

Meet-U

A meeting story

edition 2020

When protein genes are sequenced

It's Time To Fold!

<http://www.meet-u.org/>

OBJECTIVES

By realizing a project from A to Z to address a challenging open question in biology

OBJECTIVES

By realizing a project from A to Z to address a challenging open question in biology

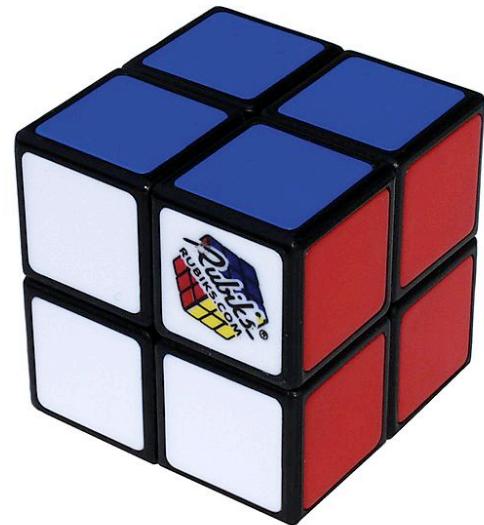


Prediction of protein 3D structures based on sequence similarity

OBJECTIVES

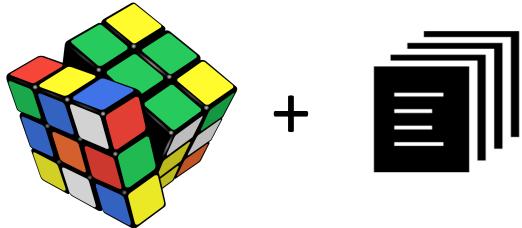


OBJECTIVES



FRAMEWORK

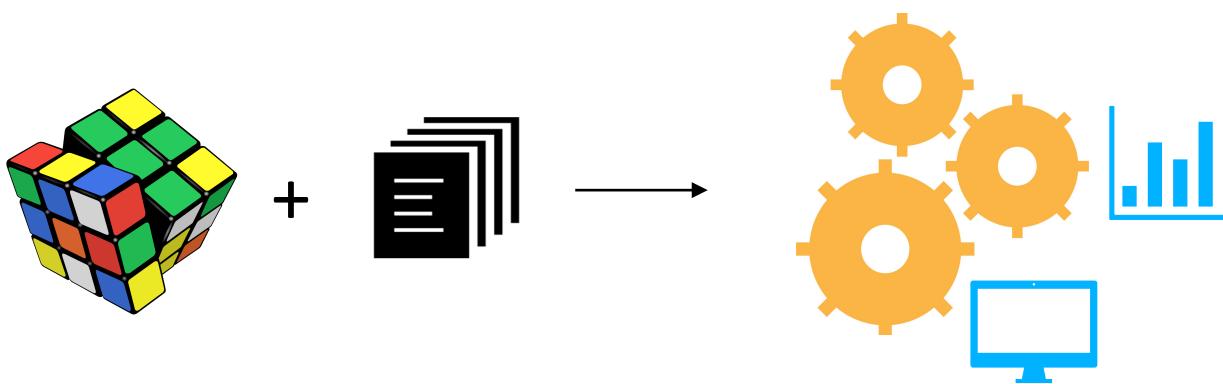
By realizing a project from **A** to **Z** to address a challenging open question in biology



Pedagogical team

FRAMEWORK

By realizing a project from **A** to **Z** to address a challenging open question in biology

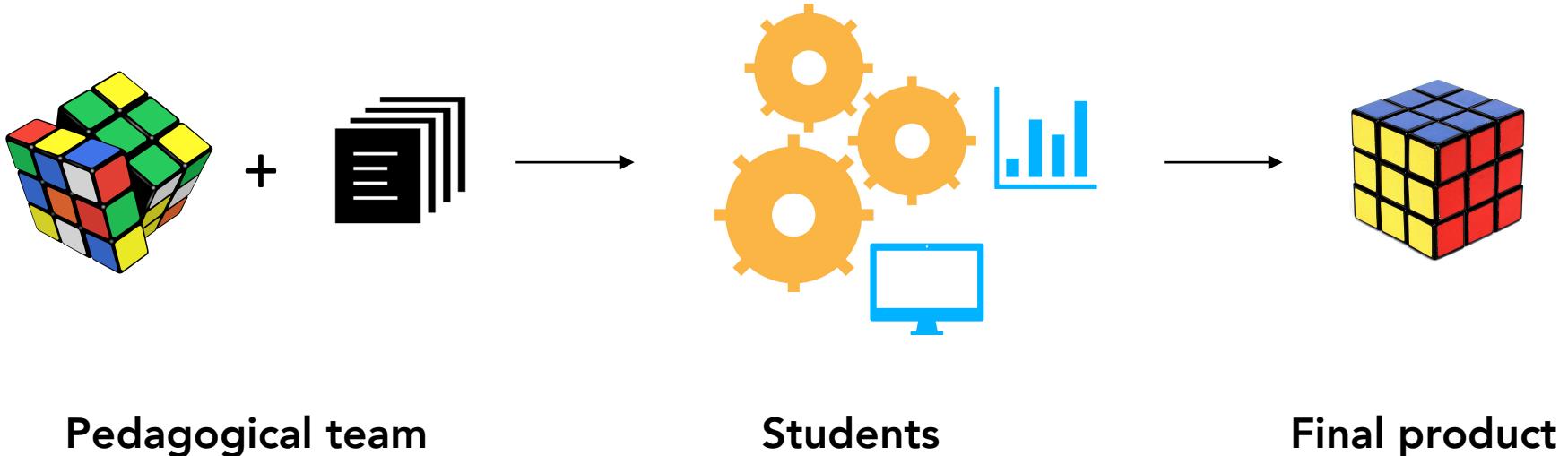


Pedagogical team

Students

FRAMEWORK

By realizing a project from **A** to **Z** to address a challenging open question in biology



WHO?

Pedagogical Teams:

- Univ. Pierre et Marie Curie (Sorbonne Université)
J. Bernardes, E. Laine, Y. Mohseni Behbahani, M. Muscat
- Univ. Paris Saclay
- G. Lelandais, A. Lopes, C. Papadopoulos

Students from:

- 2 masters of bioinformatics

Jury members from France and Europe
(invitations are being sent...)

HOW?

Team
upstream

HOW?

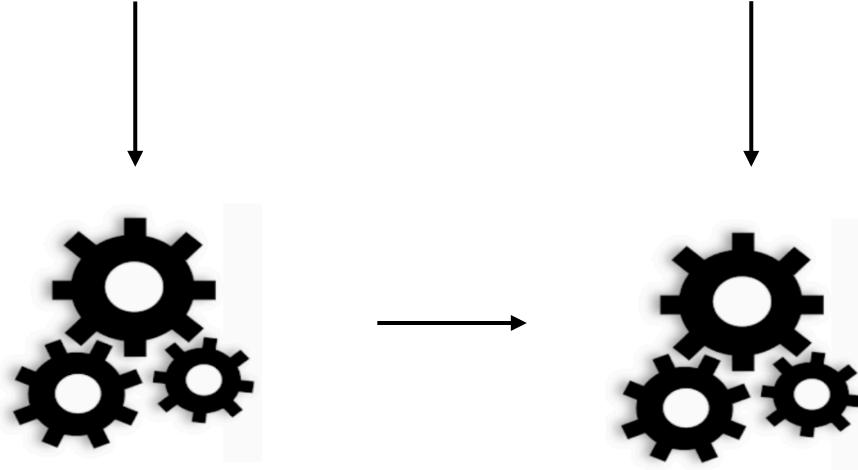
Team
upstream



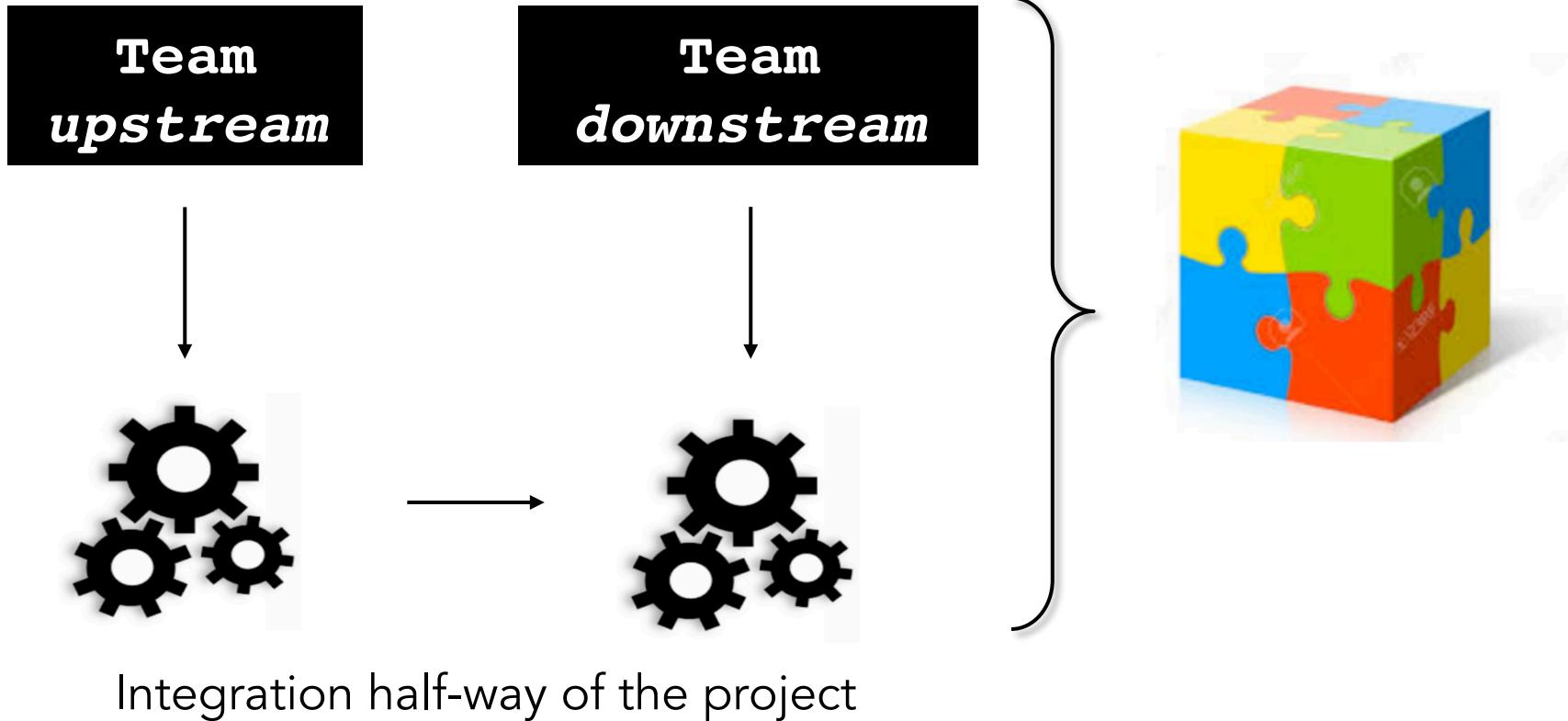
HOW?

Team
upstream

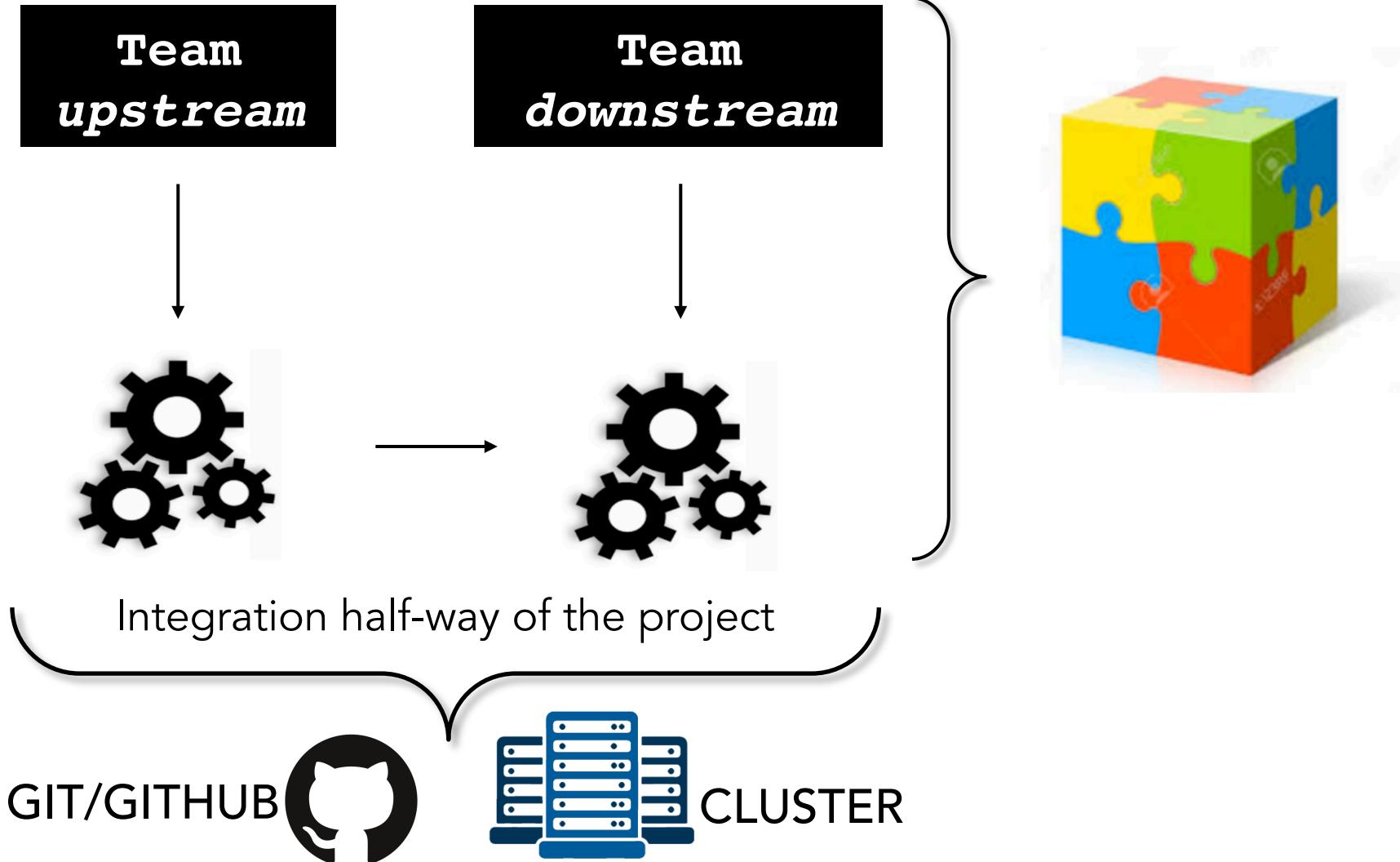
Team
downstream



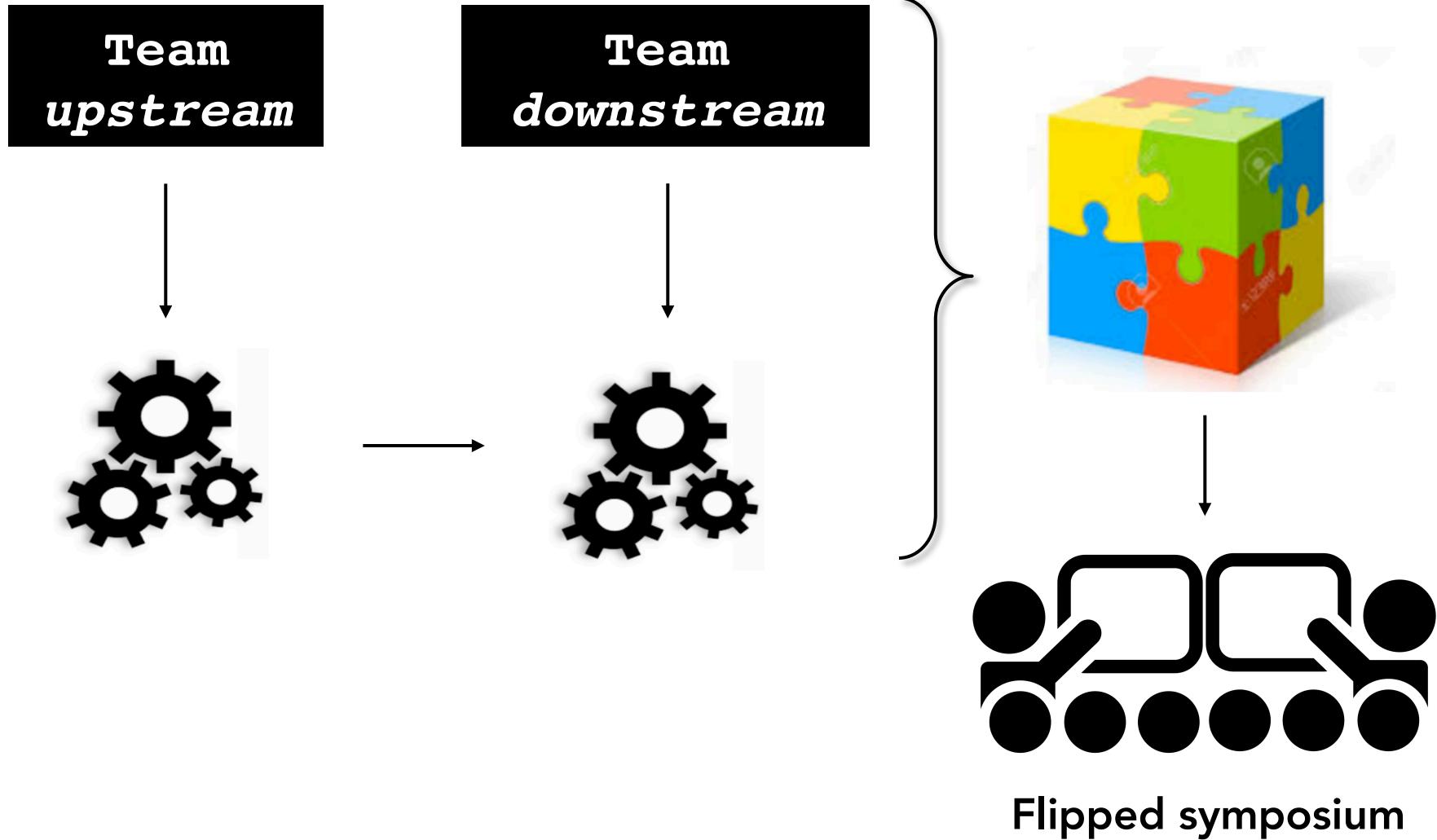
HOW?



HOW?

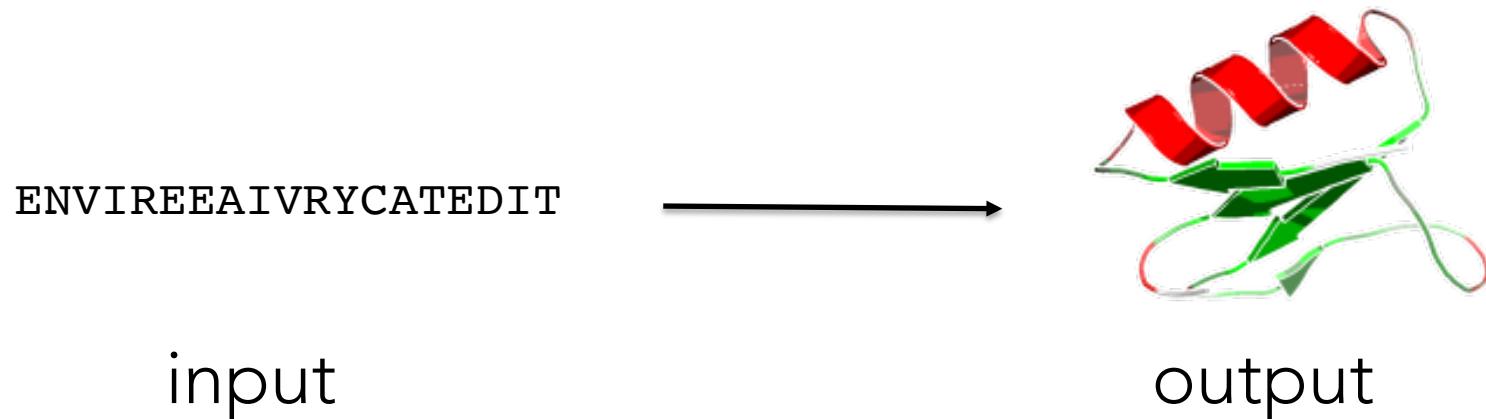


HOW?



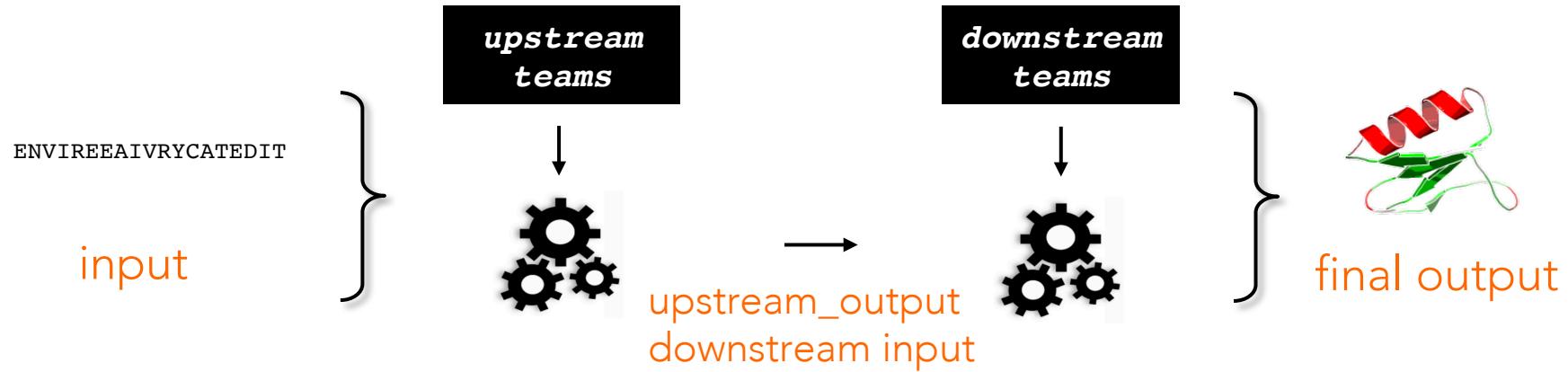
Meet-U 2020

- Prediction of the fold of a protein from its sequence



Meet-U 2020

3 constraints!



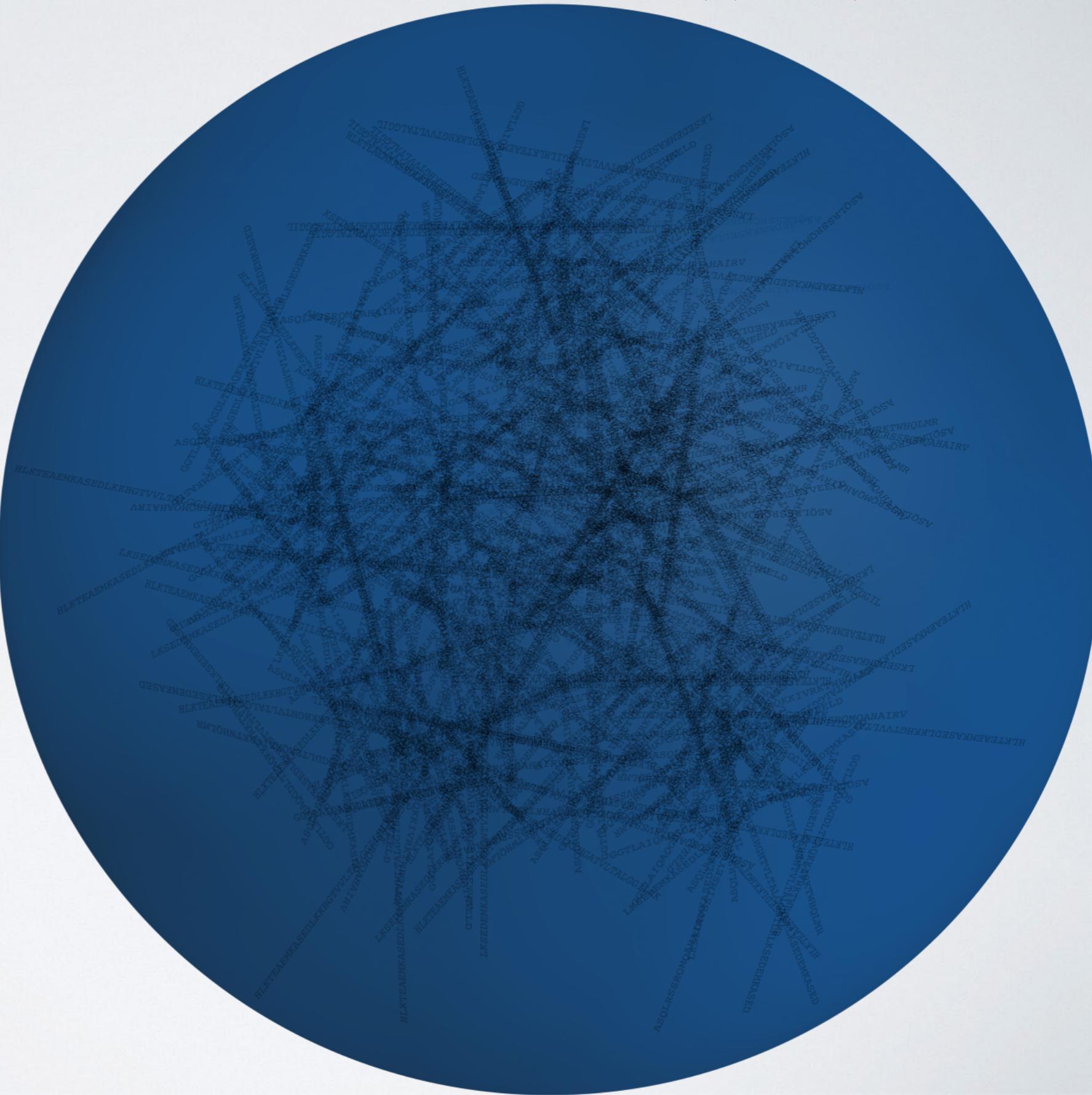
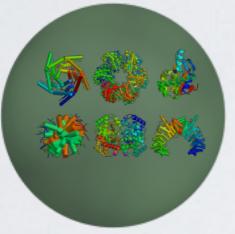
- 2 steps
 - **upstream teams:** building and querying sequence profiles against a database of protein structure families (dynamic programming)
 - **downstream teams:** assessment of profile-structure compatibility and identification of the most probable fold

~50 millions protein sequences

(UniprotKB/TrEMBL)

~100 000 protein structures

(Protein DataBank)



SEQUENCE/STRUCTURE CONSERVATION

Similar sequences

GLLT-ESQLVIKSAWEEFN...
GVSDGDWNIVLNTSGKVEA...
GDLSEL-QKKIRSTWHQLM...
GISDGEWQLVVNASGKVES...



Similar structures



Sequence similarity : 30 %

Structural similarity : ~70 %
(secondary structure conservation)

Structure is more conserved than sequence

PRINCIPLE

Protein A ● ● Protein B



Evolutive distance

PRINCIPLE

Protein A  Protein B



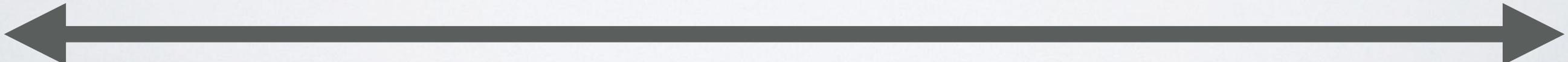
Evolutive distance

PRINCIPLE

Protein A

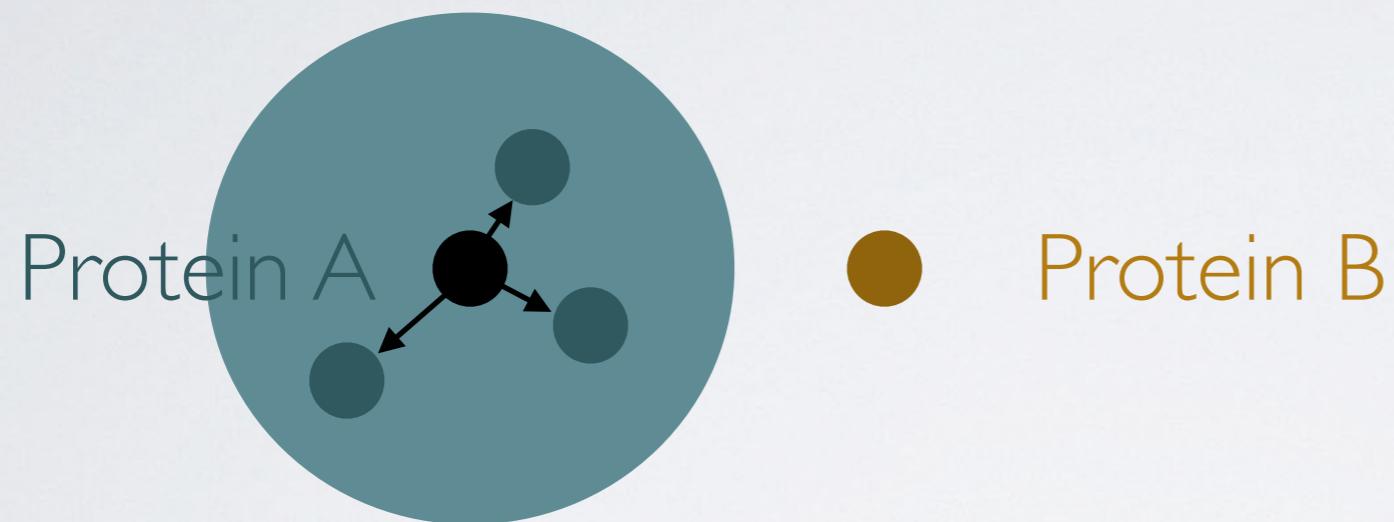


Protein B

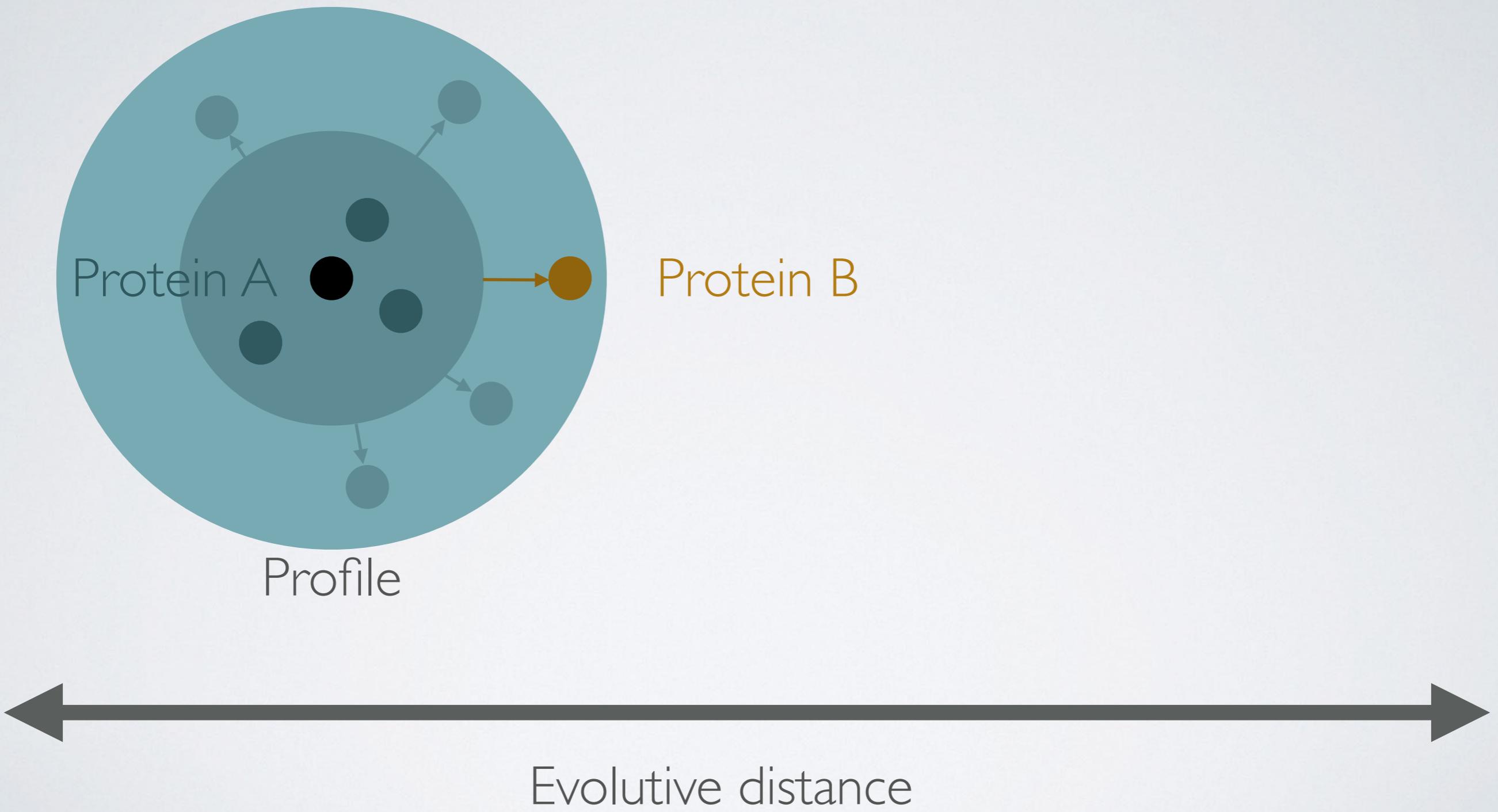


Evolutive distance

PRINCIPLE



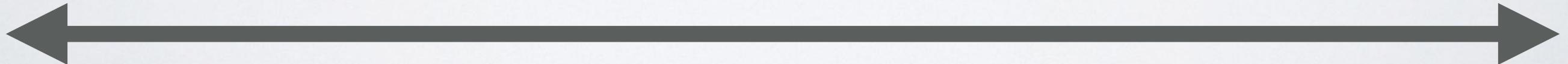
PRINCIPLE



PRINCIPLE

Protein A ●

● Protein C



Evolutive distance

PRINCIPLE

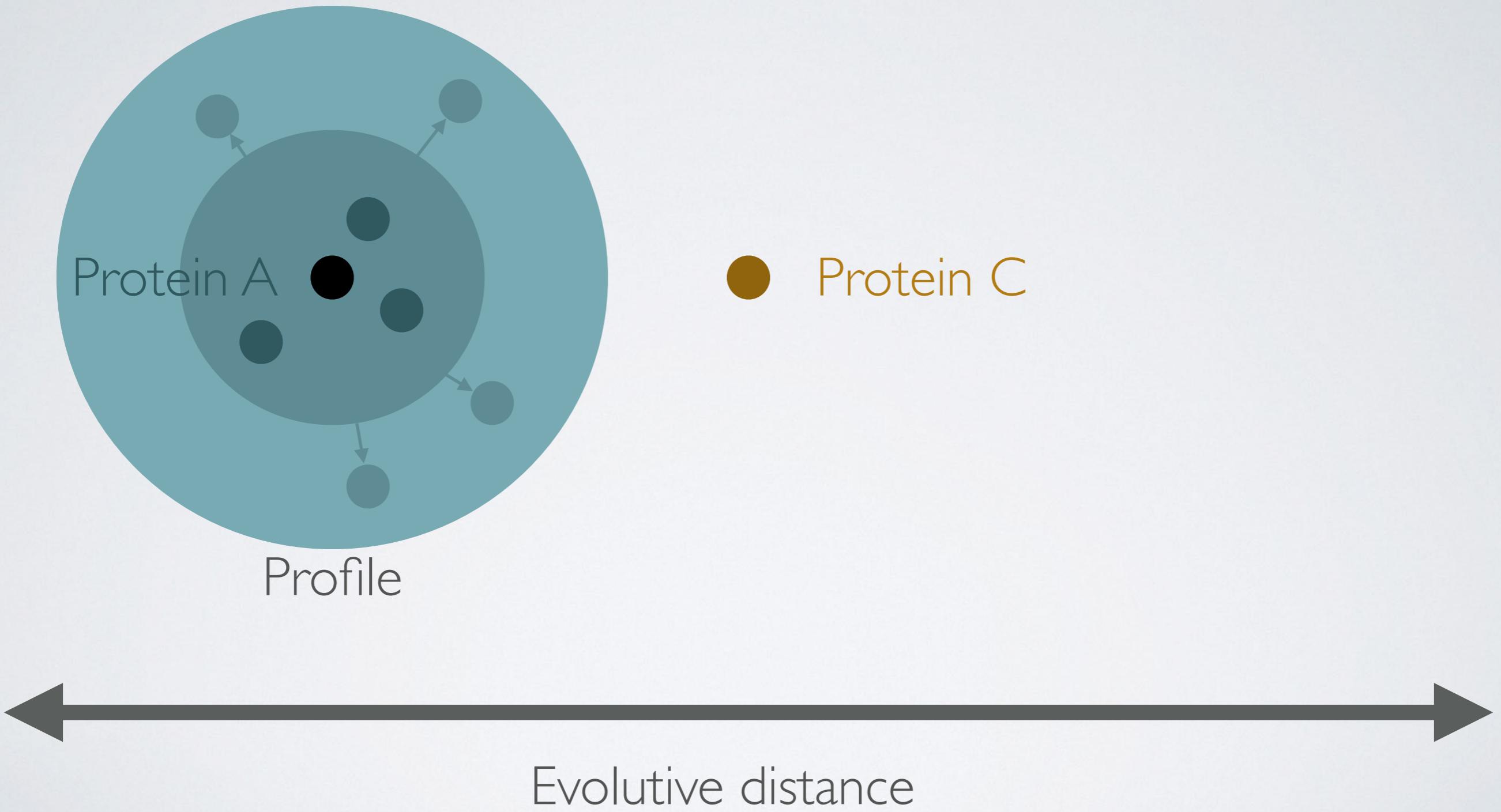


Profile

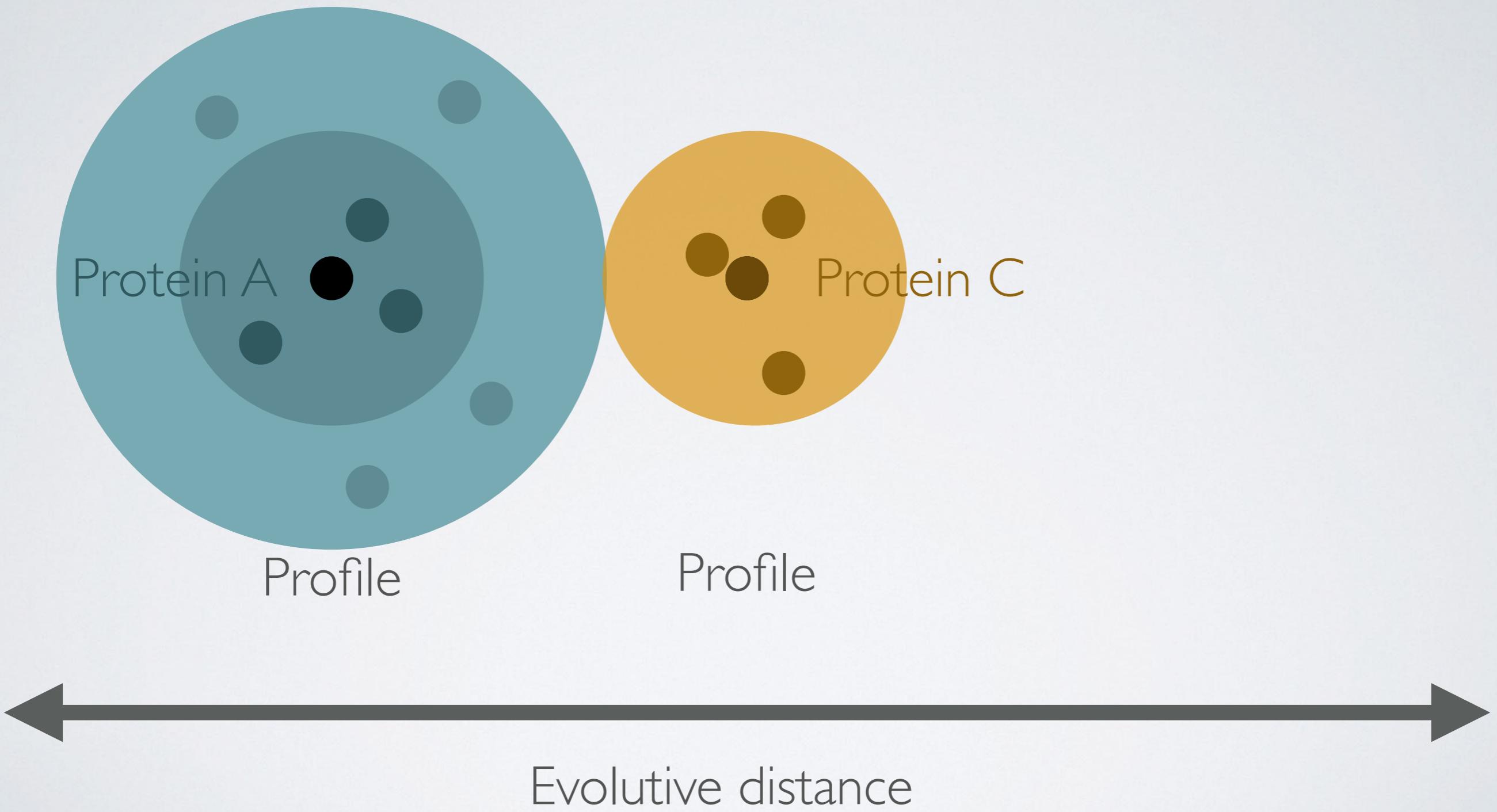


Evolutive distance

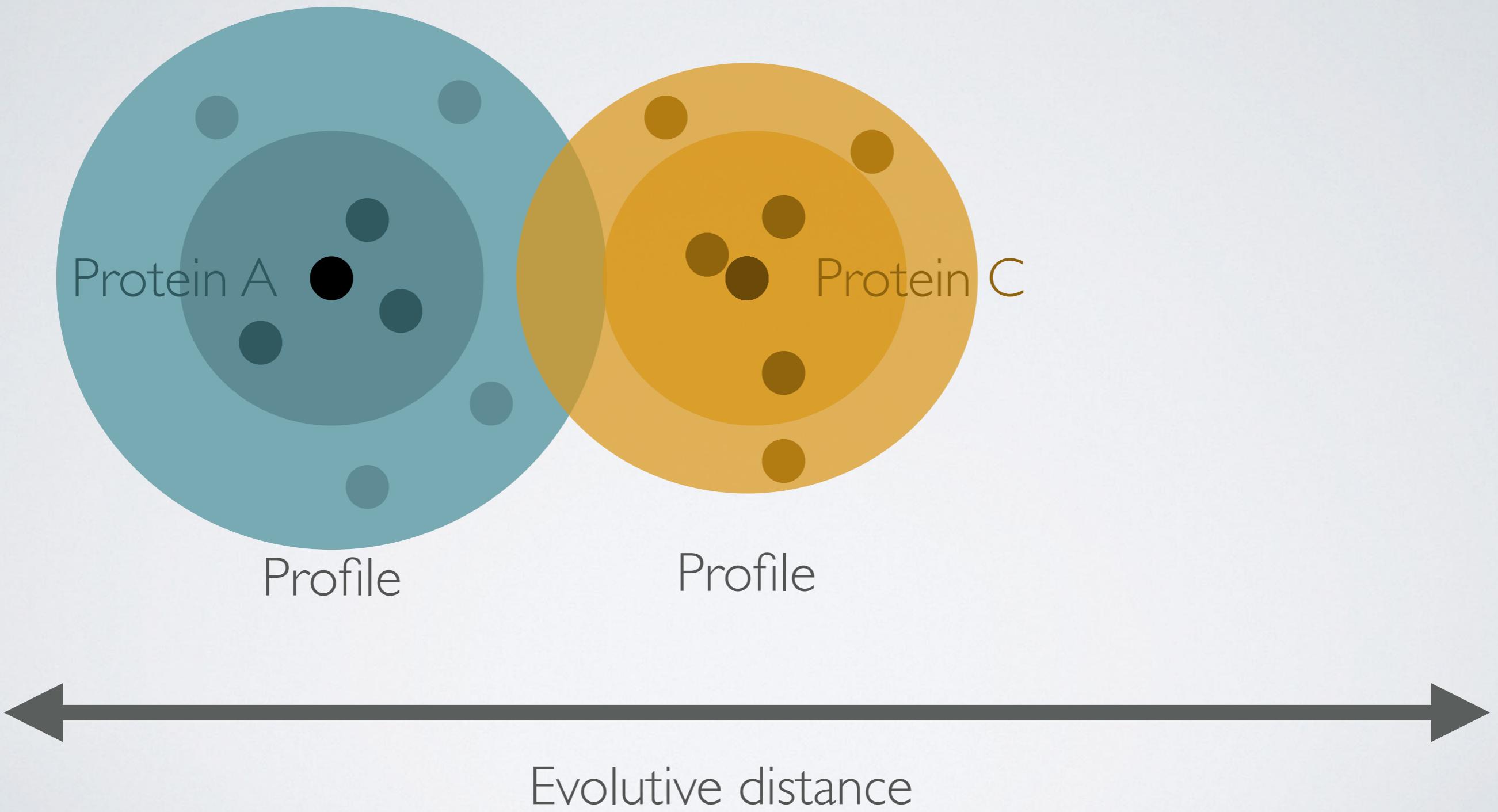
PRINCIPLE



PRINCIPLE



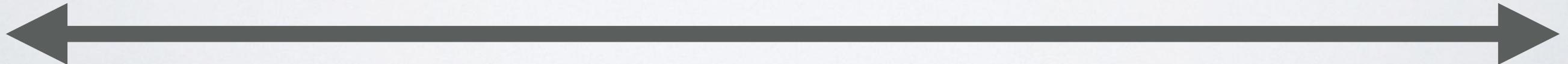
PRINCIPLE



PRINCIPLE

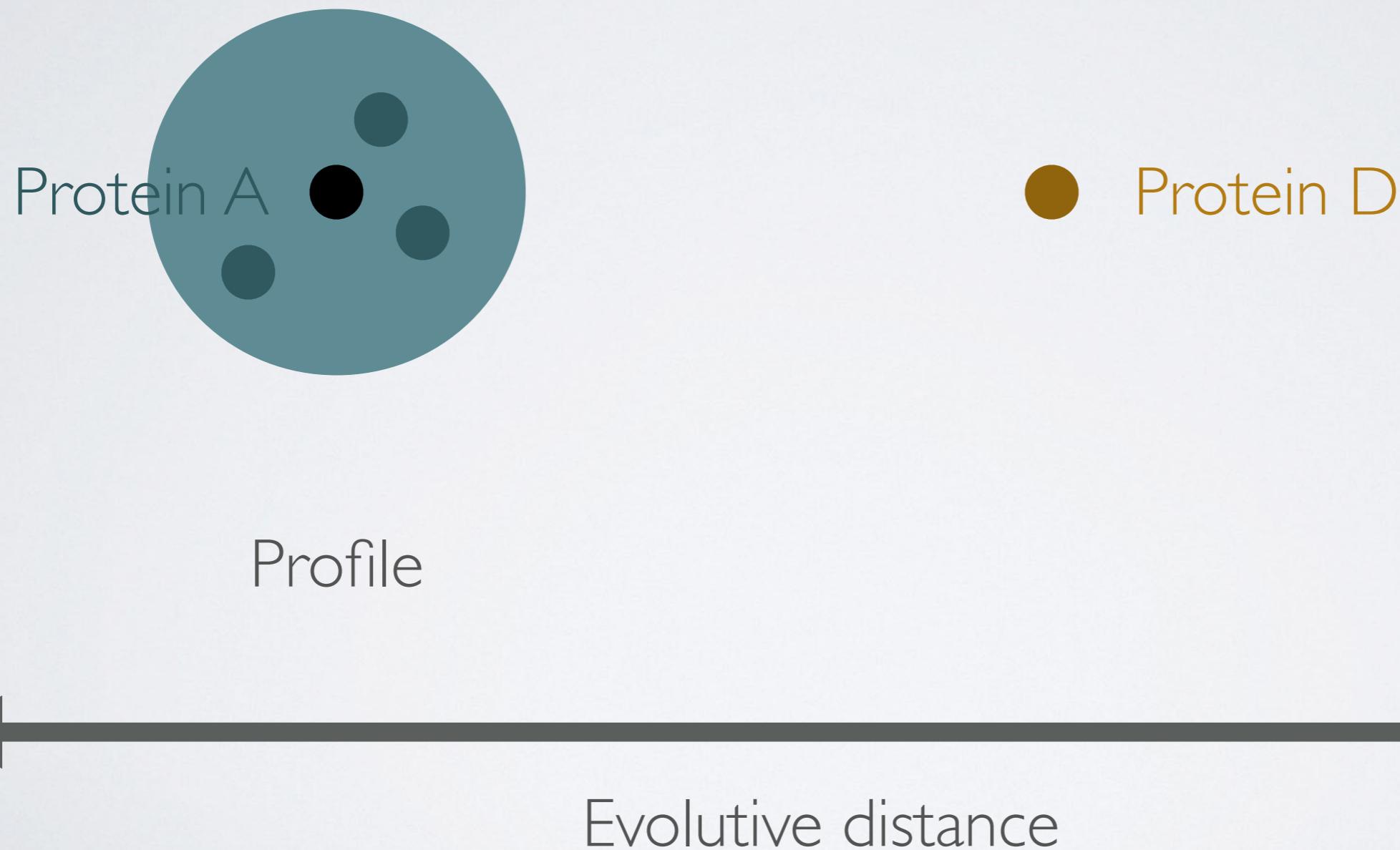
Protein A ●

● Protein D

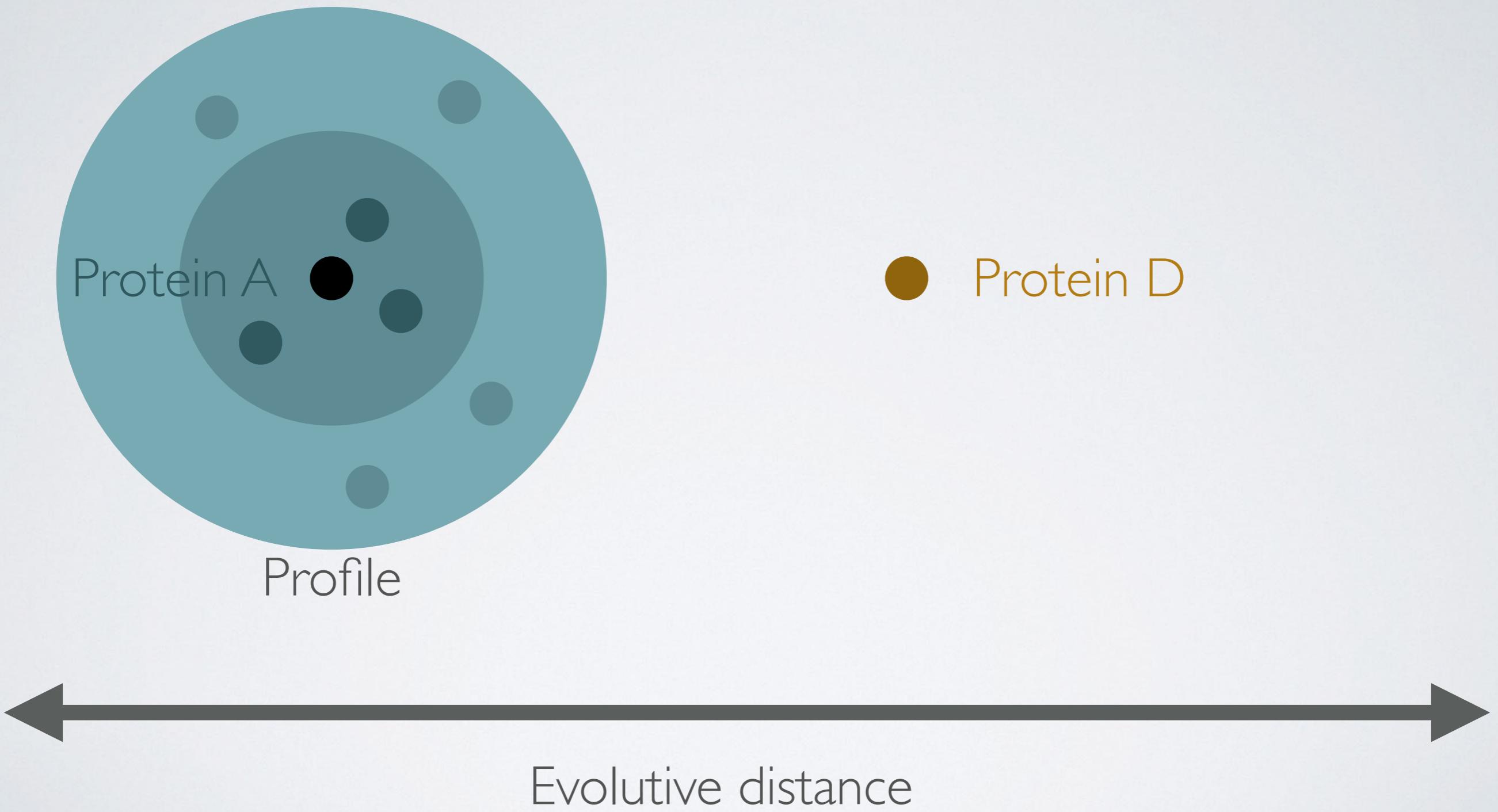


Evolutive distance

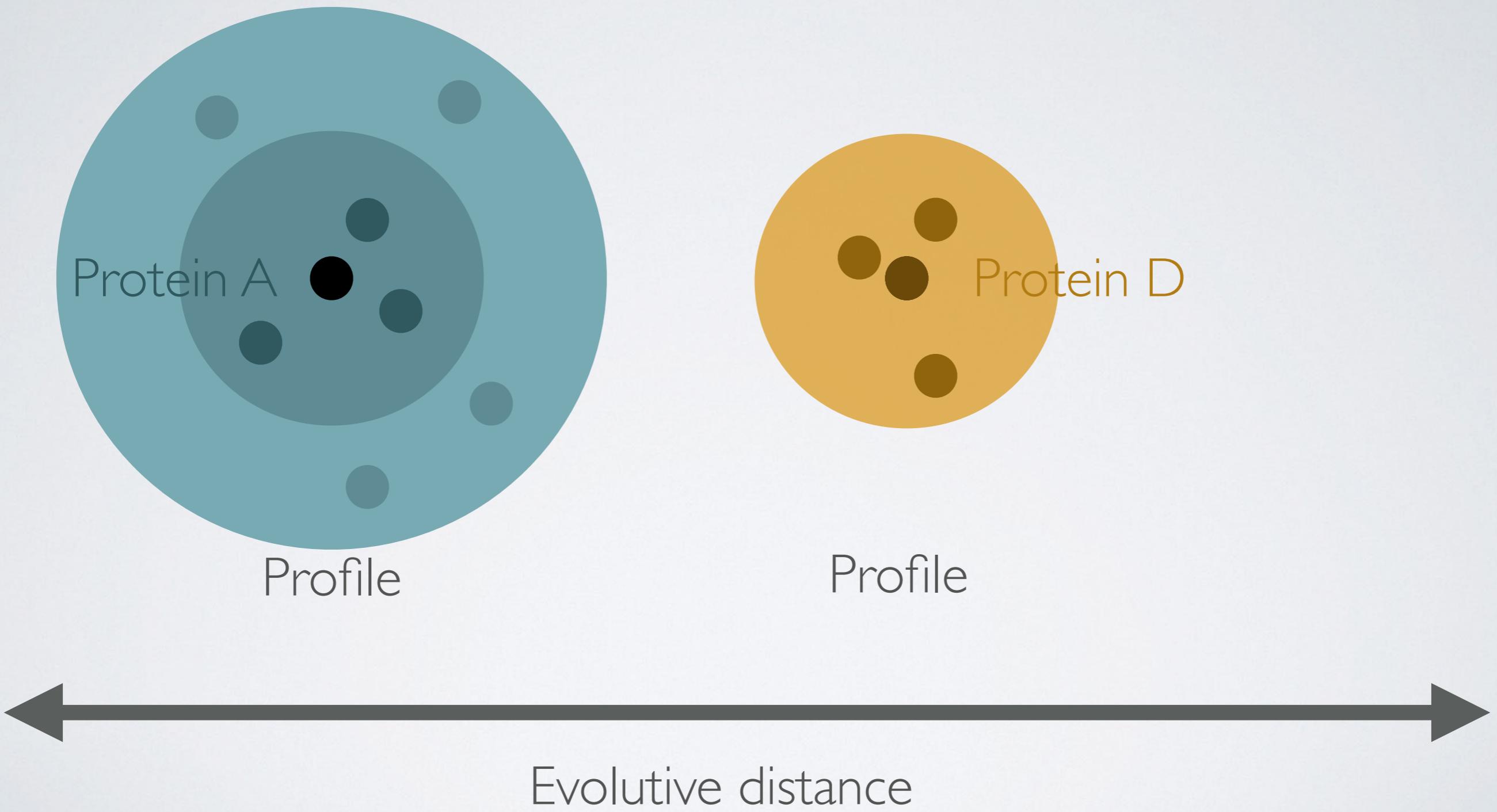
PRINCIPLE



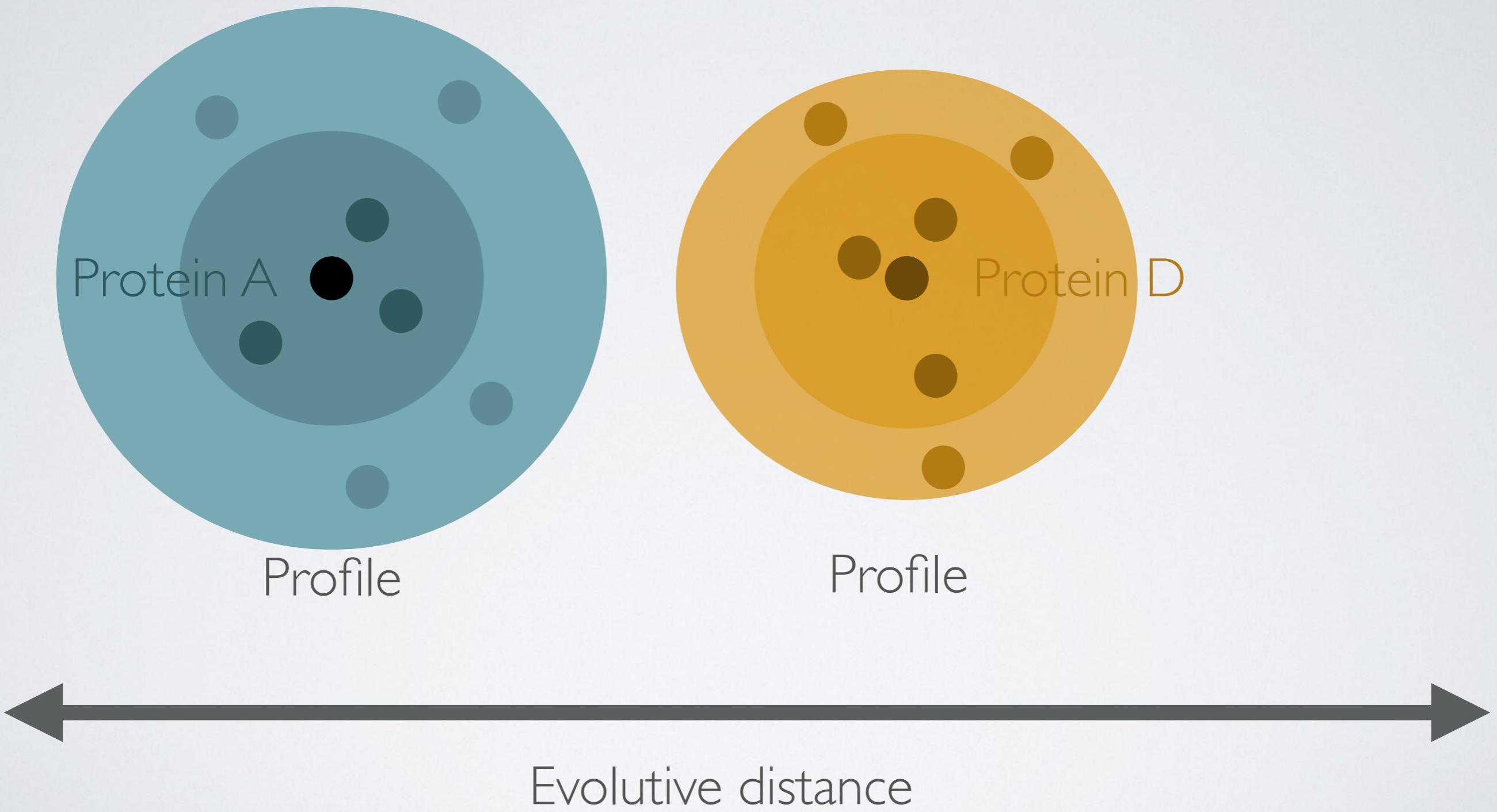
PRINCIPLE



PRINCIPLE



PRINCIPLE



HISTORY OF FOLD RECOGNITION

Query Sequence

RHPVVMGNWKLNGKEMVVDLNLNGLNAELEGVTGVDV

Template
Sequences

KTVRWMLKNADSQE

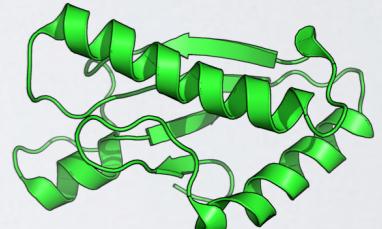
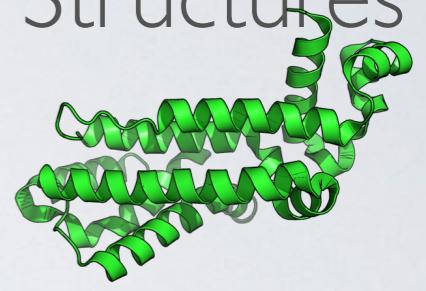
ERHKWMLKNADSQ

LHHKWMLKNADSQL

MLHKWMLKNADSQL

GFSKWMLKNADSQL

Template
Structures

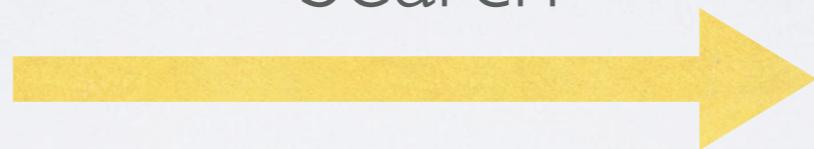


HISTORY OF FOLD RECOGNITION

Query Sequence

RHPVVMGNWKLNGKEMVVDLNLNGLNAELEGVTGVDV

Search



KTVRWMLKNADSQE

ERHKWMLKNADSQ

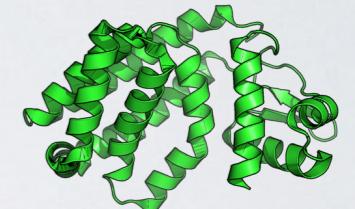
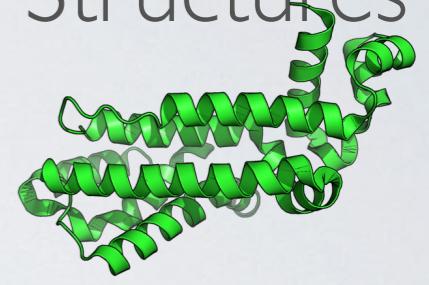
LHHKWMLKNADSQL

MLHKWMLKNADSQL

GFSKWMLKNADSQL

Template Sequences

Template Structures



Sequence-Sequence

SSEARCH (Smith & Waterman)

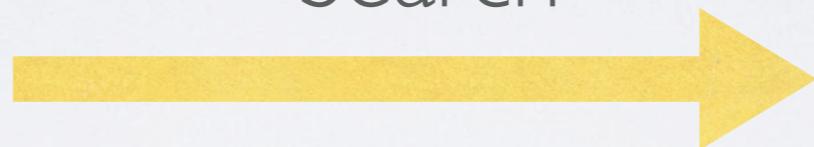
BLAST

HISTORY OF FOLD RECOGNITION

Query Sequence

RHPVVMGNWKLNGKEMVVDLNLNGLNAELEGVTGVDV

Search



KTVRWMLKNADSQE

ERHKWMLKNADSQ

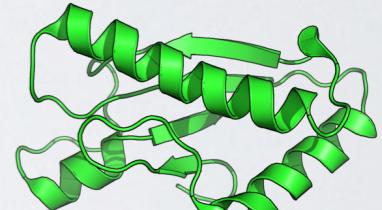
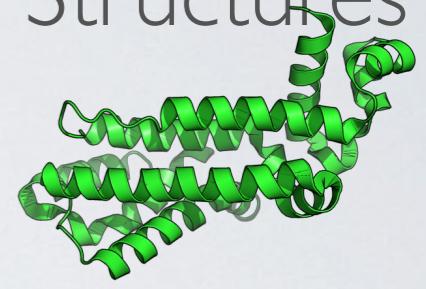
LHHKWMLKNADSQL

MLHKWMLKNADSQL

GFSKWMLKNADSQL

Template Sequences

Template Structures

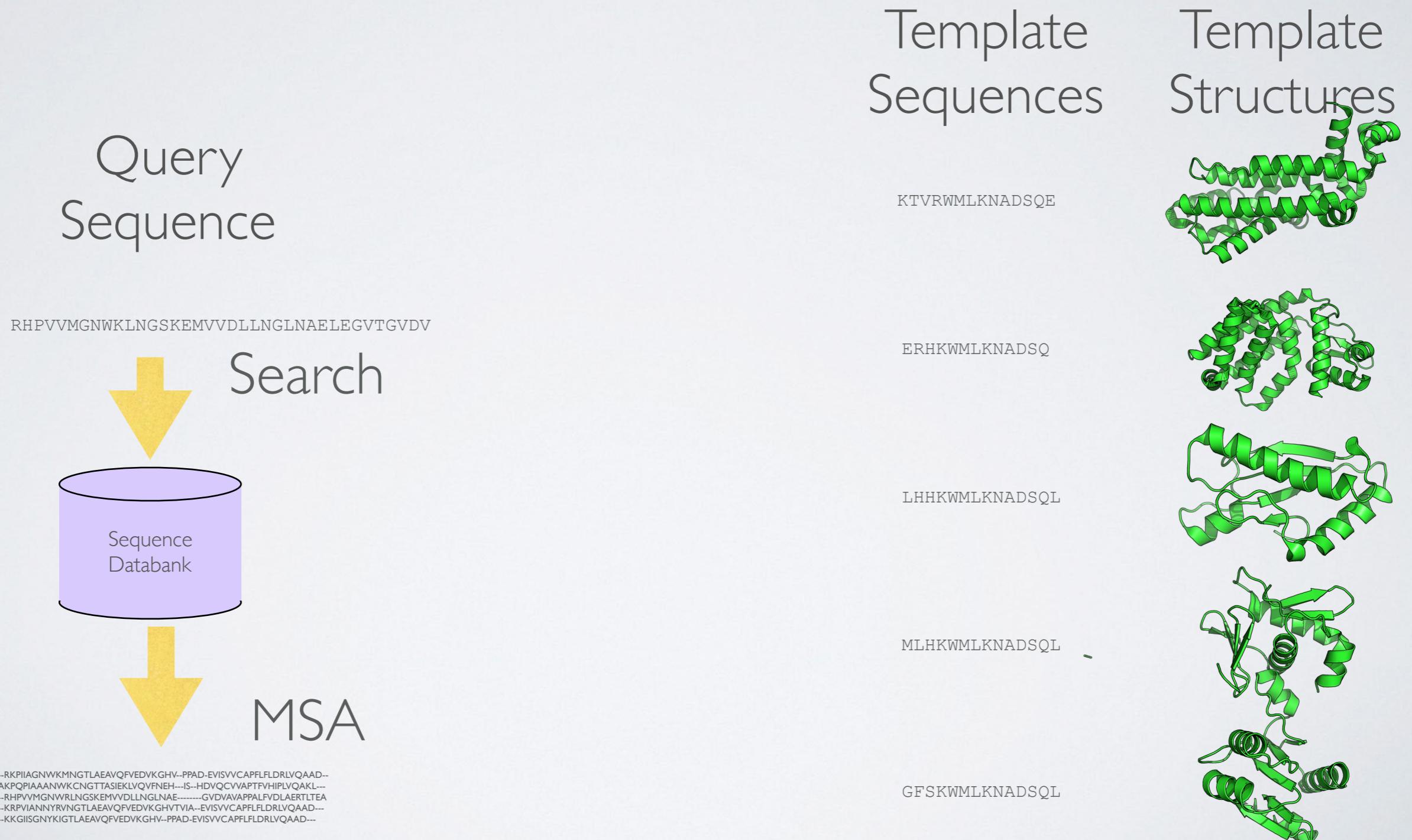


Sequence-Sequence

SSEARCH (Smith & Waterman)

BLAST

HISTORY OF SEQUENCE/SEQUENCE METHODS



HISTORY OF SEQUENCE/SEQUENCE METHODS

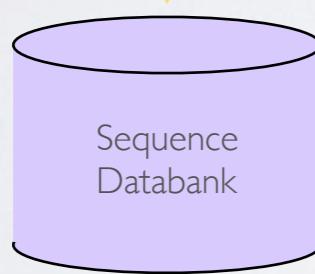


HISTORY OF SEQUENCE/SEQUENCE METHODS

Query
Sequence

RHPVVMGNWKLNGSKEMVV DLLNGLNAELEGVTGVDV

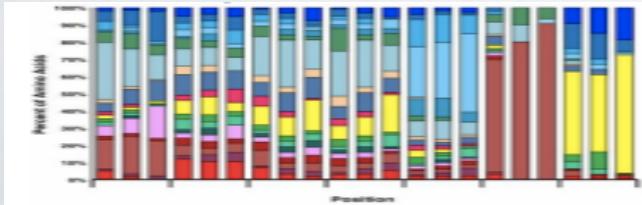
Search



Sequence
Databank



Profile



Template
Sequences

KTVRWMLKNADSQE

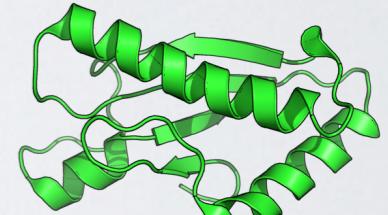
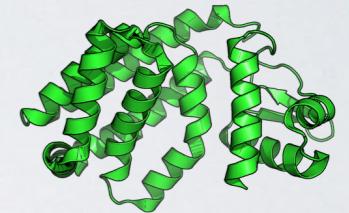
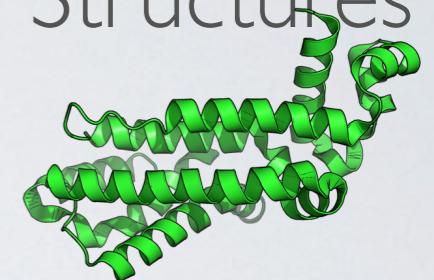
ERHKWMLKNADSQ

LHHKWMLKNADSQL

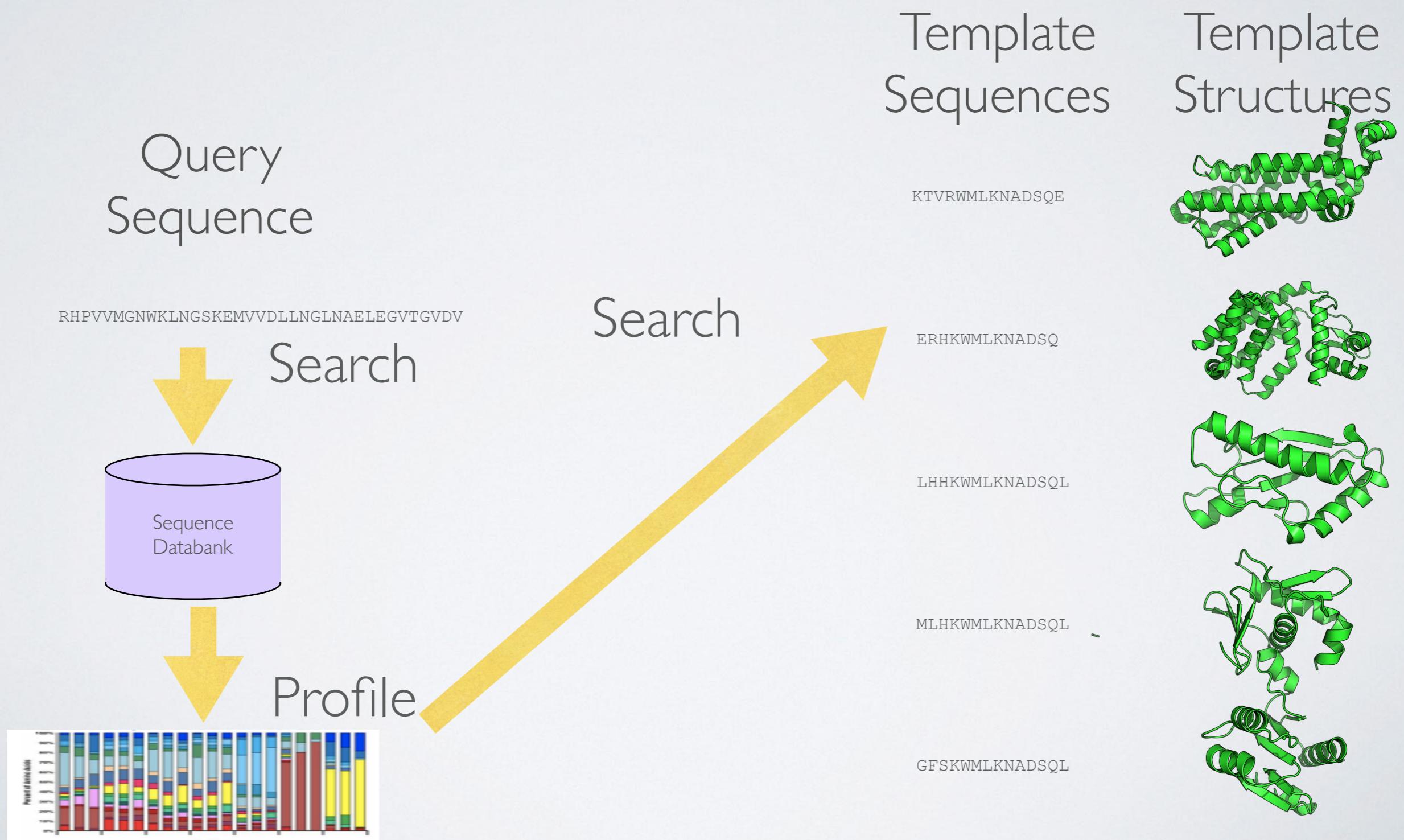
MLHKWMLKNADSQL

GFSKWMLKNADSQL

Template
Structures



HISTORY OF SEQUENCE/SEQUENCE METHODS

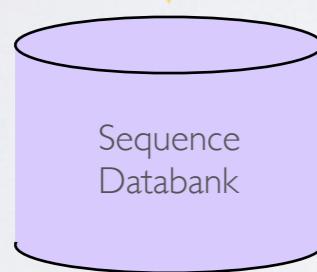


HISTORY OF SEQUENCE/SEQUENCE METHODS

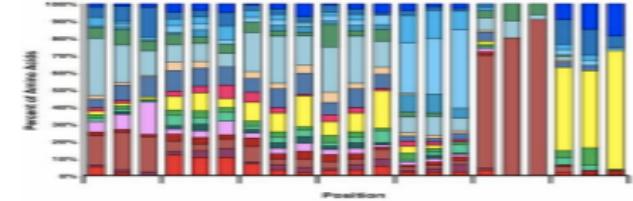
Query
Sequence

RHPVVMGNWKLNGSKEMVVDLLNGLNAELEGVTGVDV

Search



Profile



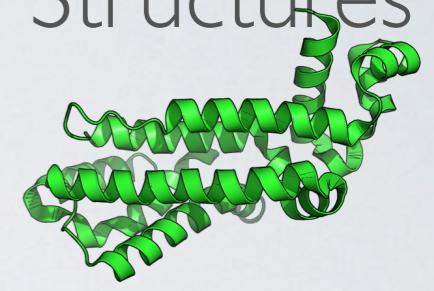
Sequence-Profile
PSI-BLAST
SAM
HMMER

Search

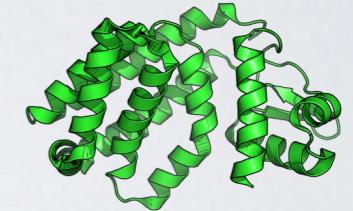
Template
Sequences

KTVRWMLKNADSQE

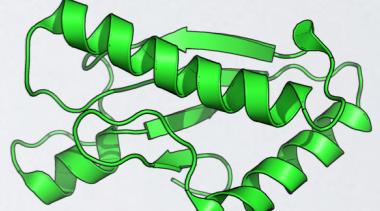
Template
Structures



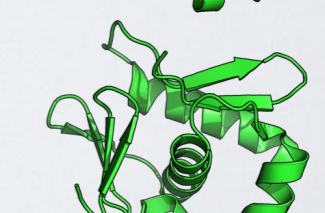
ERHKWMLKNADSQ



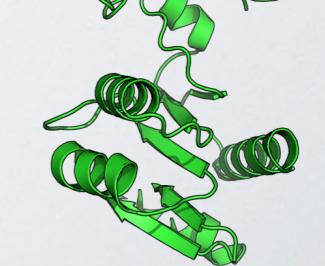
LHHKWMLKNADSQL



MLHKWMLKNADSQL



GFSKWMLKNADSQL



TRANSITIVITY



if $\text{query} \sim H$ and $H \sim \text{target}$ then
 $\text{query} \sim \text{target}$

HISTORY OF SEQUENCE/ SEQUENCE METHODS

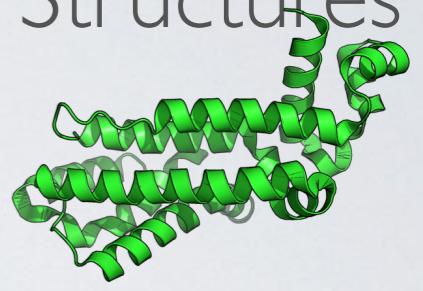
Query Sequence

RHPVVMGNWKLNGSKEMVVVDLLNGLNAELEGVTGVDV

Template
Sequences

KTVRWMLKNADSQE

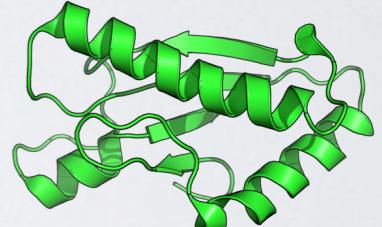
Template
Structures



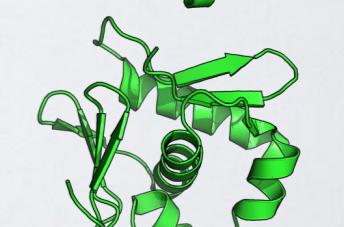
ERHKWMLKNADSQ



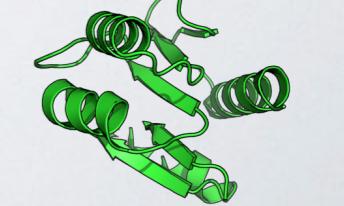
LHHKWMLKNADSQL



MLHKWMLKNADSQL



GFSKWMLKNADSQL



HISTORY OF SEQUENCE/ SEQUENCE METHODS

Query Sequence

RHPVVMGNWKLNGKEMVVDLNLNGLNAELEGVTGVDV

Template
Sequences

RTVKWMLKNAESQE
KTVKWLKNNTDSQE
KTVRWMLKNADSQE

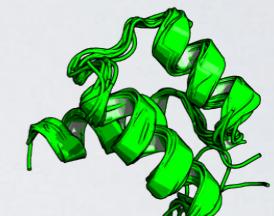
Template
Structures

ERHKWMLKNADSQL
ERHKWMLKNADSQL

LHHKWMLKNADSQL
LHHKWMLKNADSQL

MLHKWMLKNADSQL
-LHKWMLKNA----

GFSKW--KNADSQL
GFSKWMLKNADSQL



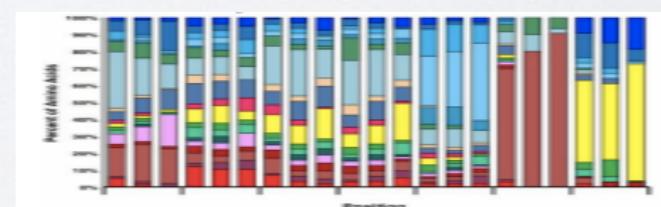
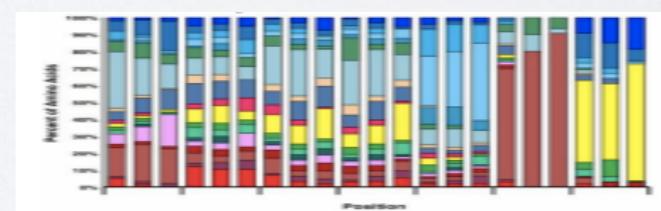
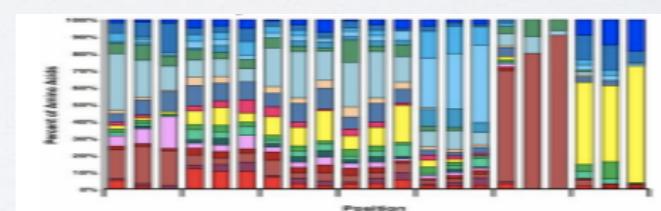
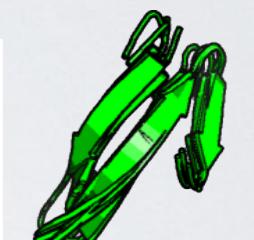
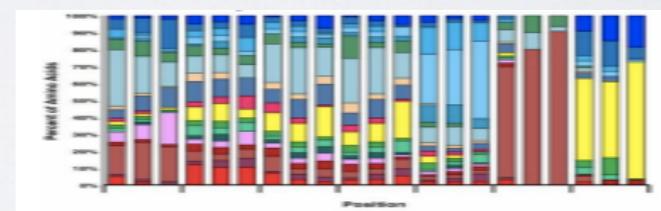
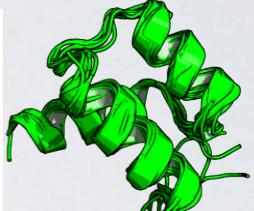
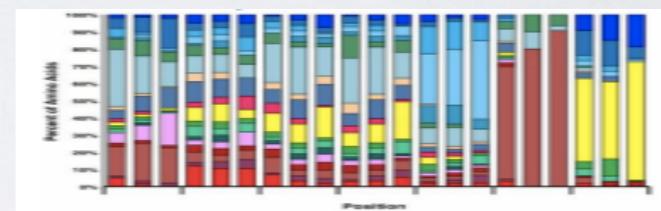
HISTORY OF SEQUENCE/ SEQUENCE METHODS

Template
Profiles

Template
Structures

Query Sequence

RHPVVMGNWKLNGSKEMVV DLLNGLNAELEGVTGVDV

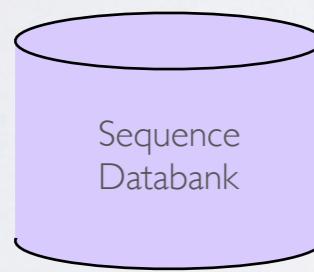


HISTORY OF SEQUENCE/SEQUENCE METHODS

Query Sequence

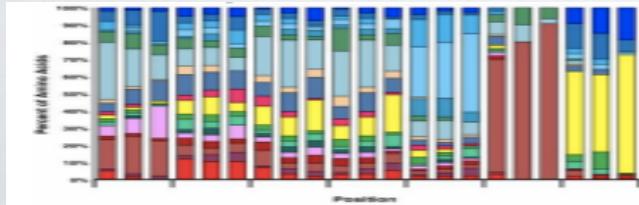
RHPVVMGNWKLNGSKEMVV DLLNLNGLN ALEGVTGV DV

Search



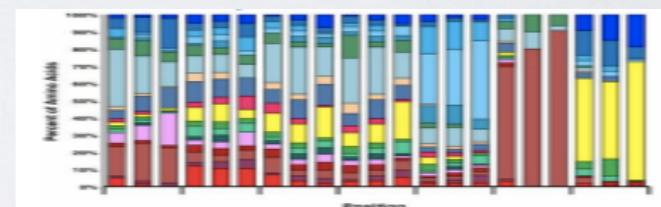
Sequence
Databank

Profile

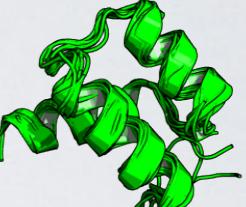


Search

Template
Profiles



Template
Structures



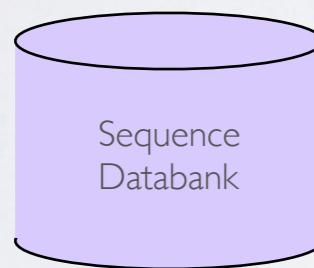
HISTORY OF SEQUENCE/SEQUENCE METHODS

Profile-Profile

Query Sequence

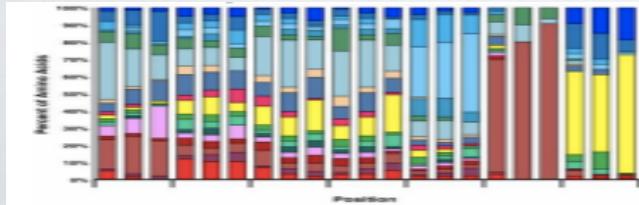
RHPVVMGNWKLNGSKEMVV DLLNLNAELEGVTGVDV

Search



Sequence
Databank

Profile



Profsim

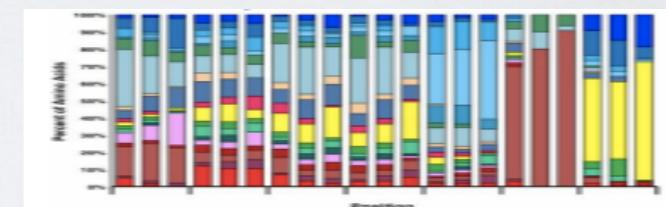
Compass/Procain

3DPSSM/Phyre/Phyre2

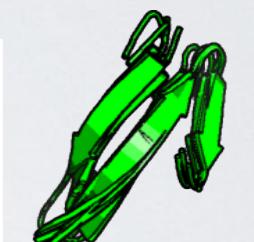
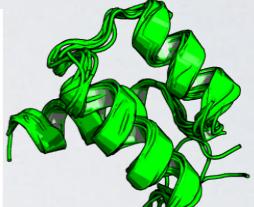
PRC

HHsearch

Template
Profiles



Template
Structures



HISTORY OF SEQUENCE/SEQUENCE METHODS

Query Sequence

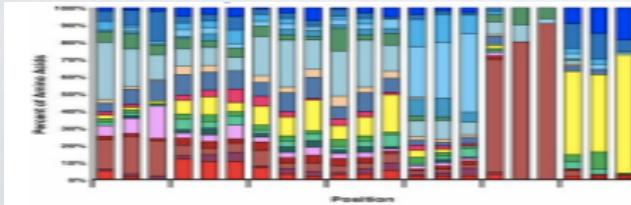
RHPVVMGNWKLNGSKEMVV DLLNGLNAELEGVTGVDV



Search

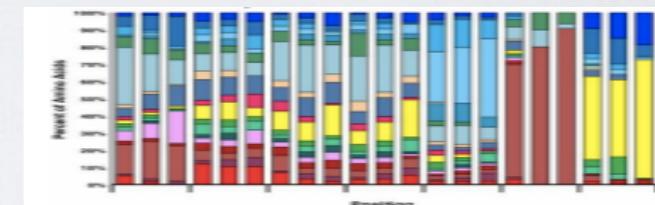
Profile-Profile
(+ local predictions: S2D and ACC prediction)
FFAS3D (Godzik lab)
Procain (Grishin lab)
Phyre2 (Sternberg Lab)
RaptorX (Xu lab)
HHsearch [+s2d] (Söding lab)
SPARX (Zhou lab)
Muster (Zhang lab)
Orion (Gelly team)

Profile

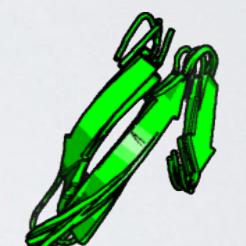
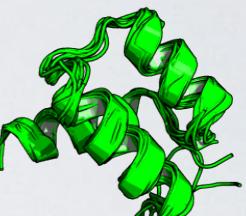


+ local structure predictions

Template
Profiles

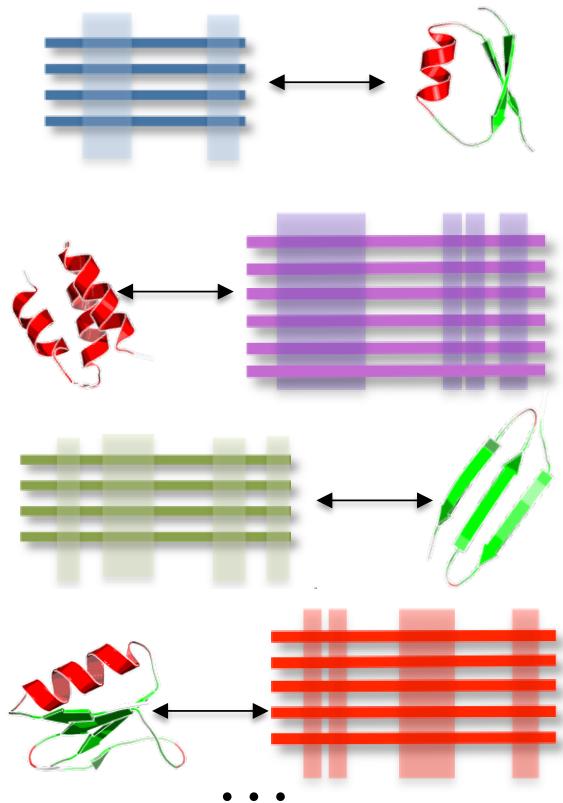


Template
Structures



Meet-U 2020

- Reference database: HOMSTRAD (*Mizuguchi K et al, Protein Science 1998*)



We provide a database of about 1000 protein families.

Each family is represented by:

- a master sequence
- a 3D structure (PDB file)
- a multiple sequence alignment
- A SCOP identifier

Meet-U 2020

- upstream teams

INPUT

a query sequence

ENVIREEAIVRYCATEDIT

Meet-U 2020

- upstream teams

INPUT
a query sequence

ENVIREEAIVRYCATEDIT

VS



searching homologs

Meet-U 2020

- upstream teams

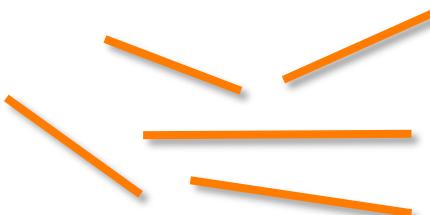
INPUT
a query sequence

ENVIREEAIVRYCATEDIT

VS



searching homologs



Meet-U 2020

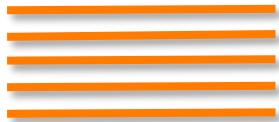
- upstream teams

INPUT

a query sequence

ENVIREEAIVRYCATEDIT

MSA (MUSCLE, mafft, T-Coffea...)



Meet-U 2020

- upstream teams

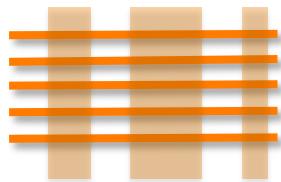
INPUT

a query sequence

ENVIREEEAIVRYCATEDIT



building profile



- Full or reduced alphabet
- Added information, like secondary structure, solvent accessibility...
- Sequences can be weighted to reduce redundancy

Meet-U 2020

- upstream teams

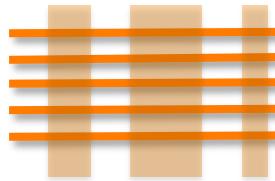
INPUT

a query sequence

ENVIREEAIVRYCATEDIT



building profile



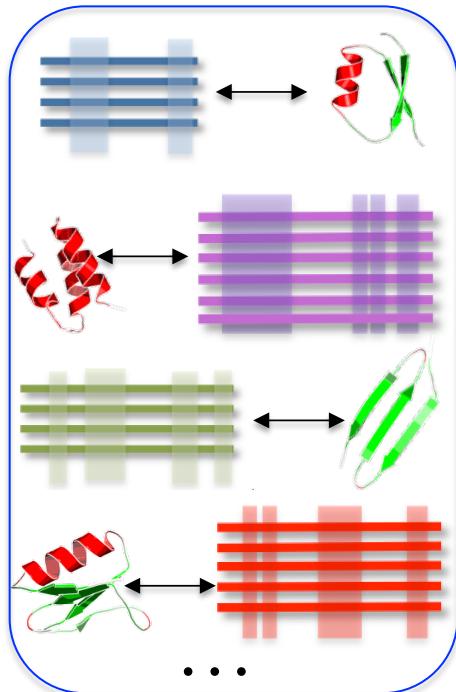
Meet-U 2020

- upstream teams

INPUT

a query sequence

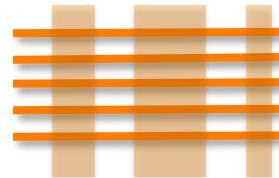
sequences-structures DB



ENVIREEAIVRYCATEDIT

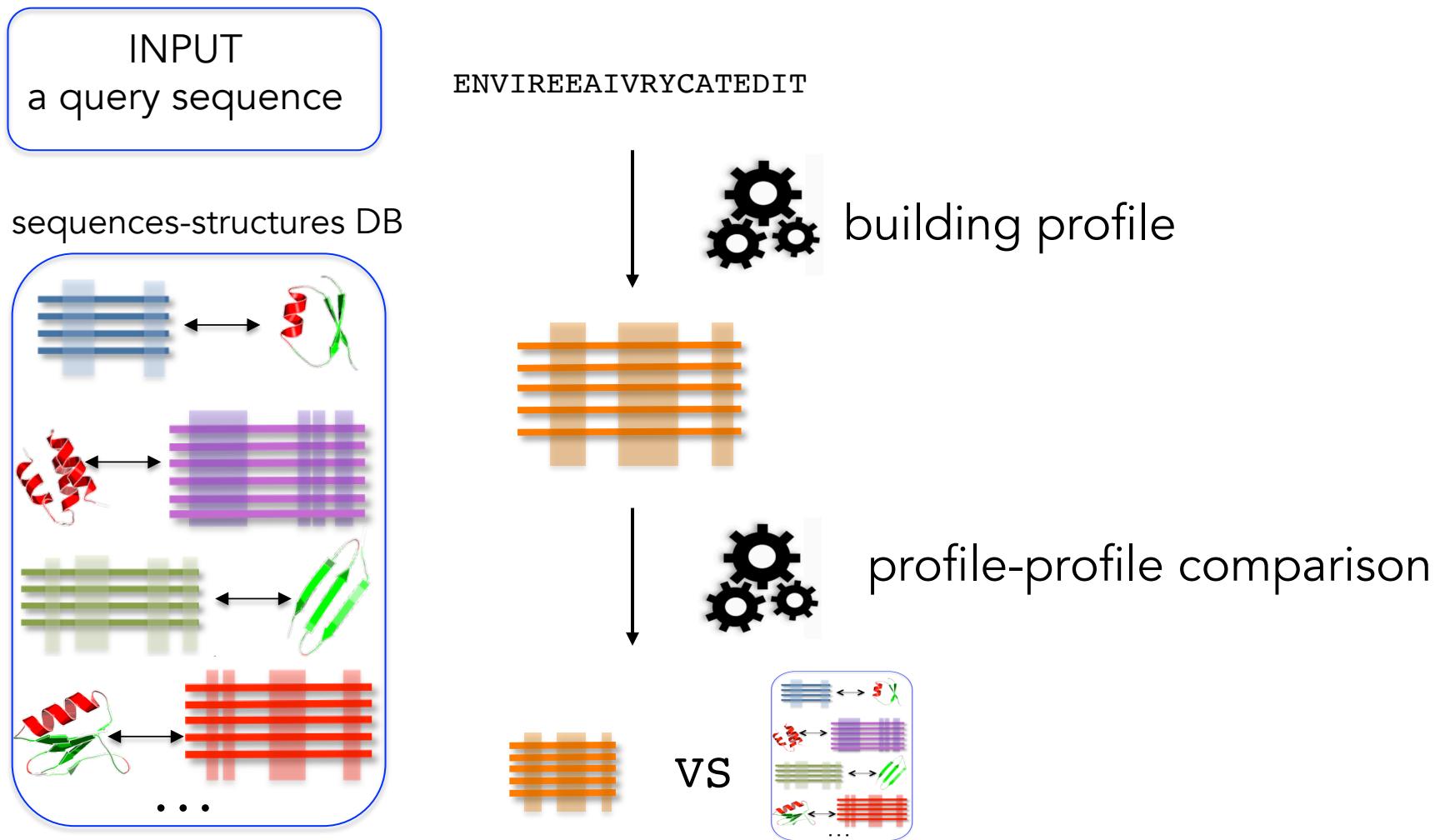


building profile



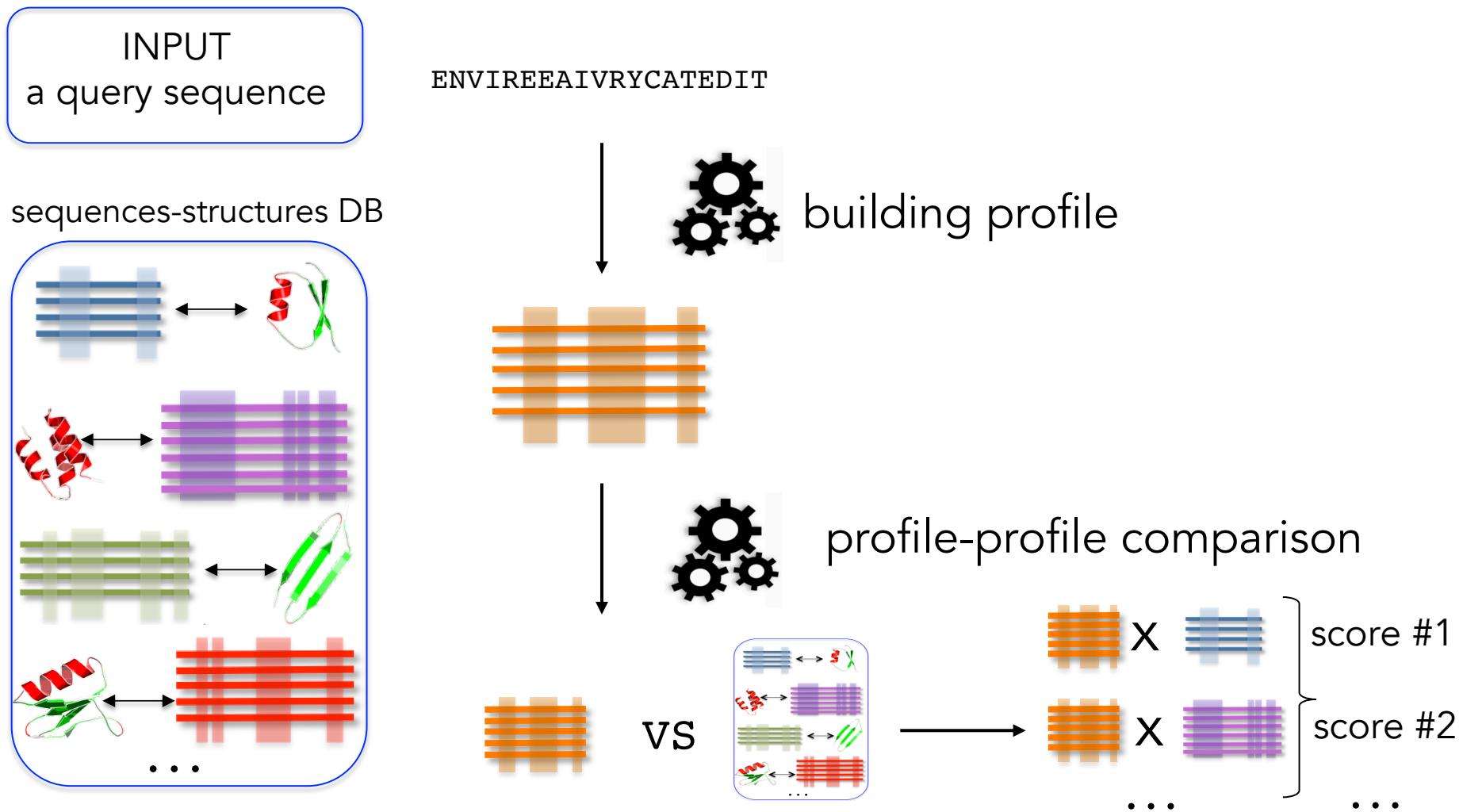
Meet-U 2020

- upstream teams



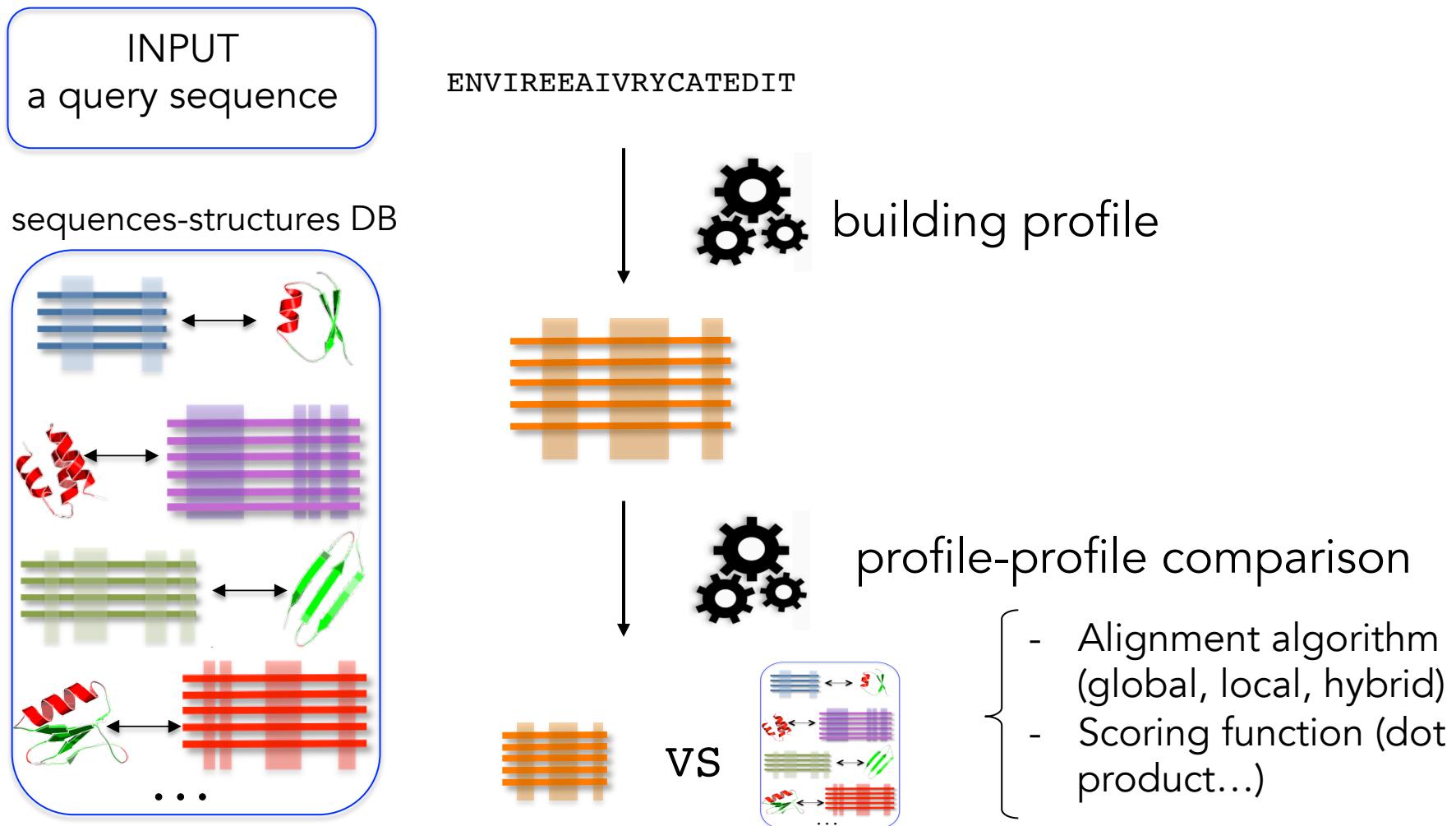
Meet-U 2020

- upstream teams



Meet-U 2020

- upstream teams



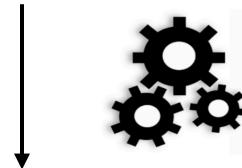
Meet-U 2020

- upstream teams

INPUT

a query sequence

ENVIREEEAIVRYCATEDIT



OUTPUT

N profil-profile
alignments and their
corresponding scores



score #1



score #2

...

Meet-U 2020

- upstream teams

INPUT

a query sequence

ENVIREEEAIVRYCATEDIT



OUTPUT

N profil-profile
alignments and their
corresponding scores

 X  score #1

 X  score #2

...

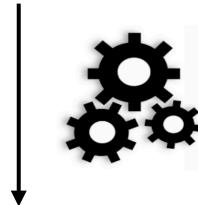
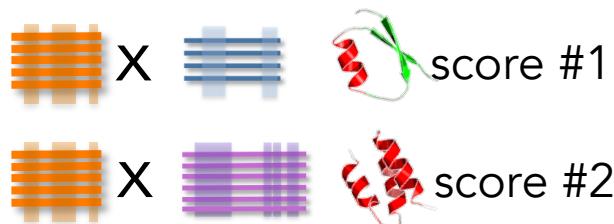


- downstream teams

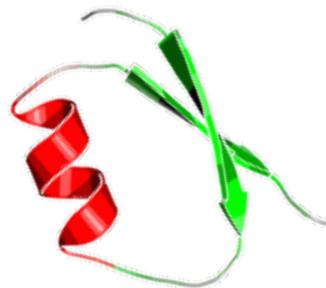
Meet-U 2020

- downstream teams

INPUT
N profile-profile
alignments and their
corresponding scores



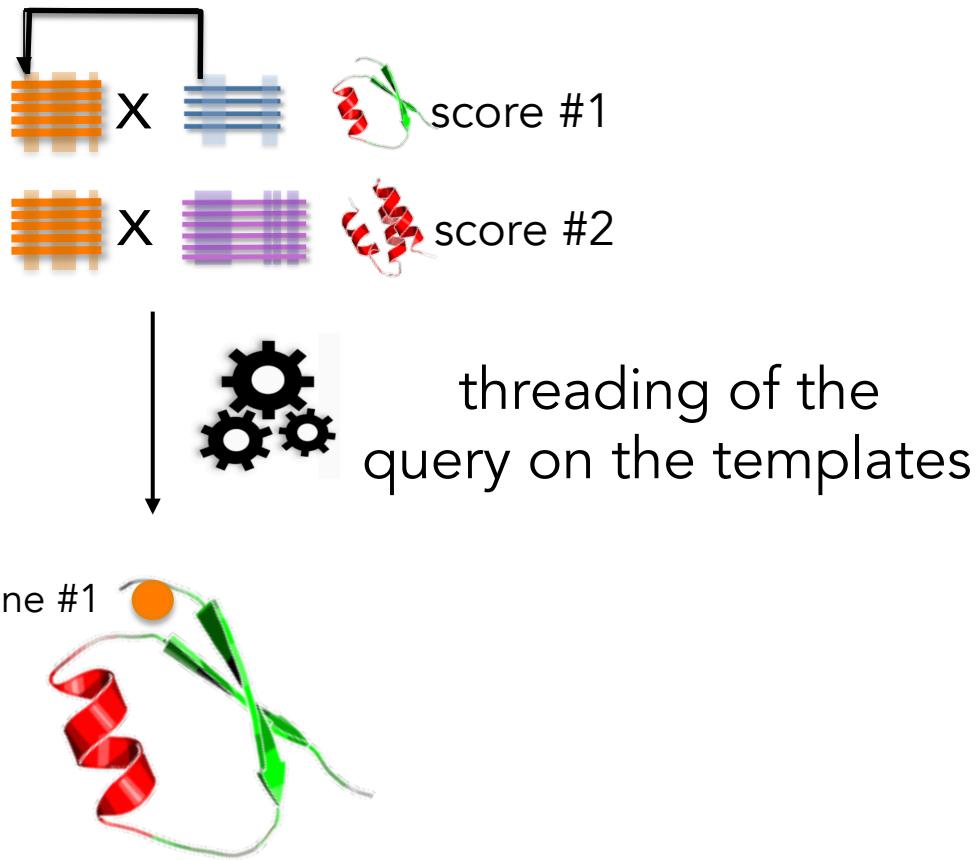
threading of the
query on the 3D structures
of each family master



Meet-U 2020

- downstream teams

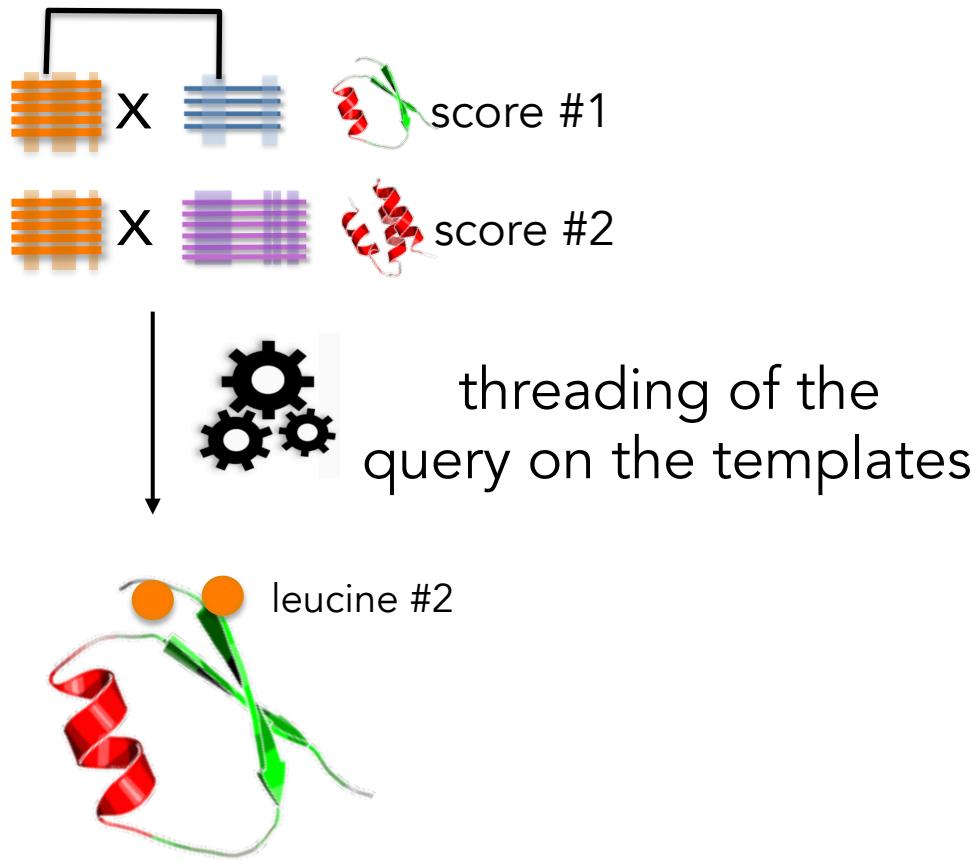
INPUT
N profile-profile
alignments and their
corresponding scores



Meet-U 2020

- downstream teams

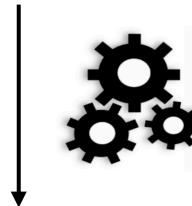
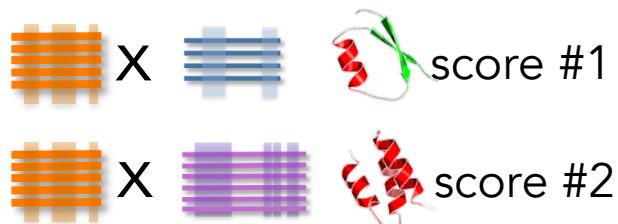
INPUT
N profile-profile
alignments and their
corresponding scores



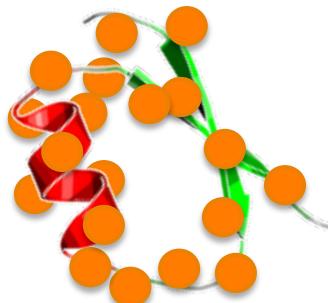
Meet-U 2020

- downstream teams

INPUT
N profile-profile
alignments and their
corresponding scores



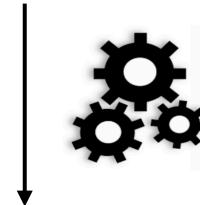
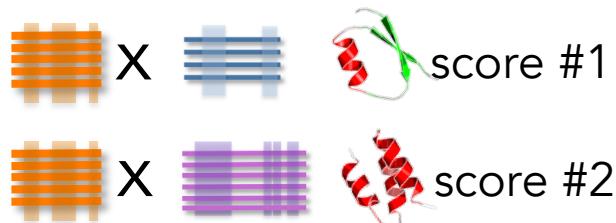
threading of the
query on the templates



Meet-U 2020

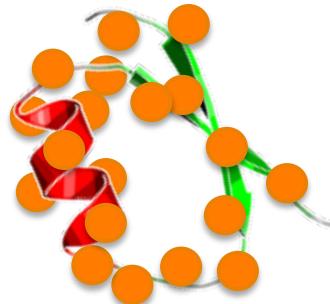
- downstream teams

INPUT
N profile-profile
alignments and their
corresponding scores

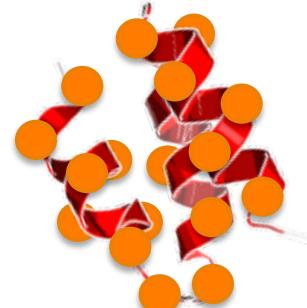


threading of the
query on the templates

N models



structure family #1

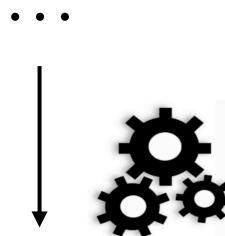


structure family #2

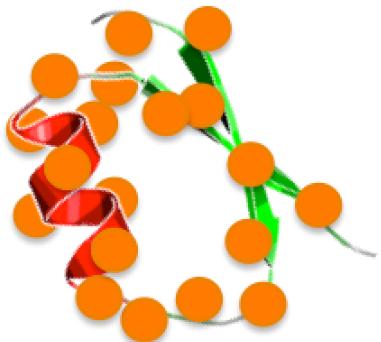
...

Meet-U 2020

- downstream teams



assessment of the quality
of the models



model #1

SCORING FUNCTION

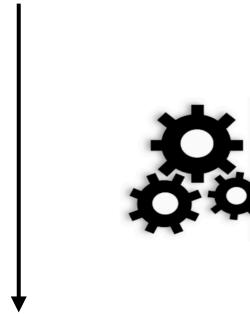
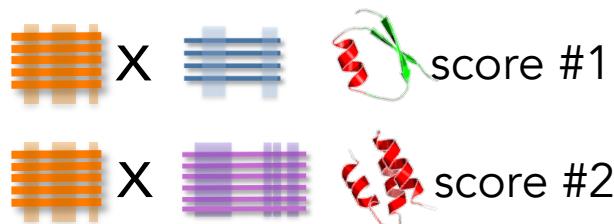
Pairwise distance-dependent DOPE score

Plus information coming from the alignment,
coevolution signals, secondary structure & solvent
accessibility consistency, physical terms...etc

Meet-U 2020

- downstream teams

INPUT
N profile-profile
alignments and their
corresponding scores



OUTPUT
X models and their
corresponding scores

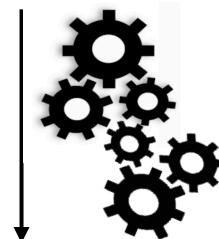


Meet-U 2020

INPUT
query sequence

ENVIREEEAIVRYCATEDIT

FASTA



- Upstream part
- Downstream part



model #46
TOP1



model #38
TOP2

...

model #72
TOPX

OUTPUT
X models and their
corresponding scores

- PDB files containing model coordinates
- One 2-column file with models scores

Meet-U 2020

- Benchmark to train(?), test/evaluate your method
 - 1000 examples with different levels of difficulties

seqtest #1 MTKRVELALIDKFPFLAVITYCESTFINI
seqtest #2 MSPRKFVACSLMSDEETLILAYLAFL
...
seqtest #N MENVIREEAPARISSAMEDIENIVREAIVRYCESTPERMIS



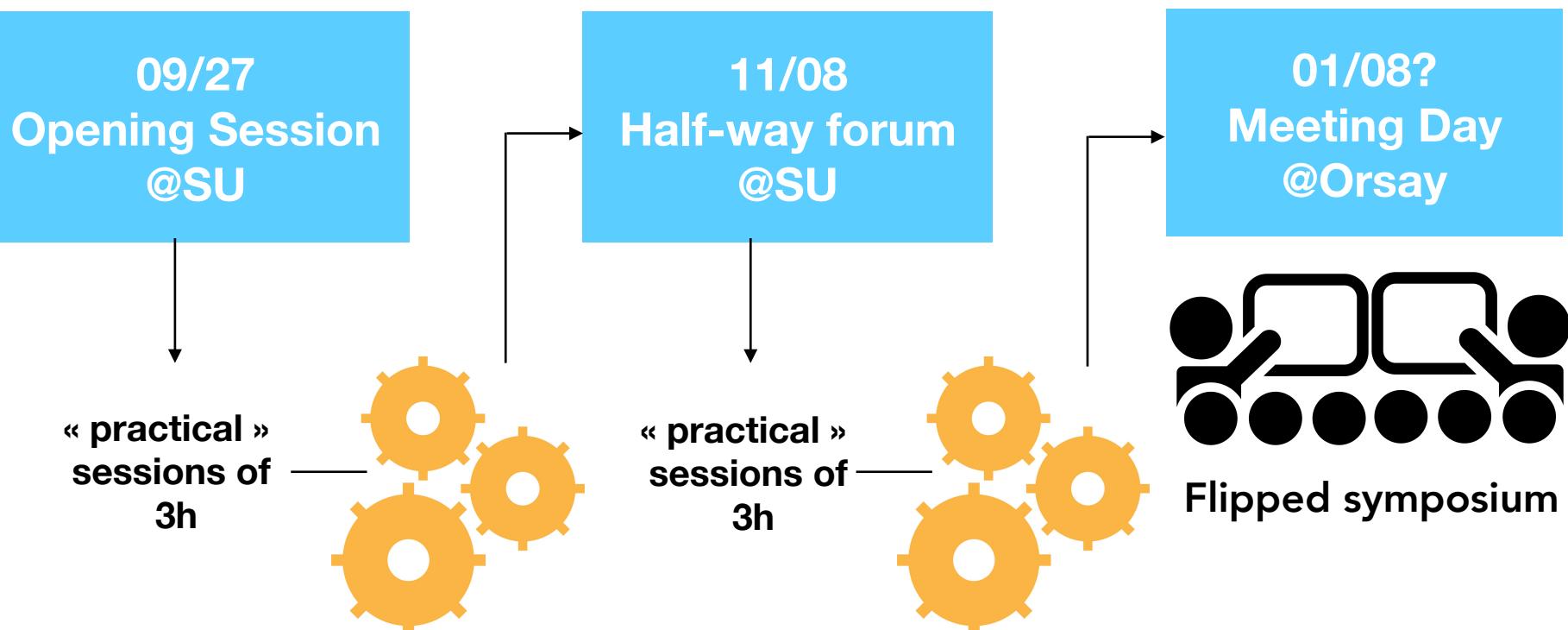
seqtest #2 MSPRKFVACSLMSDEETLILAYLAFL
seqtest #N MENVIREEAPARREAIVRYCESTPERMIS
...
seqtest #1 MTKRVELALIDKFPFLAVITYCESTFINI



sequences homologous to the different families of the database

File with the solutions (i.e. the true associations sequences-families)

WHEN&WHERE?



It's time to score!

Then, it is our turn!

Evaluations - criteria

5 criteria evaluated by the pedagogical team

- program (quality, property of the code, lisibility, usability)
- strategy: originality, relevance, risk
- result analysis: relevance, critical mind
- organization, participation, team work
- report (report 10 pages + digest 2 pages), presentations

3 criteria evaluated by the jury members

- presentation
- strategy: originality, relevance, risk
- result analysis: relevance, critical mind

Tâches principales

- **Formation des groupes (4, 5 max) en fonction des compétences de chacun**
 - qui fait quoi ? définition des rôles
- Recherche bibliographique et appropriation des concepts
 - rapport synthétique, présentation à l'équipe
- Définition de la stratégie et des choix méthodologiques
 - rapport détaillé, envisager les risques
- Implémentation de l'outil
 - code et documentation (README & commentaires)
- Validation de l'outil
 - applications à des cas tests, justification des résultats
- **Présentation de l'outil et des résultats, dissémination**
 - présentation orale devant le jury

Déroulé du projet

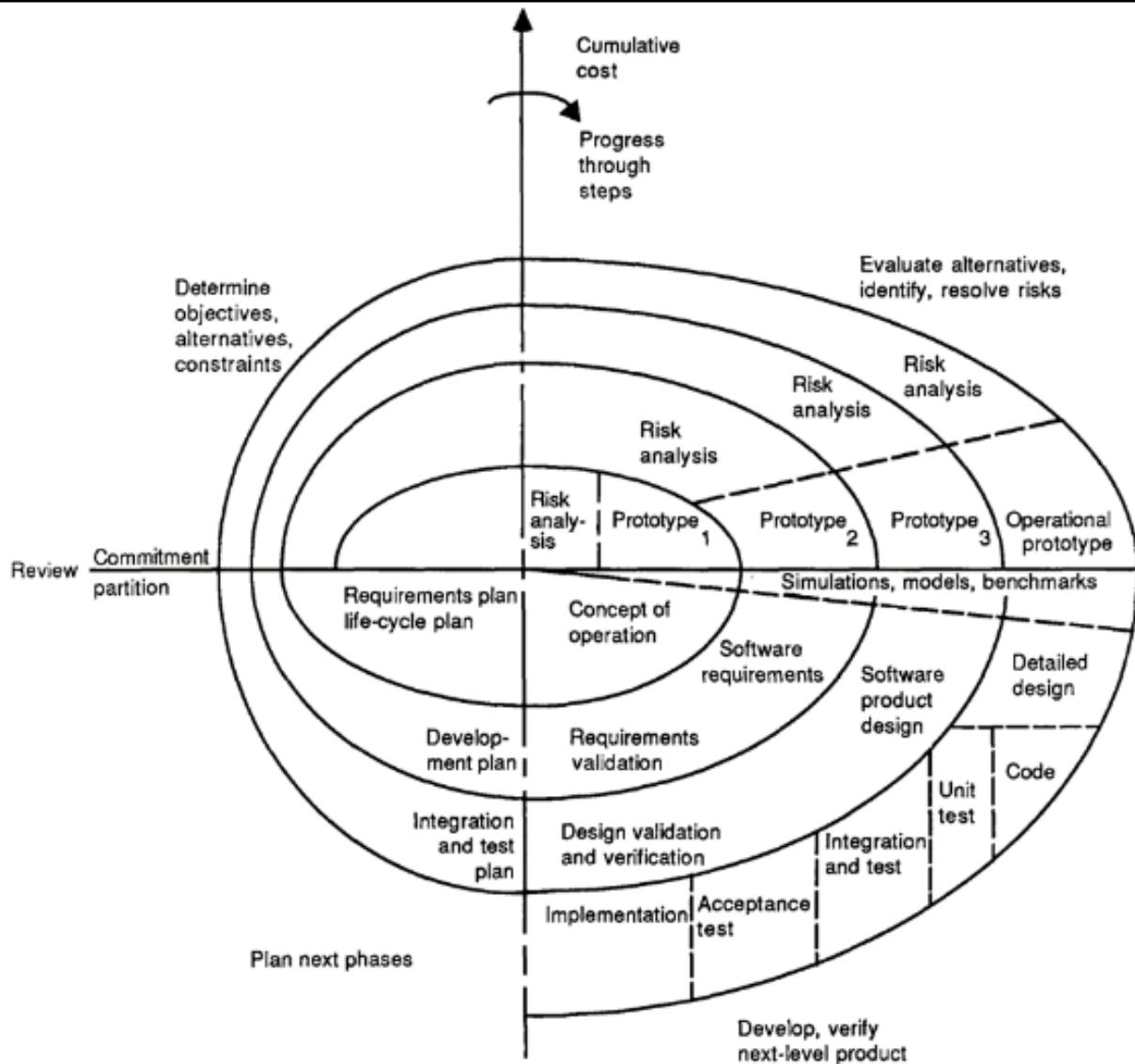


Figure 2. Spiral model of the software process.

Déroulé du projet



Figure 2. Spiral model of the software process.

Rôles dans le projet

- ✓ **Manager** : gestion des ressources et du temps, animation meetings
- ✓ **Expert(s) technique(s)** : programmation (implémentation et validation)
- ✓ **Expert(s) veille scientifique** : recherche bibliographique & propositions
- ✓ **Expert(s) deliverables** : rédaction des rapports, présentation

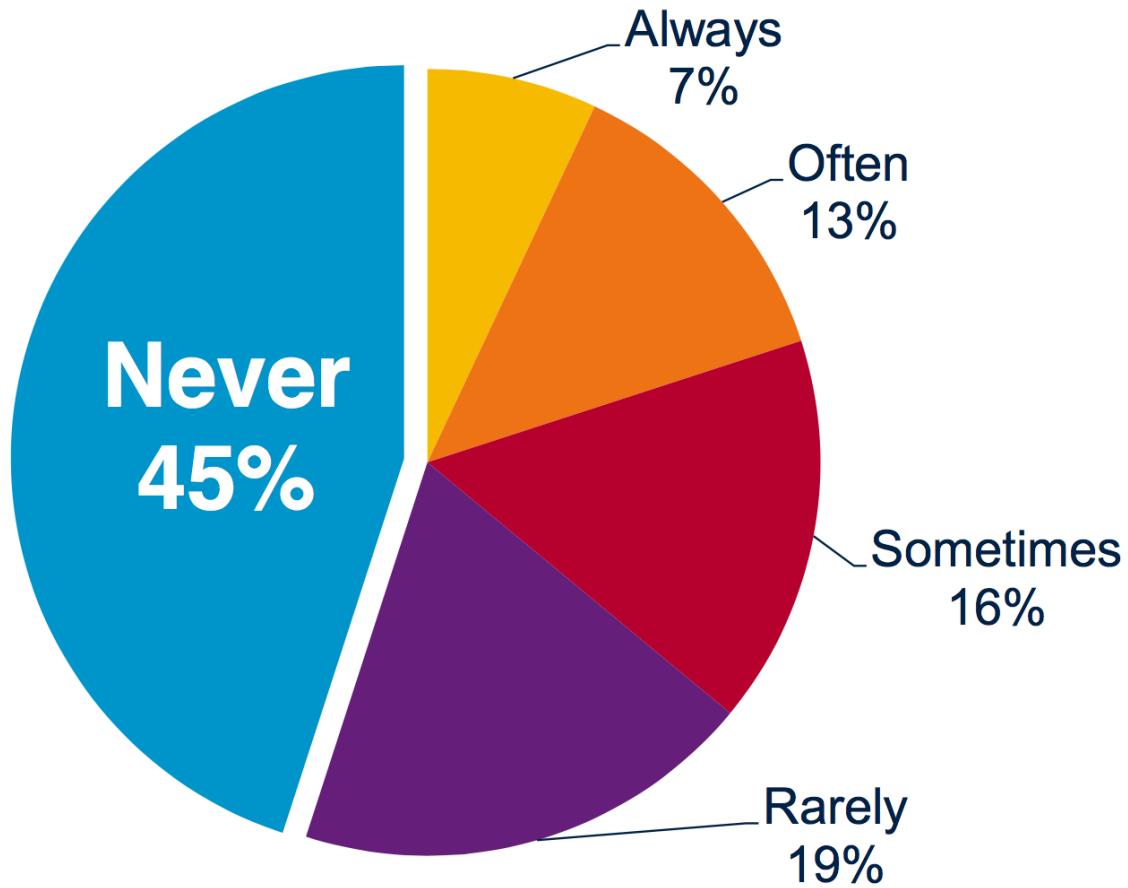
Tous ces rôles s'exerceront tout au long du projet.

Les experts doivent échanger !

Chacun peut jouer plusieurs rôles, mais il est important de bien définir les responsabilités associées.

Bonnes pratiques

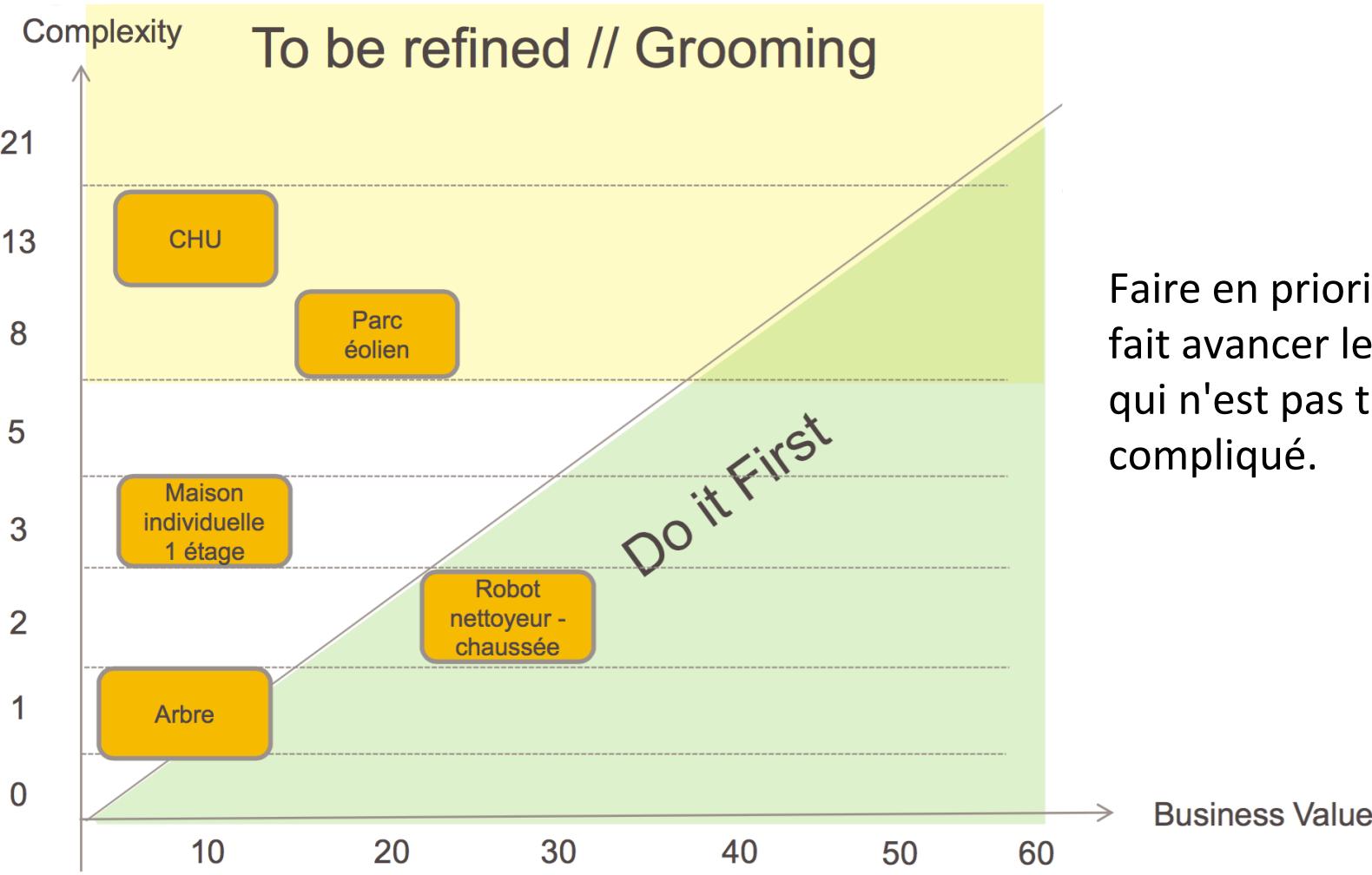
Près de la moitié des fonctions développées ne sont jamais utilisées.



Bonnes pratiques

- ✓ Il est impossible de rassembler/connaître tout ce qui sera requis/utilisé dès le début du projet.
 - N'essayez pas de concevoir le projet parfaitement et complètement avant de commencer à implémenter
- ✓ Quels que soient votre stratégie et vos choix initiaux, ils seront amenés à changer au cours du projet.
 - Soyez flexibles ! Autorisez-vous des allers-retours, implémenter vos solutions de manière modulaire, code facile à enrichir/adapter
- ✓ Il y aura toujours plus de choses à faire que le temps et les ressources ne le permettent.
 - Ne vous fixez pas d'objectifs inatteignables, soyez réalistes et faites des compromis.

Bonnes pratiques



Faire en priorité ce qui fait avancer le projet et qui n'est pas trop compliqué.

Bonnes pratiques

TO DO	DOING	DONE
		<p>Item #1</p> <p>t1.6 t1.1 t1.3 t1.5 t1.2</p>
<p>Item #2</p> <p>t2.7</p>	<p>t2.6 t2.5</p>	<p>t2.1 t2.3 t2.2 t2.4</p>
<p>Item #3</p> <p>t3.4 t3.5 t3.3 t3.2</p>	<p>t3.1</p>	
<p>Item #4</p> <p>t4.4 t4.2 t4.1 t4.5 t4.3</p>		
<p>Item #5</p> <p>t5.4 t5.3 t5.1 t5.5 t5.2</p>		

	A FAIRE	EN COURS	FAIT
CRITIQUE			Pb1
MAJEUR			Pb5 Pb4
MINEUR	Pb6	Pb7	

	PAS URGENT	URGENT
CRITIQUE		
MAJEUR		
MINEUR		

To do for next time

- creation of your GitHub account and watch the tutorial
- choice of the strategy – (and split into tasks)
- definition of the roles of each team member
- start implementing your part
- preparation of 2-3 slides (informal) to present your strategy

Conseil du jour

LA PERFECTION EST L'ENNEMI DU BIEN...