

Protomata-Learner

Learning Automata on Protein Sequences

François Coste, Goulven Kerbellec



Rencontres Bioinformatique de la Plateforme Genouest,
le 3 décembre 2007

Outline

1 Introduction

2 Characterization

3 Generalization

4 Demonstration

5 Conclusion

Introduction

Bioinformatic problem

- Biological question :
- Computer science answer :

Bioinformatic problem

- Biological question :
How to define signatures of known protein families ?
- Computer science answer :

Bioinformatic problem

- Biological question :
How to define signatures of known protein families ?
- Computer science answer :
Using machine learning algorithms !

Protein families

- Amino acid alphabet
→ {A,R,N,D,C,Q,E,G,H,I,L,K,M,F,P,S,T,W,Y,V}



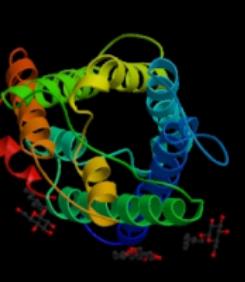
Protein families

- Amino acid alphabet
 - {A,R,N,D,C,Q,E,G,H,I,L,K,M,F,P,S,T,W,Y,V}
- Protein sequence
 - >AQP1_bovin MASEFKKKLFWRAVVAEFLAMILFIFISIG-SALGFHYPIKSQTTGAVQDNVKVSLAFGLSI...



Protein families

- Amino acid alphabet
 - {A,R,N,D,C,Q,E,G,H,I,L,K,M,F,P,S,T,W,Y,V}
- Protein sequence
 - >AQP1_bovin MASEFKKKLFWRAVVAEFLAMILFIFISIG-SALGFHYPIKSNSQTTGAVQDNVKVSLAFGLSI...
- Protein data set
 - >AQP1_bovin MASEFKKKLFWRAVVAEFLAMILFIFISIG-SALGFHYPIKSNSQTTGAVQDNVKVSLAFGLSI...
 - >AQP2_rat MWELRSIAFSRAVLAEFLATLLFVFFGLGSALQ-WASSPPSVLQIAVAFGLGIGILVQALGH...
 - >AQP3_mouse MGRQKELMNRCGEMLHIRYRLLRQALAE-CLGTLILVMFGCGSVAQVVLSRGTHGGFLT...



Protein family & Natural selection properties

- Common function
- Common topology (3D structure)
- Common signature

Pattern of the zinc finger protein family

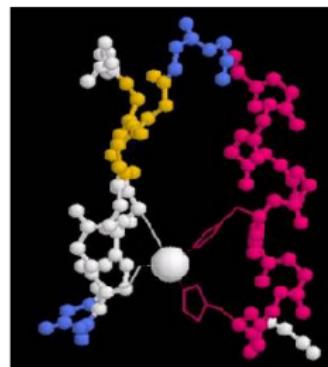
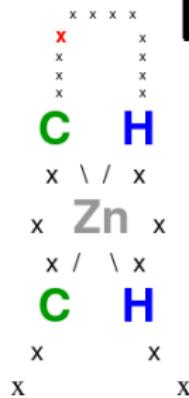
ZBT11 ...**C**si...**C**grt**L**pkl...**H**mlk...**H**...

ZBT10 ...**C**di...**C**gkl**F**trrehvkr**H**slv...**H**...

ZBT34 ...**C**kf...**C**gkk**Y**trkdqley**H**irg...**H**...

Zinc Finger Pattern

C-x(2,4)-C-x(3)-[LIVMFYW]-x(8)-H-x(3,5)-H



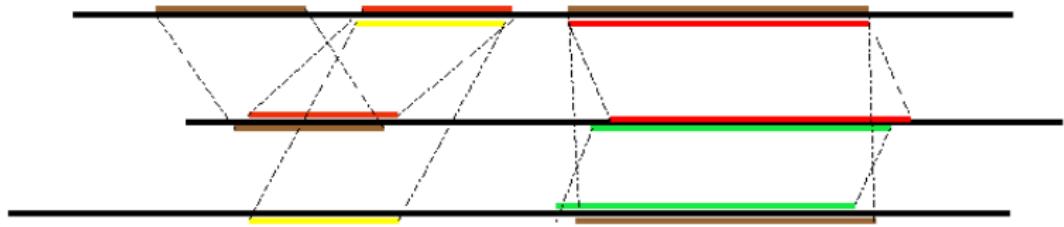
Characterization

Significantly similar areas

- The first phase of the approach is based on the similarity in protein sequences.
- The problem is how to describe and modelize the significantly similar areas.
- And how to optimize the process of identification.

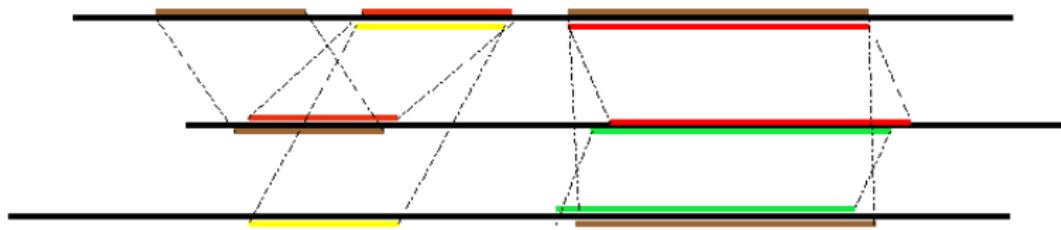
Similar Fragment Pairs Score

Significantly Similar Fragment Pair (SFP) :



Similar Fragment Pairs Score

Significantly Similar Fragment Pair (SFP) :



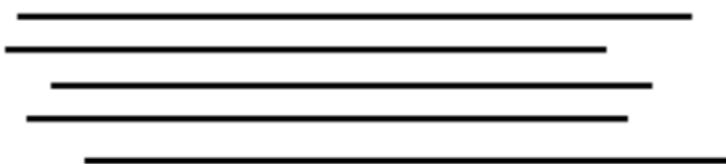
SFP : (f_1, f_2) s.t. $w(f_1, f_2) > 0$ DIALIGN [Morgenstern et al...]

$w(f_1, f_2)$: weight of a fragment pair = $-\log P(s, l)$

$P(s, l)$: probability for a random fragment pair of length l to have a similarity greater than s

s : similarity of (f_1, f_2) , l : length of f_1 and f_2

From a set of protein sequences P
... to a graph of fragments G

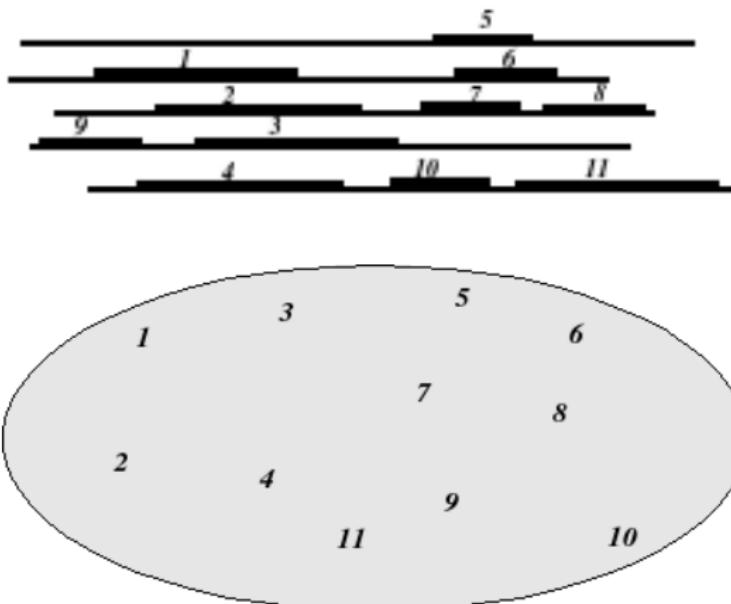


For any sequence S from P

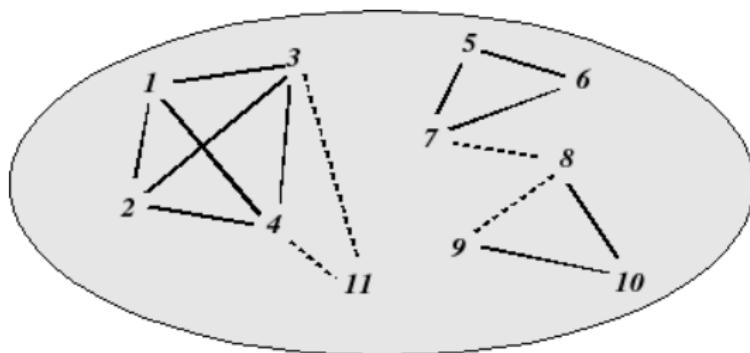
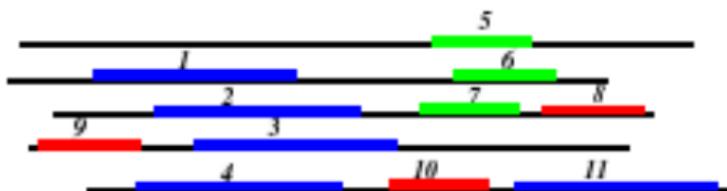
- We can break up S in $\frac{|S|(|S|+1)}{2}$ fragments of particular positions and length.



The graph G where nodes are the fragments of P



The graph G where edges shows the similarity level between fragments



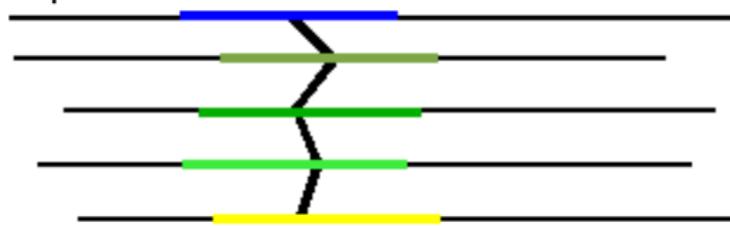
PLMA block

Definition

A PLMA (Partial Local Multiple Alignment) block is a significantly similar set of fragments that testifies a characteristic area in a set of biological sequences. It's a connected subgraph of the graph of fragments.

Weak or Strong consensus ?

- A pathwise connected involves a weak consensus by transitivity



Weak or Strong consensus ?

- A pathwise connected involves a weak consensus by transitivity

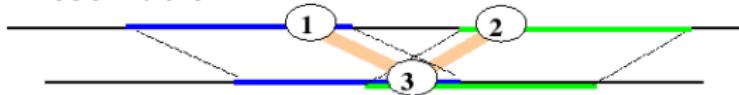


- A clique involves a strong consensus



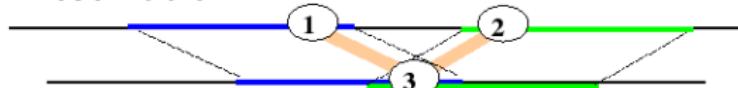
Incompatibilities between SFPs

- Preservation

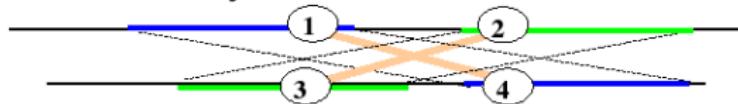


Incompatibilities between SFPs

- Preservation

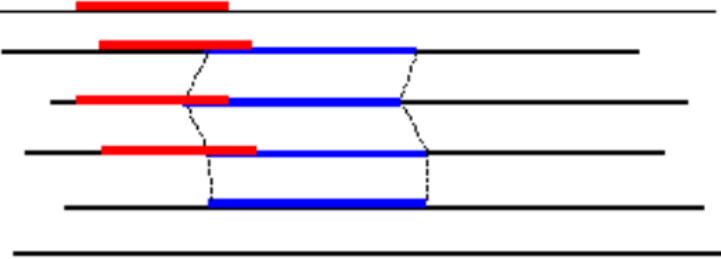


- Inconsistency



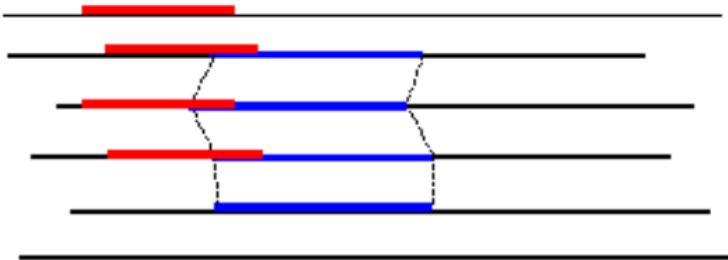
Incompatibilities between PLMA blocks

- Interference

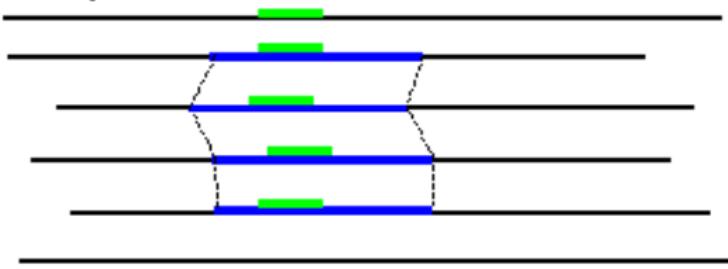


Incompatibilities between PLMA blocks

- Interference

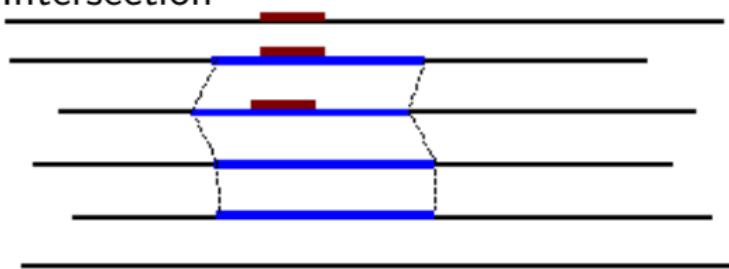


- And yet



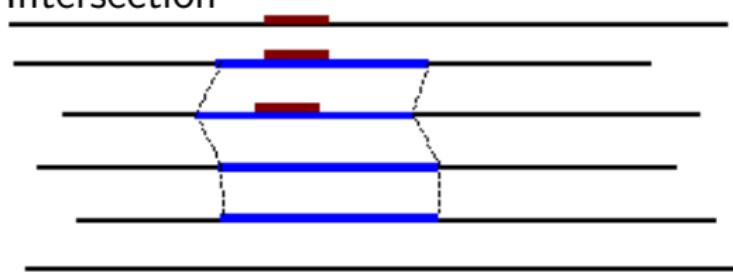
Incompatibilities between PLMA blocks

- Intersection

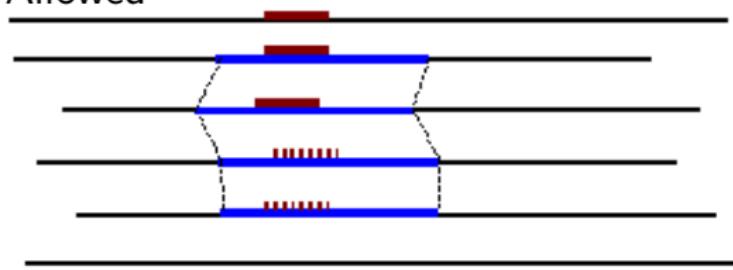


Incompatibilities between PLMA blocks

- Intersection



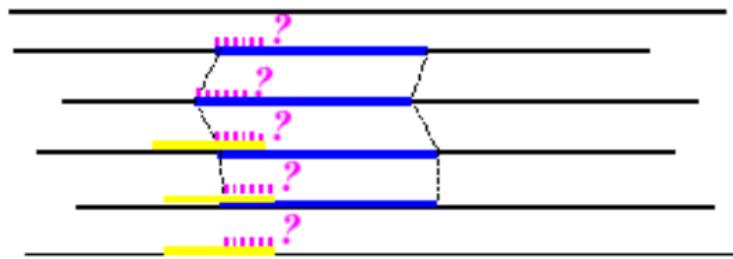
- Allowed



Preservation property of a PLMA block

Invariant

For each pair of aligned positions $P1$ and $P2$,
There exists an SFP that aligns $P1$ and $P2$.

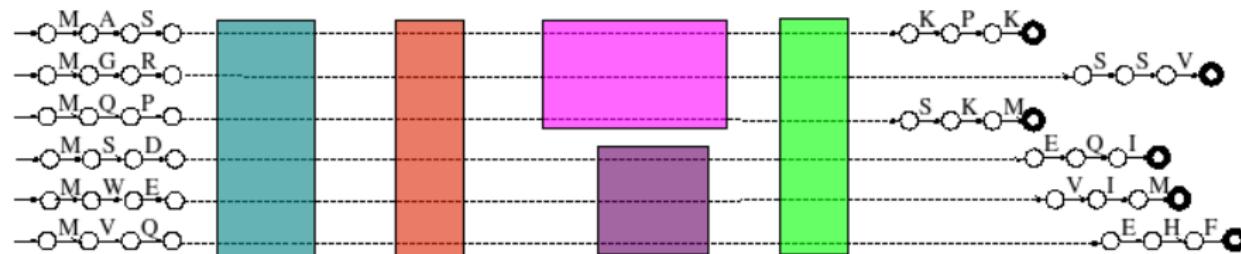


Generalization

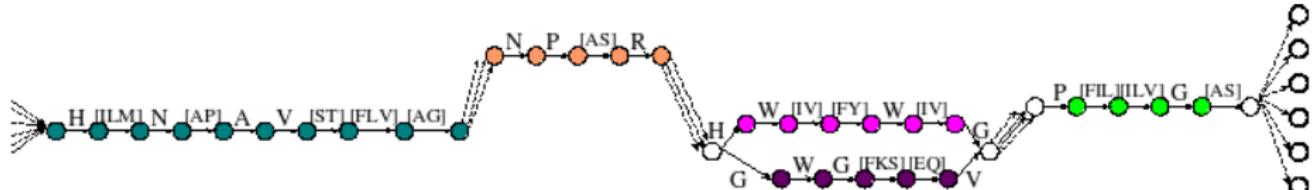
PLMA blocks merging

Characterization stage → Partial Local Multiple Alignments (PLMA) :

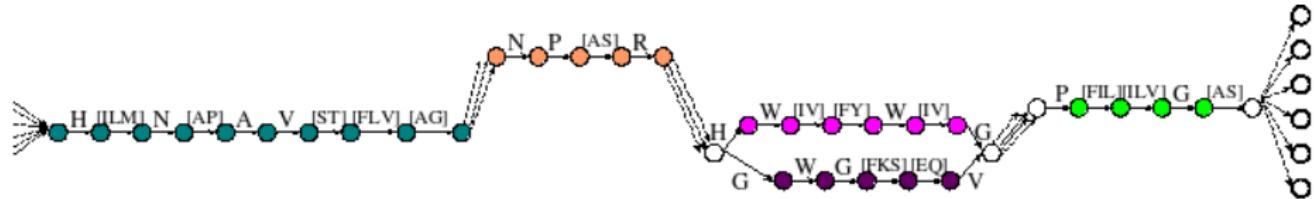
Maximal Canonical Automaton



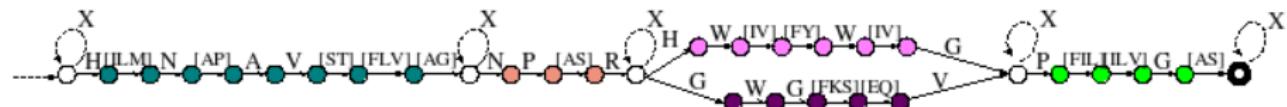
Merging PLMA blocks :



Gap generalization



Transitions used less than a given quorum are treated as gaps :
(if not in an exception path)

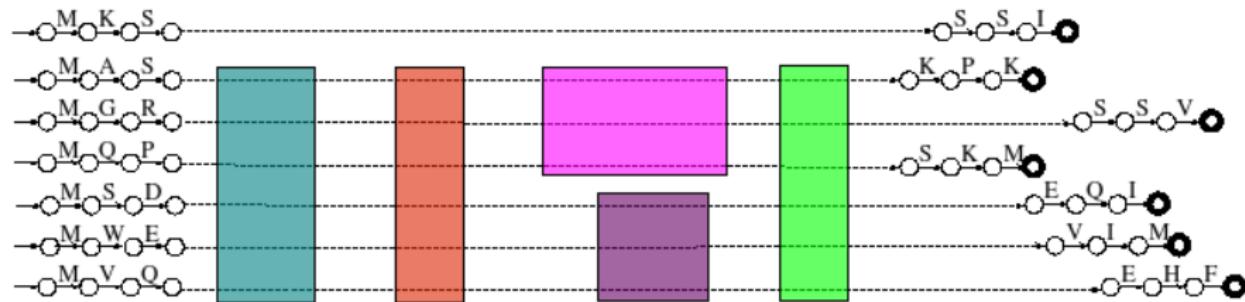


Representative segments : sub-paths used more than quorum

What if?

Characterization stage → Partial Local Multiple Alignment (PLMA) :

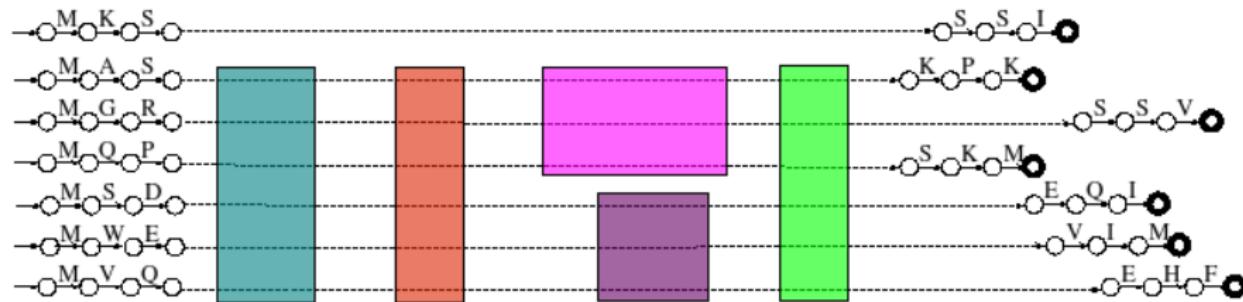
Maximal Canonical Automaton



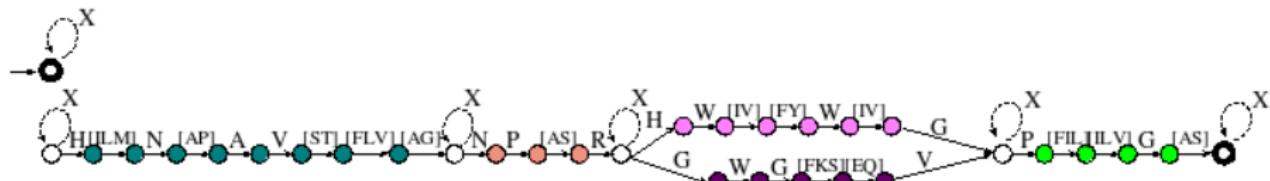
What if?

Characterization stage → Partial Local Multiple Alignment (PLMA) :

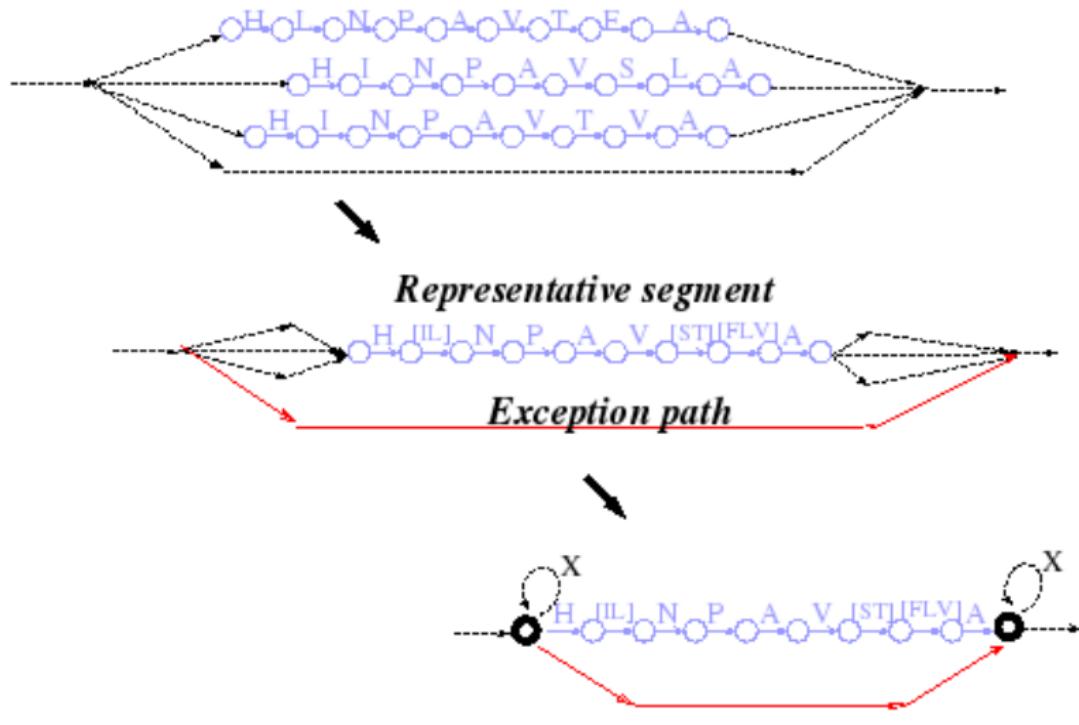
Maximal Canonical Automaton



Merging PLMAs :



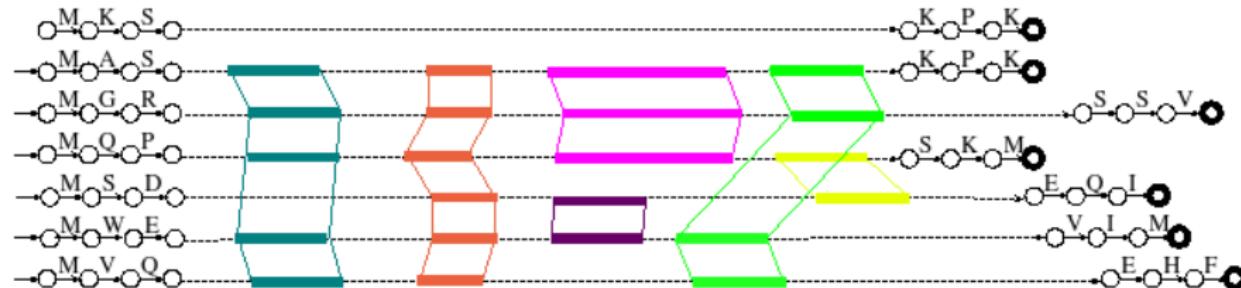
Exceptions



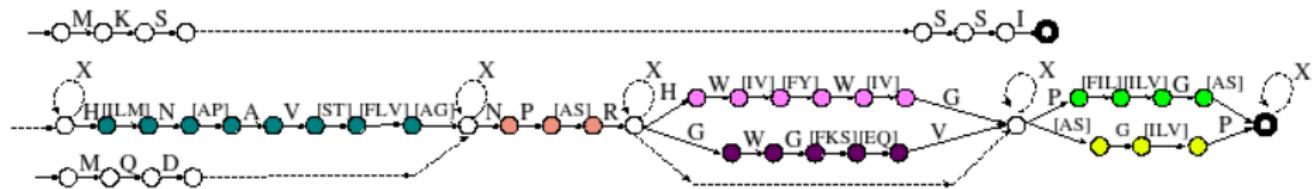
Exceptions

A more realistic result of the characterization stage :

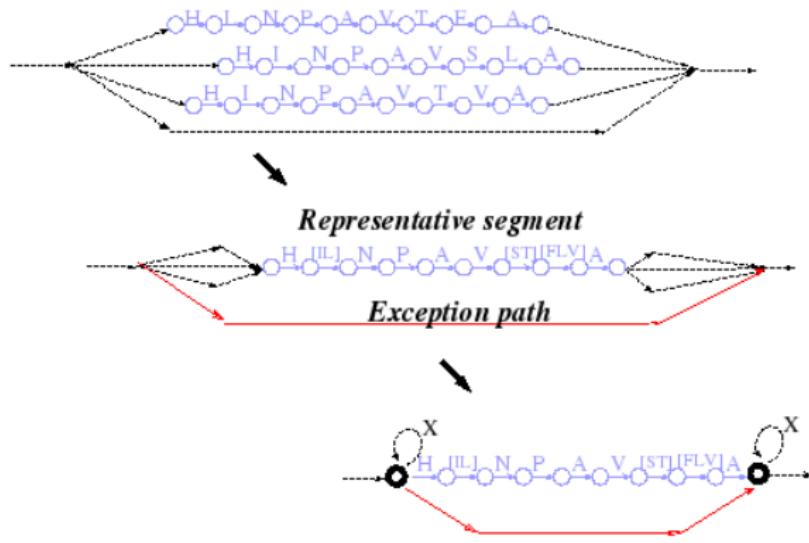
Maximal Canonical Automaton



Merging PLMAs, exception detection, gap generalisation :



Exceptions

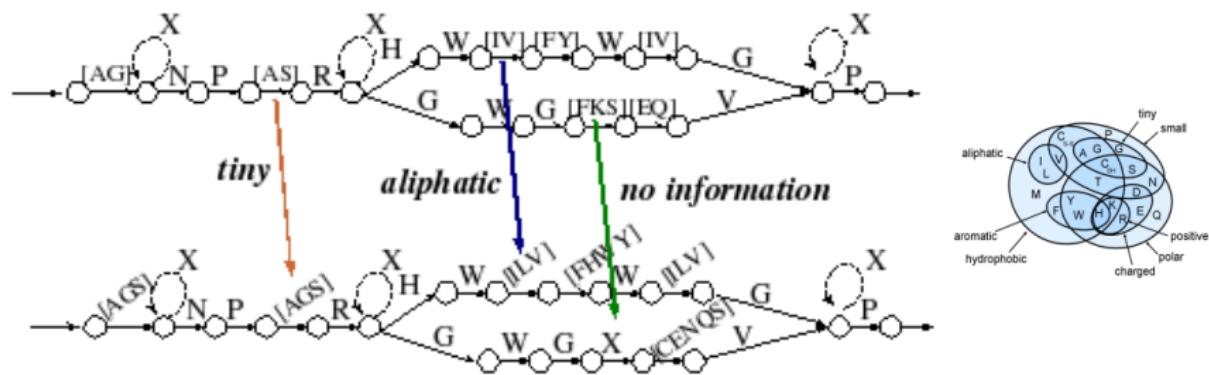


Detection

Transitive closure over accessibility graph between representative segments. Exceptions are edges between two representative segments s.t. a path of length ≥ 2 between them also exists.

Identification of informative positions

- Representative segments : average conservation of sub-sequences (BLOSUM)
- Identification of actual important *physico-chemical properties*



Likelihood ratio tests...

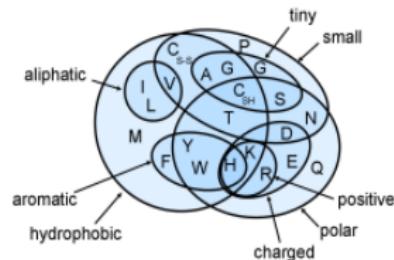
Testing expansion by likelihood ratio test

P : multiset of amino acid at a given position

G : smallest group containing aa of P

Σ : amino acids alphabet

$\lambda_G, \lambda_\Sigma$: thresholds



- Expansion $P \rightarrow G$

$$\text{Reject if } LR_{G/P} = \frac{L_G}{L_P} = \left(\frac{\sum_{a \in P} p_a}{\sum_{a \in G} p_a} \right)^n < \lambda_G$$

- or else Expansion $P \rightarrow \Sigma$

$$\text{Reject if } LR_{\Sigma/P} = \frac{L_G}{L_P} = \left(\sum_{a \in \Sigma} p_a \right)^n < \lambda_\Sigma$$

Demonstration

Protomata Learner Parameters - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils Aide

http://protomata-learner.genouest.org/ Google

Plate-forme Bio-informatique Genouest

Protomata-Learner version 0.07 (warning : Beta version)

Protomata-Learner documentation

Examples data

Description of your request no_named_request

Email (optional)

Partial Local Multiple Alignments (PLMA)

Sequences

Advice : reduce the redundancy in your set of sequences, for example use the program decrease redundancy at Expasy.

Upload your sequences file (FASTA) Parcourir...

Or paste your sequences

This screenshot shows the Protomata-Learner web application running in Mozilla Firefox. The title bar indicates the page is 'Protomata Learner Parameters'. The menu bar includes 'Fichier', 'Édition', 'Affichage', 'Historique', 'Marque-pages', 'Outils', and 'Aide'. The address bar shows the URL 'http://protomata-learner.genouest.org/'. Below the address bar is a navigation toolbar with icons for back, forward, search, and refresh. The main content area features a logo for 'Plate-forme Bio-informatique Genouest' with a blue gradient background. A banner below the logo reads 'Protomata-Learner version 0.07 (warning : Beta version)'. There are two links: 'Protomata-Learner documentation' and 'Examples data'. Below these are two input fields: one for 'Description of your request' containing 'no_named_request' and another for 'Email (optional)' which is empty. A section titled 'Partial Local Multiple Alignments (PLMA)' follows. Within this section, there is a 'Sequences' heading and a note about reducing redundancy. An 'Upload your sequences file (FASTA)' input field is present with a 'Parcourir...' button. Below it is a text area for pasting sequences with the placeholder 'Or paste your sequences'. The bottom of the page contains footer information: 'F. Coste, G. Kerbellec (IRISA)', 'Protomata-Learner', 'Rennes, 3 décembre 2007', and '31 / 53'.

Viral nucleic acid binding Clan

Carla_C4.fas

```
>Q02123|VNBP_P0PMV  
MVNNRNCNPFCDIAVISIVCRSELDFINPESLNSYKRRRARRL0  
RCVRCCFVNM[GFT]TICRDGIVTCVPSIWSWVNDVEYLIRGRVTDRETSPSTHGYGPV6  
HKT  
>E6654|VNBP_PVSP  
MKAERELMLLLCVYRLGTYLPPWDCKIKLISVAQWSVQGRSTSYCKRRARSIGRCWRCYRV  
YPPVCVNSC[CONTCG]CIPSMWPKVVVTFIRGWSN  
>Q00574|VNBP_PLV9S  
MDKRNKANALWSLC-SMFAFSRGNCIPIPIVFNIMYMRAPPKLVGRCTTAYRRRARSILRCL  
EPLV9S[PMLP]F5YKQDRCITCPV01SNTKVAQFIVWVIVTEVPHGPNF#  
>P27336|VNBP_PLV9S  
SVWGAPEWVTTGTCALKSSELIIDTQIMDEALKRRTTIVLCLLSAPPDRCIDILRRTS  
SHIVOLRSRYAANRNLQIOTCRCTCYRVYPWCGSKC[CONTCG]CPLSINTNVANVYIDH0  
VIVVIPWSPHGGQFLRPLR  
>O16871|VNBP_PLV9S  
MDKVTKWALLIARACMTSGCTVPELAFSIAECAGRPGLGGRSKYARRRAISTARCHRC  
YRLWPPZPVEETIRODKNTCYCWP01SNTVNRVAQFIDEVGTEVIPSVINRE  
>P37982|VNBP_C1V  
MDIVIVVNLILRKFVEQGNVCP1HLQVDV1KXRAFPRSVNGKRSYARRRRAILEGRCHRCY  
RVVPPLEPEISRCRDNTCPV01SNTNSKVRDYLWVGTEVIPHGPYME#  
>P22625|VNBP_C1V  
MREKRLRKQLLEDLFRKFPAFQSVGHGSDCINII1AKIKSDGEKSYARRRRAKSIAACPRC  
ARVCKQPYTITRCRDKTGTGDLARSPDULLEPIGIDLCVRSK  
>O38024|O38024_9VIRU  
MKCVQGALLVARALYLSHGTVFYLSTIRAGRPGLGGRSKYARRRAIAAGRCHRC  
YRLWPPZPVEETIKCNRSPCVP01SYNRVSAASR  
>O09491|O09491_9VIRU  
MFSKTQPRESMI[KR]YRRLRAIFKLHDKNCVNDV01IVNKIVCORTGASKYRAR  
RAKS10RCPRCPRCAPOGYTPVNCNDTCTCTP01SYNEVKVNFIW0VTM  
>O41484|O41484_9VIRU  
MKADRLATLLCVRHLGVYLTVCGVNIILSAAPIGSRSTYORRRAARSIGRCWRCYR  
VYPPICNSKCCDNTRCP01SPNVMVIRGWSN
```

u---XEmacs: Carla_C4.fas (Fundamental) ---All---

CTV_P23.fas

```
>Q66245|Q66245_9CLOS  
MDNTSGQTFSVNLSDENSTATTIVEPVSSEADRDLFLQKMNPII1DALIRKNSYQGARF  
RAR10IVCVD[CRKHKD]KALAKTERKKVKNNTUSQNEVAHMLMDPVKYLNRKRAFNSAE  
IFAI1DVLVHMTKEROLAIADAAEREKTRLARRHPSPEETPEYYKFONTAKMLPDINAV  
DVGDNEDTSSEYPSVLSVSG0VLRHEHFI  
>Q66J80|Q66J80_9CLOS  
MDNTSGQTFSVNLSDENSTASTRVERVNSSEADRFLRFLQKMNPII1DALIRKNSYQGARF  
RAR10IVCVD[CRKHKD]KALAKTERKKVKNNTUSQNEVAHMLMDPVKYLNRKRAFNSAE  
MFA1ELVLVHMTKEROLAIADAAEREKTRLARRHPSPEETPEYYKFONTAKMLPDINAV  
DVGDNEDTSSEYPSVLSVSG0VLRHEHFI  
>Q66258|Q66258_9CLOS  
MDNTSGQTFSVNLSDENSTASTETKAVSSEADRFLRFLQKMNPII1DALIRKNSYQGARF  
RAR10IVCVD[CRKHKD]KALAKTERKKVKNNTUSQNEVAHMLMDPVKYLNRKRAFNSAE  
MFA1ELVLVHMTKEROLAIADAAEREKTRLARRHPSPEETPEYYKFONTAKMLPDINAV  
DVGDNEDTSSEYPSVLSVSG0VLRHEHFI  
u---**-XEmacs: CTV_P23.fas (Fundamental) ---All---
```

Viral_NABP.fas

```
>Q04580|VNBP_SHVX  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYNRCFCDFGAYLYDDHHVCKRFTSRSNSDCLSVIRHQGPAKLYAEGAYRANSDAEQLIMND  
LLKSLKL  
>O67698|O67698_9VIRU  
MHTYVANFLACQFAPQNLPSDWRISIYMLSSASRKIGRKSQQNKPPTGTSKCAAPRRAK  
RYRCPDQGALLNTDHVCKLFTSRASTDCLHVIREOPAKLYAERKTFRKSSFAEQLILDD  
ELMLVLE  
>O67702|O67702_9VIRU  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYRCPDQGALLNTDHVCKLFTSRASTDCLHVIREOPAKLYAERKTFRKSSFAEQLILDD  
LYTAKL  
>O67696|O67696_9VIRU  
MHRDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYRCPDQGALLNTDHVCKLFTSRASTDCLHVIREOPAKLYAERKTFRKSSFAEQLILDD  
LYTAKL  
>O67694|O67694_9VIRU  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYNRCFCGAYLYDDHHVCKRFTSRSNSDCLSVIRHQGPAKLYAEGAYRANSDAEQLIMND  
QYKLFQNRKA  
>O67662|O67662_9VIRU  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYNRCFCGAYLYDDHHVCKRFTSRSNSDCLSVIRHQGPAKLYAEGAYRANSDAEQLIMND  
ELMLVLE  
>O67661|O67661_9VIRU  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYNRCFCGAYLYDDHHVCKRFTSRSNSDCLSVIRHQGPAKLYAEGAYRANSDAEQLIMND  
VTKL  
>O67660|O67660_9VIRU  
MHPDPDNLCLCPKSPNPLDNLKTLPRACETSCKLNRLNLDNPFF07SKAACKRRAK  
RYNRCFCGAYLYDDHHVCKRFTSRSNSDCLSVIRHQGPAKLYAEGAYRANSDAEQLIMND  
RCRCONCARWFHRDOPCLHORPDYSLQAPPDPHQLNSFEPILLALASVLRPLPRD  
IQITIISIACADFVHSVCASSRYLGSSRSASVKKRAARLNVYCKC00HPLVLYNKPHCTRCPG  
LCSASISERLALLREGPISLTELNPINARAHTLAHELLOPR
```

u---XEmacs: Viral_NABP.fas (Fundamental) ---All---

Protomata Learner Parameters - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils Aide

http://protomata-learner.genouest.org/ Google

Description of your request: Viral_nucleic_acid_binding

Email (optional): goulven.kerbellec@irisa.fr

Partial Local Multiple Alignments (PLMA)

Sequences

Advice : reduce the redundancy in your set of sequences, for example use the program decrease redundancy at Exasy.

Upload your sequences file (FASTA)

Or paste your sequences:

```
>Q01123 |VNP_PODMV  
MVNMRKVLAQMQVFRERYDHKCDPNFCDAVSIVCRSELDFINEPGLENYAKRRRARRLG  
RCVRCFPRVMPGPFYPTKRCDG1TCVPG1SWNYDVEYIKRGRTQDRETPSTPHGYGPVG  
HET  
>P16654 |VNP_PVSP  
MKAERLEMLLLCVYRLGYILPVDVCCIKIISVAQVEVQGRSTYSCKRRARSIGRCWRCYRV  
YPPVCNSKCDNR7CRPG1SPNPKVVFIRGWN
```

Sequence type: Proteins DNA DNA (with coding regions)

Paloma parameters

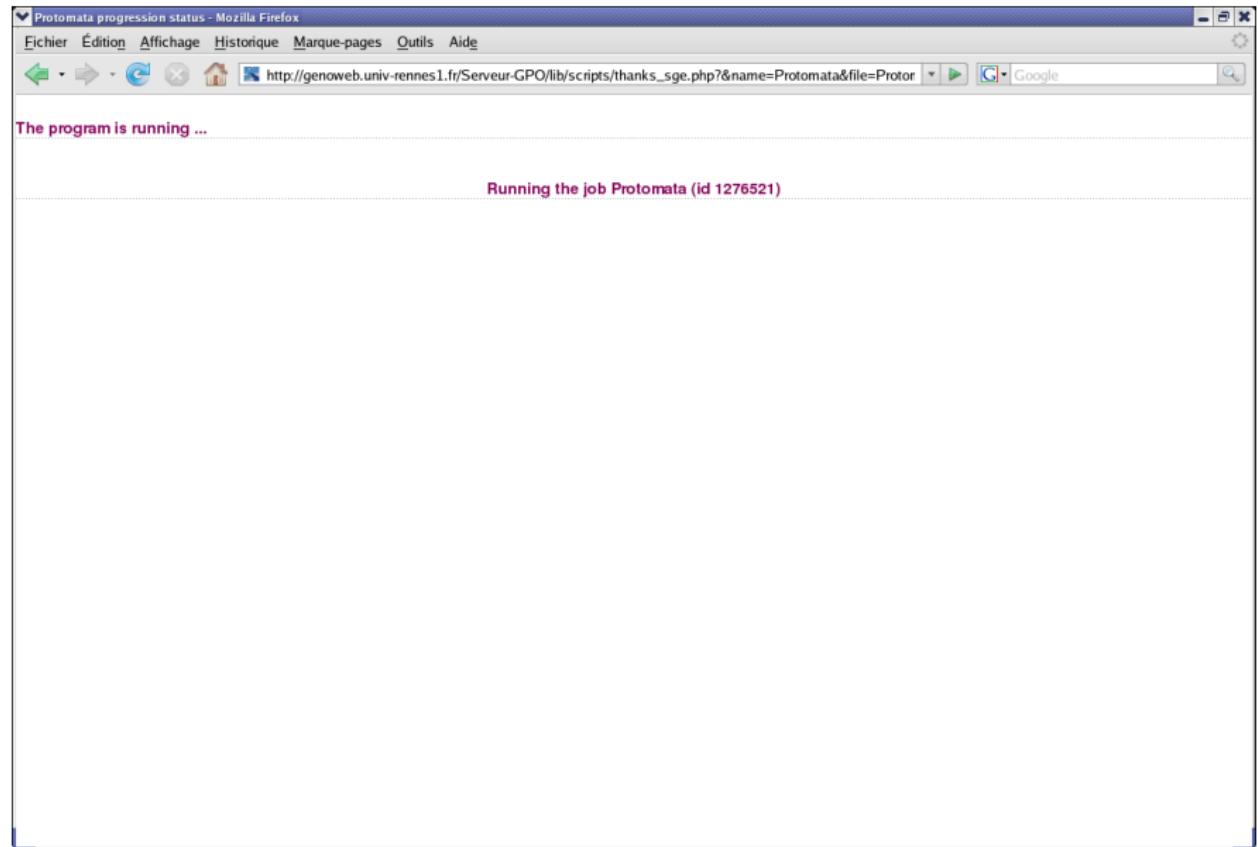
Fragments similarity threshold: 5

Consensus: weak

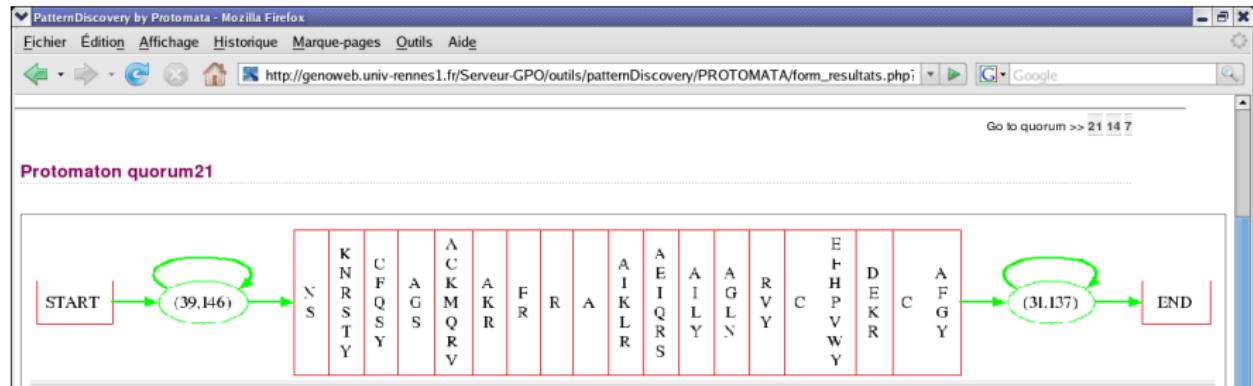
Minimum fragments size: 1

Maximal fragments size: 15

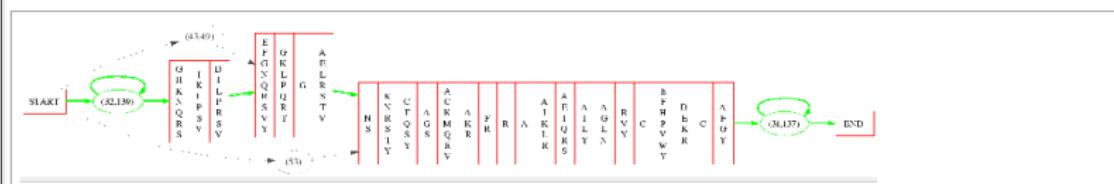
process is running



Results of protomata ordered by quorum



Protomaton quorum14

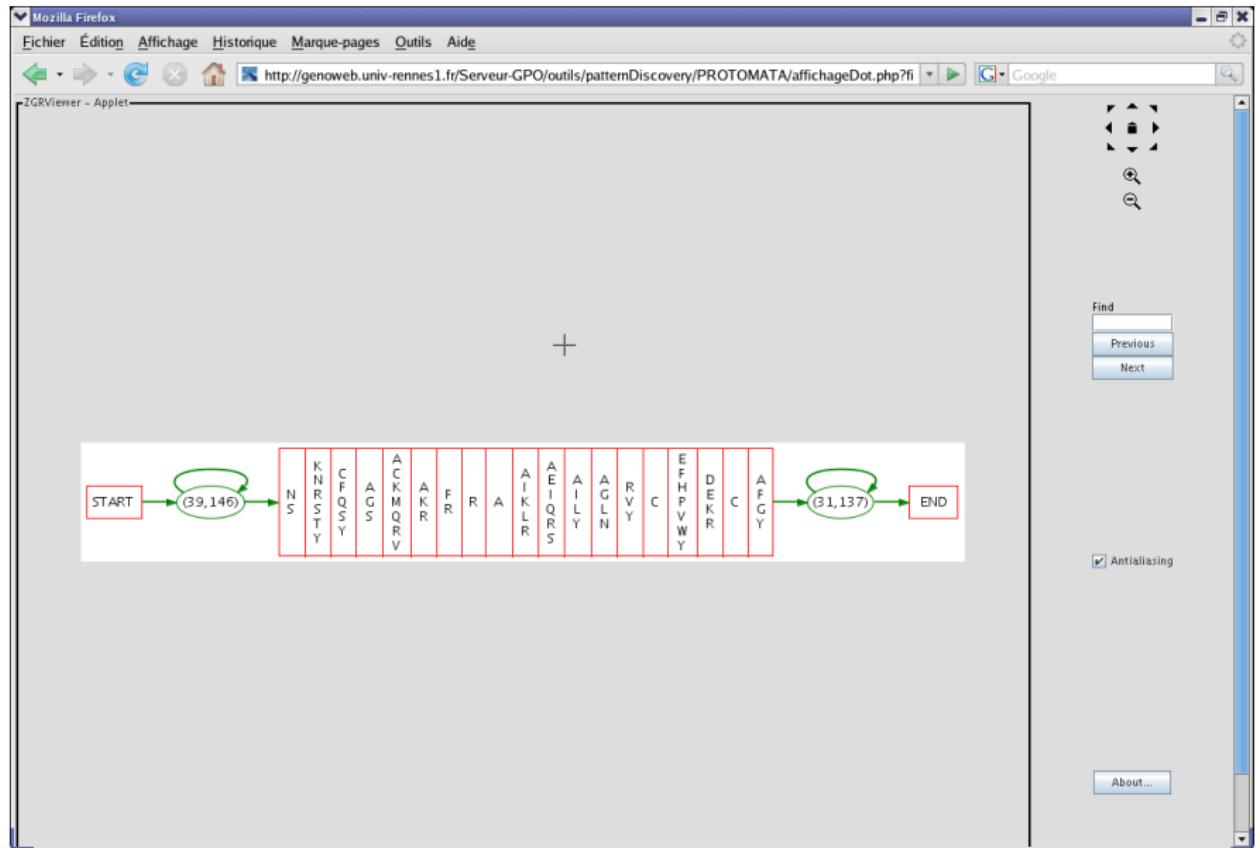


View protomaton : JAVA applet PNG format SVG format | View alignment : JAVA applet PNG format SVG format

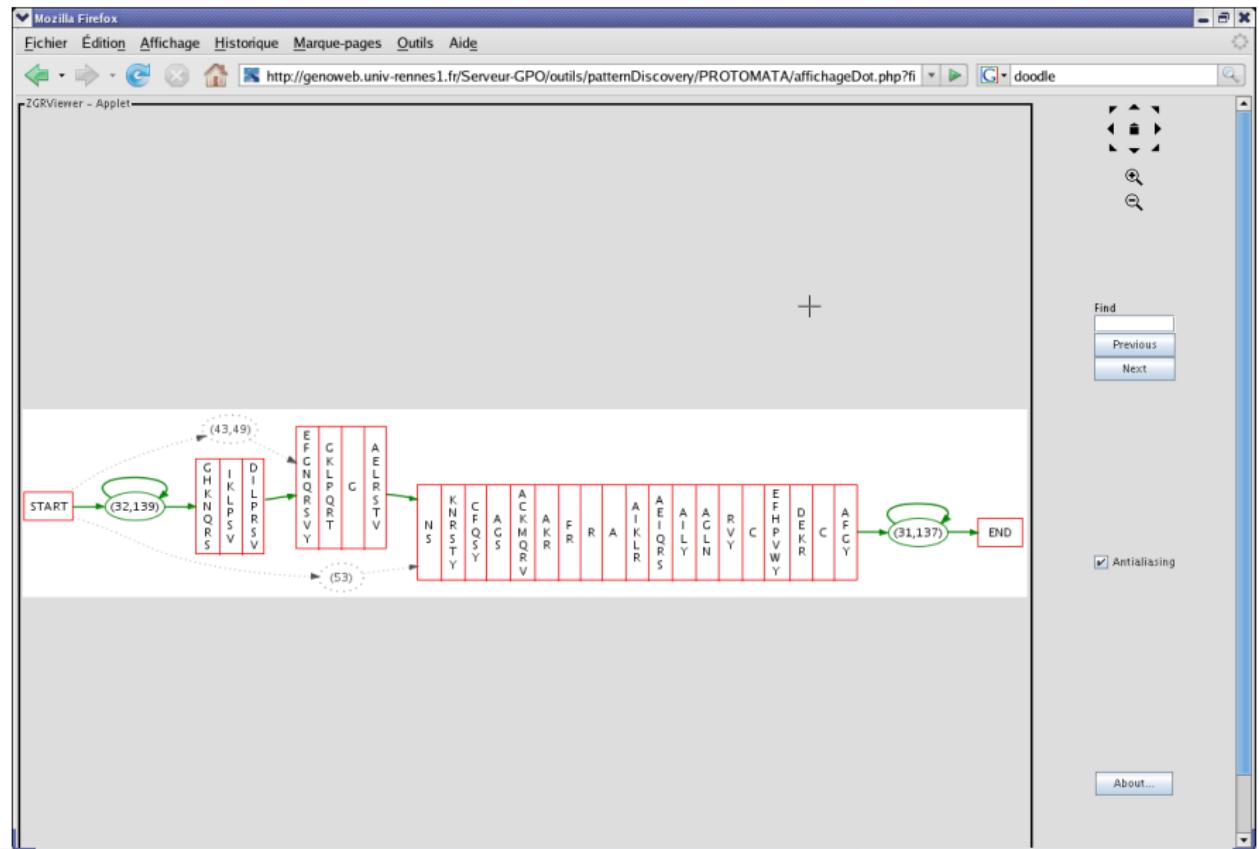
Save protomaton (XML) Scan sequences with protomaton

Protomaton quorum7

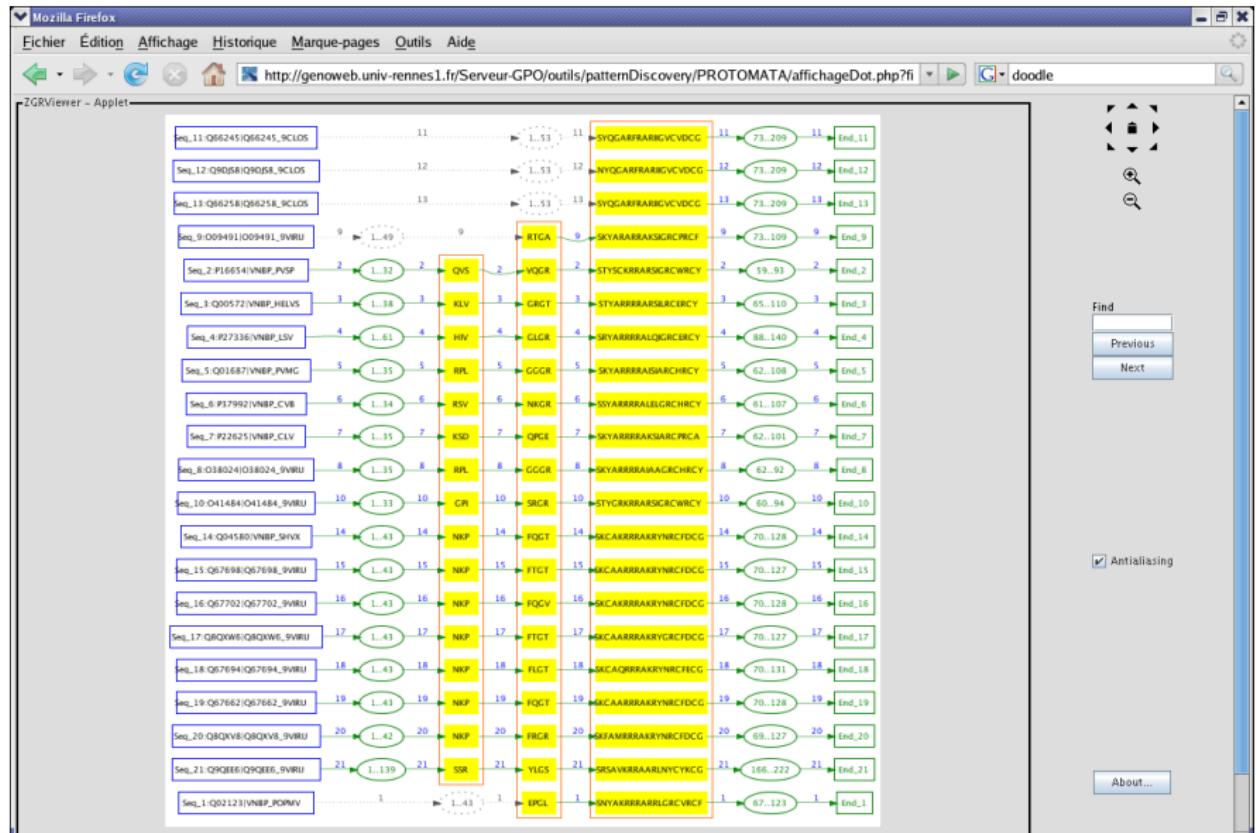
Quorum of 100%



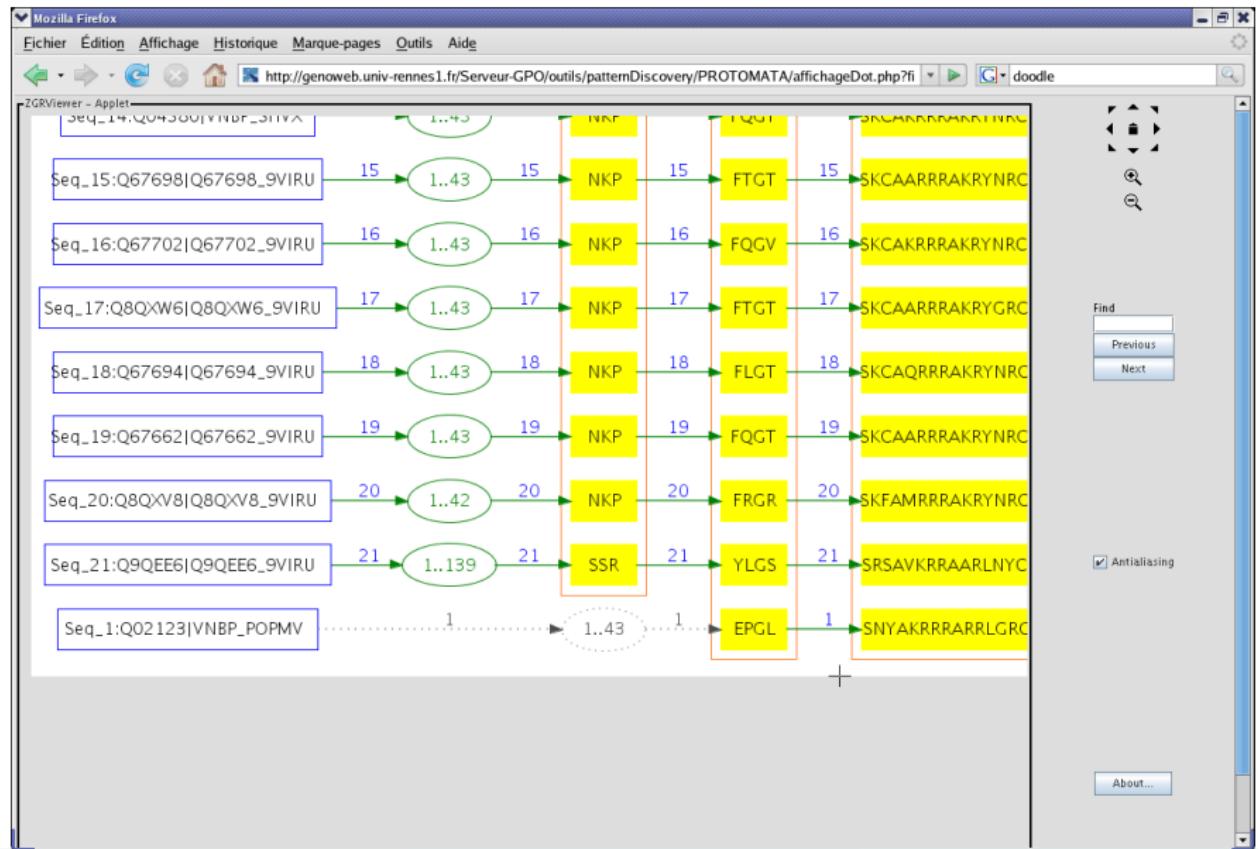
Quorum of 2/3



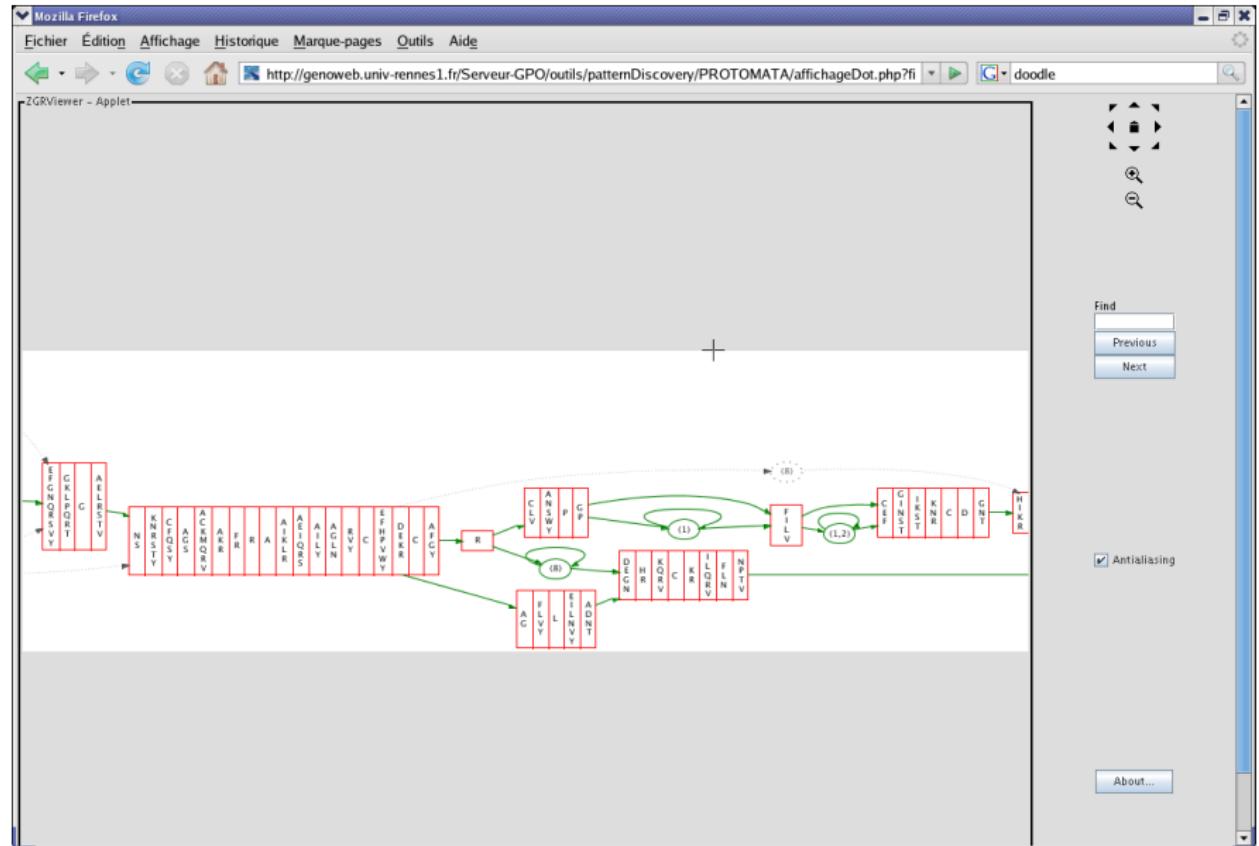
Partial Local Multiple Alignment view



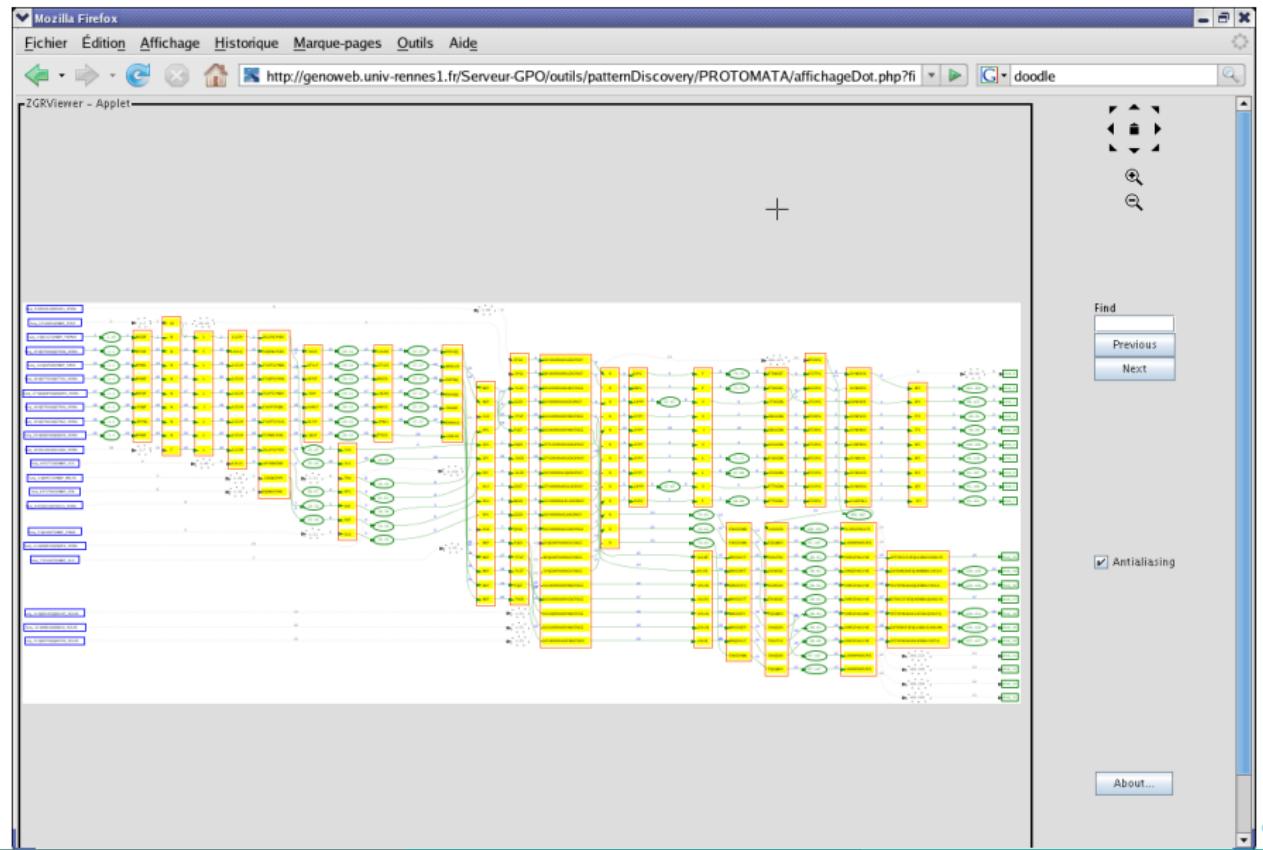
An exception



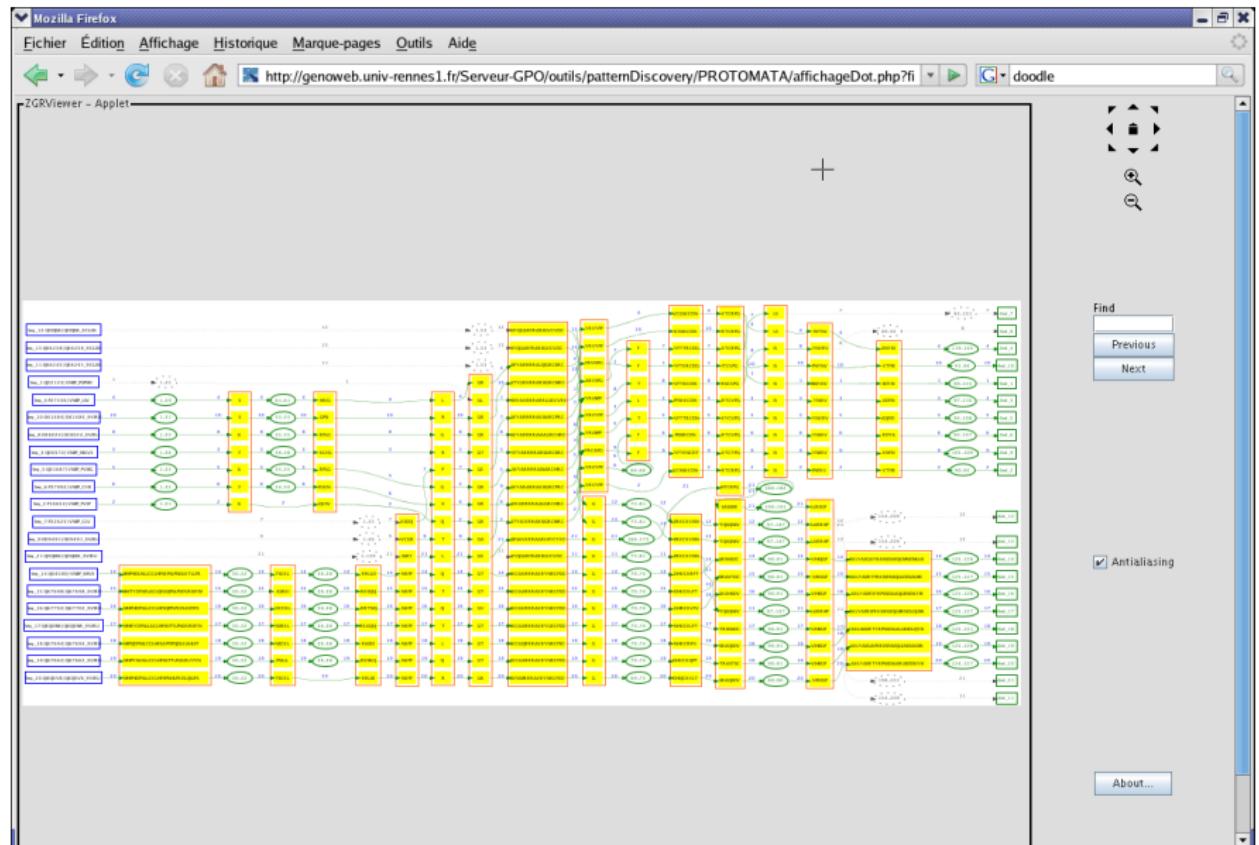
Quorum of 1/3



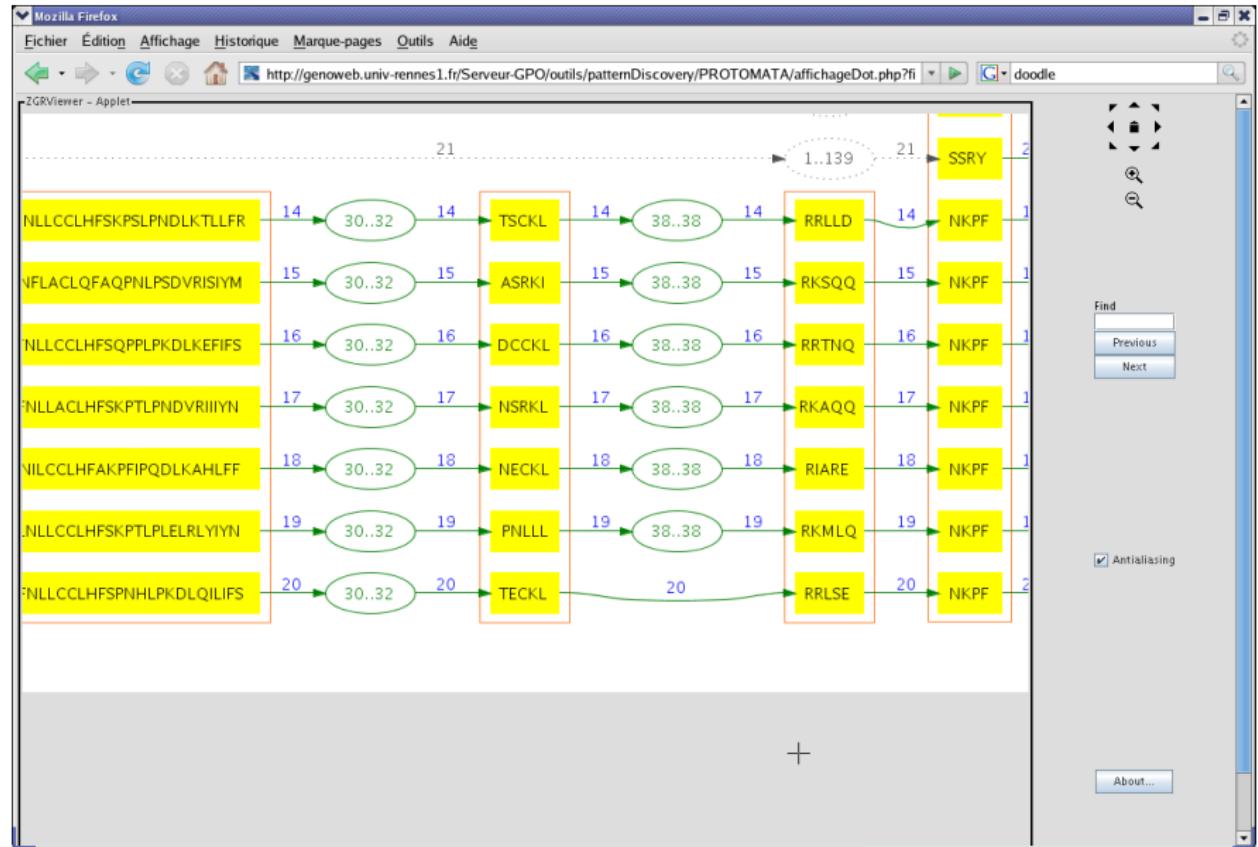
Alignment view of a weak consensus result



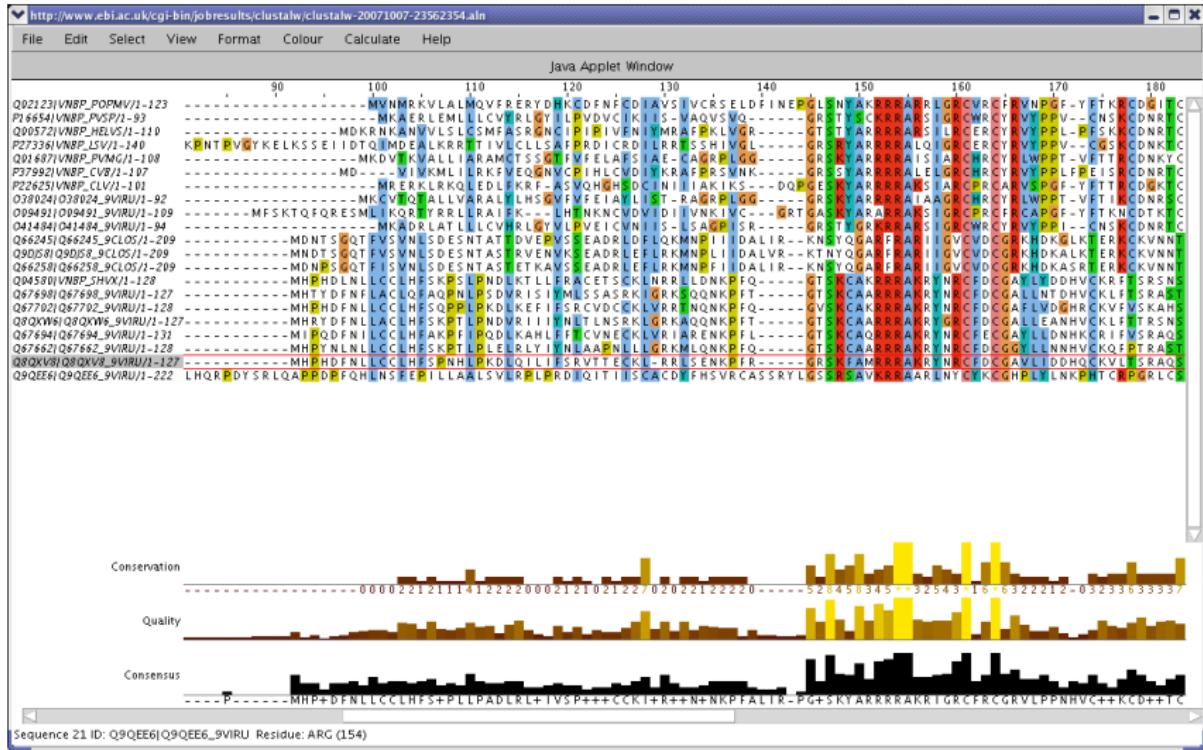
Alignment view of a strong concensus result



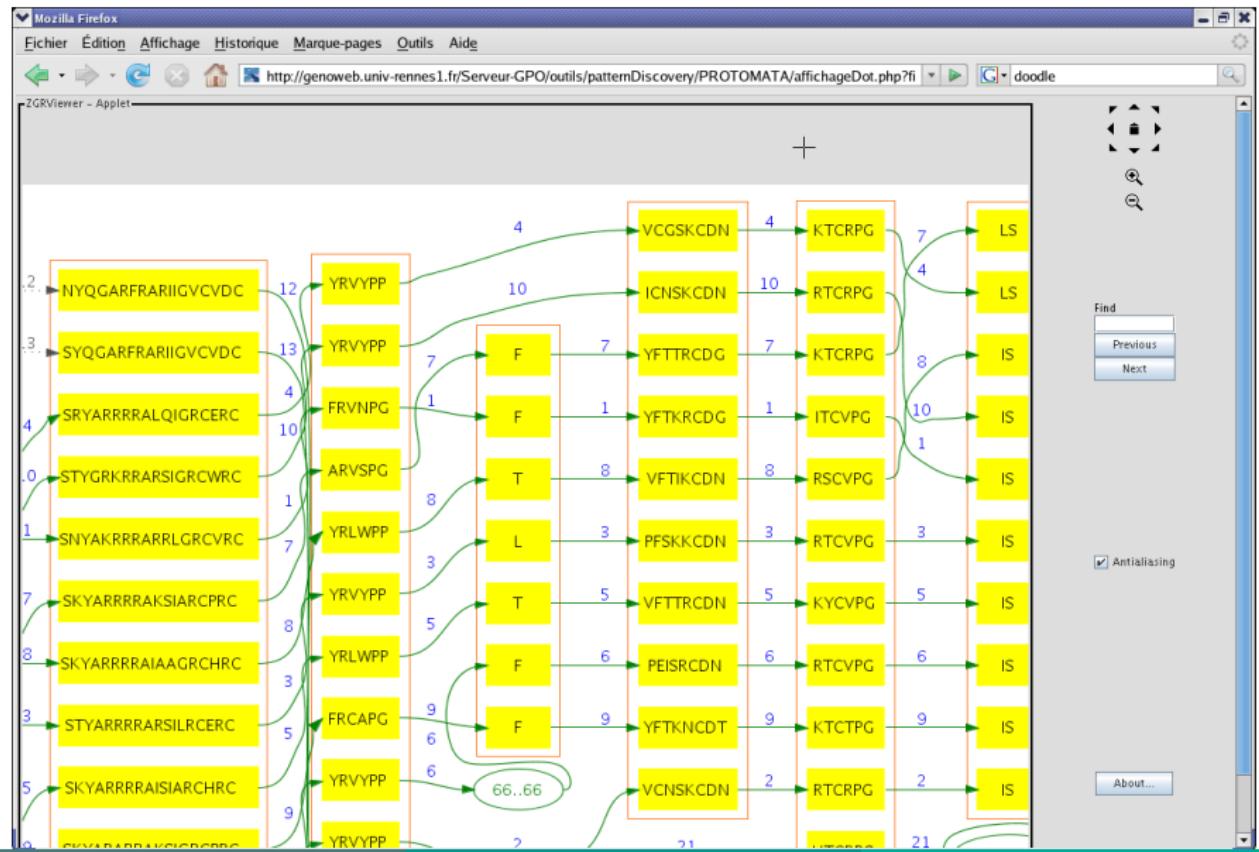
A deletion



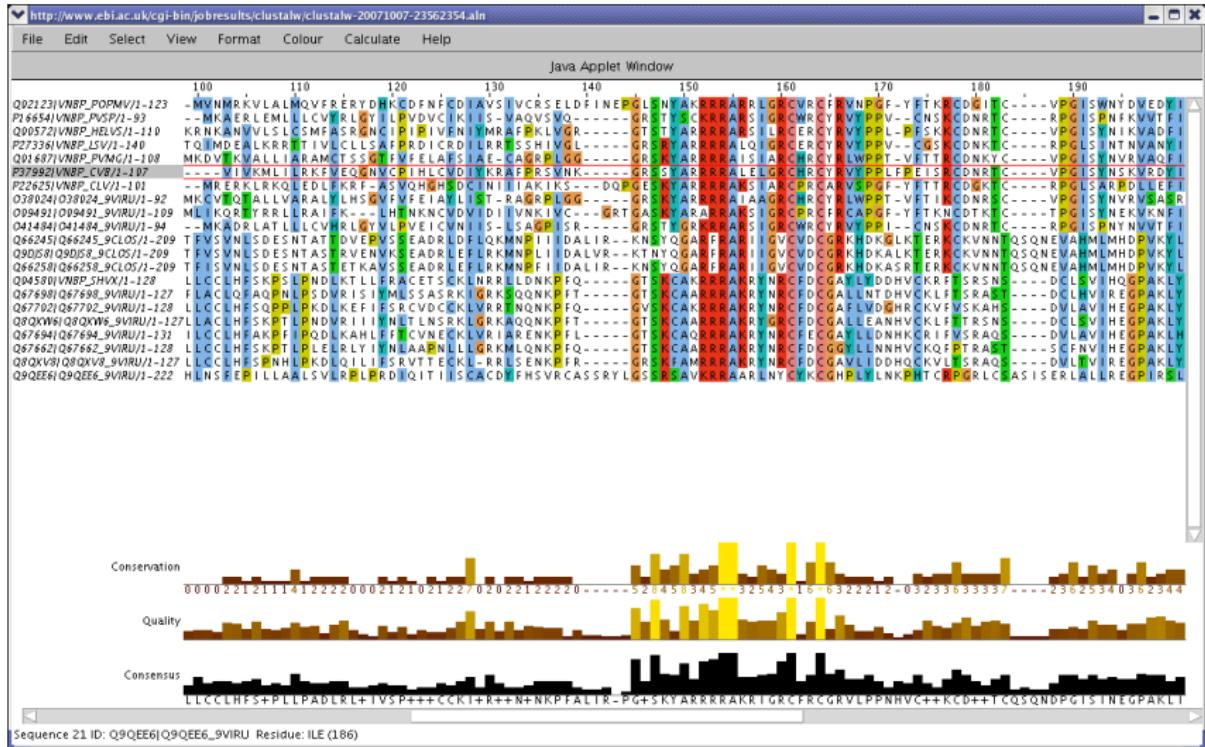
The same in Clustal-W



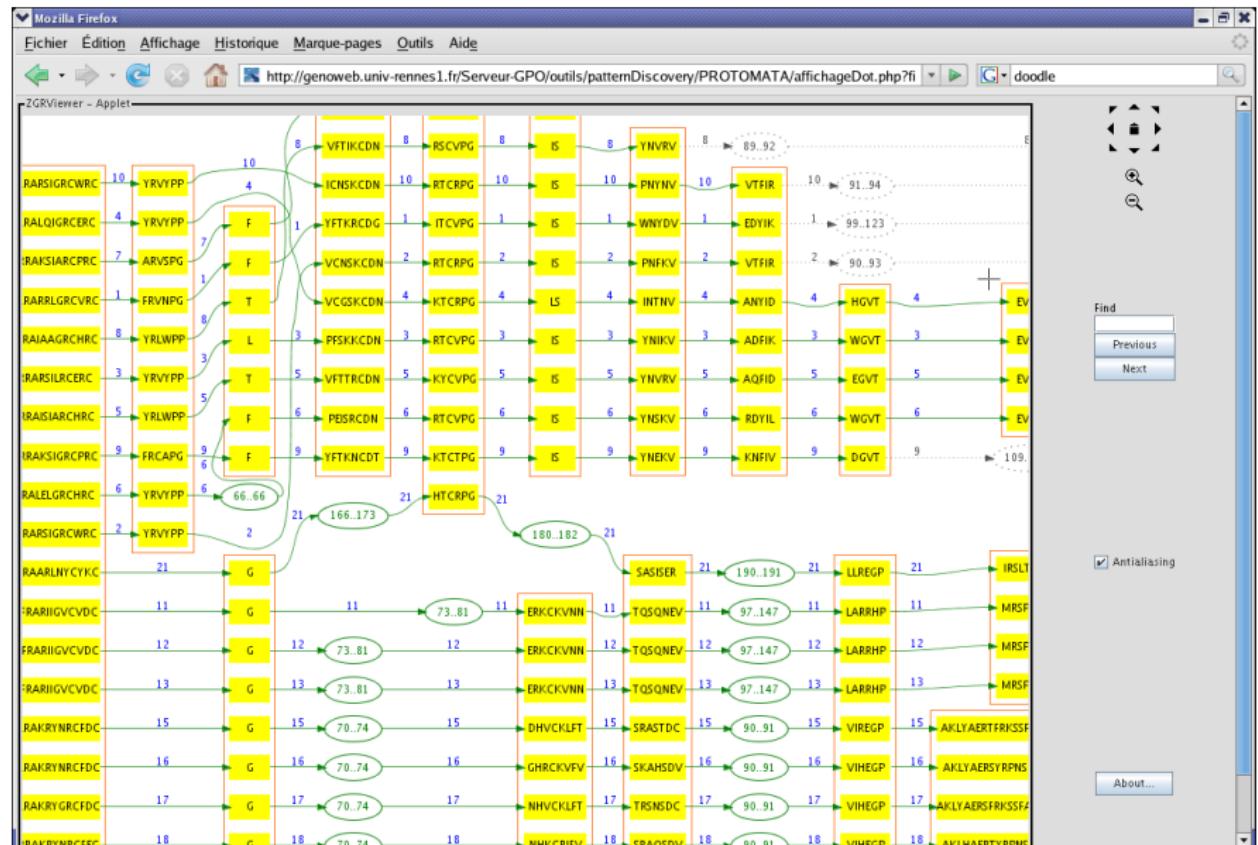
An insertion



The same in Clustal-W



Zinc finger motif



protomata-scan

PROTOSCAN - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils Aide

http://genoweb.univ-rennes1.fr/Serveur-GPO/ouutils/patternMatching/PROTOSCAN/index.php?file=/web

-- PROTOSCAN --

-- BETA VERSION --

Your email (optional) goulven.kerbellec@irisa.fr

Your XML file of protomata

Scores matrix Select a default matrix Blosum62_plus4_inv.mat Upload your matrix file

The sequences

Search within a data base or a Personnal sequences swissprot

Scan

Score threshold default WA threshold Define an other threshold

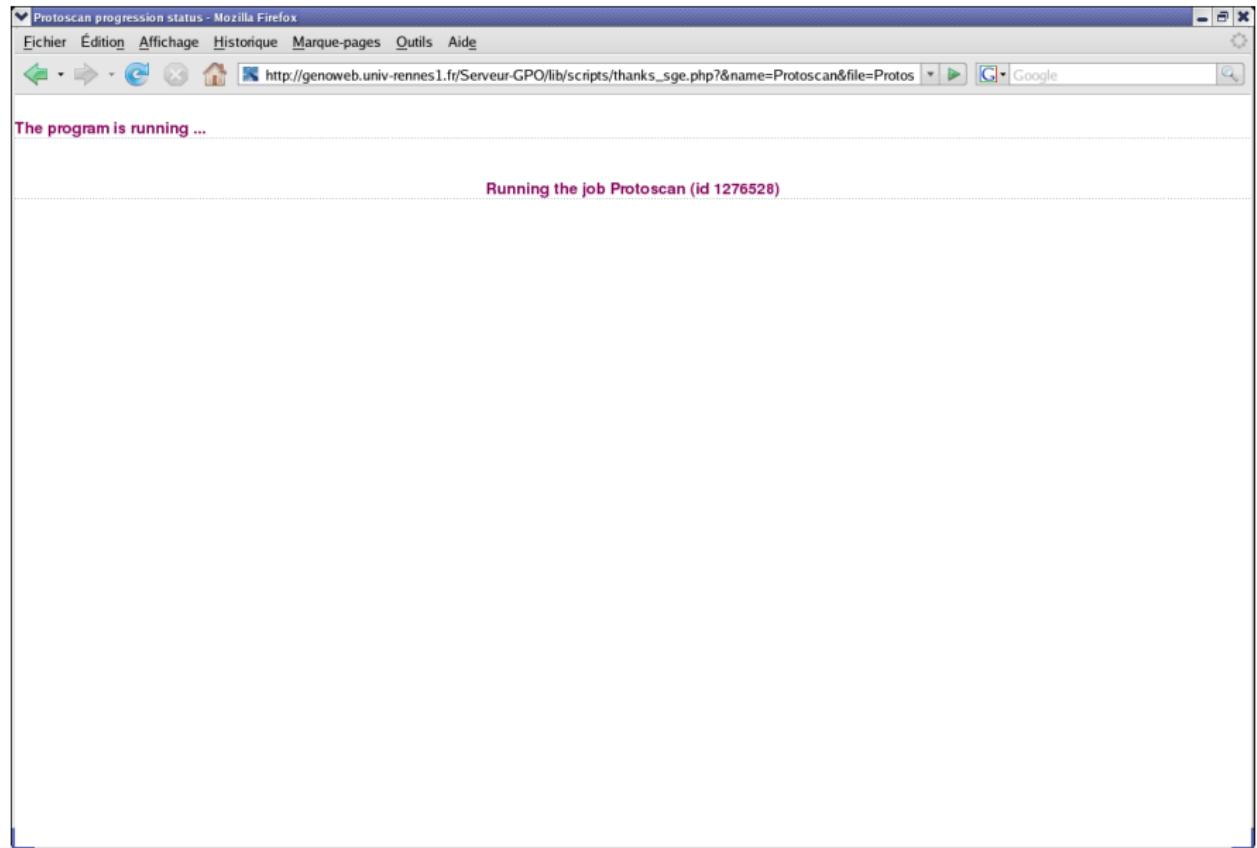
Enter your threshold 20

Start the scan

For any problem or any question don't hesitate to contact the Webmaster

XHTML

scanning databases



results

WAPAM result - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils Aide

Back Results per page : 1500 Next
 Begin Jump to : P14274|CB2A_SOLLC Chic End
 Maximum sequences length : 30

http://geneweb.univ-rennes1.fr/Serveur-GPO/lib/pattern/matching/parser_xml.php

Result 1 to 23 of 23

Num	Chromosome	Strand	F-index/www-tmp-Protocan4591810.wa				
			begin	end	cost	sequence	length
1	P14274 CB2A_SOLLC Chlorophyll a-b binding protein 1A, chloroplast precursor - Solanum lycopersicum (Tomato) (Lycopersicon esculentum)	plus	0	0	0	MAAAAMALSSPSFAG ... NNNAWAFA TNFVPGK	265
2	P14275 CB2C_SOLLC Chlorophyll a-b binding protein 1C, chloroplast precursor - Solanum lycopersicum (Tomato) (Lycopersicon esculentum)	plus	0	0	0	MAAAATMLSSPSFAG ... NNNAWAFA TNFVPGK	265
3	P14276 CB2E_SOLLC Chlorophyll a-b binding protein 3A, chloroplast precursor - Solanum lycopersicum (Tomato) (Lycopersicon esculentum)	plus	0	0	0	MAAATMLSSSTFAG ... NNNAWAFA TNFVPGK	267
4	P14277 CB2F_SOLLC Chlorophyll a-b binding protein 3B, chloroplast precursor - Solanum lycopersicum (Tomato) (Lycopersicon esculentum)	plus	0	0	0	MAAATMLSSSTFAG ... NNNAWAFA TNFVPGK	267
5	P0C2W8 CO1A1_MAMAE Collagen alpha-1(I) chain - Mammut americanum (American mastodon)	plus	0	0	0	GFGSGLDGAKXXXXXX ... GLNLGPPIGPPGPR	900
6	P0C2W2 CO1A1_TYREX Collagen alpha-1(I) chain - Tyrannosaurus rex (Tyrant lizard king)	plus	0	0	0	GATGAGPGIAAGAPGFP ... XXXXXXGVGVGLPGR	570
7	P85154 CO1A2_MAMAE Collagen alpha-1(II) chain - Mammut americanum (American mastodon)	plus	0	0	0	GSDGSVPGVPA GPN ... XXXXXGFPAGSPGVGK	830
8	P85153 CO2A1_MAMAE Collagen alpha-1(III) chain - Mammut americanum (American mastodon)	plus	0	0	0	GERGGVGPIGPPGER ... GDGGASGPSPGPAGP	669
9	Q40805 CYB_ERYTCA Cytochrome b - Eryx tataricus (Tatar sand boa)	plus	0	0	0	MPHQQMLILKGLPFLV ... NPIAGWIENNMKDN	371
10	P16274 IFEA_HELPO Non-neuronal cytoplasmic intermediate filament protein A - Helix pomatia (Roman snail) (Edible snail)	plus	0	0	0	TSKISTTYEEERGQS ... ADNEQIADMFSLGVG	551
11	Q19816 MATK_ALLCA Matsumura K - Alasmunda cathartica (Yellow alasmunda)	plus	0	0	0	MEAIIQYLOFDRSQO ... WYLDITCINDLNQK	505
12	P41413 PCSK5_RAT Proprotein convertase subtilisin/kexin type 5 precursor - Rattus norvegicus (Rat)	plus	0	0	0	MDNWGWSRCRCPGRR ... EDDOLEYDESYSYQ	1877
13	P22625 VNBP_CLV 11.6 kDa protein - Carnation latent virus (CLV)	plus	0	0	0	MERIKRLKQLEQDFK ... DLLFEGIDLCVRSK	101
14	P37992 VNBP_CVB 12.6 kDa protein - Chrysanthemum virus B (CVB)	plus	0	0	0	MDIVIKMLILRKFVE ... LWGVTEIVPHPGYNF	107
15	Q00572 VNBP_HELVS 12.6 kDa protein - Helminth virus S (HeV/S)	plus	0	0	0	MDKRKNKANVVLSCS ... KWQVTEIVPHPGFN	110
16	P27336 VNBP_LSV 16 kDa protein - Lily symptomless virus (LSV)	plus	0	0	0	MSVWGAWKPNTPVGY ... PWISPHRQOFYLRPK	140
17	Q02123 VNBP_POPMV 14 kDa protein - Poplar mosaic virus (isolate ATCC Pv275) (PMV)	plus	0	0	0	MVNMRKVLA LMQVFR ... PSTFHGYGYPVGHKT	123
18	Q01687 VNBP_PVM3 12 kDa protein - Potato virus M (strain German) (PVM)	plus	0	0	0	MKDVTKVALLIARAM ... EGTVETPVINSKRE	108
19	P16654 VNBP_PVSP 10.7 kDa protein - Potato virus S (strain Persian)	plus	0	0	0	MKAERLEMMLLCVYR ... SPNFKVVTFIRGWSH	93
20	Q04580 VNBP_SHVX 14.7 kDa protein - Shallot virus X (ShvX)	plus	0	0	0	MHPHDNLCLHF ... QLINDMILKSLKL	128
21	Q6PSC4 MATK_DOLLA Matsumura K - Dolichos lab lab (Field bean) (Lablab purpureus)	plus	0	0	2	MEQOAYLELRSSRY ... WYLDIFFFFRNDFVNHF	504
22	Q71ML1 MATK_CHIUM Matsumura K - Chimaphila umbellata (Pipeweed)	plus	0	0	10	MEFKRNLELDRSQO ... WYLDICINDLSNKE	504
23	P17530 VNBP_PVMR 10.7 kDa protein - Potato virus M (strain Russian) (PVM)	plus	0	0	19	MKDVTKVALLIARAM ... EGTVETPVINSKRE	108

Back Results per page : 1500 Next
 Begin Jump to : P14274|CB2A_SOLLC Chic End
 Maximum sequences length : 30

Result 1 to 23 of 23

Conclusion

Conclusion

- New sequence-driven algorithm
- Bridging the gap between Multiple Sequences Alignments and Pattern Discovery and beyond ...
- Web interface
- Let's play with all kind of data !

Acknowledgments

- Région Bretagne
- Inria/Irisa
- Symbiose Team
- Plateforme Genouest (Laetitia Guillot)
- Christian Delamarche (MIP) and Thierry Guillaudeux (TNF)
- Boris Idmont, Daniel Fredouille, Thi Hong Hanh Hoang, Laetitia Guillot, François Coste and Jacques Nicolas