



Evaluation and Integration of Existing Methods for Computational Prediction of Allergens

Jing Wang



What's allergen?

Allergens

Food allergens



Others



Venom/salivary



Contact allergens



Aero allergens



Hazards of allergens

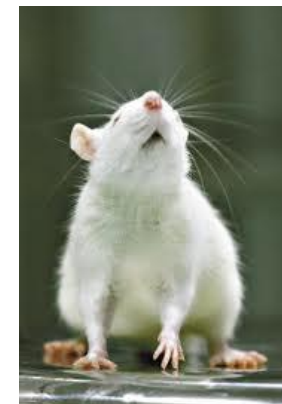
- Chronic ill health are mainly caused by allergy, affecting about 25% of the population in the world.
- Asthma and atopic dermatitis, respectively, affect 10% and 15% of the children in some countries.
- Fish allergy and general food allergy were reported in 2.3% and 4% of the US population, respectively.



.....

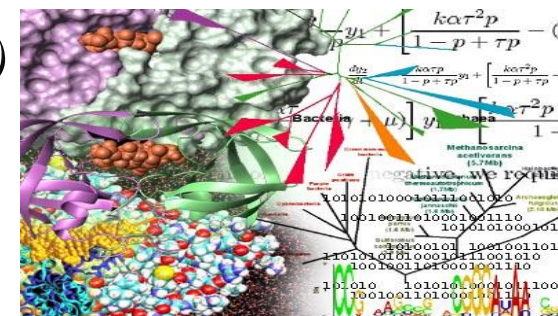
Experimental methods

- High time consumption
- High cost
- Difficulty for choosing candidate



Bioinformatics prediction methods

- Sequence-based method (*FAO/WHO, 2001*)
- Motif-based method (*Stadler, M. B. et al., FASEB J, 2003*)
- SVM-based method (*Saha S et al., NAR, 2006*)



FAO: Food and Agriculture Organization of the United Nations;
WHO: World Health Organization

Bioinformatics prediction methods

- Sequence-based method (*FAO/WHO, 2001*)

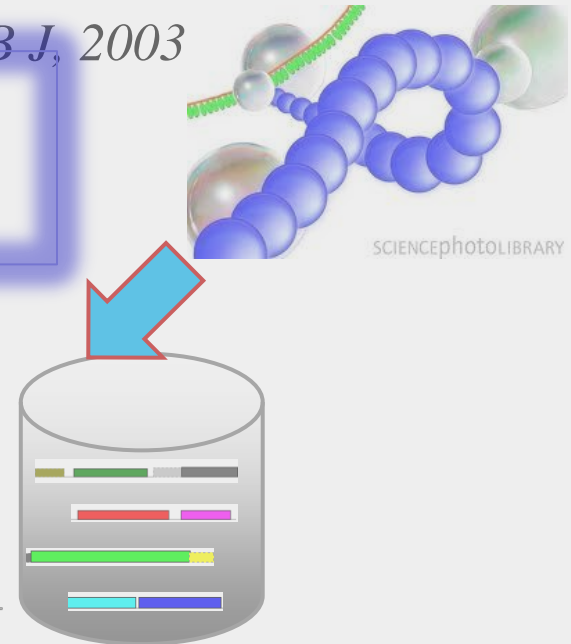
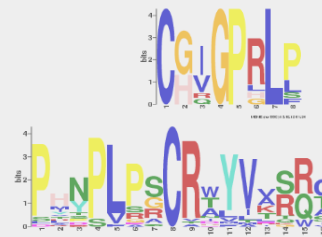
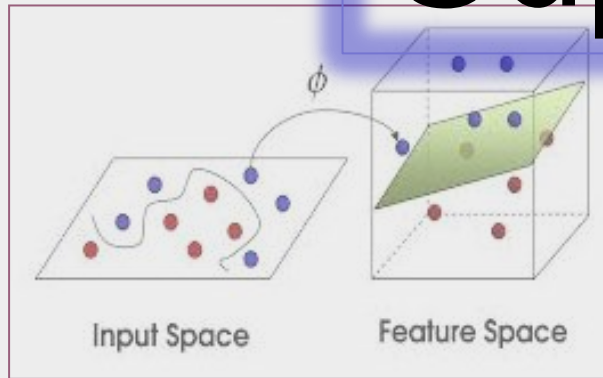
```

EDRRRGRGSRGRGNGLIEETICTASAKKNIGRMRSPDIYNPQAGSLKTAN - -
LDGRRGRGRGGGLIEETICTACVKKLIGGMRSPHIYDPR - - LFTGNCH
    
```

- Motif-based method (*Stadler, M. B. et al., FASEB J, 2003*)

- SVM-based method (*Saha S et al., NAI, 2006*)

Superior?





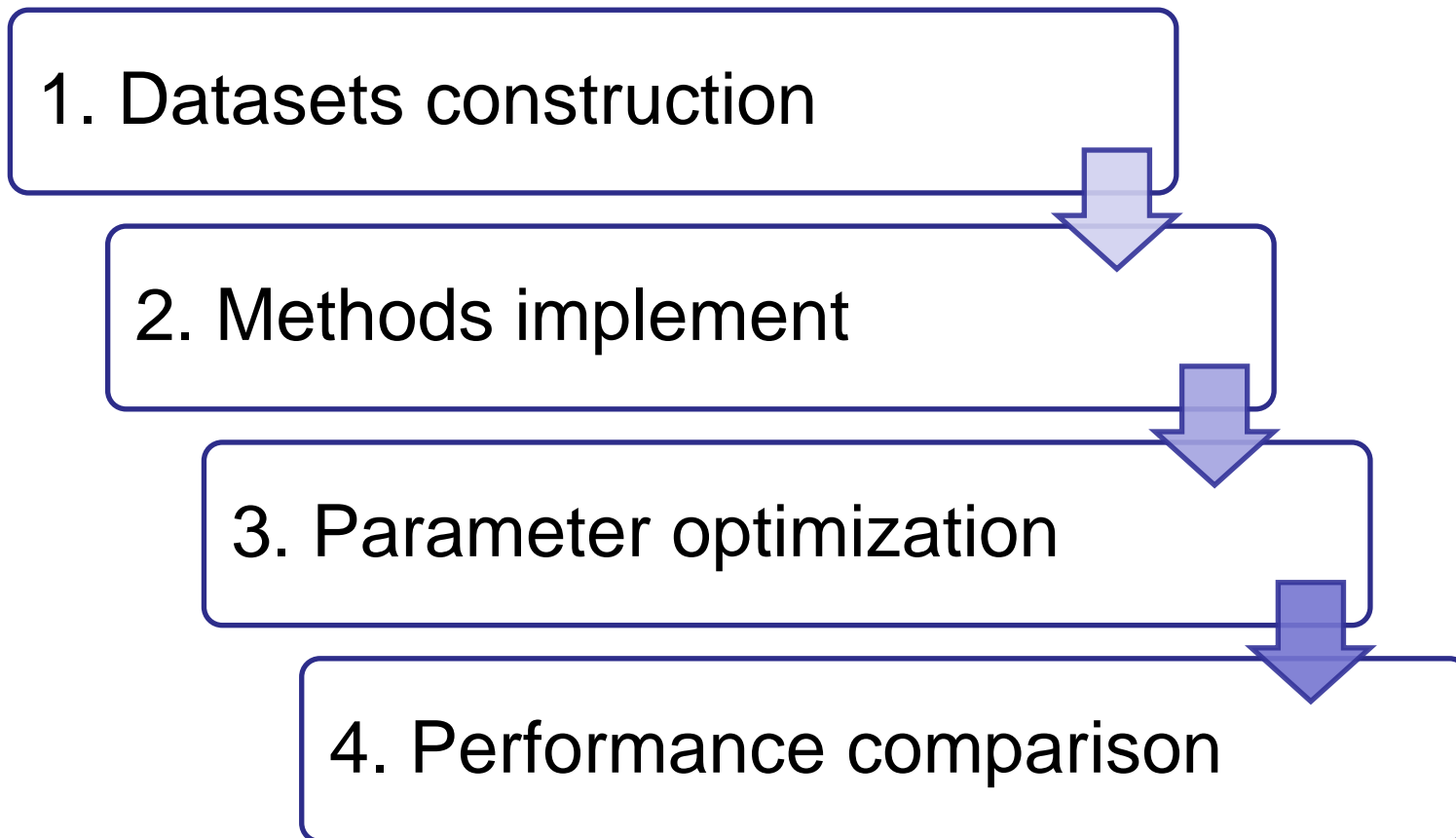
Methods Evaluation

- Compare the performance of a variety of computational methods for allergen prediction
- Find the cons and pros of each method, and perform parameter optimization

Methods integration

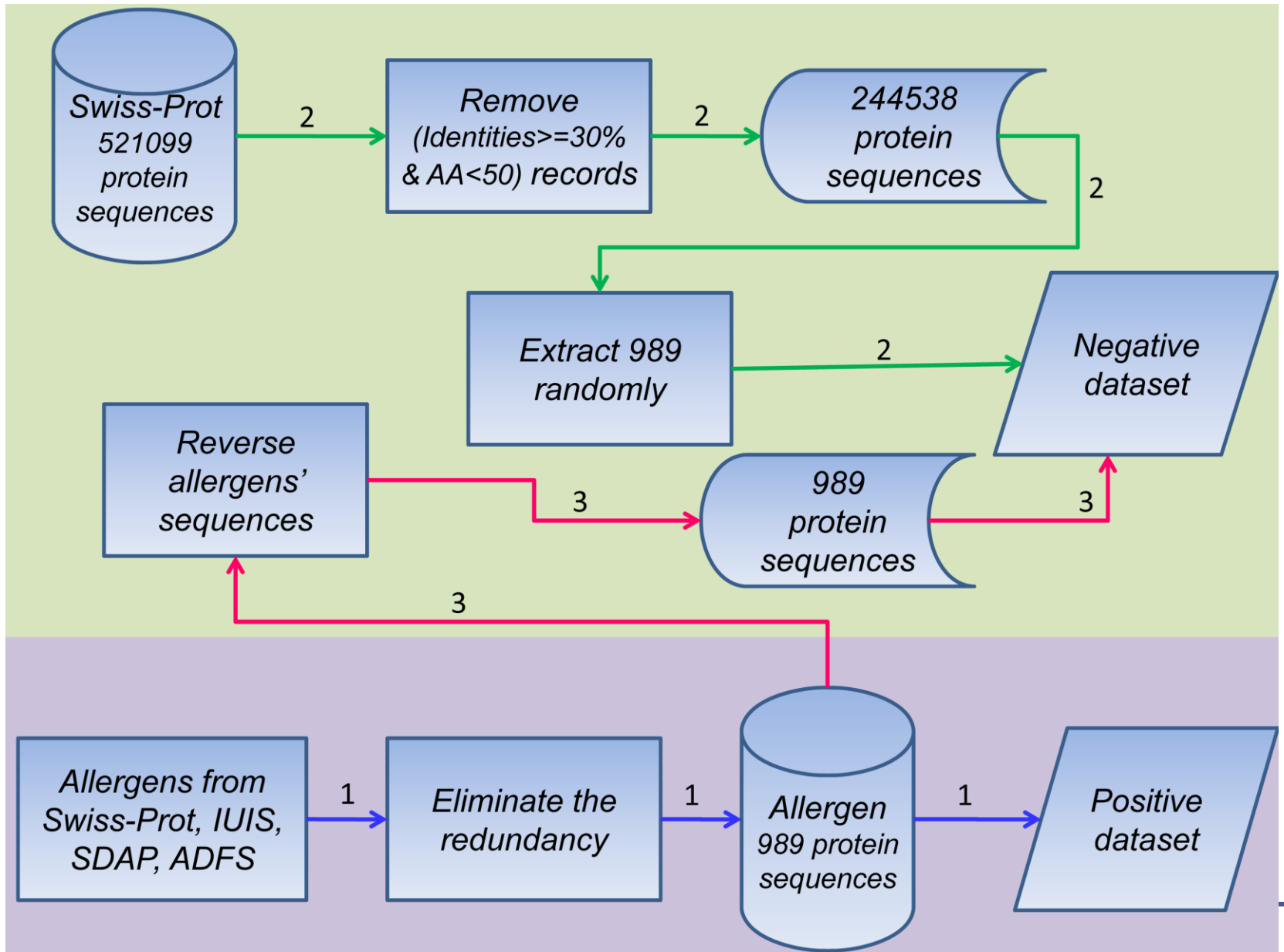
- Database search for all known allergen
- One-stop prediction for the protein allergenicity
 - Single protein prediction
 - Batch prediction

Workflow of the evaluation



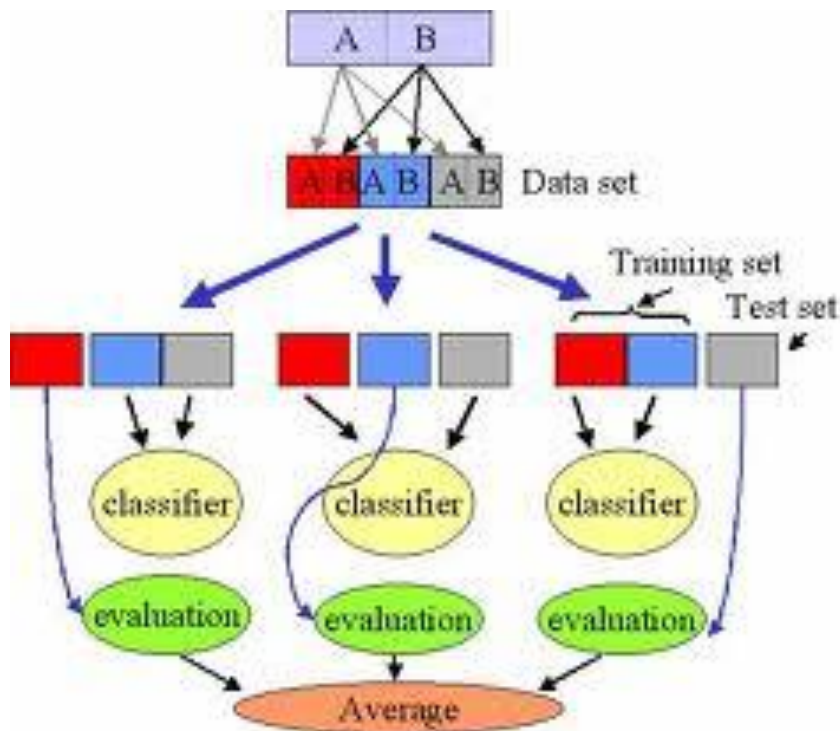


1. Datasets construction



1. Datasets construction

⊙ Ten-fold cross validation



The dataset was randomly partitioned into ten subsets, where each subset had nearly equal number of allergens and non-allergens. Of the ten subsets, a single set was retained as the validation data for testing the method, and the remaining nine subsets were used as training data. This process was then repeated 10 times with each of the ten subsets used exactly once as the validation data. The overall performance of a method was the average performance over ten subsets.

2. Methods implement

Sequence-based method (FAO/WHO criteria)

Querying:

MQTRSI **DNVVNW** SRCVHPSCVAWVIFIHF CFAKNCSILY.....

Rule 1: Cutoff=6



AYYVAAGKL **DNVVNW** SRCVHPSCVAWVIFIHF CFAKNCSILY.....

Known allergen sequences

Querying:

EQCRFRCLQKQGGGEQRECQRM CQSLHKEAPEHTSPEDV.....

Identities = 31/94 (32%),

EQCRFRC-LQKQGGGEQRE-CQRM CQSLHKE----APEHTSPEDVGGTGWEEEEEEHERE
+QC+ +C +Q+Q +Q+E C + C+ +KE EH E+ GTG +E HE
KQCKHQCKVQRQYDEQQKEQCVKECEKYKKEKKGREREREHEEEEEEWGTGGVDEPSTHE--

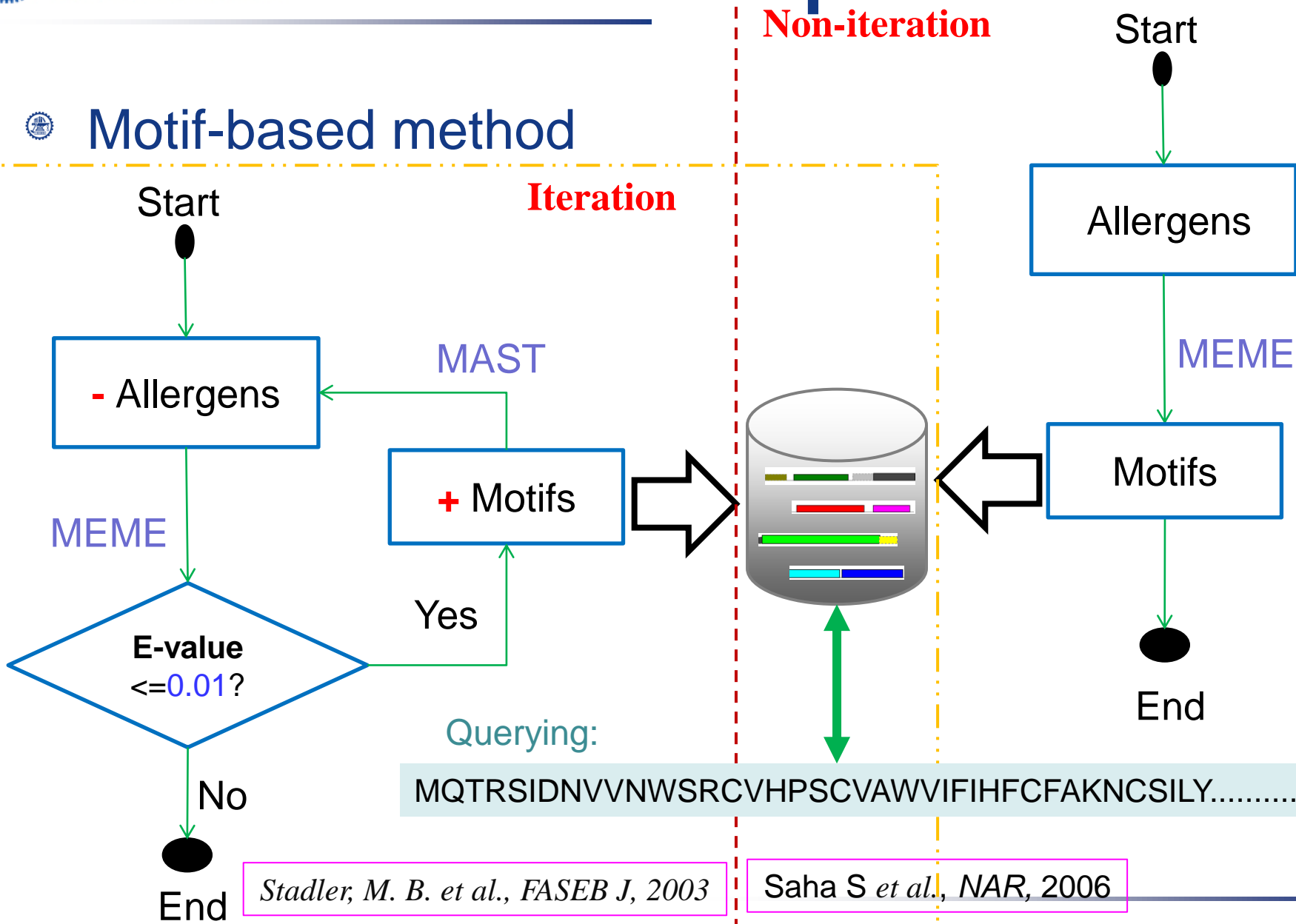
Rule 2: Cutoff=35%

KQCKHQCKVQRQYDEQQKEQCVKECEKYKKEKKGRERE.....

Known allergen blast database

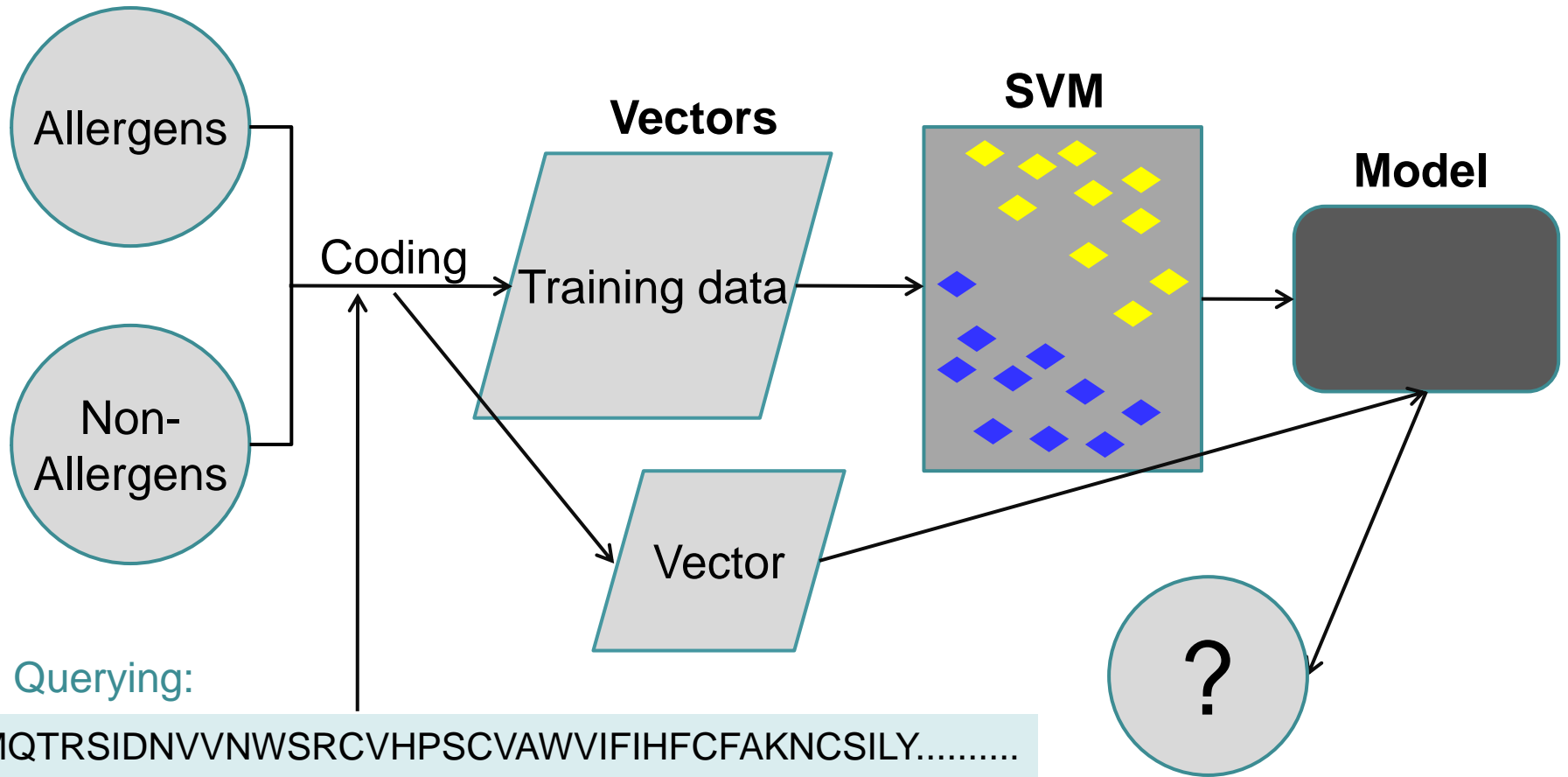
2. Methods implement

Motif-based method



2. Methods implement

SVM-based method

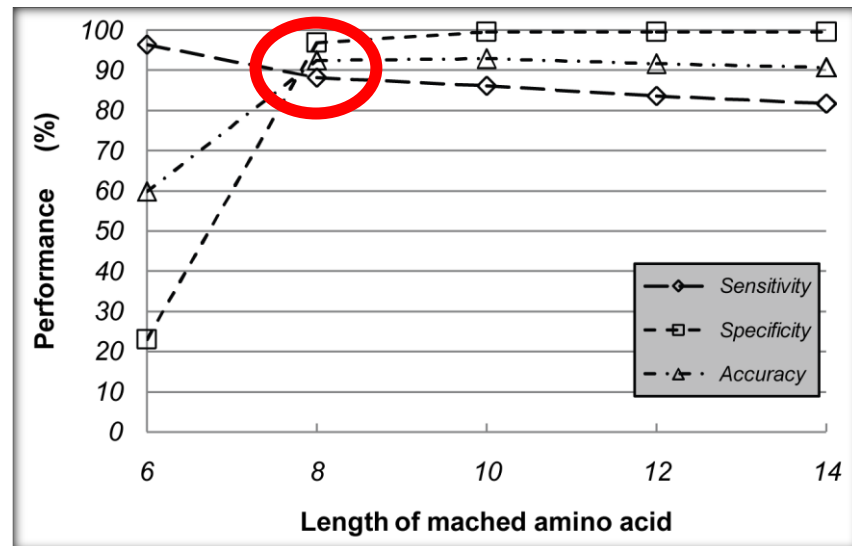
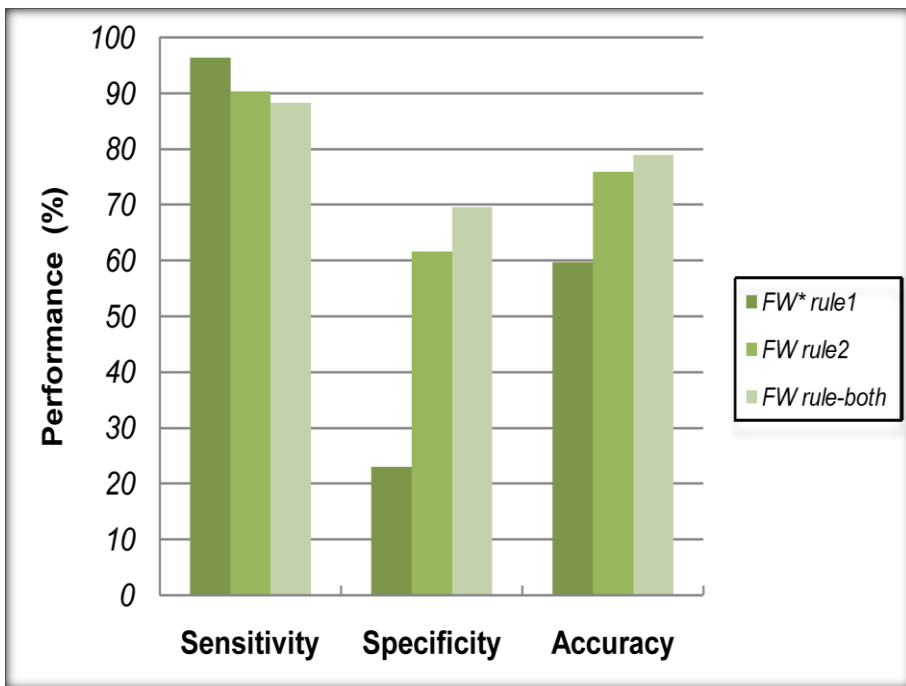


MQTRSIDNVVNWSRCVHPSCVAWVIFIHFCFAKNCSILY.....

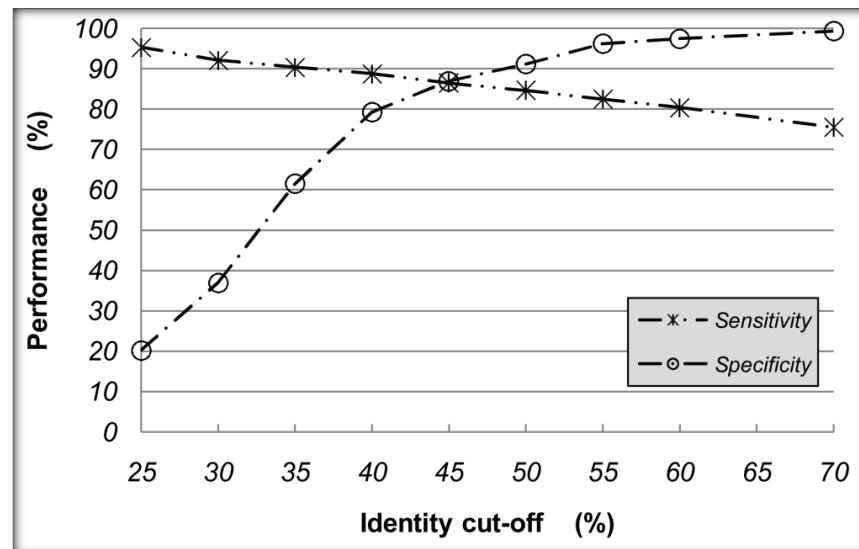
Amino Acid Composition Fraction of amino acid $i = \frac{\text{total number of amino acids } (i)}{\text{total number of amino acids in protein}}$

3. Evaluation results

FAO/WHO



Rule 1



Rule 2

3. Evaluation results

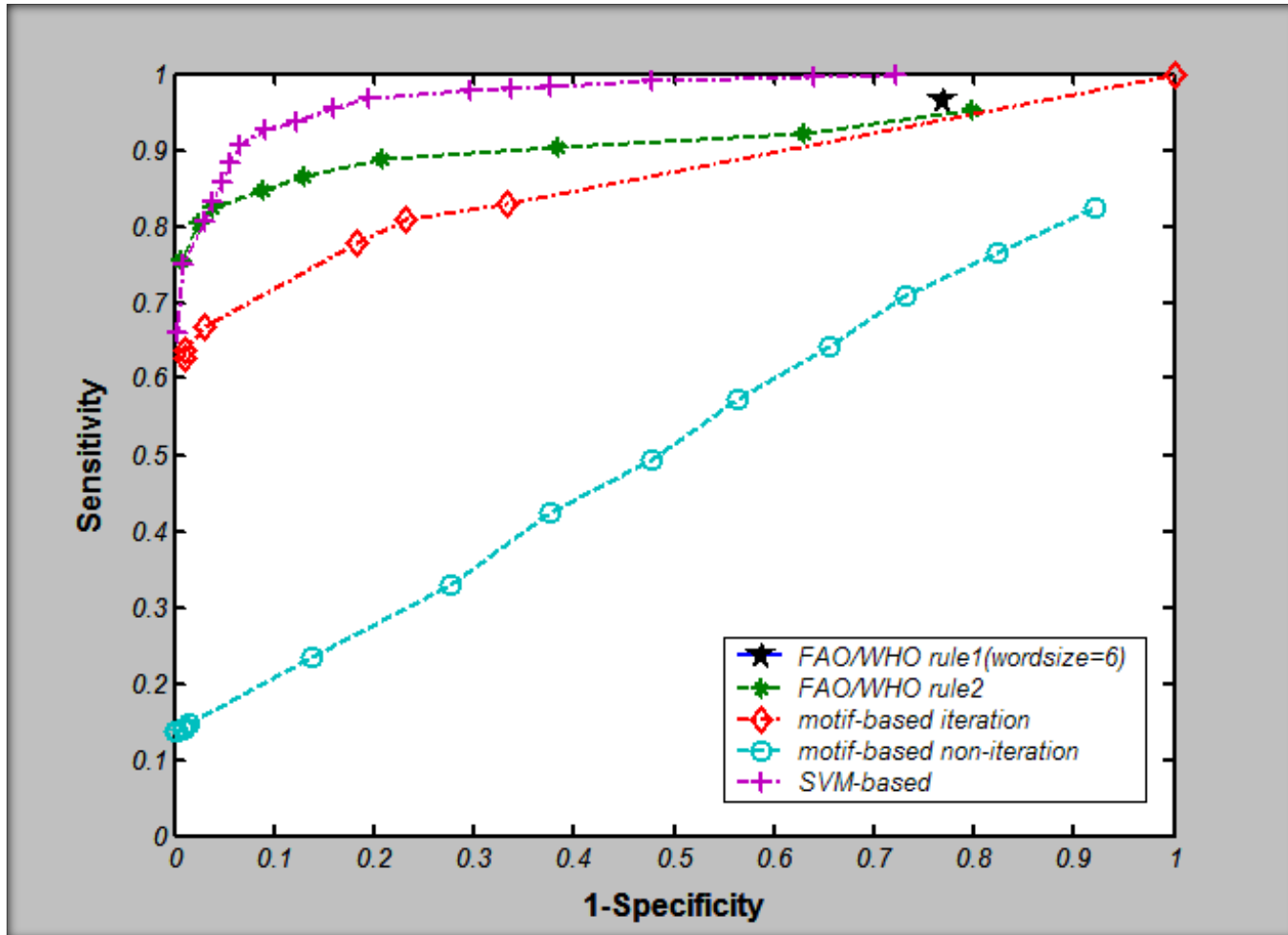
Motif-based method

MAST E-value		0.001	0.01	0.1	0.5	0.7	1	10
Iteration	Sensitivity	62.63%	63.64%	66.67%	77.78%	80.81%	82.83%	100%
	Specificity	98.99%	98.99%	96.97%	81.82%	76.77%	66.67%	0%
Non-iteration	Sensitivity	13.66%	13.66%	13.66%	13.95%	14.16%	14.77%	23.56%
	Specificity	100%	100%	99.70%	99.19%	98.89%	98.48%	86.15%



4. Methods comparison

SVM-based method and comparison



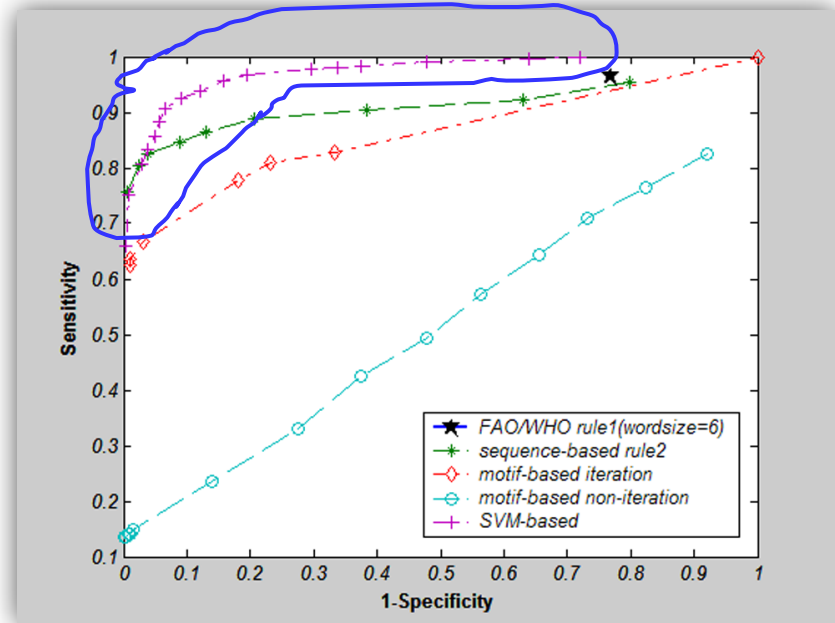
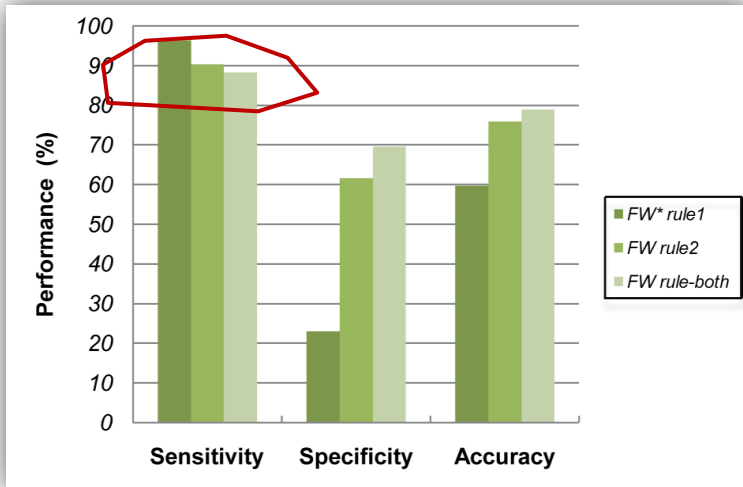
4. Methods comparison

- Methods comparison
 - Time complexity

Approaches	F/W rule 1	F/W rule2	Motif-based	SVM-based
Time (ms*)	15940	58640	87	10

ms means millisecond

Why integration?





[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- Introduction
- Allergen search
- Allergenicity prediction
- Batch prediction

Statistics

- Allergen: 1096
- Allergen category: 13
- Species: 249
- Prediction method: 4

Links

[Swiss-Prot/TrEMBL](#)
[NCBI-Entrez](#)
[SDAP-Structural Database of Allergenic Proteins](#)
[BLAST](#)
[FASTA](#)
[MEME/MAST](#)

Home

Welcome to proAP !

proAP, Protein Allergenicity Prediction, is a web-based database of allergenic proteins, providing bioinformatics tools to determine cross-reactivities between potential allergens and known allergens.



Main modules

Allergen search

Allergen search

By category. Aero Animal
 By species. Aero Animal
 List all.

Aero Animal
 Aero Fungi
 Aero Insect
 Aero Mite
 Aero Plant
 Contact
 Food Animal
 Food Fungi
 Food Plant
 Gladin
 Protozoan
 Venom/Salivary
 Worm

Search Res



Statistics

- o Allergen: 1096
- o Allergen category: 13
- o Species: 249
- o Prediction method: 4

Allergen search

By category. Aero Animal
 By species. Actinidia
 List all.

Actinidia
 Alternata
 Anisakis simplex
 Aspergillus
 Buckwheat
 Carpinus
 Cattle
 Cedar
 Chironomus
 Cockroach
 Grass
 Hazelnut
 Hevea brasiliensis
 Maize
 Malus
 Mite
 Olive
 Peanut
 Pendula
 Potato

Search Re

Allergenicity prediction

Allergenicity prediction

Requiring:

Paste or type your sequence: (One sequence only. Fasta format or pure aa sequence)

```
MYWSNQITRRLGERVQGFMSGISPQQMGEPEGSWSGKNPGTMGASRLYL  
VLVLQQRVLLGMKKRGFGAGRWNFGGGKVQEGETIEDGARRELQEESGLTV  
DALHKVGQIVFEFVGEPELMDVHVFCTDSIQGTPVESDEMPCWFQLDQIPF  
KDMWPDDSYWFPLLLQKKKFHGYFKFQGQDTILDYTLREVDTV
```

Methods:

Select method(s): (One or more methods)

- FAO/WHO: Amino acids sliding window: 80, sequence identity cutoff \geq 35 %
- FAO/WHO: Exact match for \geq 6 contiguous amino acids
- Motif-based method
- SVM-AAC method

Batch prediction

Requiring:

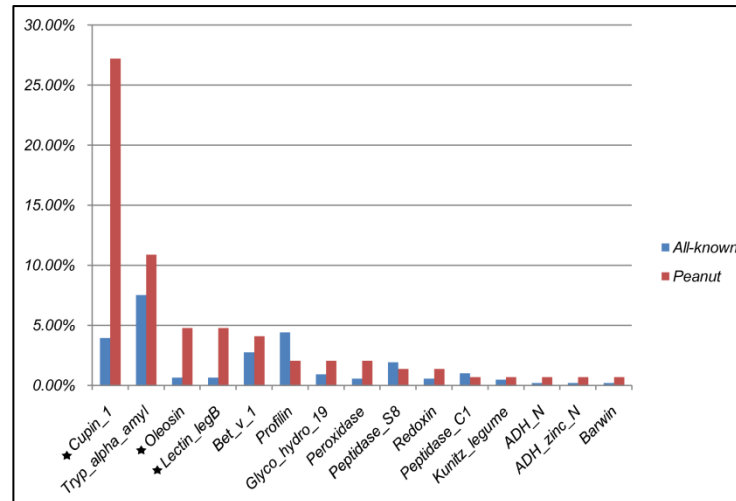
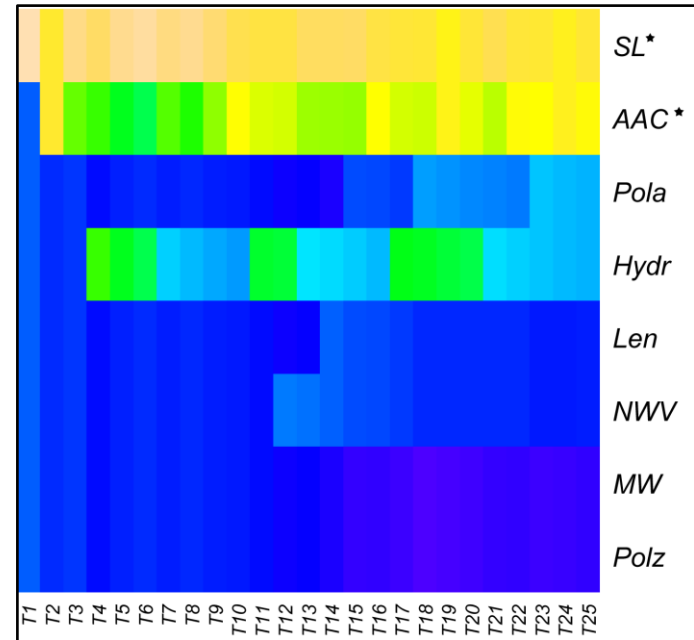
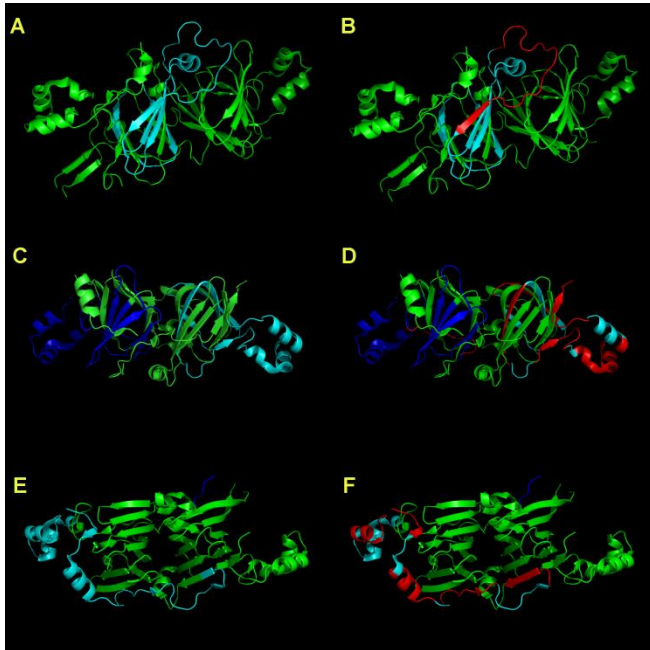
Sequences file:

 浏览...

Email address:

Next step ...

- Key features for allergenicity
- Family preference
- Specific structures



Acknowledgments



 Dr. Jing Li



 Prof. Dabing Zhang

 Yabin Yu

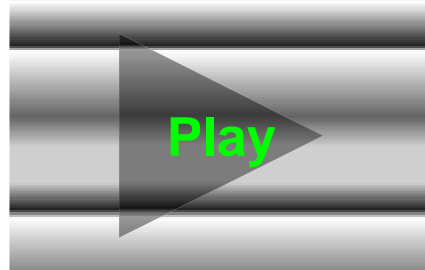
 Yunan Zhao



上海交通大學
SHANGHAI JIAO TONG UNIVERSITY

Software Demonstration

 ***proAP*** -- Protein Allergenicity ***P***rediction





上海交通大學
SHANGHAI JIAO TONG UNIVERSITY

Thank you



- 我的文档
- Zend Studio - 7.1.2
- plsvar.m
- 简介内容--学生.doc
- domainap...
- 我联系国外博士后职...
- Resume (王婧).doc
- VBA密码破解.txt
- PF00190_...
- 我的电脑
- 新建 文本文档.txt
- 甜菜蚜 半翅目_蚜科_...
- clustal-...
- Microsoft Word遇到...
- 在美国做博士后.txt
- IMG_0387...
- 148.jpg
- PF00190_...
- 网上邻居
- ipapp.doc
- 进化树.jpg
- foodplan...
- writing worksh...
- Manual_w...
- Resume(E...
- 2012xing...
- 1111111.nwk
- 回收站
- 安装和用户指南.pdf
- IMG_0901...
- OS.docx
- shanghai jiaotong...
- Tra_Gene...
- IMG0029A...
- 2012shen...
- Figure2.pdf
- ICBC在线银
- 安装和使用
- email.docx
- 百奥知LIMS
- GraphApp...
- QMEAN_pl...
- IMG_3639...
- 2012shen...
- 上海交通大学公开...

Administrator

- Internet Explorer
- 迅雷看看播放器
- 强力卸载电脑上的软件
- Mozilla Firefox
- Micro 2007
- 飞信2012
- 腾讯QQ
- Microsoft Office Excel 2007
- Microsoft Office PowerPoint 2007
- 所有程序 (E)

- 我的文档
- 我最近的文档 (D)
- 我的电脑
- 网上邻居
- 控制面板 (C)
- 网络连接
- 搜索 (S)
- 运行 (R)...

位置: C:\Program Files\Mozilla Firefox

- Model_1_...
- Group时间表.xls
- 博士等高层次人才
- 2012第五届“海内外”
- SGM-DY-J...
- 平均.xlsx
- Swiss_pd... 使用说明
- 080613_S...
- Biostat2...
- Tra_Exp...
- 111.txt
- Book1.xlsx
- 上海交通大学+王婧+
- gettemp.gif
- 暑假班车.txt
- 学生党支部信息.xls
- gettemp2...
- Doc1.docx
- 上海交通大学...
- well...i...
- game.pdf

<http://gmobl.sjtu.edu.cn/proAP/main.html>

Allergen Prediction Based on Protein Features Home

gmobl.sjtu.edu.cn/PRAL/

Allergenicity Prediction Database of Protein

gmobl.sjtu.edu.cn/proAP/main.html

Allergenicity Prediction Database of Protein

gmobl.sjtu.edu.cn/proAP/prediction.html

Allergenicity Prediction Database of Protein

gmobl.sjtu.edu.cn/proAP/batch.html

Allergenicity Prediction Database of Protein

gmobl.sjtu.edu.cn/proAP/search.html

Allergenicity Prediction Database of Protein

gmobl.sjtu.edu.cn/proAP/introduction.html

proAP Protein Allergenicity Prediction

Home

Navigation

- Introduction
- Allergen search**
- Allergenicity prediction
- Batch prediction

Statistics

- Allergen: 1096
- Allergen category: 13
- Species: 249
- Prediction method: 4

Links

[Swiss-Prot/TrEMBL](#)
[NCBI-Entrez](#)
[SDAP-Structural Database of Allergenic Proteins](#)
[BLAST](#)
[FASTA](#)
[MEME/MAST](#)

Welcome to proAP !

proAP, Protein Allergenicity Prediction, is a web-based database of allergenic proteins, providing bioinformatics tools to determine cross-reactivities between potential allergens and known allergens.



1 2 3 4 5



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- Introduction
- Allergen search
- Allergenicity prediction
- Batch prediction

Allergen search

By category.

 By species.

 List all.

- Aero Animal
- Aero Animal
- Aero Fungi
- Aero Insect
- Aero Mite
- Aero Plant
- Contact
- Food Animal
- Food Fungi
- Food Plant
- Gladin
- Protozoan
- Venom/Salivary
- Worm

沪交ICP备20111135

Reaction Laboratory in Shanghai Jiao Tong University (GMODL-SJTU)



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Search result

Total records: 17 Page: 1/2 [First](#) [Previous](#) [Next](#) [Last](#)

No.	Name	UniProtAcc	TaxonomicName	CommonName	Category	Species	Epitope	Sugar	
1	Bos d 3	Q28050	Bos taurus	Bos bovis,Bos primigenius taurus,Bos taurus,bovine,cow,domestic cow,domestic cattle,cattle,	Aero Animal	Cattle			Prote calcium A7) allergen all (Calcium in 2)(C
2	Can f 1	O18873	Canis lupus familiaris	Canis domesticus,Canis familiaris,Canis canis,Canis lupus familiaris,dogs,dog,	Aero Animal	other		sugar	Major (Allergen (Fla
3	Can f 2	O18874	Canis lupus familiaris	Canis domesticus,Canis familiaris,Canis canis,Canis lupus familiaris,dogs,dog,	Aero Animal	other		sugar	Minor (Allergen (Fla
4	Can f 3	P49822	Canis lupus familiaris	Canis domesticus,Canis familiaris,Canis canis,Canis lupus familiaris,dogs,dog,	Aero Animal	other			Serum (Fla
				Cavia aperea					

Navigation

- Introduction
- Allergen search
- Allergenicity prediction
- Batch prediction

By category.

Aero Animal

By species.

Actinidia

List all.

Search

Reset



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- [Introduction](#)
- [Allergen search](#)
- [Allergenicity prediction](#)
- [Batch prediction](#)

Allergenicity prediction

Predict methods including FAO/WHO criteria, Motif-based method and SVM-AAC method (take amino acid composition as protein features). Choose the method according to your demand. Details about predict methods see [Introduction](#) page.

Paste or type your sequence: (One sequence only. Fasta format or pure aa sequence)

Select method(s): (One or more methods)

- FAO/WHO: Amino acids sliding window: 80, sequence identity cutoff \geq %
- FAO/WHO: Exact match for \geq contiguous amino acids
- Motif-based method
- SVM-AAC method

Select category: (Only for FAO/WHO methods)

- By .
- By all.



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- [Introduction](#)
- [Allergen search](#)
- [Allergenicity prediction](#)
- [Batch prediction](#)

Allergenicity prediction

Predict methods including FAO/WHO criteria, Motif-based method and SVM-AAC method (take amino acid composition as protein features). Choose the method according to your demand. Details about predict methods see [Introduction](#) page.

Paste or type your sequence: (One sequence only. Fasta format or pure aa sequence)

```
MGIKHCCYILYFTLALVTLVQAGRLGEEVDILPSPNDTRRSLQGCEAHNIIDK
CWRCKPDWAENRQALGD
CAQGF GKATHGGKWDIYMTSDQDDDVVNPKEGTLRF GATQDRPLWII
FQRDMIIYLQQEMVVTSDKTI
DGRGAKVFI VYGGITI MNVKNVITHNIDIHDV RVI PGGRIKSNGGPAIPRHO ...
```

Select method(s): (One or more methods)

- FAO/WHO: Amino acids sliding window: 80, sequence identity cutoff \geq %
- FAO/WHO: Exact match for \geq contiguous amino acids
- Motif-based method
- SVM-AAC method

Select category: (Only for FAO/WHO methods)

- By .
- By all.



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- [Introduction](#)
- [Allergen search](#)
- [Allergenicity prediction](#)
- [Batch prediction](#)

Prediction results

Query protein:

```

MQLLLLVGLALICGLQAQEGNHEEPQGGLEELSGRWHSVALASNKSDLIKPWGHRVFI
HSMSAKDGNLHGDILIPQDGGQCEKVSMTAFKTATSNKFDLEYWGHNDLYLAEVDPKSYLI
LYMINQYNDTSLVAHLMVRDLSRQQDFLPAFESVCEDIGLHKDQIVVLSDDDDRCQGSRD

```

Prediction result

FAO/WHO: sequence alignment	FAO/WHO: amino acids match	Motif-based	SVM-AAC
allergen	--	--	--

FAO/WHO: Sequence identity cutoff \geq 45% (allergen)

Matched records:

No.	Name	UniProtAcc	TaxonomicName	CommonName	Category	Species	Epitope
				Canis			

FAO/WHO: Sequence identity cutoff >= 45% (allergen)

Matched records:

No.	Name	UniProtAcc	TaxonomicName	CommonName	Category	Species	Epitope	Sugar	Description
1	Can f 2	O18874	Canis lupus familiaris	Canis domesticus,Canis familiaris,Canis canis,Canis lupus familiaris,dogs,dog,	Aero Animal	other		sugar	Minor allergen Can f 2 (Allergen Dog 2) (Can f 2) (Flags: Precursor)



Alignments detail:

```

query:(1--80)
>sp|O18874|ALL2_CANFA Minor allergen Can f 2 OS=Canis familiaris PE=1 SV=1
    Length = 180

    Score = 166 bits (421), Expect = 4e-44, Method: Compositional matrix adjust.
    Identities = 80/80 (100%), Similarity = 100%, Positives = 80/80 (100%)

Query: 1  MQLLLLVGLALICGLQAQEGNHEEPQGGLEELSGRWSVALASNKSDLIKPWGHFRVFI 60
          MQLLLLVGLALICGLQAQEGNHEEPQGGLEELSGRWSVALASNKSDLIKPWGHFRVFI
Sbjct: 1  MQLLLLVGLALICGLQAQEGNHEEPQGGLEELSGRWSVALASNKSDLIKPWGHFRVFI 60

Query: 61  HSMSAKDGNLHGDILIPQDG 80
          HSMSAKDGNLHGDILIPQDG
Sbjct: 61  HSMSAKDGNLHGDILIPQDG 80

query:(2--81)
>sp|O18874|ALL2_CANFA Minor allergen Can f 2 OS=Canis familiaris PE=1 SV=1
    Length = 180

    Score = 166 bits (420), Expect = 5e-44, Method: Compositional matrix adjust.
    Identities = 80/80 (100%), Similarity = 100%, Positives = 80/80 (100%)

Query: 1  QLLLLTVGLALICGLQAQEGNHEEPQGGLEELSGRWSVALASNKSDLIKPWGHFRVFIH 60
    
```





Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- [Introduction](#)
- [Allergen search](#)
- [Allergenicity prediction](#)
- [Batch prediction](#)

Allergenicity prediction

Predict methods including FAO/WHO criteria, Motif-based method and SVM-AAC method (take amino acid composition as protein features). Choose the method according to your demand. Details about predict methods see [Introduction](#) page.

Paste or type your sequence: (One sequence only. Fasta format or pure aa sequence)

```
LDPAIIGYVAEQENMSASDVVNALNKKSGMLALTGASDMRDVFAKPQENA  
VAIKMYVNRV  
ADYIAKYLNLQLEGNIDGLVFTGGIGENASDCVELFINAVKSLGFATDLKLFVK  
YGDSCVV  
STPQSKYKIYRVRTNEELMIVEDSIRLTQK
```

Select method(s): (One or more methods)

FAO/WHO: Amino acids sliding window: 80, sequence identity cutoff \geq %

FAO/WHO: Exact match for \geq contiguous amino acids

Motif-based method

SVM-AAC method

Select category: (Only for FAO/WHO methods)

By .

By all.

--	allergen	allergen	non-allergen
----	----------	----------	--------------

FAO/WHO: Exact match 6 amino acids (allergen)

Matched records:

No.	Name	UniProtAcc	TaxonomicName	CommonName	Category	Species	Epitope	Sugar
1	Pha a 5	P56166	Phalaris aquatica	Phalaris aquatica L., Phalaris tuberosa, Phalaris aquatica, canary grass,	Aero Plant	Grass		
2	Par c ?	A2V734	Paralithodes camtschaticus	Paralithodes camtschaticus, Kamchatka crab, red king crab,	Food Animal	other		

Hits Subjects Detail:

```
>sp|P56166|MPA53_PHAHQ Major pollen allergen Pha a 5.3 OS=Phalaris aquatica PE=1 SV=1
MAVQKYTVLFLAMALVAGPAASYAADAGTPPTPATPAVPGAAAGKATTHEQKLIEDINA
AFKWWPASAPPADKYKTFETAFSKANIAGASTKGLDAAYSVVYNTAAGATPEAKYDSFVT
ALTEALRIMAGTLEVHAVKPATEVPSAKAKILRANSRSSTRSSRFKIAATVATPLSHSTA
ANSAPANDKFTVFEGAFNKAIKERHGGPTETYKFIPLSLEAAVKQAYGATVARAPEVKYAV
FEAGLTKAITAMSEAQKVAKPVRLSPQPPQVPLAAGGAATVAAASDSRGGYKV
```

```
>tr|A2V734|A2V734_PARCM Tropomyosin slow-tonic isoform OS=Paralithodes camtschaticus GN=Tm-Pane-tonic PE=2 SV=1
MDAIIKKMQAMKLEKDNAMDKADTLEQQNKEANNRAEKTEEEIRLTPQKMQQVENEVDVA
QEQLSLANTKLEEKALQNAEGEVAALNRRIQLLEEDLERSEERLNTATTKLAEASQAA
DESERMRKVLNRSLSDEERMDALENQLKEARFLAEEADRKYDEVARKLAMVEADLERAE
ERAETGESKIVELEELRVVGNLKSLEVSEEKANQREEAYKEQIKTLANKLKAAEARAE
FAERSVQKLQKEVDRLLEDELVNEKEKYKNIADEMDQAFSELSGF
```

hits subjects Detail:

>sp|P56166|MPA53_PHAHQ Major pollen allergen Pha a 5.3 OS=Phalaris aquatica PE=1 SV=1

MAVQKYTVLFLAMALVAGPAASYAADAGTPPTPATPAVPGAAAGKATTHEQKLIEDINA
 AFKWWPASAPPADKYKTFETAFSKANIAGASTKGLDAAYSVVYNTAAGATPEAKYDSFVT
 ALTEALRIMAGTLEVHAVKPAT**EVPSAK**AKILRANSRSRSTRSRFKIAATVATPLSHSTA
 ANSAPANDKFTVFEGAFNKAIKERHGGPTETYKFI PSLEAAVQAYGATVARAPEVKYAV
 FEAGLTKAITAMSEAQKVAKPVRLSPQPPQVLPLAAGGAATVAAAASDSRGGYKV

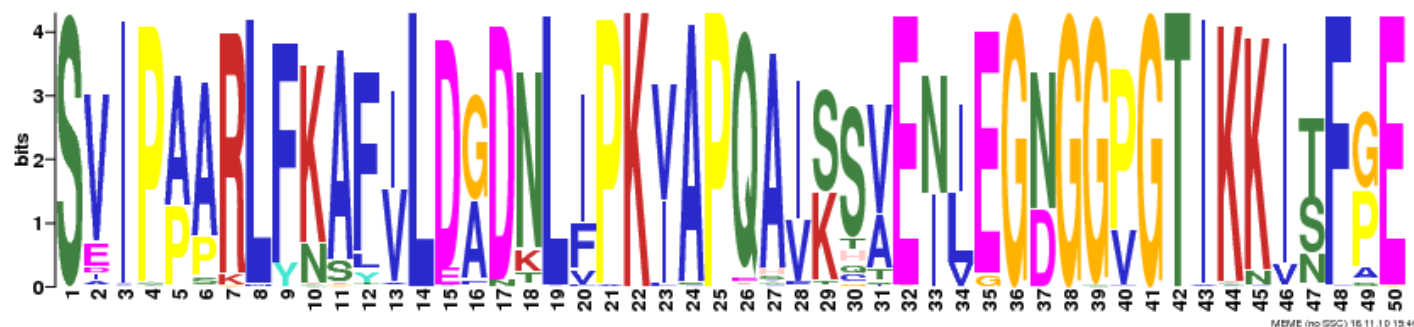
>tr|A2V734|A2V734_PARCM Tropomyosin slow-tonic isoform OS=Paralithodes camtschaticus GN=Tm-Pane-tonic PE=1

MDAIKKMQAMKLEKDNAMDKADTLEQQNKEANNRAEKTEEE**IRLTOQ**KMQQVENELDVA
 QEQLSLANTKLEEKALQNAEGEVAALNRRIQLLEEDLERSEERLNTATTKLAEASQAA
 DESERMRKVLENRSLSDEERMDALENQLKEARFLAEEADRKYDEVARKLAMVEADLERAE
 ERAETGESKIVELEEELRVVGNLKSLEVSEEKANQREEAYKEQIKTLANKLKAAEARAE
 FAERSVQKLQKEVDRLLEDELVNEKEYKNIADEMDQAFSELSGF

Motif-based result (allergen)

Your querying is predicted as allergen. Matched motif is described as:

SEQUENCE NAME	DESCRIPTION	E-VALUE	LENGTH
name		0.00048	390

**SVM-AAC result**

Your query is predicted as non-allergen (probability= 0.975783).



Protein Allergenicity Prediction

[Home](#) | [Site map](#) | [Contact us](#)

Navigation

- [Introduction](#)
- [Allergen search](#)
- [Allergenicity prediction](#)
- [Batch prediction](#)

Batch prediction

Upload a FASTA-format file containing multiple protein sequences to be predicted for allergenicity. Results of the prediction will be returned to you at the email address that you specify. Please check the [notes](#) below for the restrictions on uploaded sequence files.

Sequences file:

 浏览... 

Prediction method:

- FAO/WHO: Amino acids sliding window: 80, sequence identity cutoff \geq %
- FAO/WHO: Exact match for \geq contiguous amino acids
- Motif-based method
- SVM-AAC method

Email address:

Thank you