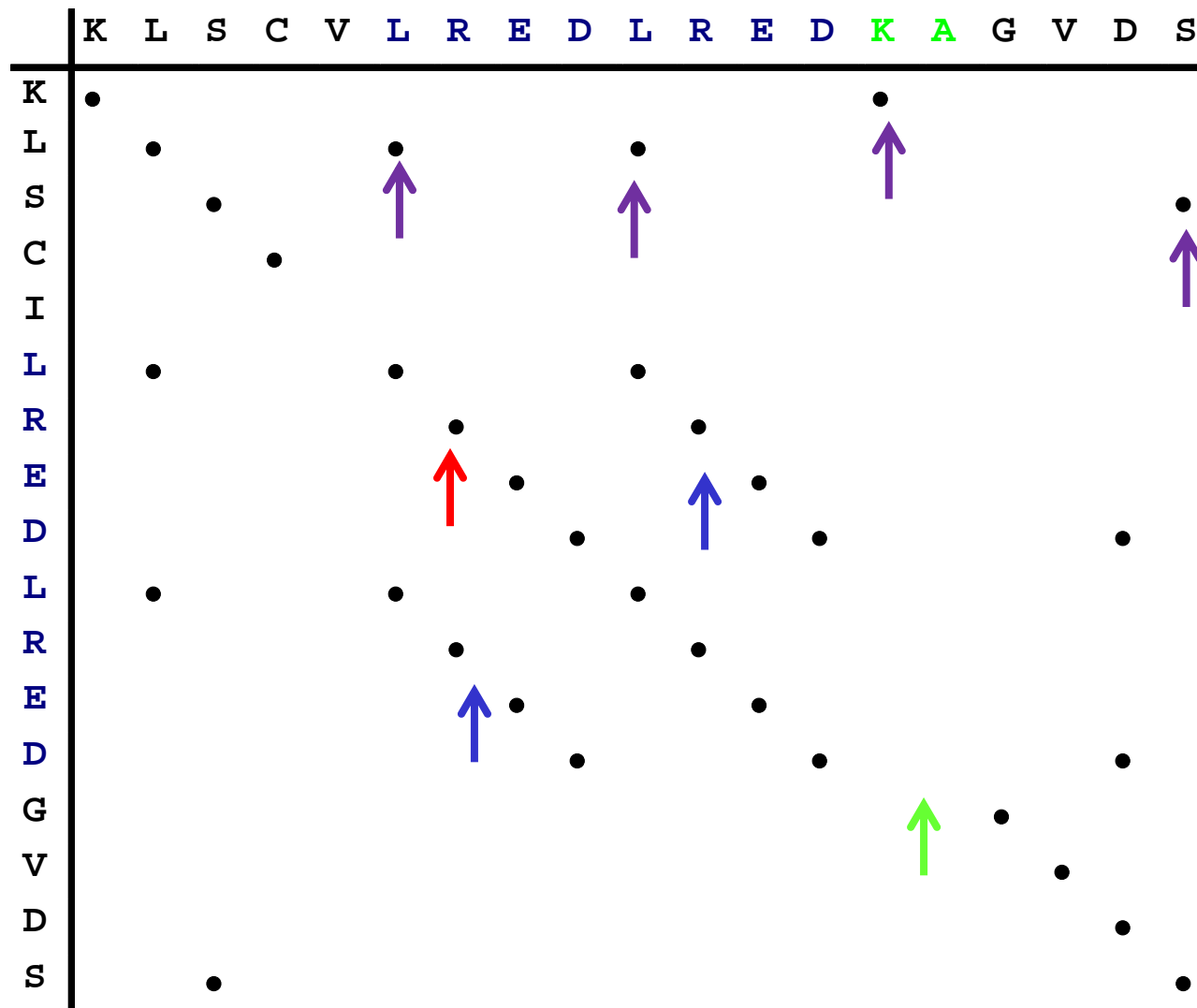
DOT PLOT

Gibbs AJ and McIntyre GA (1970) The diagram, a method for comparing sequences. Its use with amino acid and nucleotide sequences. *European Journal of Biochemistry* 16:1-11.



Najbolj osnovna metoda. Dve zaporedji primerjamo med seboj vizuelno, "na oko"

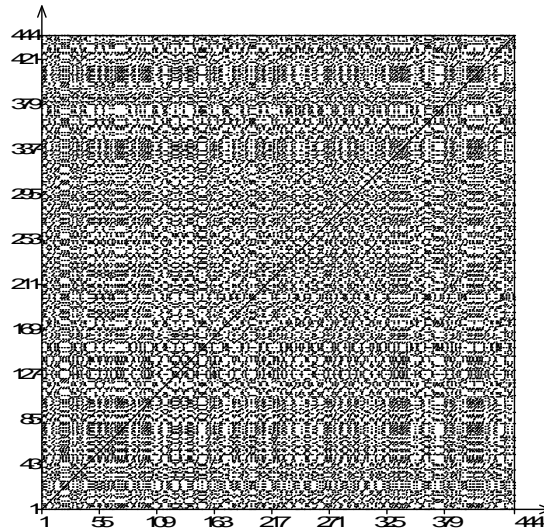
Pripravimo matriko tako, da na eno os nanesemo elemente enega zaporedja in na drugo elemente drugega zaporedja. Položaj z identičnimi elementi označimo.



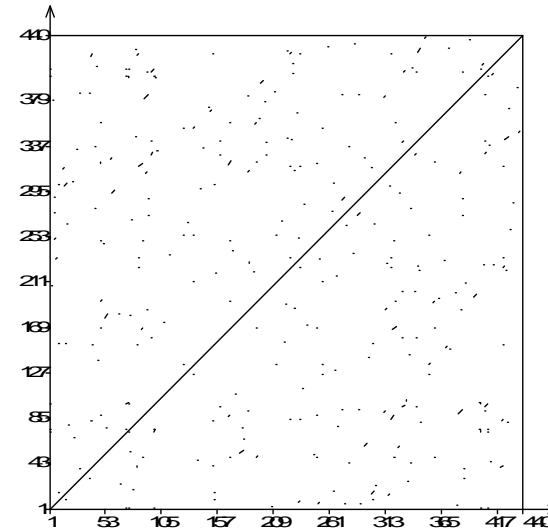
Identične aminokisliline predstavlja diagonalna črta.

Lepo so vidne **ponovitve** znotraj zaporedja in **delecije**.

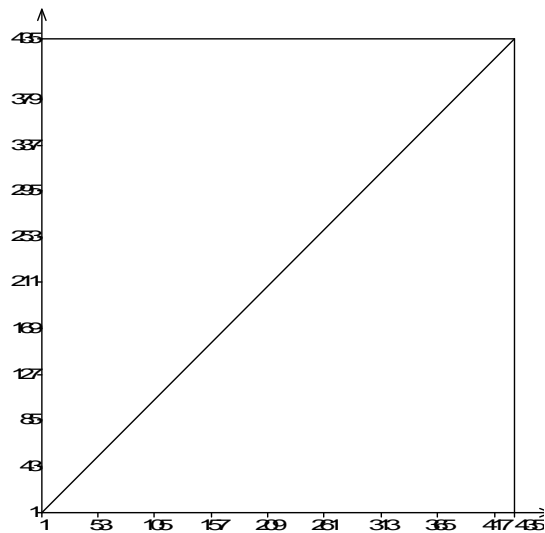
Šum zaradi zadetkov, ki se pojavijo naključno



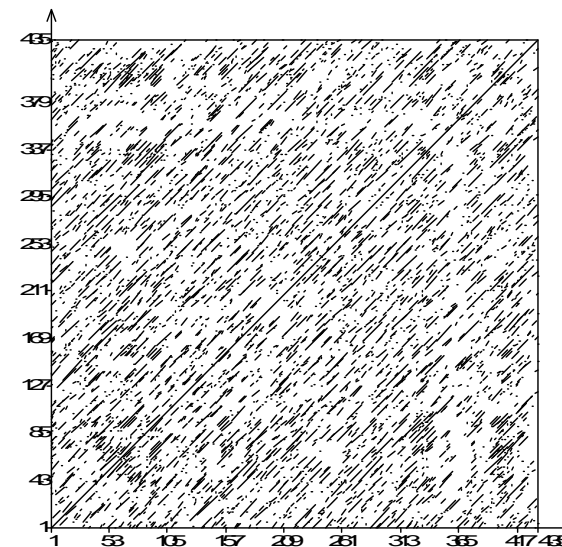
Ujemanje 100%
Okno 1



Ujemanje 100%
Okno 5



Ujemanje 100%
Okno 10



Ujemanje 50%
Okno 10

Uporaba filtrov za zmanjšanje šuma. Okna, ujemanje.

Maizel JV Jr and Lenk RP (1981) Enhanced graphic matrix analysis of nucleic acid and protein sequences. *Proceedings of National Academy of Sciences USA* 78:7665-7669.

Uporaba matrik za računanje rezultata v posameznem oknu, npr. kemijska podobnost.

Staden R. (1982) An interactive graphics program for comparing and aligning nucleic acid and amino acid sequences. *Nucleic Acids Research* 10:2951-2961.

>gi|122615|sp|P02023|HBB_HUMAN HEMOGLOBIN BETA CHAIN
MVHLTPEEKSAVTALWGKVNVDVEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLGAF
SDGLAHLNLDNLKGTFFATLSELHCDKLHVDPENFRLLG NVLVCVLAH HFGKEFTPPVQAAYQKVVAGVANALA
HKYH

>gi|7428621|pir||HBCZ hemoglobin beta chain - chimpanzee
MVHLTPEEKSAVTALWGKVNVDVEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLGA
FSDGLAHLNLDNLKGTFFATLSELHCDKLHVDPENFRLLG NVLVCVLAH HFGKEFTPPVQAAYQKVVAGVANAL
AHKYH

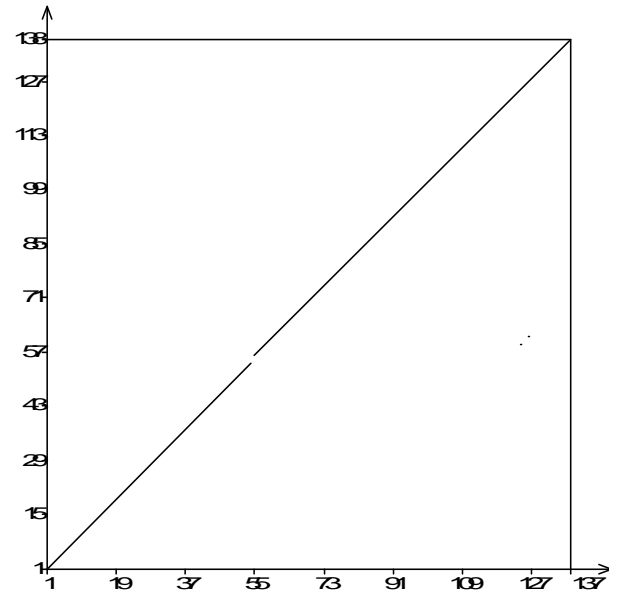
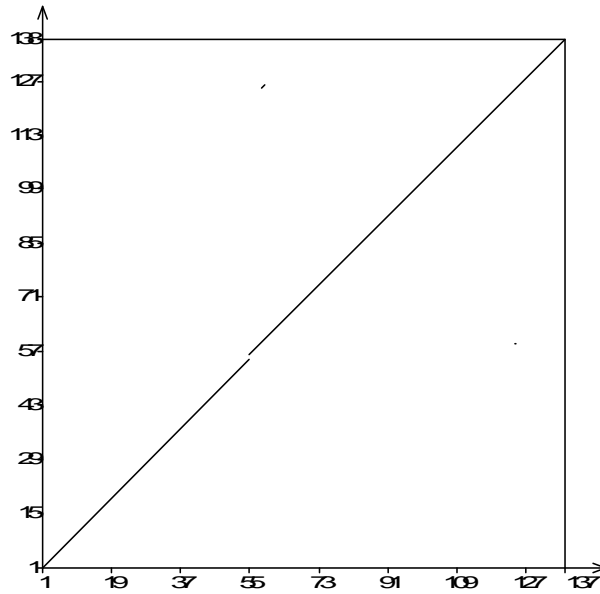
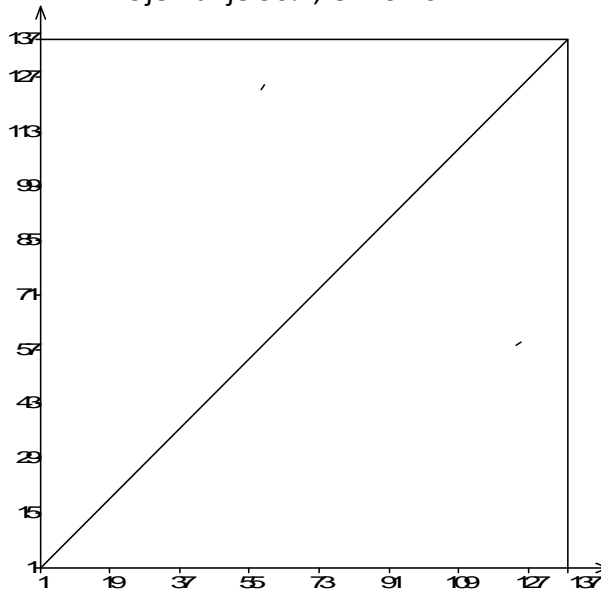
>gi|2144722|pir||HBRB hemoglobin beta chain - rabbit
MVHLSSEEKSAVTALWGKVNVEEVGGEALGRLLVVYPWTQRFFESFGDLSSANAVMNNPKVKAHGKKVLAA
FSEGLSHLDNLKGTFAKLSELHCDKLHVDPENFRLLG NVLVIVLSH HFGKEFTPPVQAAYQKVVAGVANAL
AHKYH

>gi|70509|pir||HBGY hemoglobin beta chain - goldfish
VEWTD AERSAI IGLWGKLN PDELGPQALARCLIVYPWTQRYFATFGNLSSPAAIMGNPKVAAHGRTVMGGL
ERA IKNMDNIKATYAPLSVMHSEKLVHVD PDNFRL LADCITVCAAMKFGPSGFNADVQEAWQKFLSVVVSAL
CRQYH

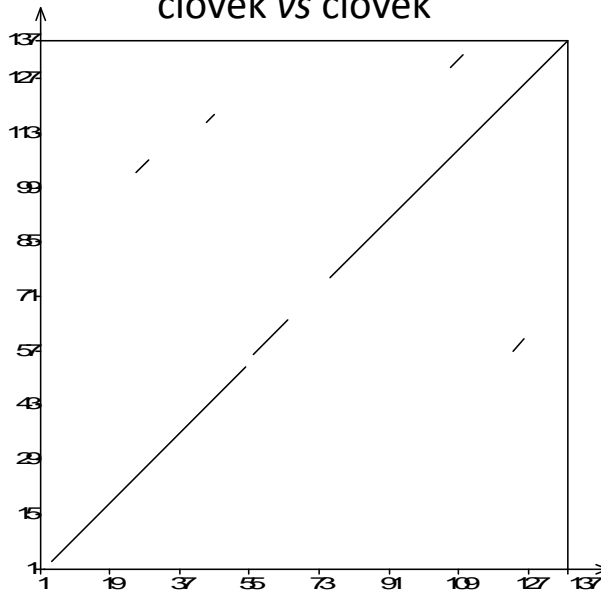
>gi|86549|pir||B24625 hemoglobin beta chain - Eurasian tree sparrow
VQWTAEEKQLITGLWGKVNVAECGGEALARLLIVYPWTQRFFASFGNLSSPTAVLGNPKVQAHGKKVLTSF
GEAVKNLDSIKNTFSQLSELHCDKLHVDPENFRL LGDILVVVLA AHFGKDFTPDCQAAWQKLV RVVAHALA
RKYH

>gi|70497|pir||HBAK hemoglobin beta chain - Nile crocodile
ASFD PHEKQLIGDLWHKVDVAHCGGEALS RMLIVYPWKRRYFENFGDISNAQAIMHNEKVQAHGKKVLASF
GEAVCHLDGIRAHFANLSKLHCEKLVHVDPENFKLLGDI I I IIVLAAHYPKDFGLECHAAYQKLV RQVAAALA
AEYH

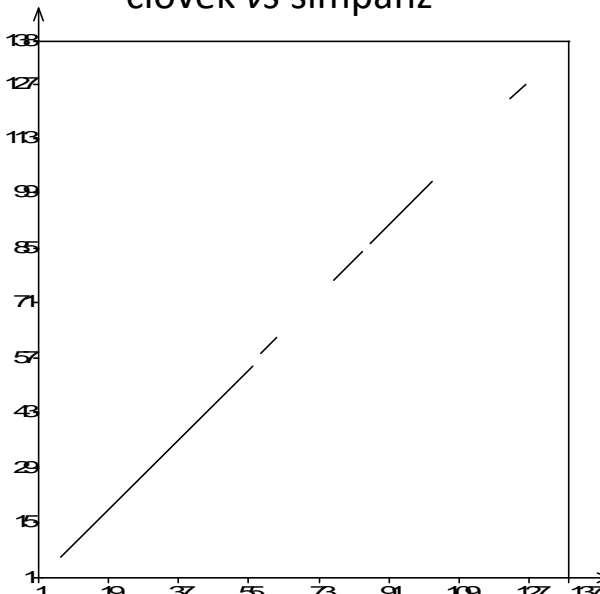
Ujemanje 50%; Okno 10



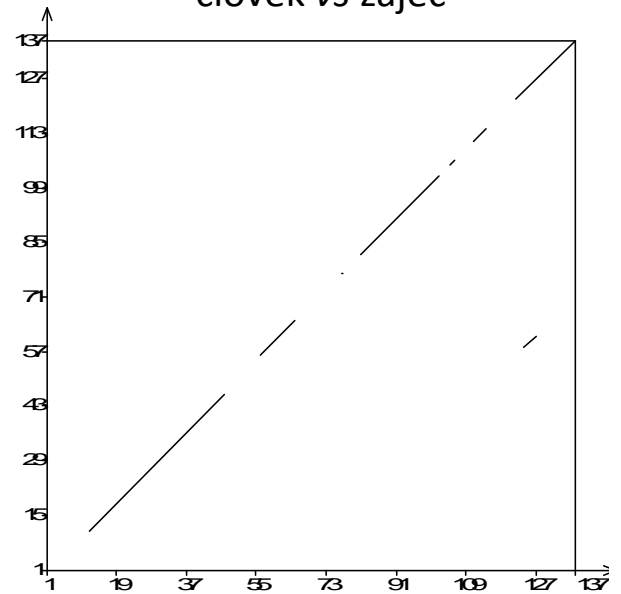
človek vs človek



človek vs šimpanz



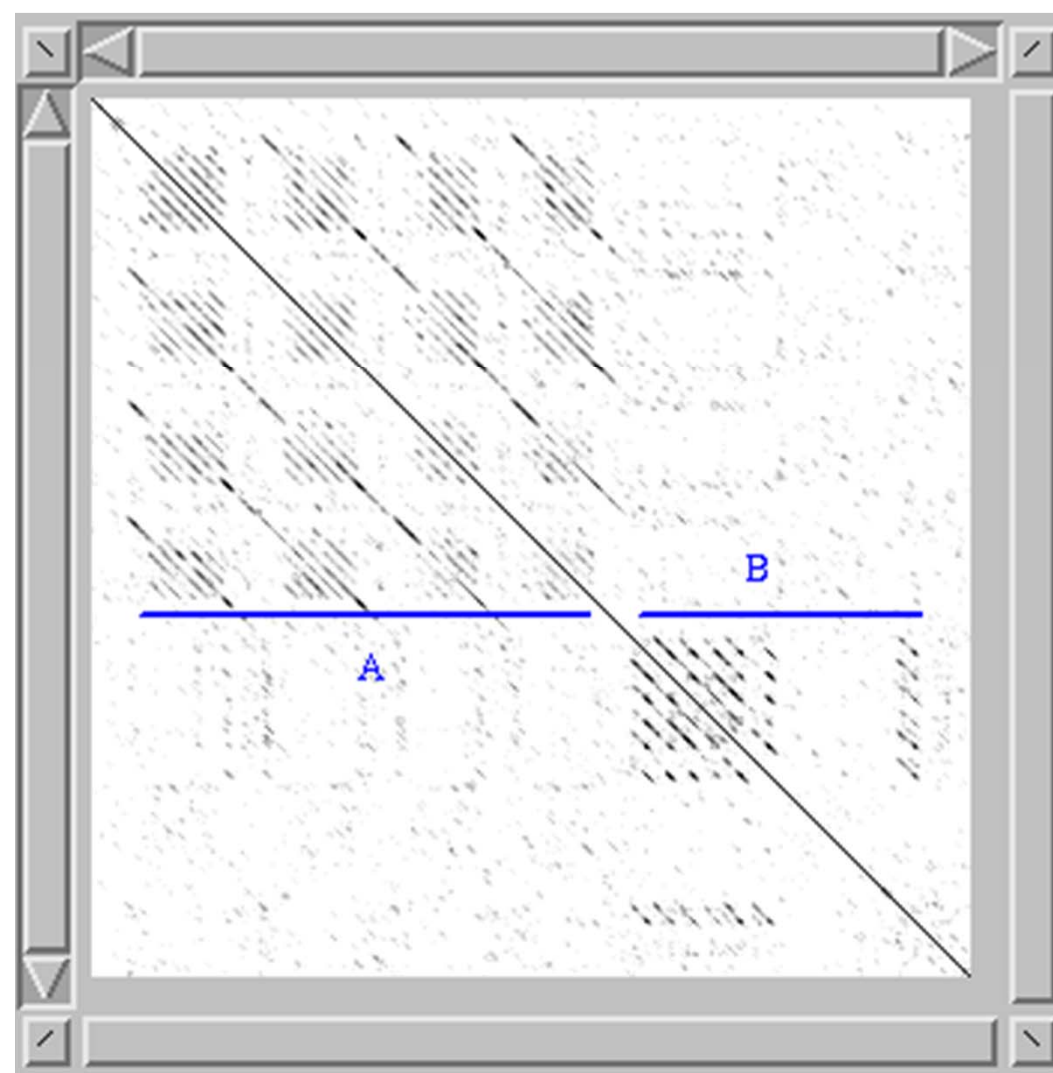
človek vs zajec



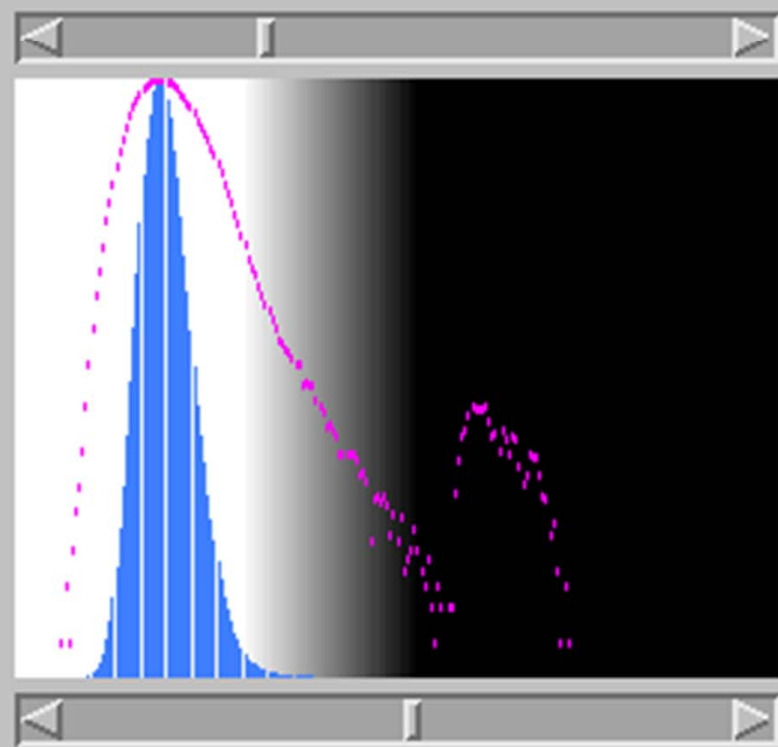
človek vs lastovka

človek vs zlata ribica

človek vs krokodil

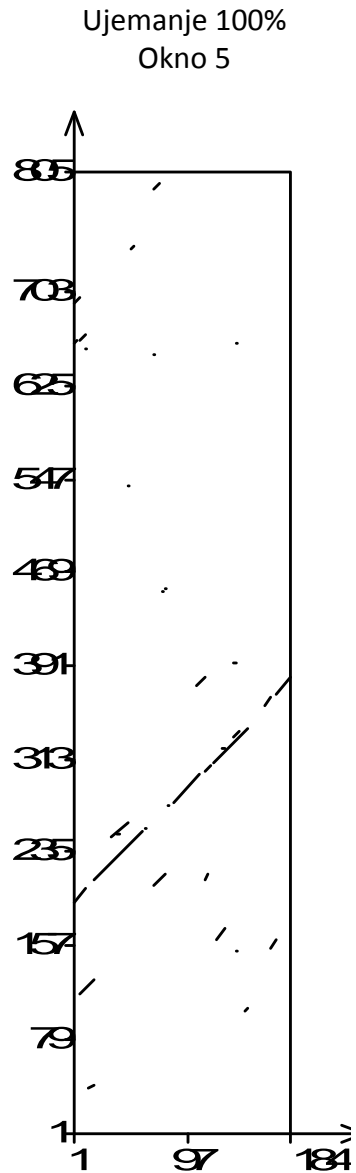


horizontal: SLIT_DROME
vertical: SLIT_DROME
matrix: Blosum62
sliding window: 15
zoom: 1:5
gray scale: 53% - 30%



Ohranjene domene

MS2 površinski celični antigen
vs adamalysin II,
metaloproteinaza iz strupa
Crotalus adamanteus (Eastern
diamondback rattlesnake).
Oba vsebujeta Zn-proteazno
domeno



```
>MS2_HUMAN (P78325)
MRGLGLWLLGAMMLPAIAPSRPWALMEQYEVVLPRLPG
PRVRRALPSHLGLHPPERVSIVLGGATGHNFLLHLRKNRDL
LGSGYTETTYTAANGSEVTEQPRGQDHCLYQGHVEGYPDS
AASLSTCAGLRGFFQVGSDDLHLIEPLDEGEGGRHAVYQ
AEHLLQTACTCGVSDDSLGSLLGPRTAAVFRPRPGDSL
SRETRYVELYVVVDNAEFQMLGSEAAVRHRVLEVVNHVD
KLYQKLNFRVVLVGLIWNISQDRFHVSPDPSVTLENLLT
WQARQRTTRRHLHDNVQLITGVDFGTGTVGFARVSAMCSH
SSGAVNQDHSKNPVGACTMAHEMGNLGMDDHENVQGC
RCQERFEAGRCIMAGSIGSSFPRMFSDCSQAYLESFLER
PQSVCLANAPDLSHLVGGPVCGNLFVERGEQCDCGPPED
CRNRCCNSTTCQLAEGAQCAHGTCCQECKVKPAGELCRP
KKDMCDLEEFCDGRHPECPEDAFQENGTPCSGGYCYNGA
CPTLAQQCQAFWGPGGQAAEESCFYSYDILPGCKASRYRA
DMCGVLQCKGGQQLGRAICIVDVCHALTTEDGTAYEPV
PEGTRCGPEKVCWKGRQCQLHVVYRSSNCSAQCHNHGVCN
HKQECHCHAGWAPPHCAKLLTEVHAASGLPVLVVVVVLV
LLAVVLTLAGIIVYRKARSRIILSRNVAPKTTMGRSNPL
FHQAASRVPAKGGAPAPSRGPQELVPTTHPGQPARHPAS
SVALKRPPPAPPVTVSSPFPVYVYTRQAPKQVIKPTFA
PPVPPVKPGAGAAANPGPAEGAVGPKVALKPPIQRKQGAG
APTAP
>ADAM_CROAD (P34179)
QQNLFPQRYIELVVVADRRVFMKYNSDLNIIIRTVHEIVN
IINGFYRSLNIDVSLVNLEIWSGQDPLTIQSSSNTLNS
EGLWREKVLNKKKKDQAQLLTAIEFKCETLKGAYLNSM
CNPRSSVGIVKDHSPINLLVAVTMAHELGHNLGMEHDGK
DCLRGASLCIMRPGLTGGRSYEFSDDSMGYYQKFLNQYK
PQCILNKP
```



Eksoni, introni

Preveden
Emericella
(Aspergillus)
nidulans
calmodulin gen vs
protein

TOČKOVNI GRAFI NA INTERNETU

DOTLET

<http://myhits.isb-sib.ch/cgi-bin/dotlet>

Pagni M and Junier T. Swiss Institute of Bioinformatics, Epalinges, Switzerland.
Točkovni graf dveh zaporedij, določiš okno, matriko.

DNADOT

<http://arbl.cvmbs.colostate.edu/molkit/dnadot/>

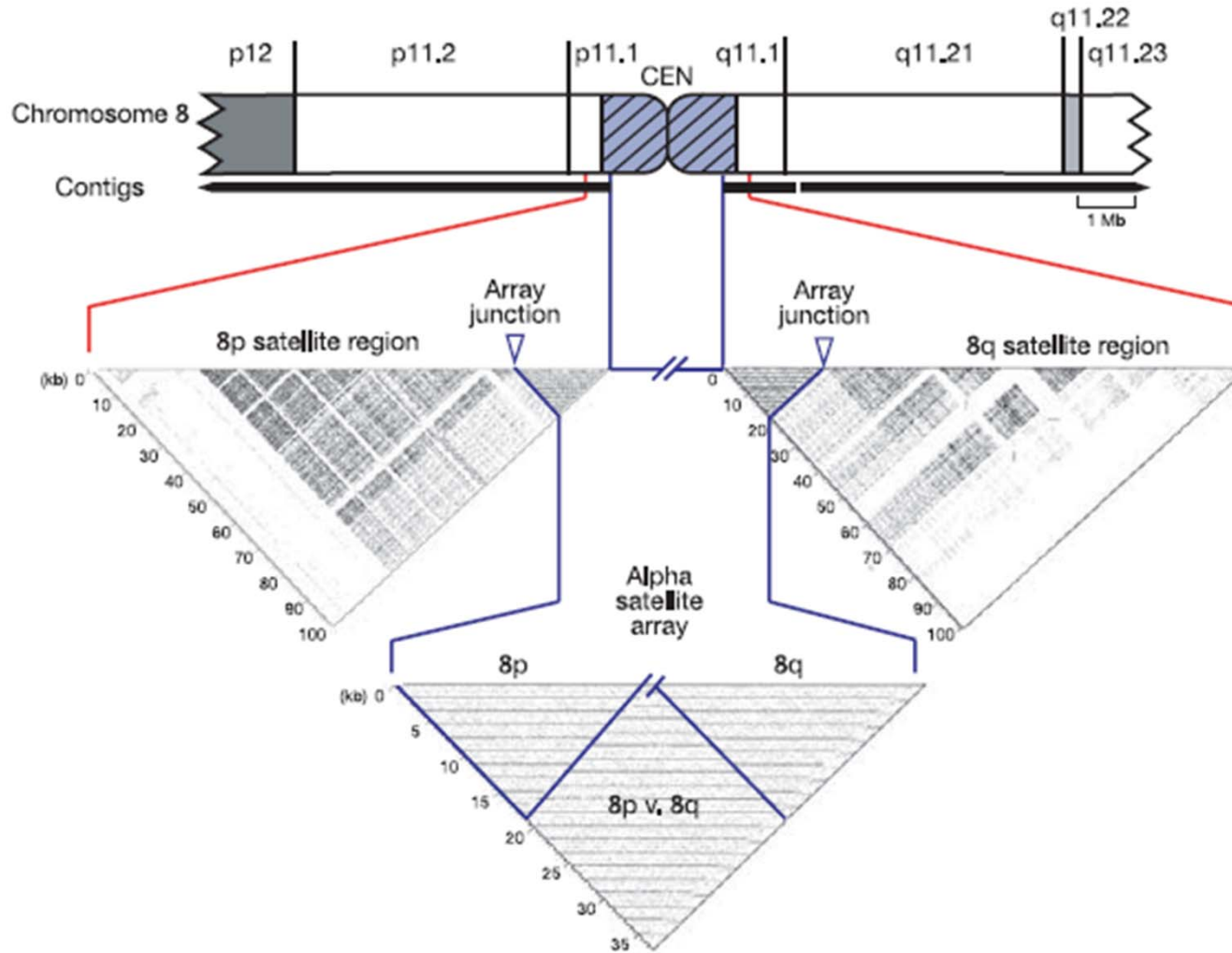
V Molecular Tool Kit. Na Colorado State University, USA.
Za primerjavo DNA zaporedij. Določiš okno in ujemanje.

DOTTER

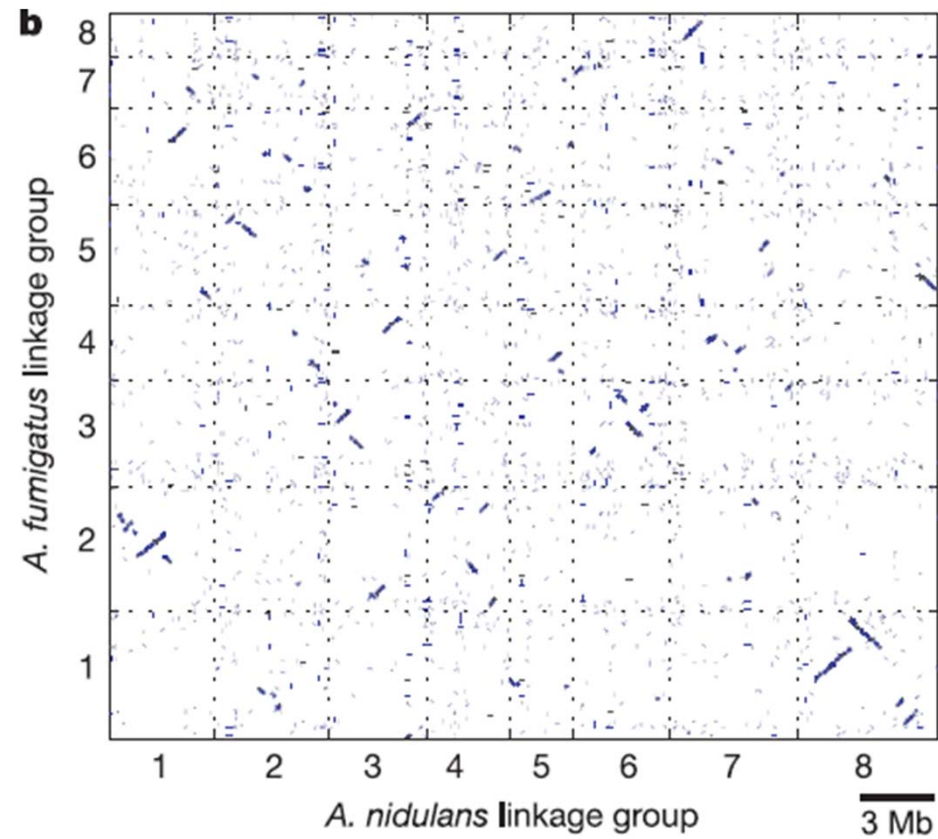
<http://sonnhammer.sbc.su.se/Dotter.html>

Karolinska Institutet, Center for Genomic Research, Švedska
Točkovni graf dveh zaporedij, določiš okno, matriko. Zaporedji primerja med seboj v dveh dimenzijah, v tretji dimenziji poda velikost rezultata kot vrh.

TOČKOVNI GRAFI V GENOMIKI



Nusbaum C et al. (2006) DNA sequence and analysis of human chromosome 8. Nature 439, 331-335.



Galagan JE et al. (2005) Sequencing of *Aspergillus nidulans* and comparative analysis with *A. fumigatus* and *A. oryzae*. *Nature* 438, 1105-1115.