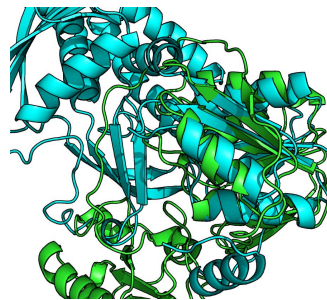


Upstream team Group 2



Daniel DE MURAT, Madeleine DE SOUSA VIOLANTE, Mei-shiue KUO, Sonia TIEO

Introduction

- **Problematic:** Find the most probable/stable 3D fold adopted by the protein in solution, given a protein primary sequence
- The project has **two steps**:
 - (i) domain annotation based on profile-profile comparison
 - (ii) identification of the most stable 3D fold by threading.



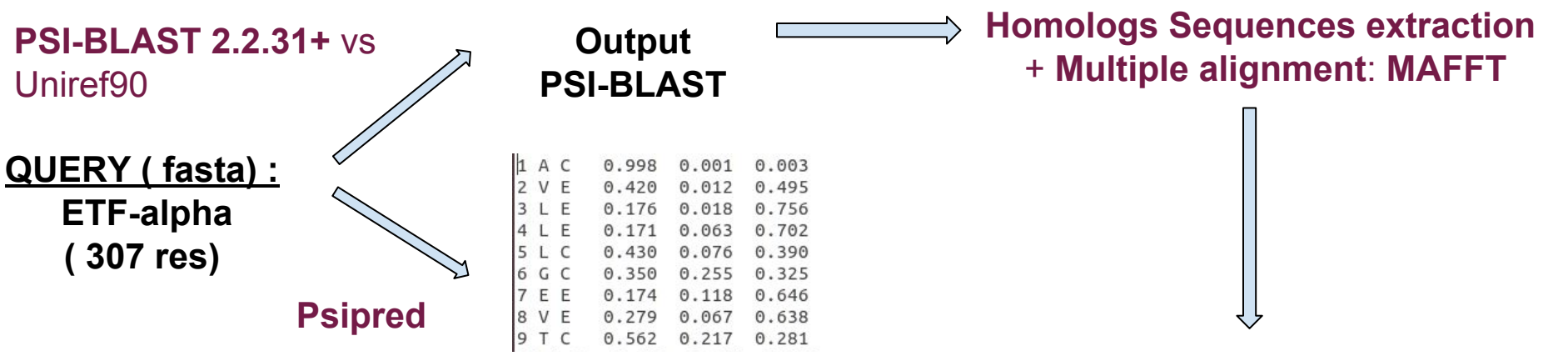
Psi-blast + Mafft alignment

PSI-BLAST 2.2.31+ vs
Uniref90

QUERY (fasta) :
ETF-alpha
(307 res)

Psipred

Output
PSI-BLAST



1	A	C	0.998	0.001	0.003
2	V	E	0.420	0.012	0.495
3	L	E	0.176	0.018	0.756
4	L	E	0.171	0.063	0.702
5	L	C	0.430	0.076	0.390
6	G	C	0.350	0.255	0.325
7	E	E	0.174	0.118	0.646
8	V	E	0.279	0.067	0.638
9	T	C	0.562	0.217	0.281

Homologs Sequences extraction
+ Multiple alignment: MAFFT

Psi-blast + Mafft alignment

PSI-BLAST 2.2.31+ vs
Uniref90

Output
PSI-BLAST

Homologs Sequences extraction
+ Multiple alignment: MAFFT

QUERY (fasta) :

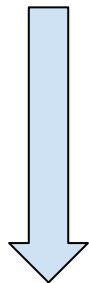
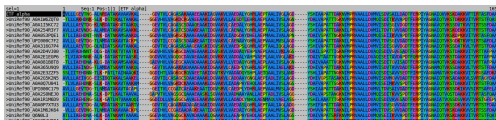
ETF-alpha
(307 res)

Psipred

sel=1	1	Seq:1 Pos:1 1 [ETF alpha]	165
ETF_alpha	AVLLLEGEVTDG-ALNRDATAKAVAAVKAL	GDVTVLCAGASAKAAEEAAKIGVAKVLAEDALYGHRLAEPAAALIVSLAGD	YSHIAAPATTDAKNVMPRAVALLDVMVLSIDVSAILDADTFERPIYAGNAIQVKSDDAKKVITRTSTFQAA
>UniRef90 A0A1W6Z0T0	-TLLAKHDNK-SLK-DSTOKALTAAKAL	GGEVHVLVAGKDCRGVAEEAAKLDGVAKVLLAEDSLTLEHRLAEPAAALIVSLAGN	YDALVAPATTGKNVMPRAVALLDVMVLSIDITKVVAPDTERPIYAGNAIQVKSDDKKVITVRTSTFQAS
>UniRef90 A0A1I5KC72	AVLLLAETVTDG-ALNRDATAKGVTAAKQL	GEVTVLCAGASCSDAAEAAKIDGVAKVLCADALYGHRLAEPAAALIVSLAGD	YSHIAAPATTDAKNIMPRVALLDVMVLSIDVSGIVDGDTERPVIYAGNAVQTVKSSDVTKVITRTSTSFDA
>UniRef90 A0A254R3Y7	AVLLLAETVTDG-ALNRDATAKAVSAAKKL	GDVTVLCAGATCSDAAEAAKIDGVAKVLCENALYGHRLAEPAAALIVSLAGD	YSHIVAPATTDAKNILPRVALLDVMILDTGVGDADTFERPIYAGNAIQTVKSKDEKVVITRTSTFQAA
>UniRef90 A0A0S3PQE1	ATLLIAEHDNT-NVK-DATNKAMSAAKDL	GGDVHVLVAGKGRCAAAEAAKIDGVAKVLIADGDYDHELAEPMALIVLAPN	YDAIVSPATTNGKNFMPRVALLDVMVLSIDITKVVAPDTERPIYAGNAIQTVKSKDPKKVITVRTSTFQAA
>UniRef90 UPI000C7F2	AVLLLEGEVTDG-TLNRDATAKAVTACKPL	GDVTVLCAGASCSDAAEAAKIDGVAKVLCENALYGHRLAEPAAALIVSLAGD	YSHIAAPATTDAKNILPRVALLDVMVLDVSGIVDADTFERPIYAGNAIQTVKSGDKTVKFSIRTSFADLA
>UniRef90 A0A316G7P4	AVLLLAETVTDG-QLNRDATSKAVTAAKAL	GDVTVLCAGASASAAEAAKIDGVSQVLAEDPSLGHRLAEPAAALIVSLAGD	YSHIVAPATTDAKNVMPRAVALLDVMVLSIDVSGVVDGDTERPIYAGNAIQTVKSSDATKVLVTRTSTFQAA
>UniRef90 A0A2D4VJ00	-TLVFADHNNT-DLG-DATAKTVTAAATKI	GGDVHVLVAGKGCDAVAAEAAKLDGVSQVLAEDAQYEHVPAEPLAALLVSLAGS	YDAILSPATTIGKNFMPRVALLDVAQISDVIQVESAADTFERPIYAGNAIQTVKSTDAKKVITVRTTAFPA
>UniRef90 A0A2E1EGP6	-TLVIADHNK-ALN-DATAKTVTAAASKL	GGDVHVLVAGKGCDAVAAEAAKLDGVSQVLAEDALYAHVPAEPLAALLVSLAGD	YDAIVTAATTIGKNFMPRVALLDVAQISDITAVDADPTFRPIYAGNAIQTVKSTDAKKVITVRTTAFPA
>UniRef90 A0A081B8T8	-TLLIAEHDNA-KLG-DATAKMSAAKAL	GADVHVLVAGEGCAAVADAANKLDGAQVLLVEDAQYGHRLAEPAAALIVSLAGN	YDAIVAAATTIGKNVMPRAVALLDVMVLSIDITGVISADTFERPIYAGNAMQTVKSKDAKKVITVRTTAFPA
>UniRef90 A0A365U9Q9	AVLLIAEINDG-TLAMDATAKMTAAKQL	GDVTVLAAGAKAADAANKLDGAQVLLVEDAQYGHRLAEPAAALIVSLAGD	YDHIVAPATTDAKNILPRVALLDVMVLSIDVTAVDADTFERPIYAGNAMQTVKSSDAKKVITVRTSTFQAA
>UniRef90 A0A2E3ZZF5	SSLIIITEHDNN-TLK-PATLITAINAAQKI	GGDIHTLVAGKGRCAVAAEAAKIDGVTKYVVDADQYSHPLAETQAAALIVSLAND	YSHLIAPATTTGKNIMPRVALLDVAQISDIAVVSSEDFTIRPIYAGNAMATVKSDDAKKVITVRGSAFDKA
>UniRef90 A0A265K2N5	AVLLLAETVGG-EISIDQAKALSAAKEL	GDVTVLCASAGCSDAAEAAKLDGVAKVLCADDAAYGNLAEPIADLIVSLAGD	YDHIVGPSTASAKNILPRVALLDVMVLSIDVMDVVDADTFERPIYAGNAIQTVKSADAKKVLVSRVAGYDAA
>UniRef90 A0A0Q7U041	AVLLLAETVTDG-ALNRDATAKSVTAARTV	GGDVHVLVAGKDCCKAADAANKLDGVAKVLIADAPAYEHQLEPAAALIVSLAGS	YDAFVAPATTSGKNVMPRAVALLDVMVLSIDITKVVAPDTERPIYAGNAIQTVKSTDAKKVITVRTSTFQAA
>UniRef90 UPI000C179	AVLLLEGEVTDG-TLNRDATAKAVTACAPL	GEITVLCGATCAEAAKIDGVSRLVCAEDPLDGHRLAEPAAALIVSLAGD	YSHIAAPATTDAKNILPRVALLDVMVLDVSAIVDADTFERPIYAGNAVQTVKSSDKIKVVSIRTSFAELA
>UniRef90 A0A250NEJ0	ATLLIAEHDNA-SLK-DPTLKALTAVAL	GAPVTVLVAGSGCQAAEAAKLSGVAKVLAEDNAAYANLLAEPTADLIVSLAGG	YDALVAPSTANGKNIMPRVALLDVMVLSIDITKVVSPDTERPIYAGNAIQTVQSTDAKKVITVRTSTFQAA
>UniRef90 A0A1R1M6D9	SILVIAEHNG-SLK-GATLNTVAAQQOI	GGDIDLLVAGSGCGAVAAEAAKLVNGVAKVLLADADCYNHQLAENMAQLVLAELAG	YSHILAAATTTGKNFMPRVAAVLVDAQISDVIQVESAADTFKRPIYAGNAIATVQSSDSVKVITVRTSTFEPV
>UniRef90 A0A0P7XTU3	AVLLLAETVTDG-ALNRDATAKSVTAARTV	GGDVTVLCAGARAADAGAAKIDGVAKVLLAEDASLGHRLAEPAAALIVSLAGD	YSHIFAPATTDAKNILPRVALLDVMVLSIDVSGVVDADTFERPIYAGNAIQTVKSSDKVITVRTSTFQAA
>UniRef90 A0A1M6JK64	-VLLLEGEVTDG-TLNRDATAKAVAAKCP	GDVHVLVAGGCAEADAGKAAKIDGVAKVLLAEDALYGHRLAEPAAALIVSLAGD	YDHVAPATTDAKNILPRVALLDVMVLDVSGVVDADTFERPIYAGNAIQTVKSKDAKKVITVRTSTFQAA
>UniRef90 Q6N0L3	ATLLIAEHDNA-HLK-DATNKAMTAAAL	GGEVHVLVAGKGCDAVAAEAAKLDGAQVLLAEDAPAYEHQLEPAAALIVSLAGS	YDAIVAPATSRFKNVMPRAVALLDVMVLSIDVSEITKVVAPDTERPIYAGNAIQTVKSKDAKKVITVRTSTFQAT
>UniRef90 A0A2R4H0Q4	ATLLIAEHDNA-HLK-DATNKAMTAAAL	GGEVHVLVAGKGCDAVAAEAAKLDGAQVLLAEDAPAYEHQLEPAAALIVSLAGS	YDAIVAPATSRFKNVMPRAVALLDVMVLSIDVSEITKVVAPDTERPIYAGNAIQTVKSKDAKKVITVRTSTFQAT

Matrix Construction

Multiple Alignment

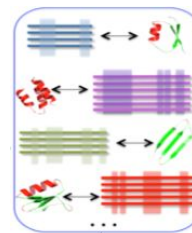


Profil

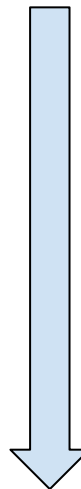


Profil creation

Profil : Matrix length of the input x 24 columns (20 amino acids + 1 gap + 3 secondary structure prediction : C, H et E)

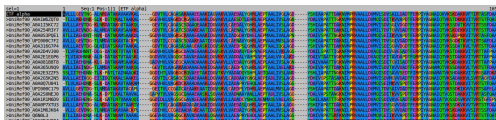


**HOMSTRAD
DATABASE**

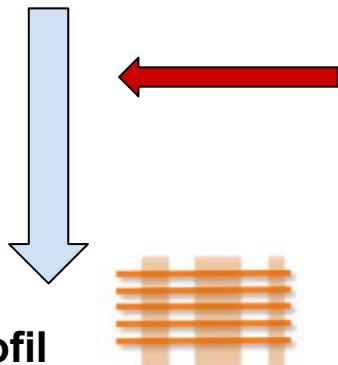


Matrix Construction

Multiple Alignment

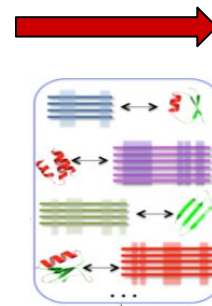


Profil



Profil creation

Profil : Matrix length of the input x 24 columns (20 amino acids + 1 gap + 3 secondary structure prediction : C, H et E)



HOMSTRAD DATABASE

[illegible]

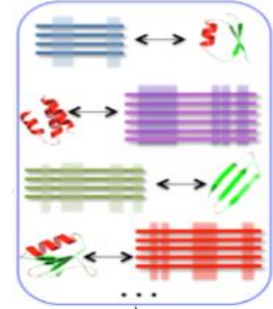
Profile-Profile Comparison



Profile-Profile Comparison



profil-profil alignment
semi-global (dot product &
Pearson correlation)

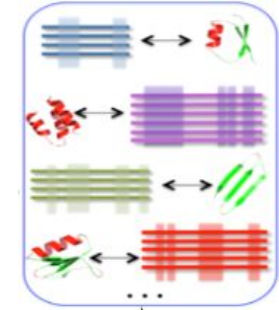


File: .foldrec

Profile-Profile Comparison



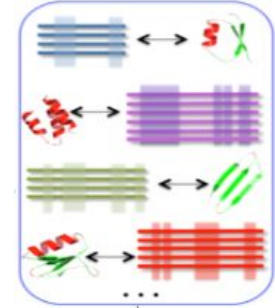
profil-profil alignment
semi-global (dot product & Pearson correlation)



File: .foldrec

*** HITS RANKED ***												
SEQUENCE QUERY FILE : ETF_alpha, 307												
#	Score	Ungaped_score	Pvalue_Q	Pscore	PQTscore	P-Value_T						
1	183.7905	X	X	X	X	X	X	X	X	X	X	X Arginosuc synth
2	182.5134	X	X	X	X	X	X	X	X	X	X	X PGI
3	181.479	X	X	X	X	X	X	X	X	X	X	X Glu_syn_central
4	181.2363	X	X	X	X	X	X	X	X	X	X	X aldedh
5	180.5341	X	X	X	X	X	X	X	X	X	X	X chorismate_bind
6	180.4503	X	X	X	X	X	X	X	X	X	X	X psA_psaB
7	179.8043	X	X	X	X	X	X	X	X	X	X	X Trypan_glycop
8	178.3679	X	X	X	X	X	X	X	X	X	X	X DNA_photolyase
9	177.9143	X	X	X	X	X	X	X	X	X	X	X ATP-synt
10	177.5459	X	X	X	X	X	X	X	X	X	X	X Ald_Xan_dh_2
11	177.4527	X	X	X	X	X	X	X	X	X	X	X COX1
12	177.0703	X	X	X	X	X	X	X	X	X	X	X Sec23_NC
13	176.9332	X	X	X	X	X	X	X	X	X	X	X cytochrome_b
14	176.9118	X	X	X	X	X	X	X	X	X	X	X CODH
15	176.3796	X	X	X	X	X	X	X	X	X	X	X Sec1
16	175.7833	X	X	X	X	X	X	X	X	X	X	X Gly_radical
17	175.3281	X	X	X	X	X	X	X	X	X	X	X GPGD
18	175.2631	X	X	X	X	X	X	X	X	X	X	X PK
19	173.8668	X	X	X	X	X	X	X	X	X	X	X lipoxygenase
20	173.6621	X	X	X	X	X	X	X	X	X	X	X AFOR
21	173.2912	X	X	X	X	X	X	X	X	X	X	X RNA_dep_RNA_pol
22	173.2841	X	X	X	X	X	X	X	X	X	X	X alk_phosphatase
23	173.2703	X	X	X	X	X	X	X	X	X	X	X PAP2
24	172.8111	X	X	X	X	X	X	X	X	X	X	X pgk
25	171.8797	X	X	X	X	X	X	X	X	X	X	X lyase_1

 X  score #1
 X  score #2



File: **.foldrec**

[illegible]

Query - First Template Alignment

ETF-alpha

4 LLGEVTNGALNRDATAKAVAAVKALGDVTVLCAGASAKAAAEAAKIAGVAKVLVAEDALYGHRLAEPTAALIVGLAGDYSHIAAPATTDAKNVMPRVAALLDVM
VLSDVSAILDADTFERPIYAGNAIQVVKSKDAKKVF-----TIRTASFDAAGEGGTAPVTETAAAADP
G-----LSSWVADEVAESDRPELTSARRVVS GGRGL-----GSKESFAIIEELADKLGAAVGASRAAVDSGYAPNDWQVGQTGKVVVA
PELYVAVGISGAIQHLAG-----MKDSKVIVAINKDE--EAPIFQIADYGLVGDLFSVPELTGKL 307

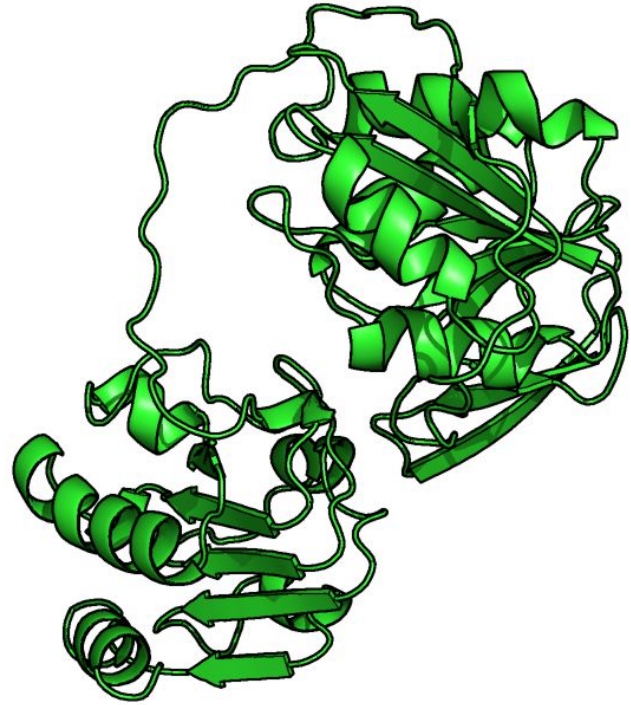
Arginosuc_synth

1 MKIVLAYSGGLDTSIILKWLKETYRAEVIAFTADIGQGEEVEEAREKALRTGASKAIALDLKEEFVRDFVFP--MMRAGAVYEGYYLLGTSIARPLIAKHLVRIA
EEEGAEIAHGATGKGNDQVR FELTAYALKPDIKVIAPWREWSFQGRKEMIAYAEAHGIPVPPYSMDANLLHISYEGGVLEDPWAEPPKGMFRMTQDPEEAPDAP
EYVEVEFFEGDPVAVNGERLSPAALLQRLNEIGGRHGVGRVDIVENRFVGMKSRGVYETPGGTILYHARRAVESLTL DREVLHQRDMLSPKYAELVYYGFWYAPE
REALQAYFDHVARSVTGVARLKLYKGNVYVGRKAPKSLYRGYDQKDAEGFIKIQUALRLRVRALVER 380

Structures of Query and First Template



Arginosuc_synth



ETF-alpha

Structures of Query and First Template



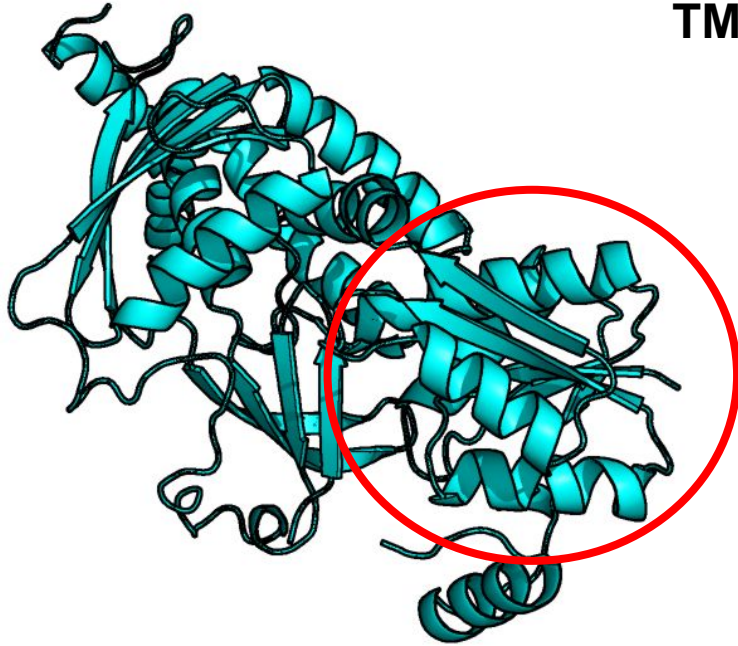
Arginosuc_synth



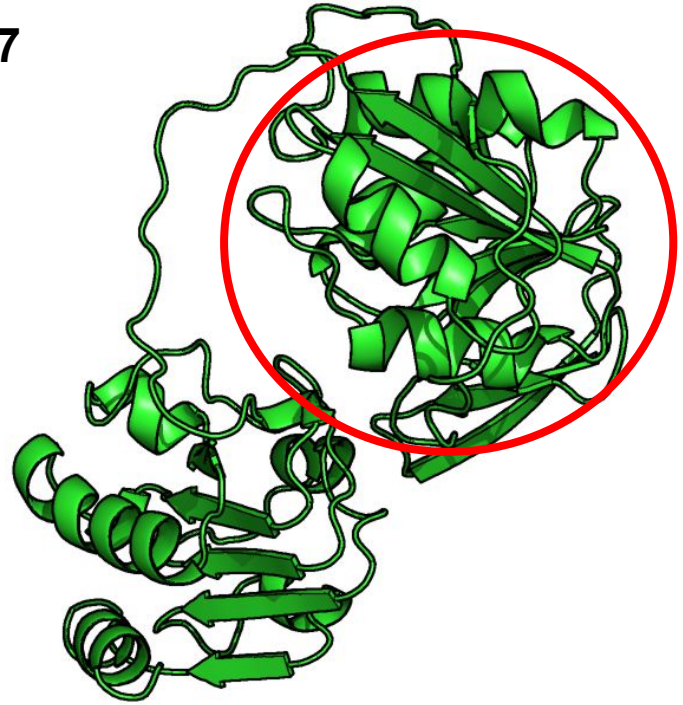
ETF-alpha

Structures of Query and First Template

TM-score = 0.37

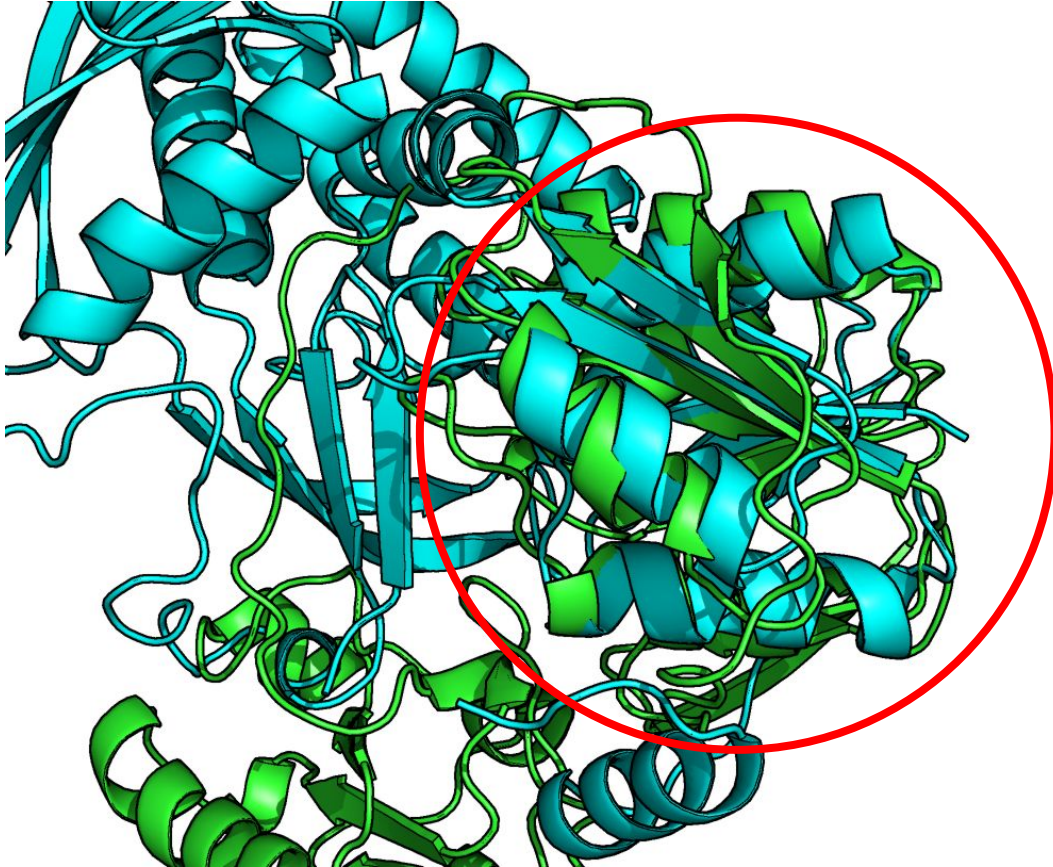


Arginosuc_synth



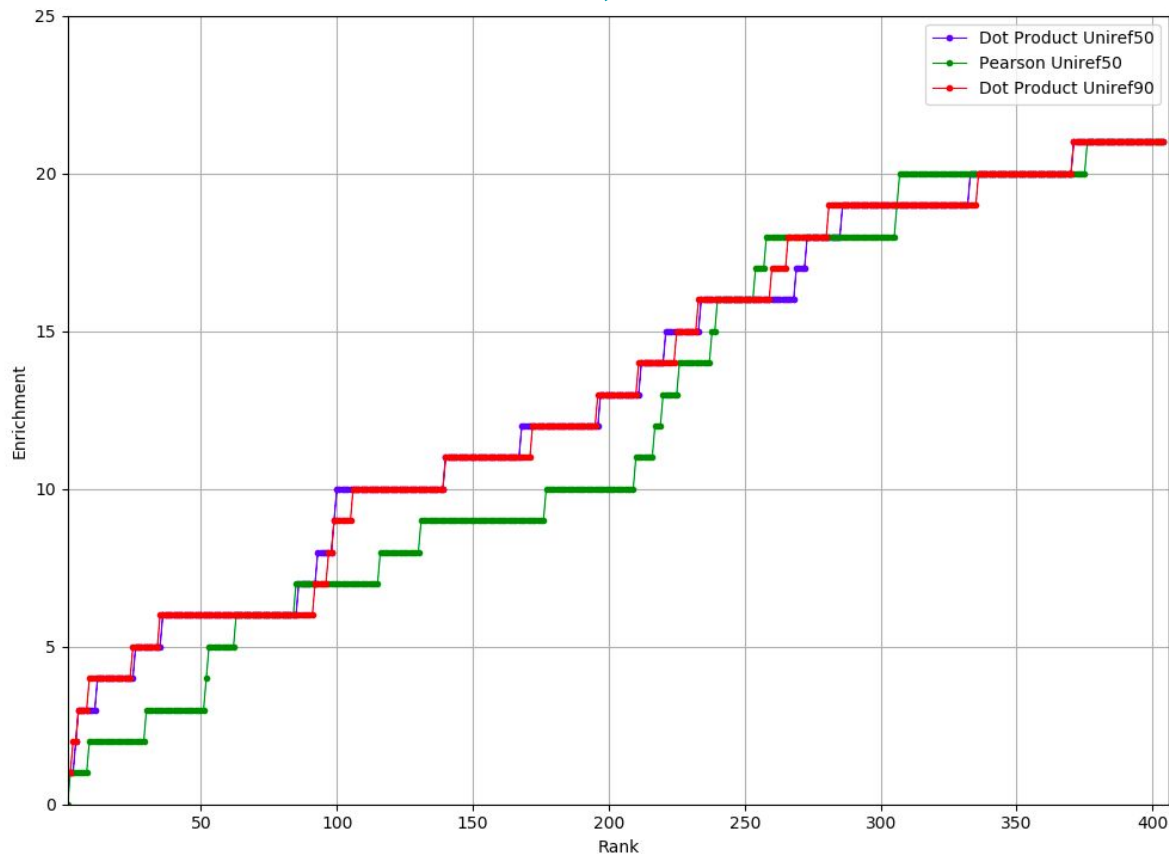
ETF-alpha

Structures of Query and First Template: Superposition



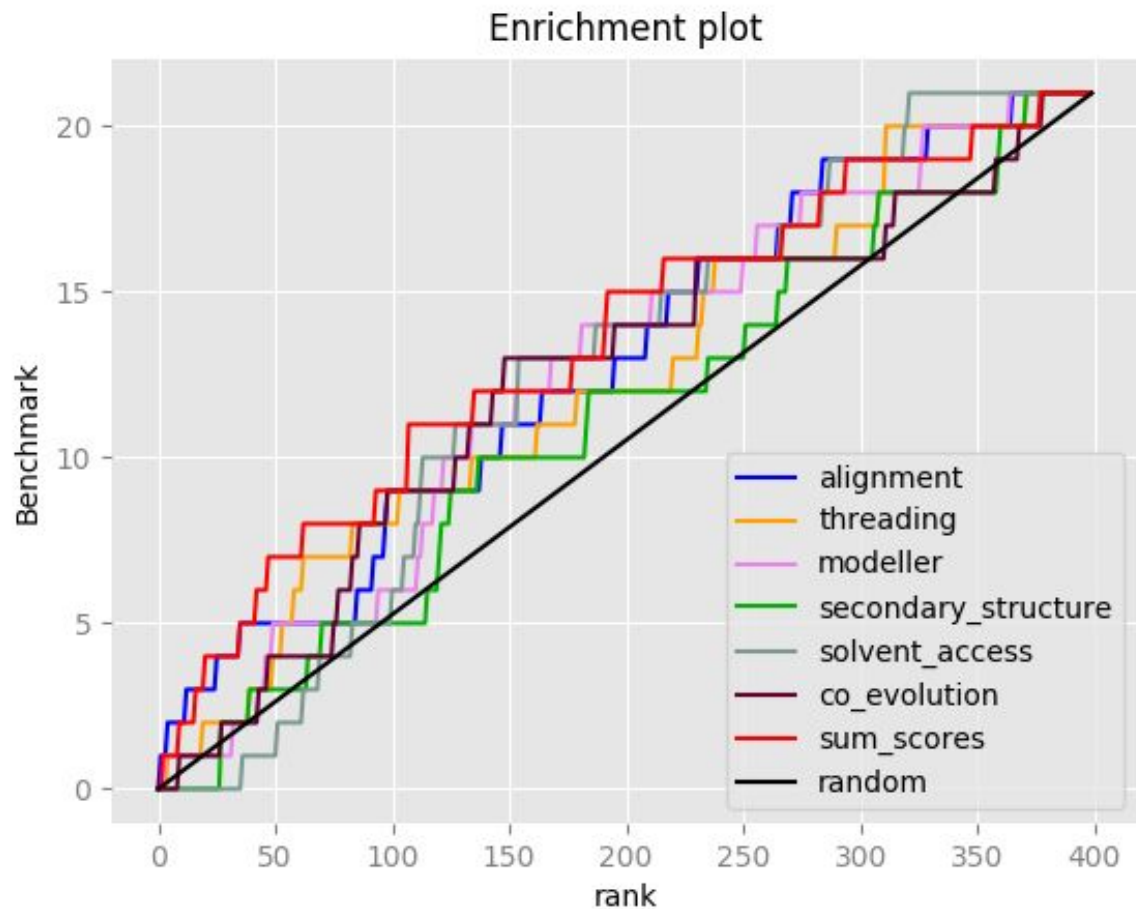
TM-score = 0.58

Results of 21 sequences (benchmark: upstream)

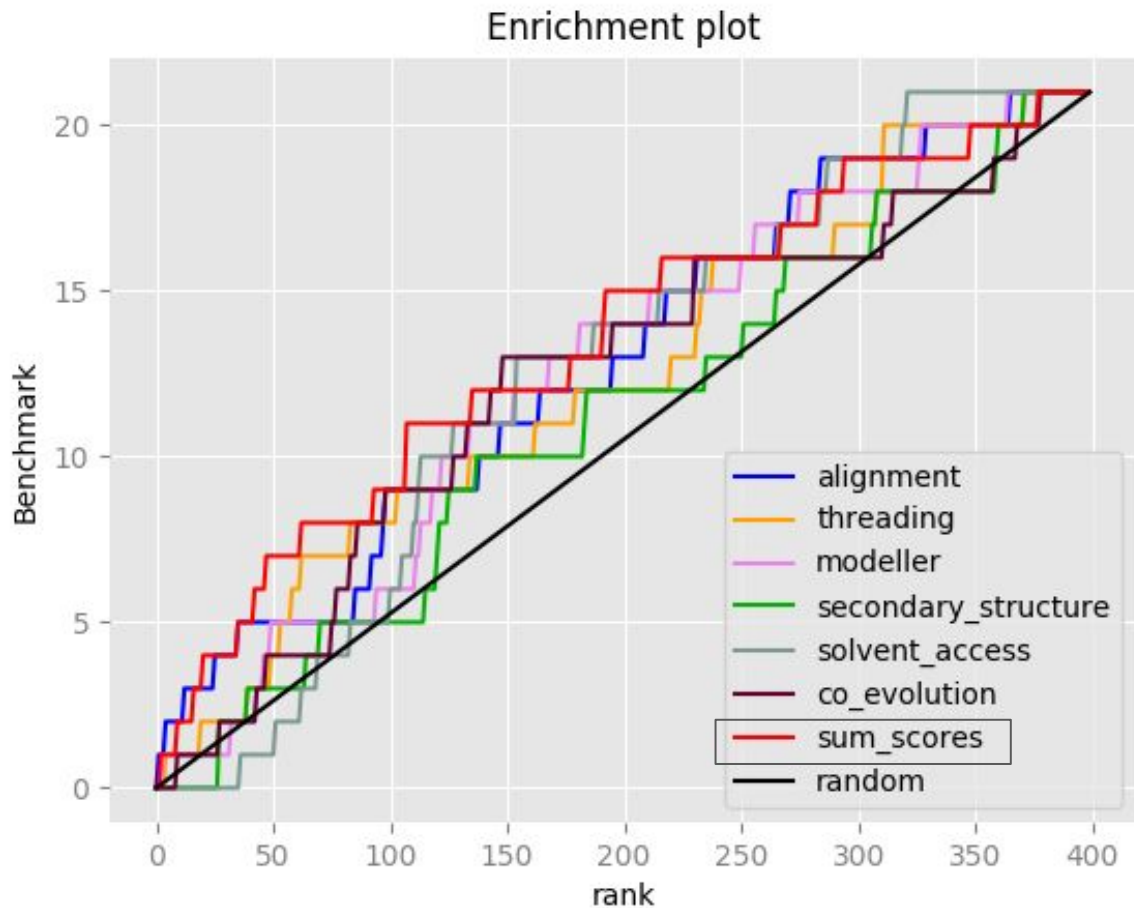


Enrichment Scores

Results of 21 sequences (upstream + downstream)



Results of 21 sequences (upstream + downstream)



Families are found at better positions

Exception, ex:
ETF-alpha (10th)

Results mysterious sequences

The whole pipeline was used on the **11 mysterious sequences**.

We noticed that :

- The results were **really divergent** amongst the teams
- **9 families** were found
- Amongst them :
 - **COX1** and **COX3** are subunits of the cytochrome c complexe, component of the mitochondrial respiratory chain.
 - **Cytochrome b** is the component of the ubiquinol-cytochrome c reductase complex, also part of the respiratory chain.

Discussion

PSIPRED	Prediction on multiple alignment of Query and hits
DSSP	On templates HOMSTRAD
Query Infos	Solvent Accessibility
Scoring Function (Profil-Profil)	PICASSO3Q

Thank you for your attention

