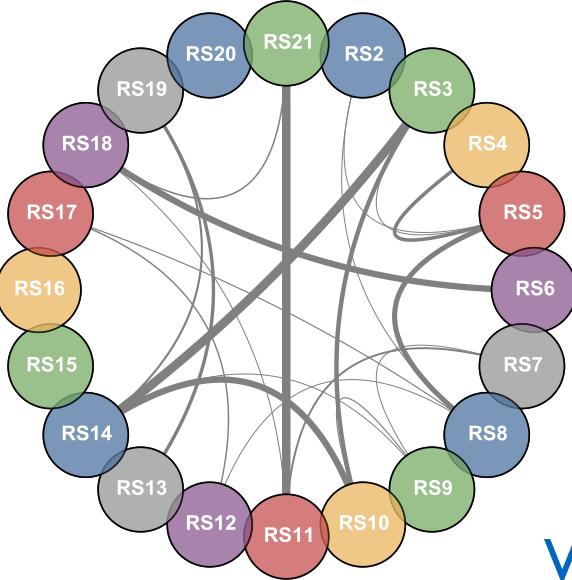


Interprotein coevolution: bridging scales from residues to genomes

Martin Weigt

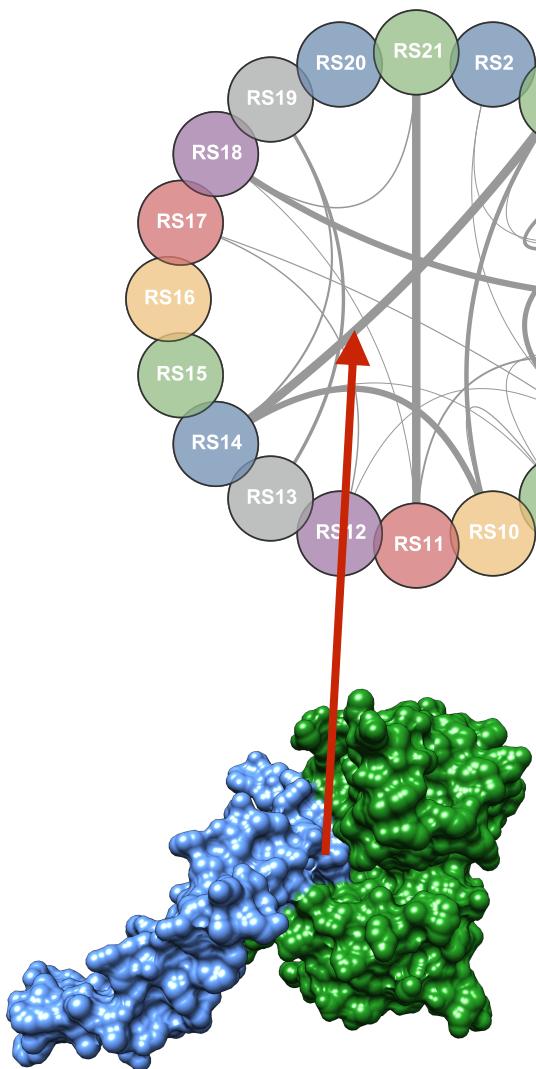
Laboratoire de Biologie Computationnelle et Quantitative
Université Pierre & Marie Curie Paris

The different scales in protein-protein interaction



Who with whom?
protein-protein interaction networks

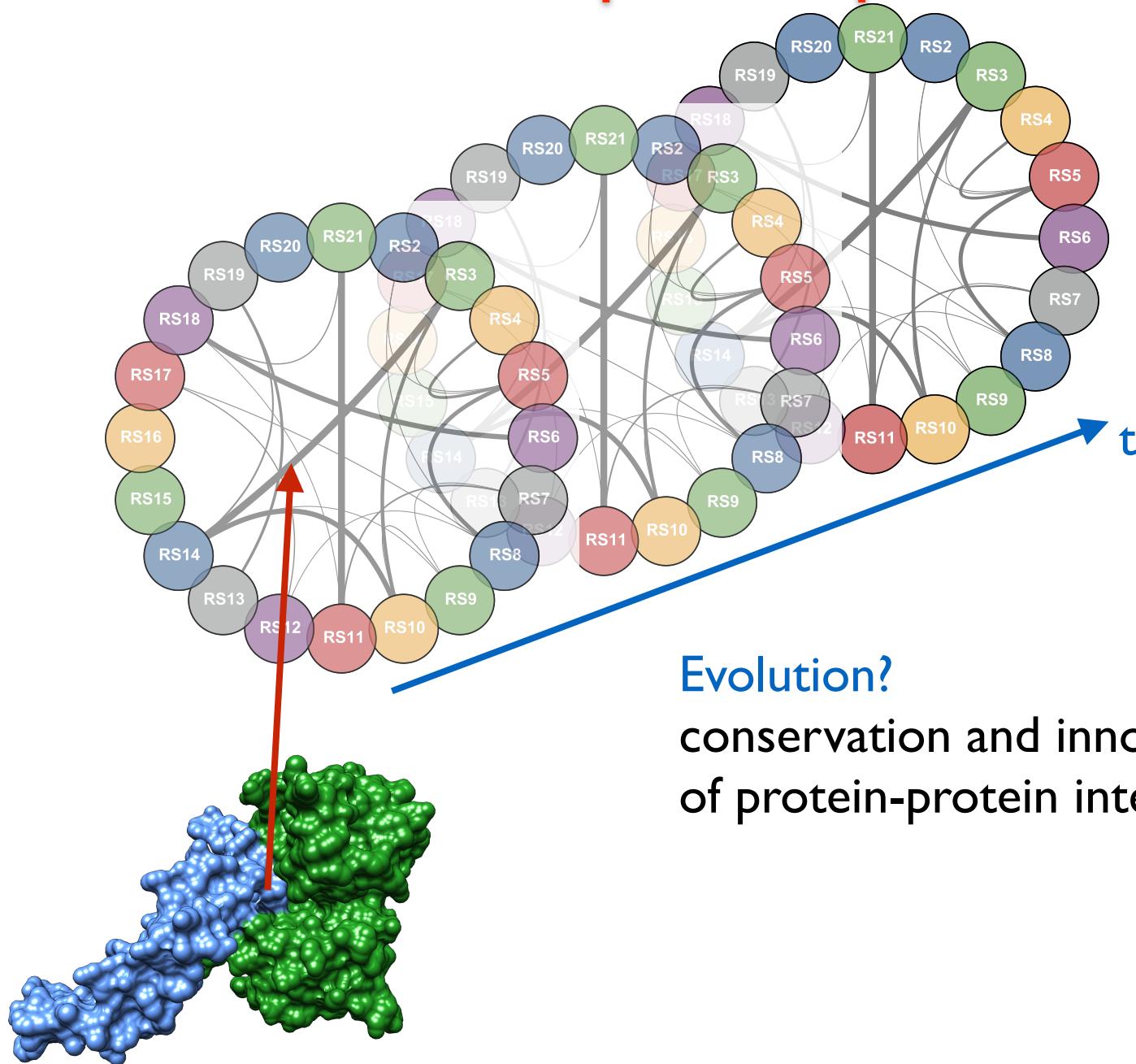
The different scales in protein-protein interaction



How?

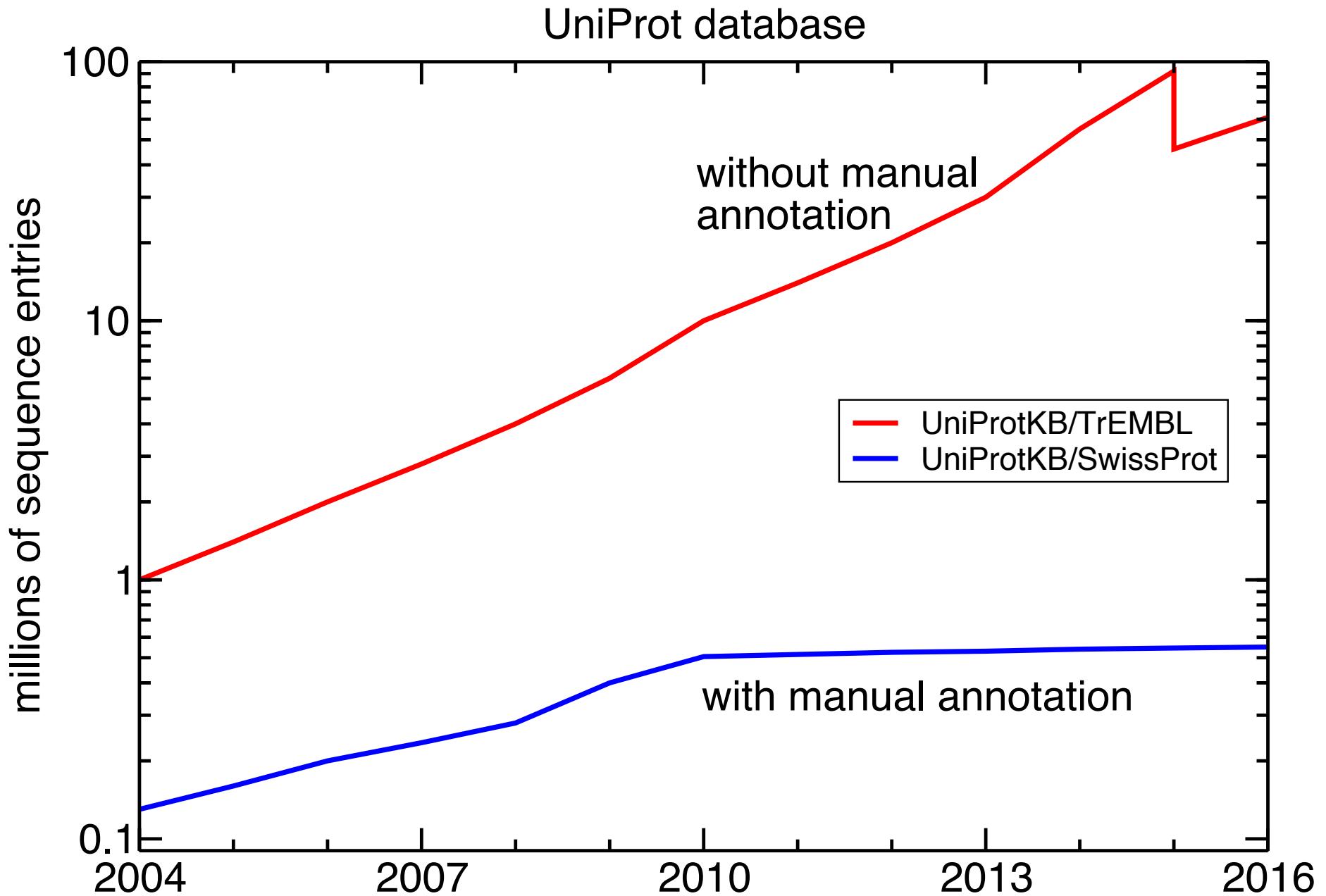
protein-protein interfaces
inter-protein residue contacts

The different scales in protein-protein interaction



Evolution?
conservation and innovation
of protein-protein interactions

Protein sequence data are accumulating...



...and are classified into homologous protein families

[HOME](#) | [SEARCH](#) | [BROWSE](#) | [FTP](#) | [HELP](#) | [ABOUT](#)



Pfam 30.0 (June 2016, 16306 entries)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

QUICK LINKS YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...

- SEQUENCE SEARCH** Analyze your protein sequence for Pfam matches
- VIEW A PFAM ENTRY** View Pfam annotation and alignments
- VIEW A CLAN** See groups of related entries
- VIEW A SEQUENCE** Look at the domain organisation of a protein sequence
- VIEW A STRUCTURE** Find the domains on a PDB structure
- KEYWORD SEARCH** Query Pfam by keywords

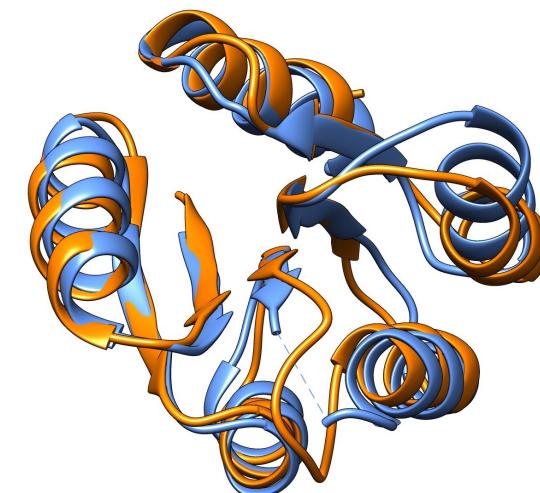
JUMP TO

[Go](#) [Example](#)

Enter any type of accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc.

Homologous proteins

- frequently 10^3 – 10^6 proteins per family
- common evolutionary ancestry
- conserved 3D structure and biological function
- diverged amino-acid sequences (~20-30% sequence identity)
- ▶ **sequence variability contains information about structure and function**
- >5000 families without example structures



Statistical physics

From models over data to thermodynamic observables:

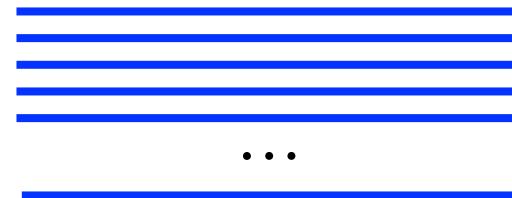
$$P(\bar{S}) \sim e^{-\beta \mathcal{H}(\bar{S})}$$

$$\mathcal{H}(\bar{S}_1) = - \sum_{i < j} J_{ij} S_i S_j - \sum_i h_i S_i$$



sample from model

$$\{\bar{S}^\mu\}_{\mu=1,\dots,M}$$



$$\langle \mathcal{O}_a(\bar{S}) \rangle_P \simeq \frac{1}{M} \sum_{\mu} \mathcal{O}_a(\bar{S}^\mu)$$



e.g. $\langle S_i \rangle_P, \langle S_i S_j \rangle_P$

Inverse statistical physics

From data over observables to models

$$P(\bar{S}) \sim e^{-\beta \mathcal{H}(\bar{S})}$$

$$\mathcal{H}(\bar{S}_1) = - \sum_{i < j} J_{ij} S_i S_j - \sum_i h_i S_i$$

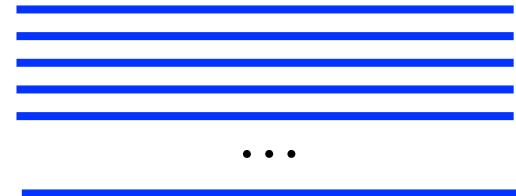
↑

$$\langle \mathcal{O}_a(\bar{S}) \rangle_P \simeq \frac{1}{M} \sum_{\mu} \mathcal{O}_a(\bar{S}^{\mu})$$

e.g. $\langle S_i \rangle_P, \langle S_i S_j \rangle_P$

Data:

$$\{\bar{S}^{\mu}\}_{\mu=1,\dots,M}$$



Inverse statistical physics

How to construct $P(\bar{S}) \sim e^{-\beta \mathcal{H}(\bar{S})}$ from data?

- coherence with data

$$\langle \mathcal{O}_a(\bar{S}) \rangle_P = \frac{1}{M} \sum_{\mu} \mathcal{O}_a(\bar{S}^{\mu})$$

- maximum entropy principle (least constrained model)

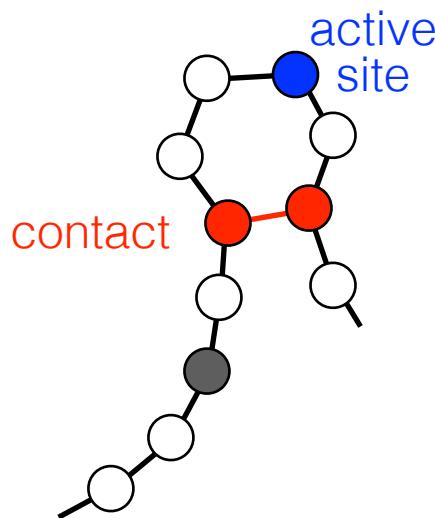
$$-\sum_{\bar{S}} P(\bar{S}) \log P(\bar{S}) \rightarrow \max$$

→ analytical form of model

$$\mathcal{H}(\bar{S}) = - \sum_a \lambda_a(\bar{S}) \mathcal{O}_a(\bar{S})$$

selection of observables
requires priori biological knowledge

Conservation and coevolution in proteins



variable residue	conserved residue
R	I
I	D
D	H
H	R
R	L
L	K
K	H
H	N
N	D
D	T
F	L
L	N
N	G
G	R
R	L
L	R
R	H
H	D
D	D
T	T
T	E
E	R
R	Q
Q	E
E	T
T	G
G	H
H	E
E	K
K	L
L	K
K	Y
Y	R
R	T
T	R
R	L
L	T
T	H
H	D
D	D
A	M
M	E
E	V
V	G
G	H
H	N
N	K
K	A
A	T
T	Q
Q	K
K	E
E	L
L	A
A	H
H	N
N	K
K	G
G	G

Profile model

$$P(a_1, \dots, a_L) \sim \exp \left\{ \sum_i h_i(a_i) \right\}$$

statistical modeling

Direct Coupling Analysis (DCA)

$$\sim \exp \left\{ \sum_{i < j} J_{ij}(a_i, a_j) + \sum_i h_i(a_i) \right\}$$

strong couplings -> residue contacts

[Weigt et al, PNAS '09]

[Morcos et al, PNAS '11]

Interactions between protein families

Family 1

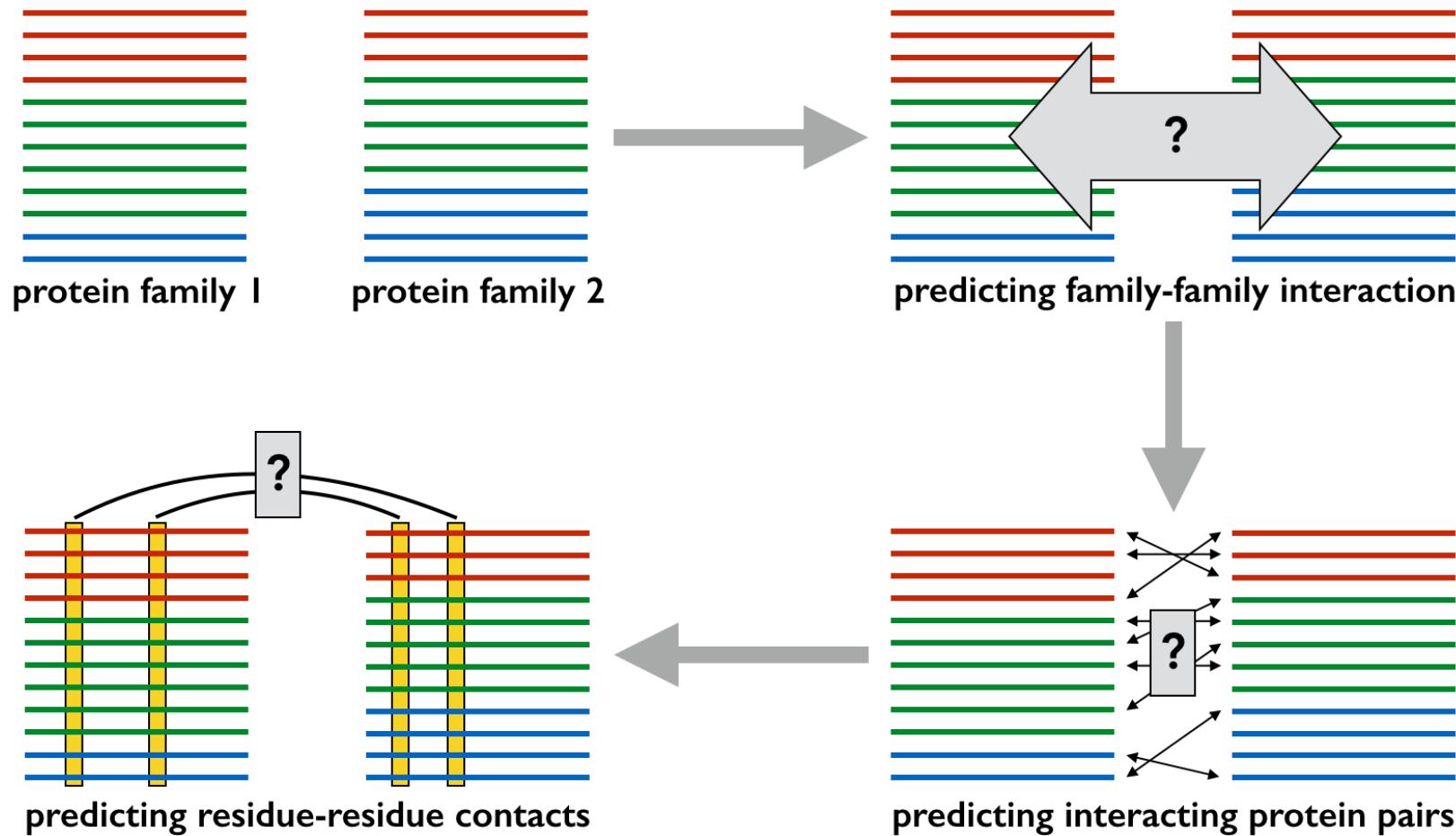
>F7XUK6_MIDMI/129-211
LAQQLEKRISFRKAAKRLIQNAM.R.....M.G..AEGIKIKISGRIG.G.AEIARDQQ
YNEGRVPL..HTLRRMMIDYGTAEAH..TTYGRIGVKVWV
>B3SEY6_TRIAD/119-201
VAEQLEKKVSFRKAVKRAISMAM.K.....M.G..AKGIKISVSGRLG.G.AEIARTEW
YKEGRVPL..HTLRAIVKYDMAEAH..TIYGLIGVKVWV
>RS3_0RITB/122-204
IAQQLERRQSFKVMKKAIHASM.K.....Q.G..AKGIKIICSGRLG.G.VEIARSES
YKEGRVPL..QTIRADIRYAFAAEI..TTYGVIGVKVWV
>RS3_RICPR/123-205
IAAQLEKRVSFRKAMKTAQASF.K.....Q.G..GQGIRVSCSGRLG.G.AEIARTEW
YIEGRMPL..HTLRADIDYSTAEAI..TTYGVIGVKVWI
>E1X0L6_HALMS/119-201
IASQLEKRVAFRRAMKKVMQSAF.R.....A.G..VKGIRVRTAGRLG.G.AEMARAEG
YSERKVPL..HTLRADIDYSTAEAH..TTYGVIGVKVWV
>I7HEJ8_9HELI/120-202
IATQLEKRVAFRRAMKKVMQAA.M.K.....A.G..AKGIKVKVSGRLA.G.AEMARTEW
YMEGRVPL..HTLRAKIDYGFTEAM..TTYGIIGVKVWI
>M4VDL1_9DELT/120-202
IAMQLEKRISSRALKKAAAT.K.....G.G..VRGIKVRVSGRLD.G.AEIARSEW
YNEKSVPL..HTLRADIDYGTAEAL..TAYGIIGMKVWI
>RS3_HYPNA/120-202
IARQLERRASFRRAMKRSIQSAM.R.....L.G..AEGVKVVVSGRLG.G.AEIARTEK
YAEGSVPL..HTLRADIDYGTAEAT..TTYGIIGVKVWV
>C0QW02_BRAHW/94-176
VARQLEMRAFRRAMKSVITQAM.K.....K.G..AKGIKVMCSGRLA.G.ADIARTEQ
YKNGSVPL..HTLRAKIDYGTAEAL..TTFGIIGIKVWI
>J9Z1W5_9PROT/119-201
IARQLEKRVAFRKAMKKSGQSAI.K.....L.G..AKGIKIVCGGRLG.G.AEIARSEK
FSEGSVPL..HTLRADIDYATARAL..TTYGIIGIKVWL
>RS3_MARMM/120-202
IAQQLERRVAFRRAMKRSQSAM.R.....M.G..AKGCKIVCGGRLG.G.AEIARTEQ
YNEGKSVPL..HTLRADIDYGTCEAK..TAMGIIGIKVWI
>G0GFA5_SPITZ/122-204
IAGQLEHRASFRRVMKLAVANAM.K.....A.G..VQGIKVRVSGRLG.G.AEIARSEV
QMAGRVPL..HTLRADIDYGFTEAM..TTYGVIGVKVWI
>V6DFZ5_9DELT/122-204
ISEQLEKRGSKKAMKRAALDM.K.....SG..AKGVKIRCAGRLG.G.AEIARDEW
IRVGSTPL..HTLRSIDYGFVEAH..TTYGVIGIKVWI
>RS3_NE0SM/120-203
IAFQLEKRSSFRRVIKKAIATVM.R.....ESD..VKGVKVACSGRLS.G.AEIARTEV
FKEGSIPL..HTMRADIDYVAEAH..TTYGVIGVKVWI
>I0III3_PHYMF/124-207
IAEQLAKRASFRRVMKMAEAAM.N.....CGV..CKGVKIMLSGRLG.G.HEMRSEV
VSLGSIPL..ATLQANVDYGFRAISK..TTYGTIGVKVWI
>F0SJ92_RUBBR/120-202
IAQQLGKRGSKRALKRSMEQVM.D.....A.G..AHGVKIELSGRLG.G.AEMSRKEK
GSRGSIPL..STLQRHVGYTTAR..TAQGIIGIKVWI



Family 2

>RS14_NE0SM/47-100
KLNQLPRNSSPARSKNRCISITGR..PRGYY..RKFGI..SRIQLVLANWGKLPGVVKSS
>I0AI30_IGNAJ/35-88
ALQKLPRNSSVTRLKNRCMFTGR..ARAYY..RKFGV..SRLVLREMALRGEIPGLKKSS
>I6YSF0_MELRP/36-88
.LQLLPRNSAPTRAHNRCLISGR..PRGYY..RKFGI..SRLVLREMALRGEIPGLKKSS
>I0IIH6_PHYMF/34-87
ALSQLPRDASPTRLVTQCAITGR..TRAVY..RKFNV..SRIVLRELALQKIPGMKKAS
>RS14_CHLT3/35-88
ALRKLPRDSSPTRLKNRCISITGR..AKGVY..KKFGL..CRHILRKYALEGKIPGMKKAS
>RS14_PROA2/35-88
ALSKLPRNSSATRVRNRCVLTGR..GRGVY..EKFGL..CRHMFRKLALEGKIPGVKKAS
>D6XYV1_BACIE/35-88
ALSKLPRDSAPSRLTRRCKATGR..PRGVL..RKFEL..SRIKFRELAHKQIPGVRKAS
>I0JIY2_HALH3/35-88
ALRKLPRDSSPTRVNRRCELSQR..PRGYM..RKFD..SRIAFARELAKHQIPGVKKAS
>RS14_EXIS2/36-88
.LSKLPVNSSAVRLHNRCISITGR..PHGYI..GKFGI..SRIKFRLAHKGQIPGVKKAS
>RS14_STR6/36-88
.LSKLPVNASPTRLHNRCRTVGR..PHSVY..RKFGL..SRIAFARELAKHQIPGVTKAS
>G0VNI1_MEGET/35-88
ALSQLPANASPVRHLNRCKVTGR..PHGYM..RKFGI..CRITFRELAYKGQIPGVKKAS
>R7PS46_9FIRM/35-88
ALSKLPRNASPTRLHNRCKLTGR..PHGYL..RKFGV..CRNFRELAHYKGQIPGVKKAS
>F8L373_SIMNZ/47-100
KLNQLPKNSSPIRRRNRCMTGR..CRGYL..RKFGI..SRLCFREMANDGSIPGVKKAS
>F8L0V7_PARAV/47-100
ALNKMPRDSSPIRLRNRCQLTGR..XRGYL..RKFKL..SRLTFREMALAGLLPGVTKSS
>D6YVK9_WADCW/47-100
QLNKMRRDTSPVRLRNRCQITGR..CRGYL..SKFKV..SRLVFREMASIGMIPGVTKSS
>L7VJR0_9FLA0/35-88
ALQKLPKNSCTVRLRNRCQLTGR..SRGYM..RKFGV..SRISFRNLVNFGLIPGVKKSS
>C7NDL0_LEPBD/41-94
ELSKLPVNASPTRLVRNRCQINGR..PRGYM..REFGI..SRVMFRQLAGEGVIPGVKKSS
>RS14_FUSUN/41-94
ELNKLKDSSAVRKNRCQLDGR..PRGYM..REFGI..SRVKFRQLAGAGVIPGVKKSS
>K0P015_9BACT/35-88
ALDKLPKNSSPVRLRNRCNITGR..ARGYI..RRFGI..SRLVFRKWALEGKIPGIRKAS
>RS14_AM0A5/35-88
ALDKLPKNASPVRVRNRCKITGR..ARGYM..RKFGI..SRIVFREWAAQGKIPGVIKAS
>I4ALV0_FLELS/42-94
.LDKLKDSSPVRLHNRCRLTGR..PRGYM..RRFGI..CRVVFREMANDGKIPGVTKSS
>RS14_SALRD/35-88
ELQKLPRDSSPVRLRNRCQLCGR..QRGYL..RKFGV..CRICFRELALEGKIPGIRKAS
>C7PU84_CHIPD/35-88
ELDQLPRNASPVRHLNRQLSGR..PKGYM..RHFGM..CRNMFRDLALAGKIPGVRKAS
>F4KWV6_HALH1/35-88
ELDKLPNSNPIRMHNRCQLTGR..PKGYM..RQFGL..CRVKFREMALYKGKIPGITKSS
.

Interactions between protein families

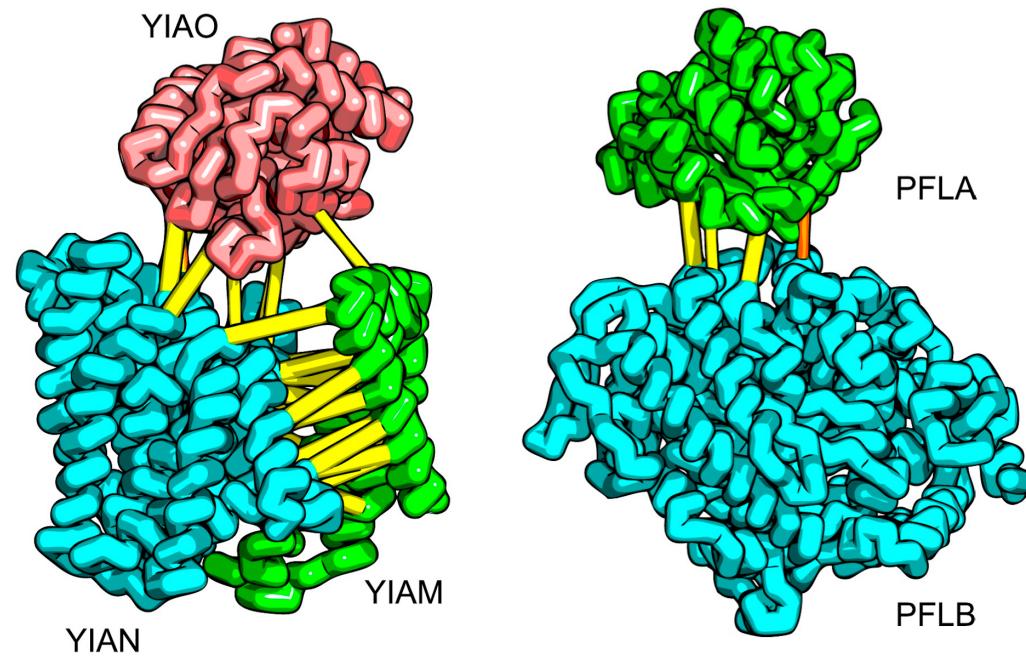
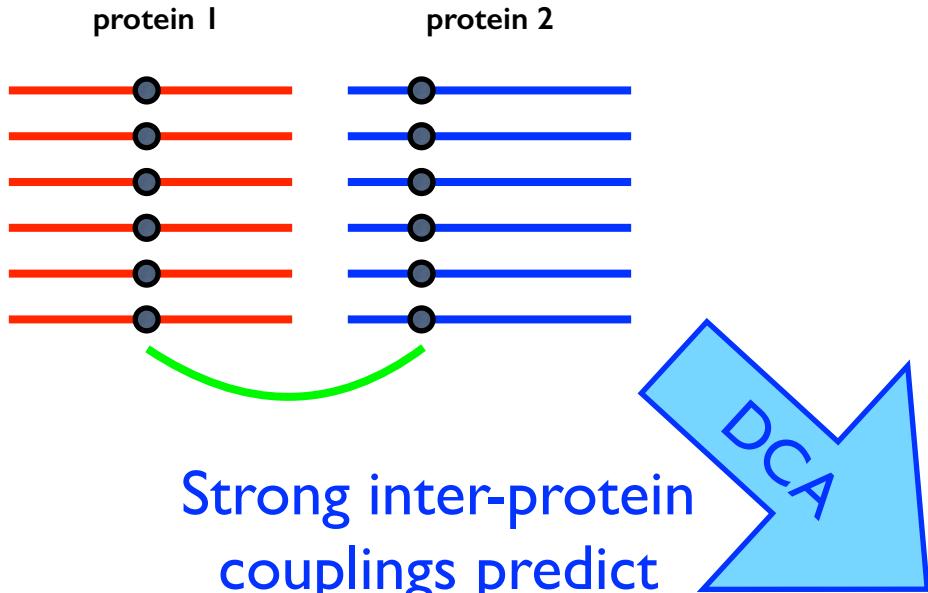


What can we learn from the empirical sequence variability:

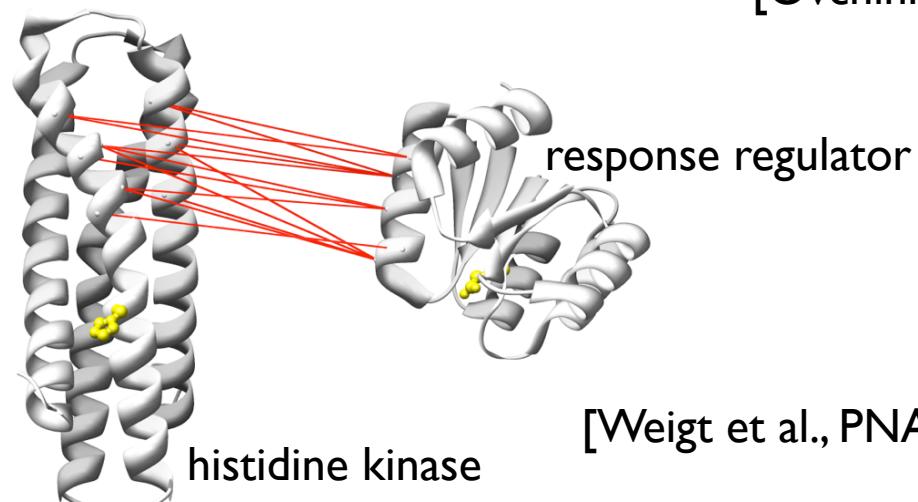
- do the families interact?
- which specific proteins interact?
- which residues are in contact?
- relation between protein structure/function and evolution

Prediction of inter-protein residue contacts

joint MSA of protein families

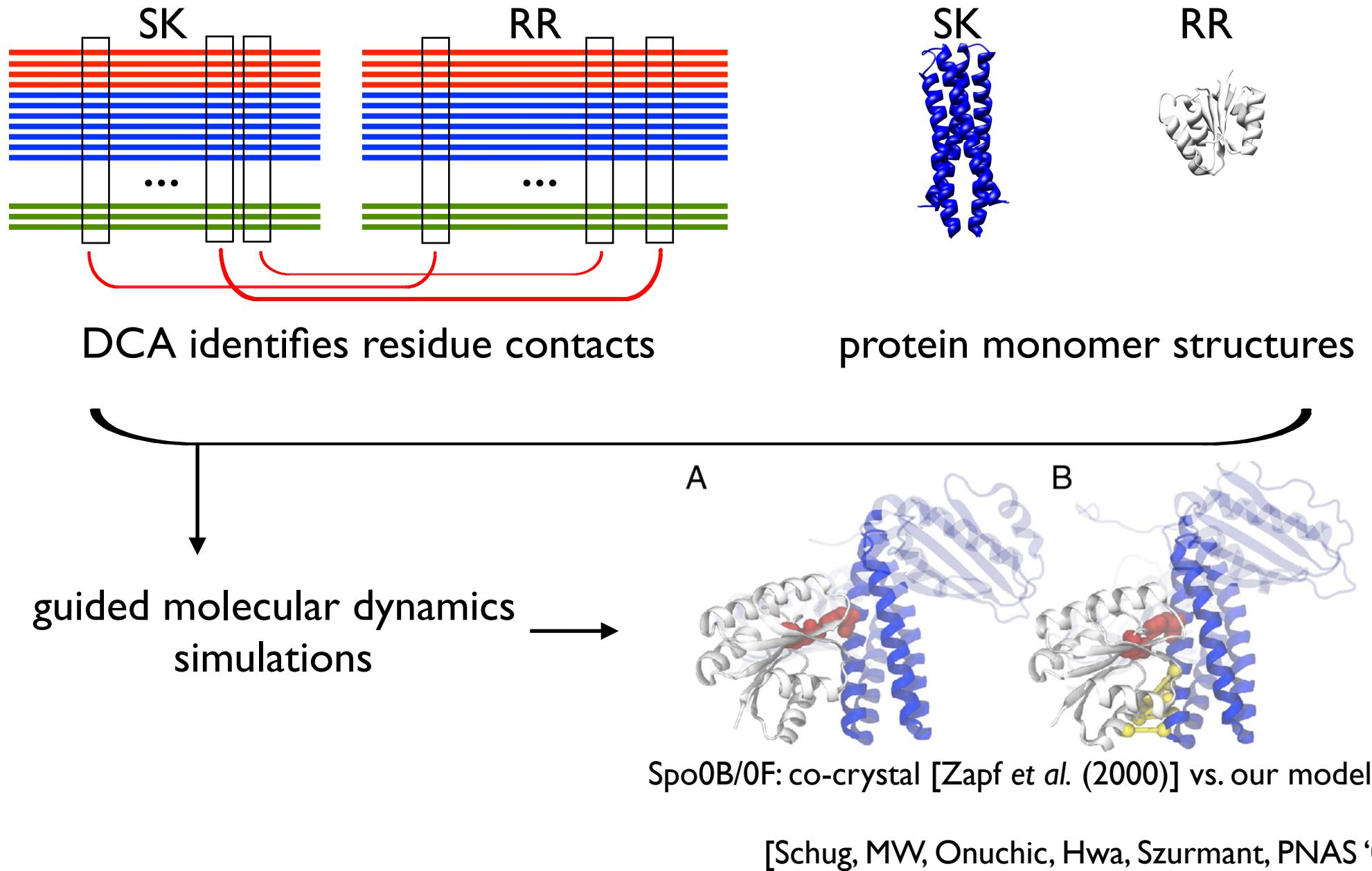


[Ovchinnikov et al., eLife '14]

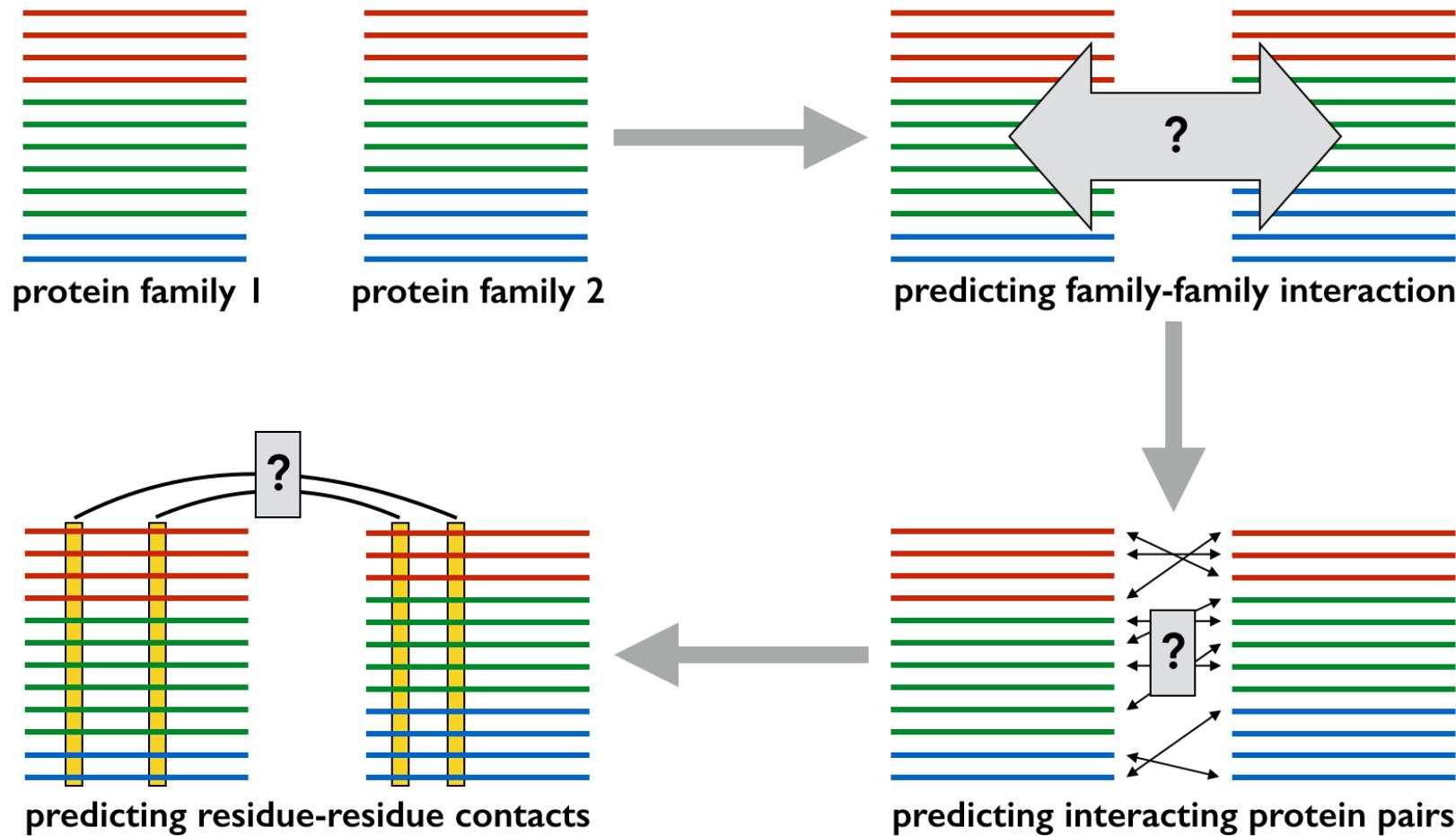


[Weigt et al., PNAS '09]

In silico prediction of high-resolution structures of transient protein complexes



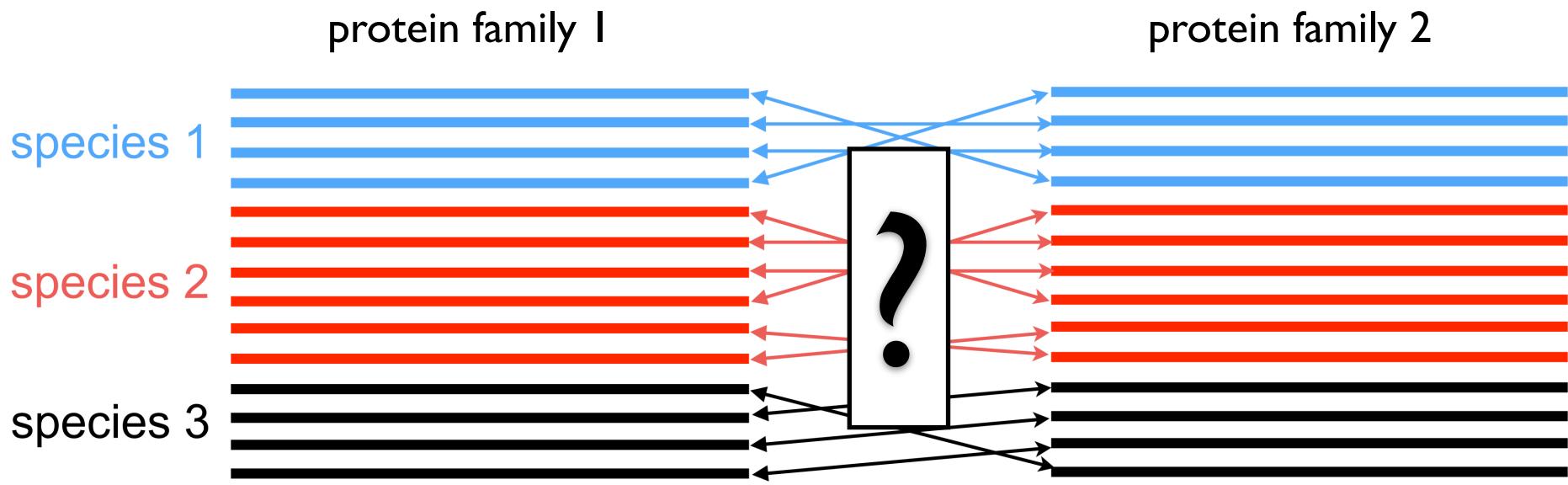
Interactions between protein families



What can we learn from the empirical sequence variability:

- do the families interact?
- which specific proteins interact?
- which residues are in contact?
- relation between protein structure/function and evolution

Specific interactions and paralog matching

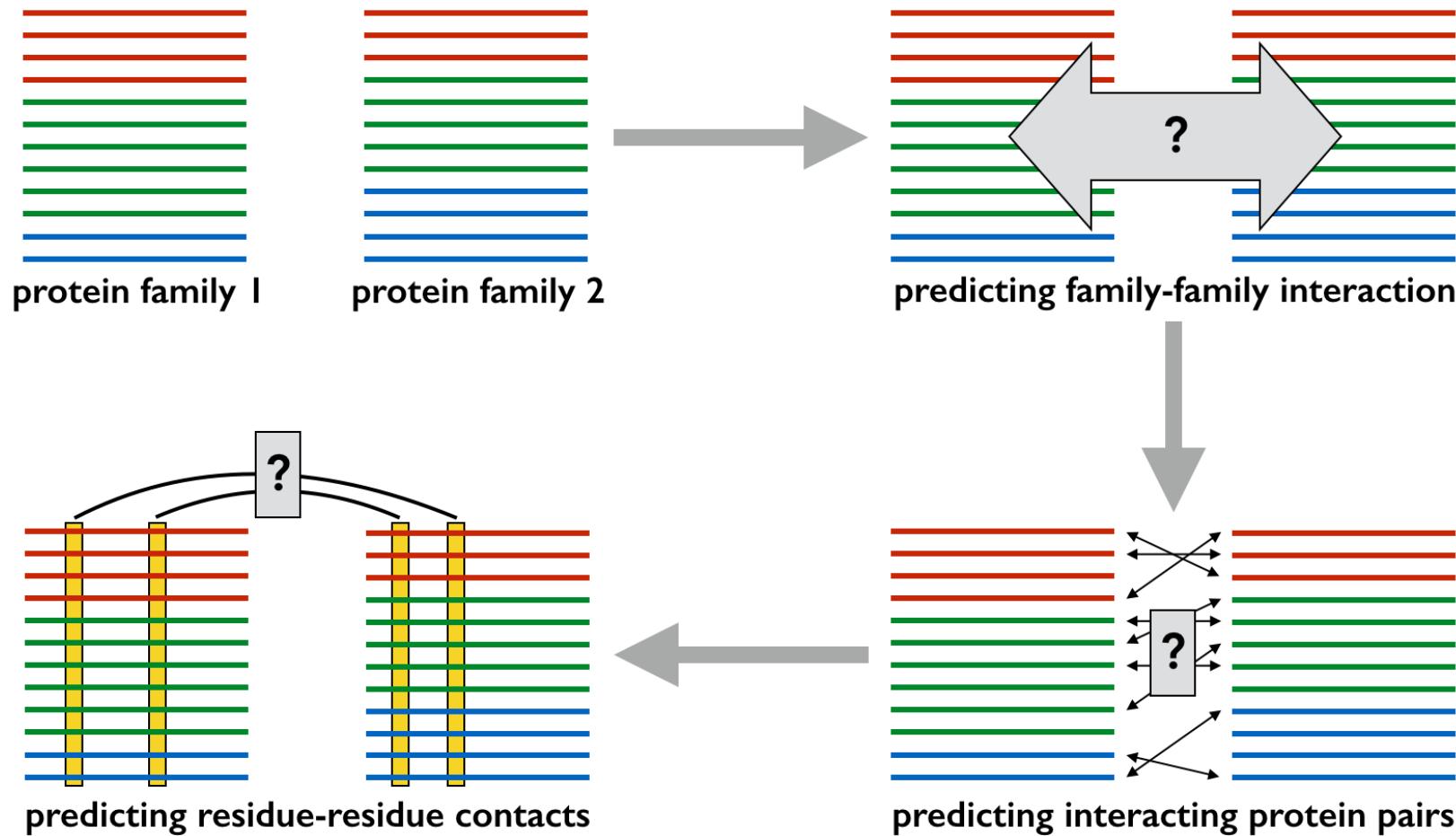


General idea:

- correct matching shows inter-protein covariation
- random matching has no inter-protein covariation
- **maximise inter-protein covariation computationally**
- reach 80-90% of accuracy in test cases
- simultaneous prediction of interacting paralogs and inter-protein contacts

[Gueudré, Baldassi, Zamparo, MW, Pagnani, PNAS '16]
[Bitbol, Dwyer, Colwell, Wingreen, PNAS '16]

Interactions between protein families

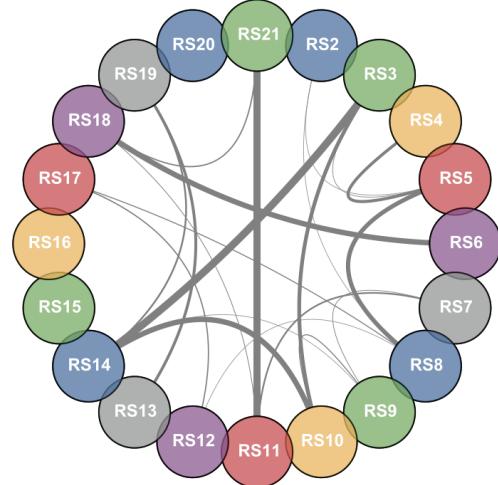


What can we learn from the empirical sequence variability:

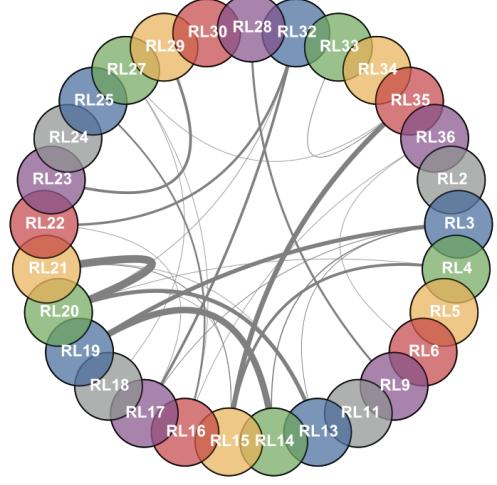
- do the families interact?
- which specific proteins interact?
- which residues are in contact?
- relation between protein structure/function and evolution

Inference of protein-protein interaction networks

SMALL RIBOSOMAL SUBUNIT



LARGE RIBOSOMAL SUBUNIT



Bacterial ribosomal proteins

Small ribosomal subunit

- 20 proteins
- 21 interactions (11% of 190 pairs)

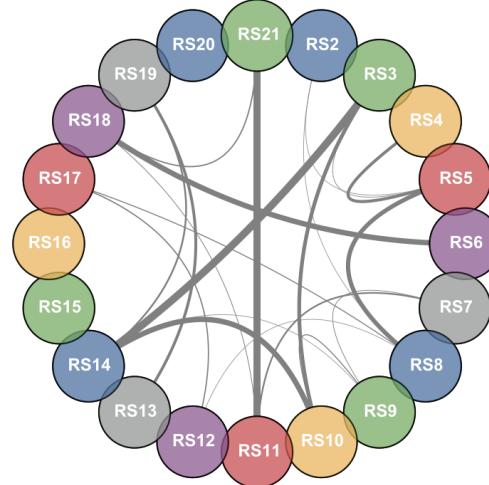
Large ribosomal subunit

- 29 proteins
- 29 interactions (7% of 406 pairs)

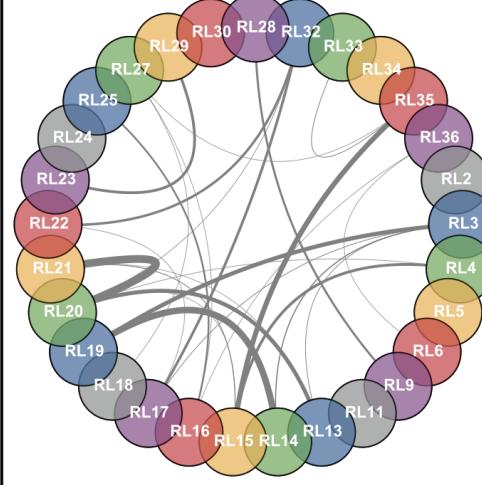
► sparse interaction network

Inference of protein-protein interaction networks

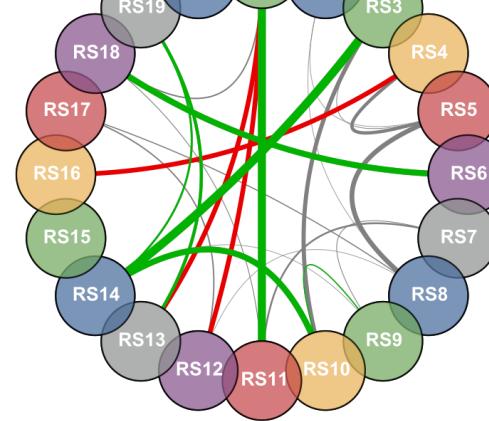
SMALL RIBOSOMAL SUBUNIT



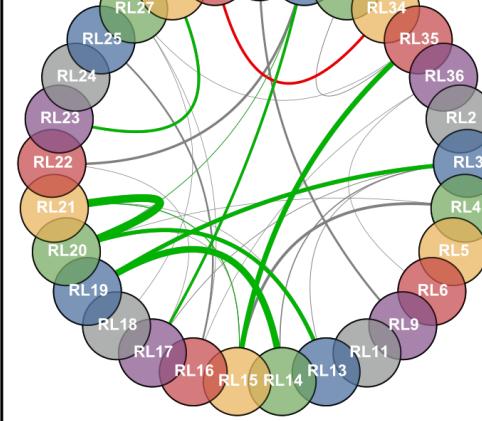
LARGE RIBOSOMAL SUBUNIT



SMALL RIBOSOMAL SUBUNIT



LARGE RIBOSOMAL SUBUNIT



Bacterial ribosomal proteins

Pairwise DCA (1000-3000 seqs.)

Top 10 predictions for each subunit

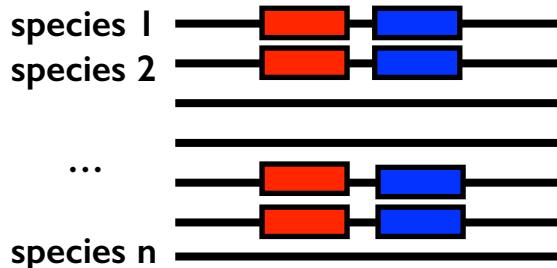
- 16 true positive interactions
(80% TP vs. 8% in random prediction)
- find most large interfaces
- fail to detect small interfaces
- false predictions appear in smaller alignments

► larger alignments needed

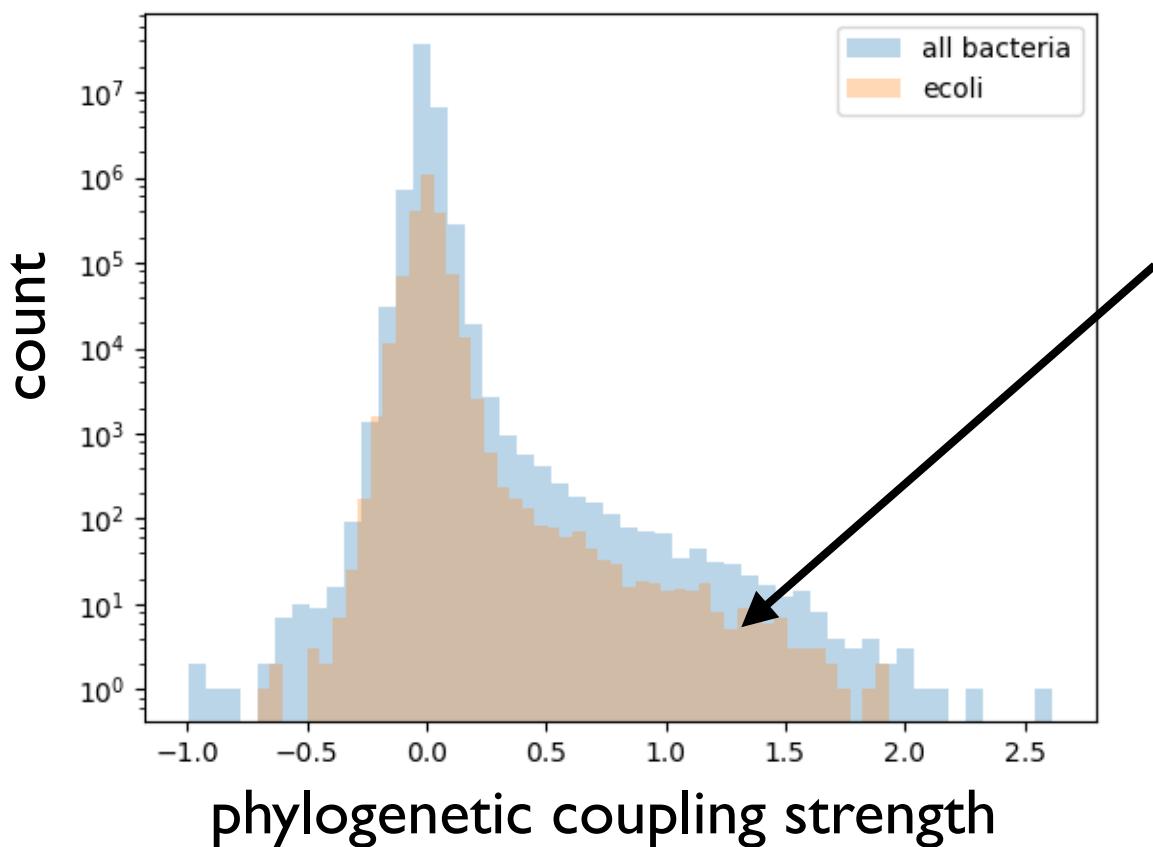
[Feinauer, Szurmant, MW, Pagnani, *PLoS ONE* '16]

cf. also [Uguzzoni, Lovis, Oteri, Schug, Szurmant, MW, *PNAS* '17]

Exploring genomic scales

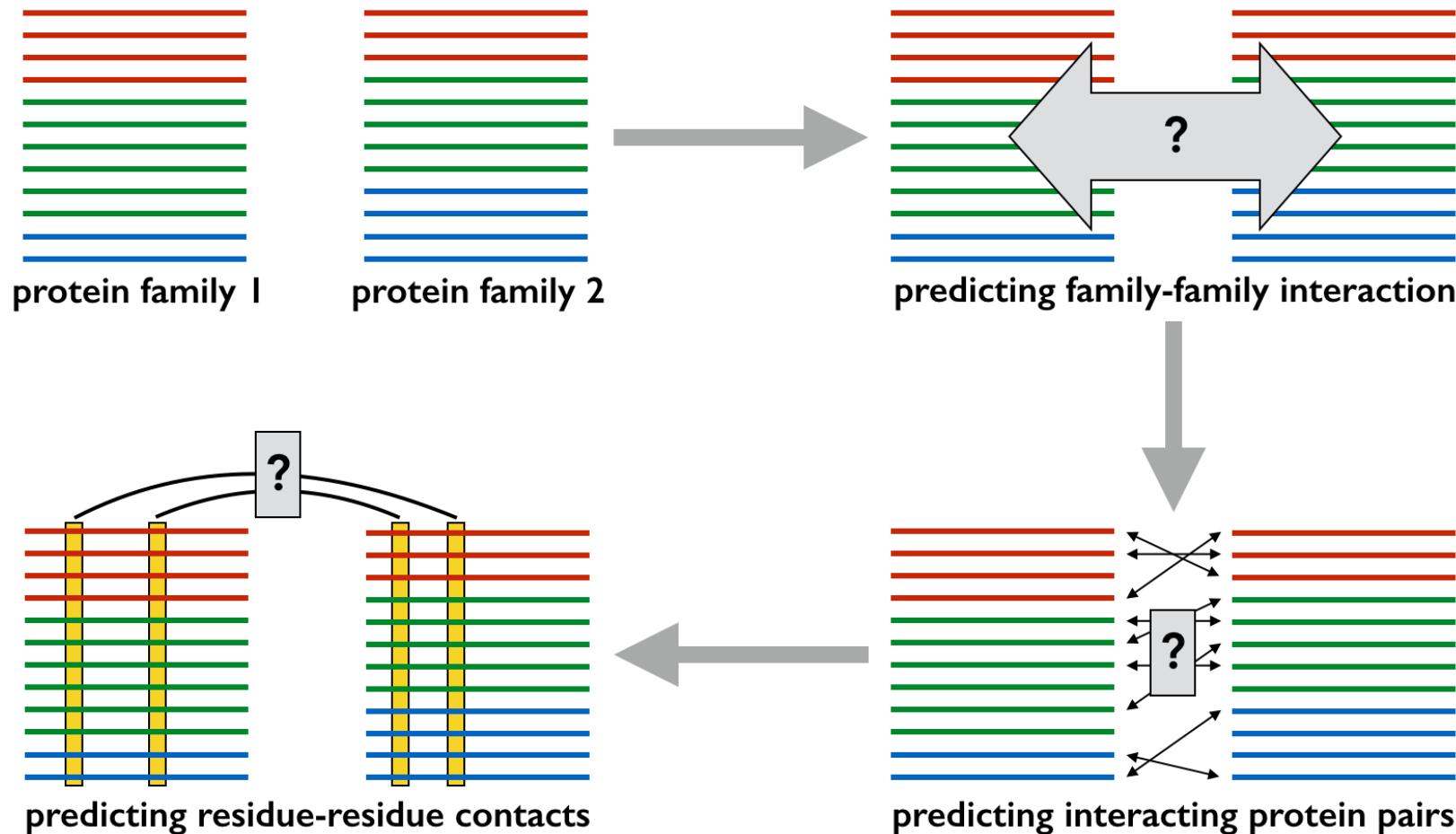


- correlated presence / absence of interacting proteins
 - phylogenetic profiles [Pellegrini et al. 1999]
 - correlated phylogenetic trees [Pazos et al. 2001]
 - phylogenetic coupling analysis [Croce et al., in prep]



- Tail of ~ 1000 strong couplings
 - 80% known relations
(interaction, colocalisation)
 - 20% new predictions

Interactions between protein families



What can we learn from the empirical sequence variability:

- do the families interact?
 - which specific proteins interact?
 - which residues are in contact?
- towards a structurally resolved & evolutionary conserved interactome

Thanks to:

The group in Paris:

Juliana Bernardes
Pierre Barrat-Charlaix
Giancarlo Croce
Kai Shimagaki
Edwin Rodriguez
Francesco Oteri

Alumni:

Eleonora de Leonardis
Guido Uguzzoni
Alice Coucke
Matteo Figliuzzi
Christoph Feinauer

Funding:



Collaborators:

Terry Hwa (UC San Diego)
Hendrik Szurmant (Western U LA)
Alexander Schug (KIT Karlsruhe)
Jose Onuchic (Rice U, Austin)
Faruck Morcos (UT Dallas)
Angel E. Dago (Scripps La Jolla)
Joanna Sulkowska (U Warsaw)
Erik Aurell (KTH Stockholm)
Andrea Pagnani (Politecnico Torino)
Thomas Gueudré (IIGM Torino)
Carlo Baldassi (U Bocconi Milano)
Rémi Monasson (ENS)
Simona Cocco (ENS)
Olivier Tenaillon (Inserm Paris)



Horizon 2020
European Union funding
for Research & Innovation