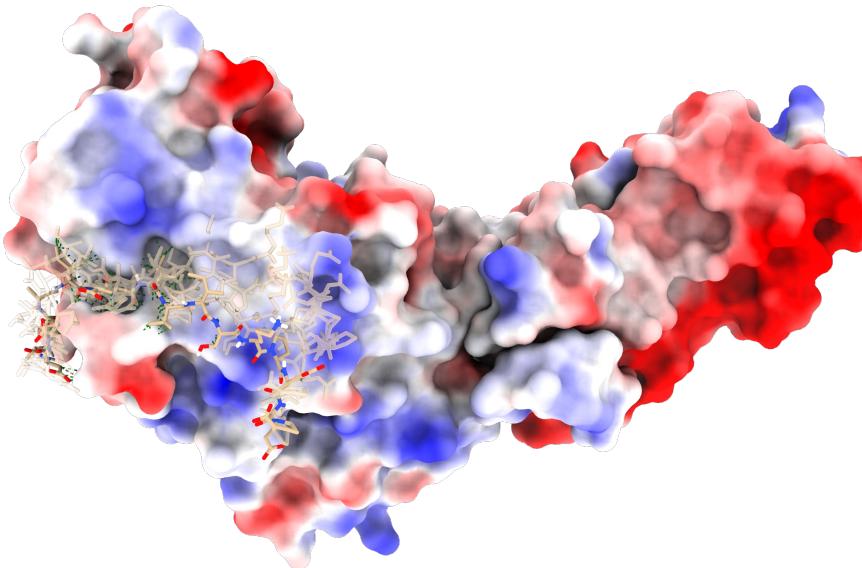


Harnessing Sequence Generative Models for Inhibitory Peptide Design: a Case Study



Dr. Jérôme Tubiana

BeVAS 18/04/2023



Joint work with:

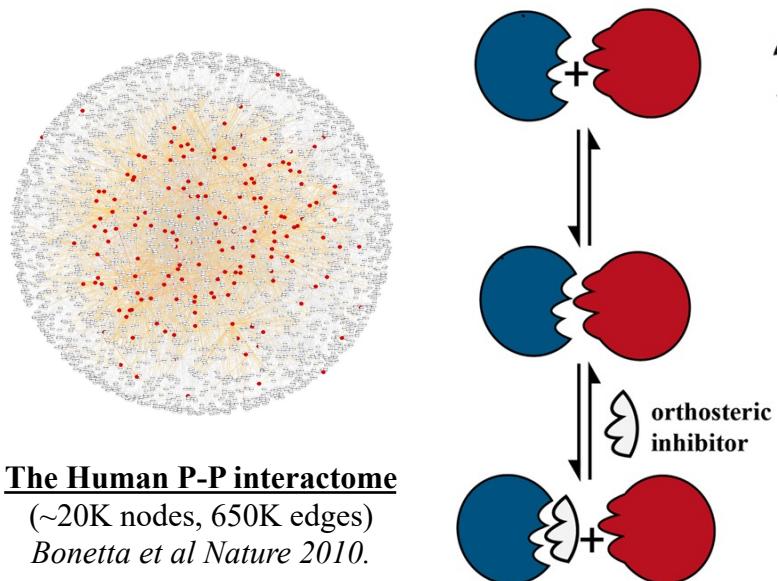
Haim J. Wolfson lab

Maayan Gal lab

Lucia Adriana-Lifshits

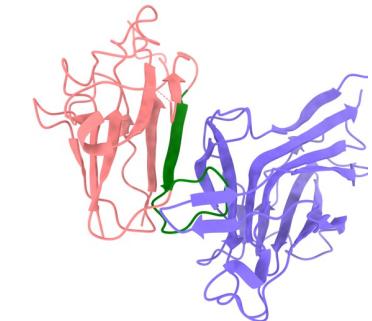
Protein-protein Interaction (PPIs) Inhibitors

- Aberrant PPIs are associated with various diseases.
- PPIs are attractive targets for basic research, therapeutic & pesticide purpose.
- Interfering with PPIs with small molecules is challenging, due to their physio-chemical properties.
 - PPI interfaces are larger, flatter, hydrophobic
 - 40% of PPIs involve a disordered partner
- mAb widely successful, but only for extracellular targets

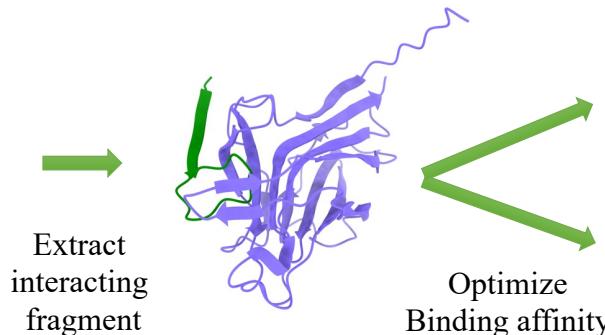


Credit: Lu et al. *Nature Signal Transduction and Targeted Therapy* 2020

Inhibiting PPIs with peptides: principle, benefits and challenges



EphB4-EphB2 complex



Display experiments: Mutagenesis & selection



Raveh, London, Schueler-Furman
PLOS CB 2010

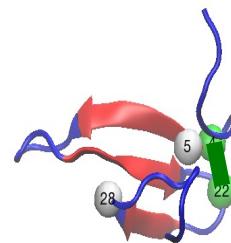
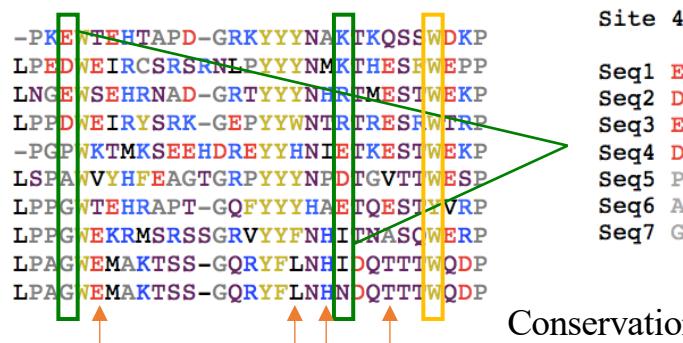
In-silico docking & binding energy optimization



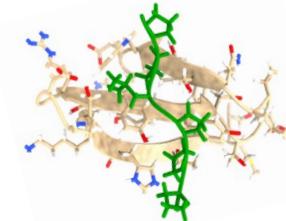
Donsky and Wolfson
Bioinformatics 2011

- ✓ Peptides are suitable for binding protein-protein interfaces
- ✓ Native protein-peptides interactions are highly specific
- ✓ Initial peptide can be derived from native substrate in >50% cases (London et al. *Proteins 2011*)
- ✗ Flexibility makes computational modeling of protein-peptide interactions challenging (sampling & scoring)
- ✗ Display experiments can be difficult to setup
- ✗ Unclear how to efficiently explore the vast search space (20^L peptides of length L)
- ✗ Unclear how to select for specificity and other desirable properties (bioavailability, immunogenicity)

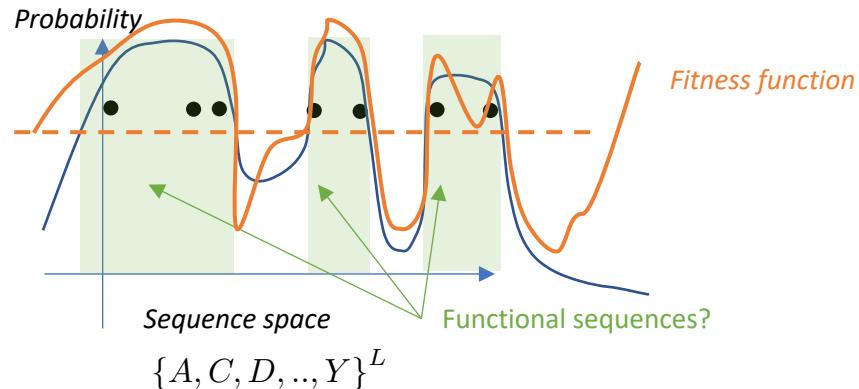
Evolutionary-based sequence generative models for protein design



Coevolution



High-order coevolution



- Non-exhaustive list of successful SGM-based design experiments:
1. Chorismate Mutase (Russ et al. Science 2020)
 2. Luciferases (Hawkins-Hooker et al. PLOS CB 2021)
 3. Malate Dehydrogenase (Repecka et al. Nat. Mach. Int. 2021)
 4. Nanobody libraries (Shin et al. Com. 2021)
 5. GFP (Biswas et al. Nat. Methods 2021)
 6. SH3 domains (Lian et al. BiorXiv 2022)
 7. Copper Superoxide Dismutase (Johnson et al. BiorXiv 2023)
 8. Cas9 PAM-interacting domain (Malbranque et al. BiorXiv 2023)

Evolutionary-based generative models for peptide design

Table 1 Peptide generation studies using deep generative models. Abbreviations: NML, neural language model; VAE, variational autoencoder; GAN, generative adversarial network; AMP, antimicrobial peptide; ACP, anticancer peptide; CPP, cell-penetrating peptide; PMO, phosphorodiimidate morpholino oligomer

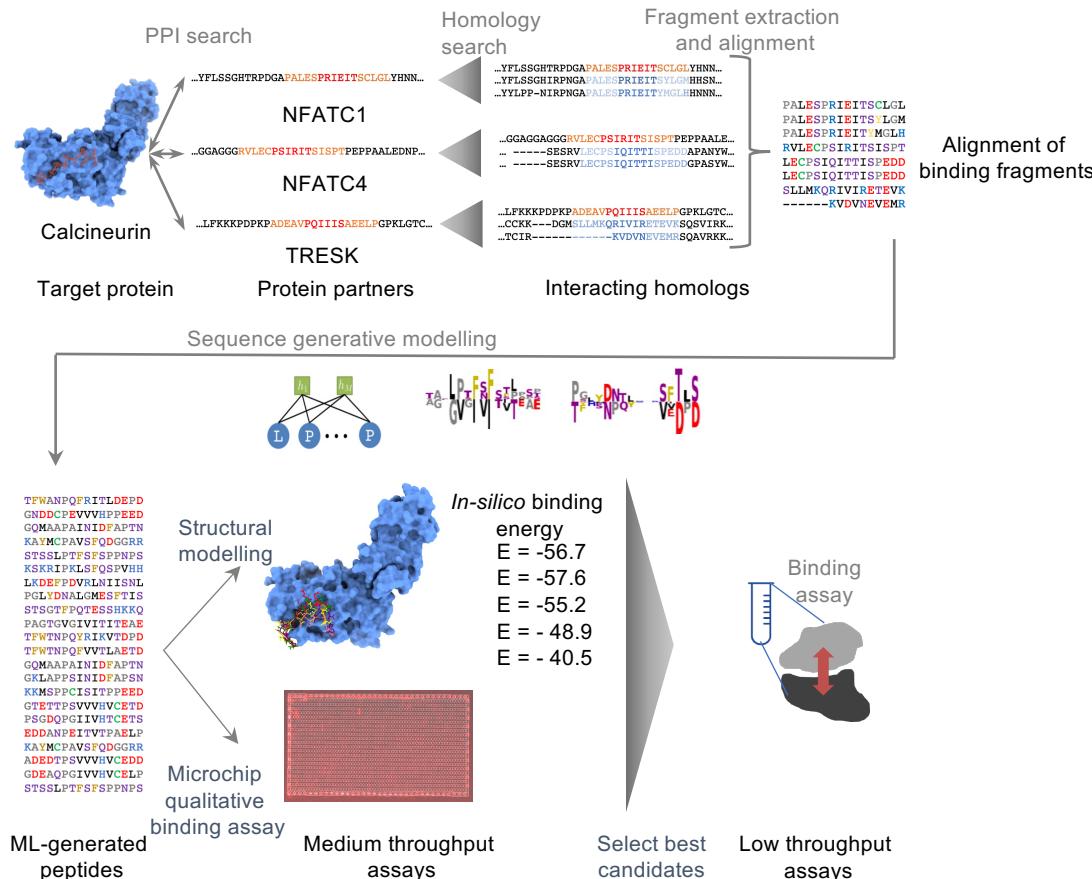
Method	Feature Representation	Application	Citation	Year
NML	One-hot	AMP generation	Müller <i>et al.</i> ⁵⁶	2018
NLM	Character sequence	AMP generation	Nagarajan <i>et al.</i> ⁵²	2018
NLM	Character sequence	ACP generation	Grisoni <i>et al.</i> ⁴⁶	2018
NLM	Learned representation using one-hot	Signal peptide generation	Wu <i>et al.</i> ⁵⁵	2020
NLM	Learned representation using structural and evolutionary data	AMP generation	Caceres-Delpiano <i>et al.</i> ⁵⁴	2020
NLM	One-hot	AMP generation	Wang <i>et al.</i> ⁴¹	2021
NLM	Character sequence	CPP generation	Tran <i>et al.</i> ⁵³	2021
NLM	One-hot	AMP generation	Capecchi <i>et al.</i> ⁴²	2021
NLM	Fingerprint, one-hot	PMO delivery peptide generation	Schissel <i>et al.</i> ⁵⁷	2021
VAE	Learned representation using character sequence	AMP generation	Das <i>et al.</i> ³⁸	2018
VAE	Learned representation using one-hot	AMP generation	Dean <i>et al.</i> ⁴⁴	2020
VAE	Learned representation using character sequence	AMP generation	Das <i>et al.</i> ⁴⁵	2021
GAN	Character sequence	AMP generation	Tucs <i>et al.</i> ⁵⁰	2020
GAN	Character sequence/PDB structure	ACP generation	Rossetto <i>et al.</i> ⁴⁸	2020
GAN	Learned representation using character sequence	AMP generation	Ferrell <i>et al.</i> ³⁹	2020
GAN	Character sequence	AMP generation	Oort <i>et al.</i> ⁴⁰	2021
GAN	Sequence of amino acid property vectors	Immunogenic peptide generation	Li <i>et al.</i> ⁵¹	2021
GAN	Character sequence	AMP generation	Surana <i>et al.</i> ⁴⁹	2021

Table reproduced from *Deep generative models for peptide design*
F. Wan, D. Kontogiorgos-Heintz and C. de la Fuente-Nunez, Digital Discovery 2022

Adapting to PPI inhibitor design?

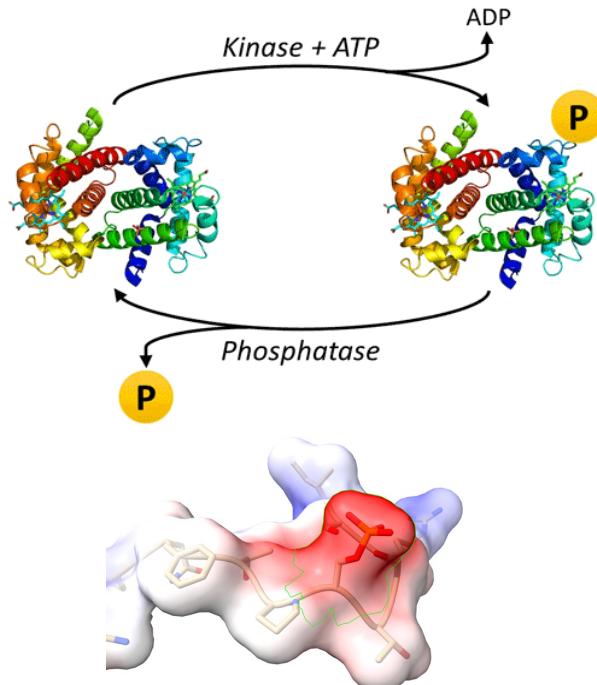
- 1/ How to gather sufficiently diverse MSA?
- 2/ How to go beyond binding affinity of natural peptides?

An integrative PPI inhibitor peptide design protocol

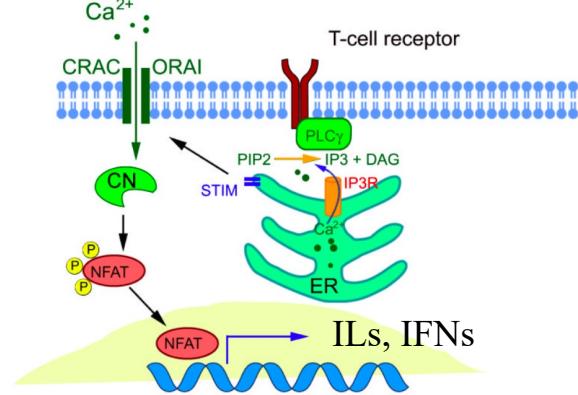


The Calcineurin (Cn) signaling pathway

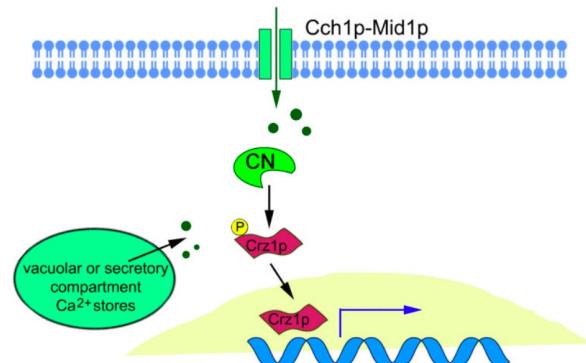
Calcineurin (Cn) is a **calcium-dependent phosphatase** involved in multiple health & disease pathways.



Phosphoserine (PDB: 1t29)



NFAT-mediated T-cell activation
(Conserved in vertebrates)

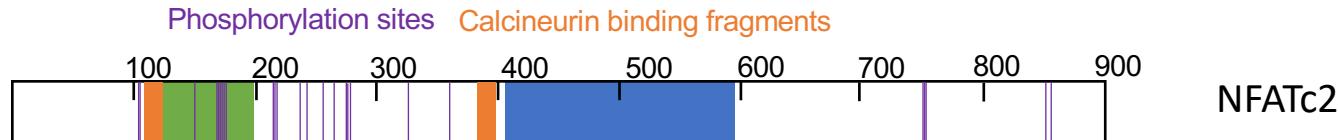


Crz1 environmental stress response
pathway (Lower eukaryotes)

(Li et al. Trends Cell. Biology 2011)

Structural basis of Calcineurin function

Most Cn substrates are disordered. They bind via two conserved short linear motifs (SLiMS)

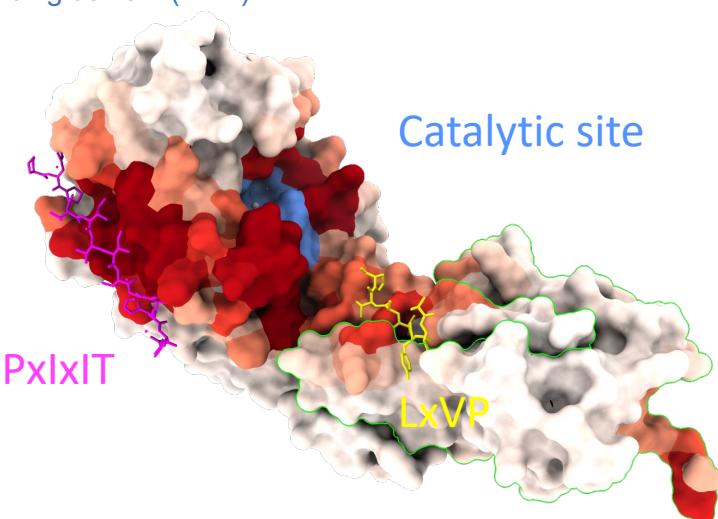


Trans-activation domain (TAD-A) DNA-binding domain (RHD)

PxIxIT

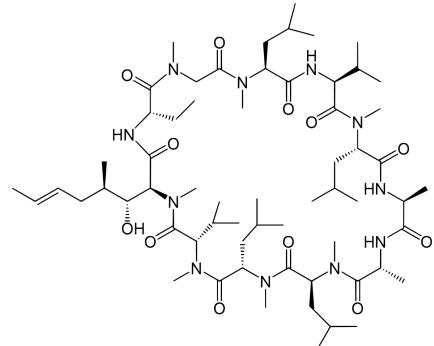
NFATc1	PALES PRIEIT SCLGL
NFATc2	ASGLS PRIEIT PSHEL
NFATc3	KPFEC PSIQIT SISPN
NFATc4	RVLEC PSIRITS SISPT
TRESK	ADEAV PQIIIS AEELP
AKAP79	KRME PIAIIT DTEIS
RCN1	GTTNT PSVIV HVCEDD
RIPOR2	SNSTN PEITIT PAEFN
CAPN11	TFWTNPQFK ISLPEGD
SFB3	RAYAN PKQFTY DSSV
CRZ1	AAPVT PIISI QEFNEG

Alignment of PxIxIT-containing
fragments from substrates

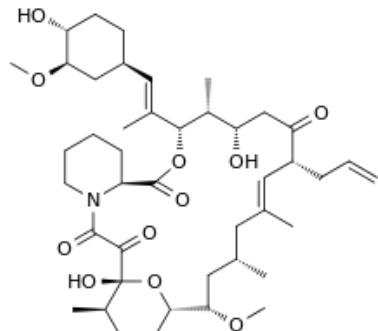


Calcineurin in **active conformation** bound to
representative substrate fragments

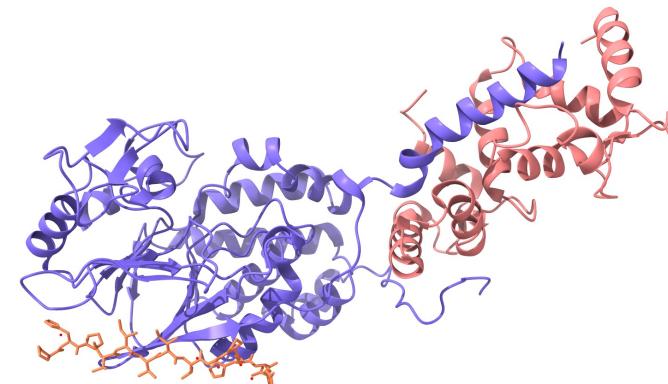
Calcineurin inhibitors



Cyclosporine A (CsA)



Tacrolimus (FK506)



The **PVIVIT** peptide
Aramburu et al. 1998 Science
(Combinatorial library + Display experiment)

Catalytic site inhibitors

Prescribed for transplantations since 80's

Protein Interaction inhibitor

Mice preclinical studies

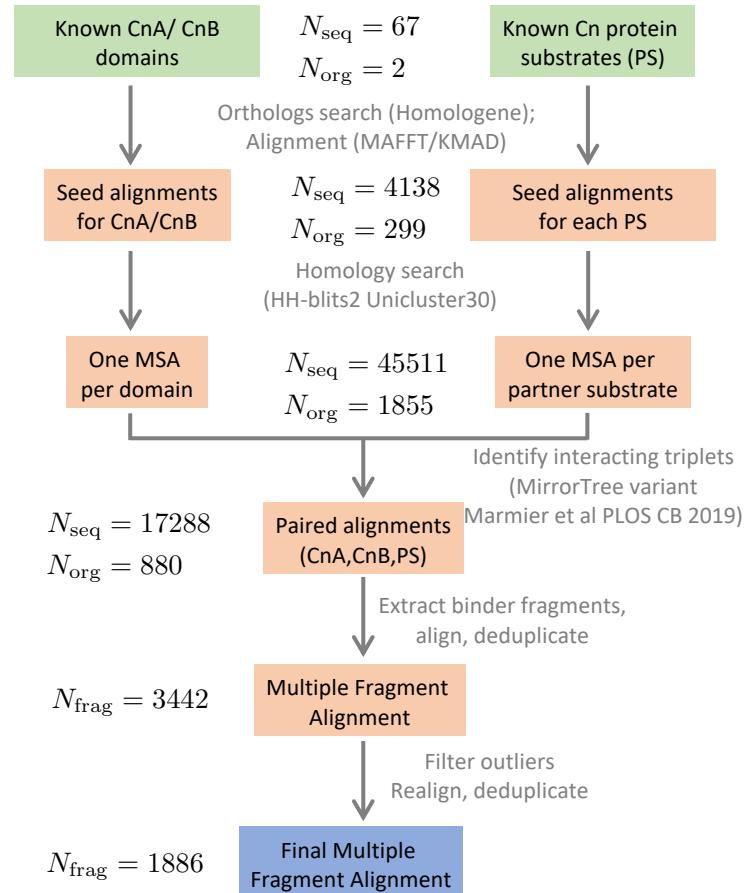
(Noguchi et al. Nature Medicine 2004)

Step 1: Construction of an alignment of putative binding fragments

Input: a list of 67 Calcineurin substrates.
Sources: integrative high-throughput experiments

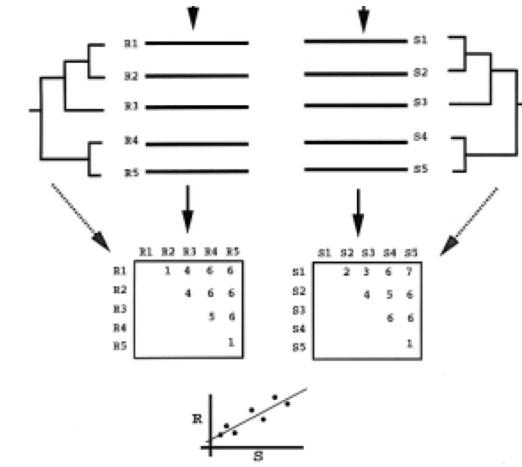
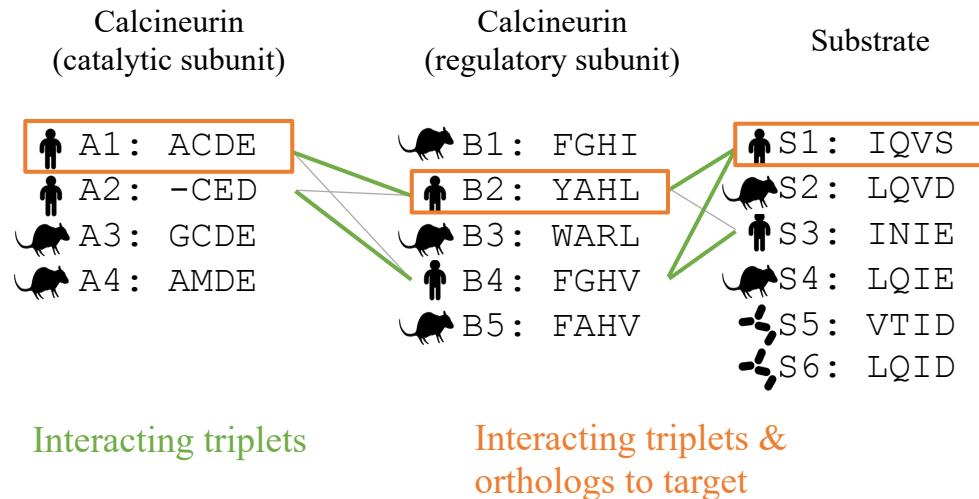
- Goldman A, Roy J, Bodenmiller B, Wanka S, Landry CR, Aebersold R, et al. The calcineurin signaling network evolves via conserved kinase-phosphatase modules that transcend substrate identity. *Mol Cell.* 2014
- Wigington CP, Roy J, Damle NP, Yadav VK, Blikstad C, Resch E, et al. Systematic Discovery of Short Linear Motifs Decodes Calcineurin Phosphatase Signaling. *Mol Cell.* 2020 (Cyert lab)

Gene	Organism	SLIM
NFATC1	Homo Sapiens	PRIEIT
NFATC2	Homo Sapiens	PRIEIT
NFATC3	Homo Sapiens	PSIQIT
NFATC4	Homo Sapiens	PSIRIT
TRESK	Homo Sapiens	PQIIIS
CRZ1	Saccharomyces Cerevisiae	PIISIQ
RCN1	Saccharomyces Cerevisiae	GAITID
SFB3	Saccharomyces Cerevisiae	PKFQFT
RGA2	Saccharomyces Cerevisiae	PQVLVS
ROD1	Saccharomyces Cerevisiae	PQIKIE
STE12	Saccharomyces Cerevisiae	PALSFS
RTS1	Saccharomyces Cerevisiae	PVLTVT
SLM1	Saccharomyces Cerevisiae	PNIYIQ
SLM2	Saccharomyces Cerevisiae	PEFYIE
RPL4A	Saccharomyces Cerevisiae	PQVTVH
RCN2	Saccharomyces Cerevisiae	PSITVN
DIG2	Saccharomyces Cerevisiae	PALNFS



The MSA pairing problem

C_n-substrate interactions are not systematically conserved across homologs



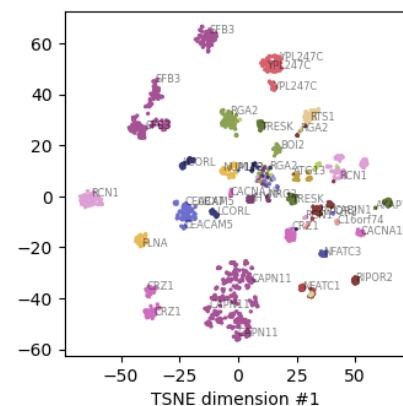
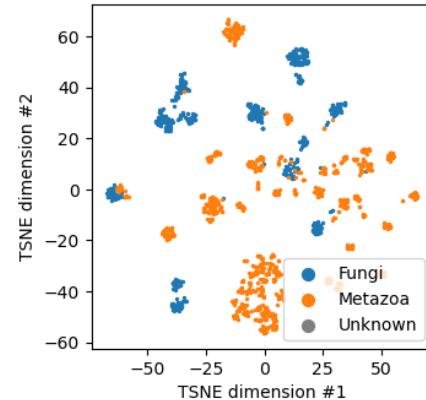
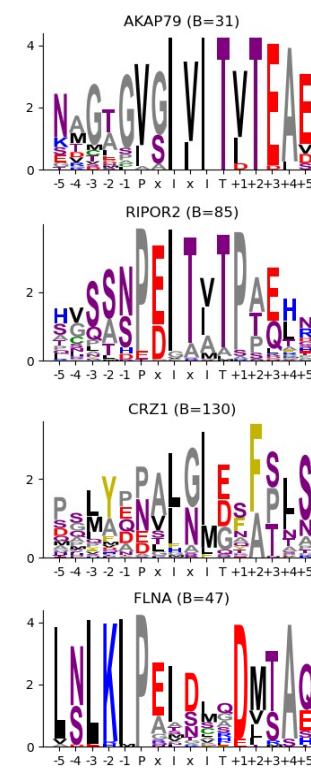
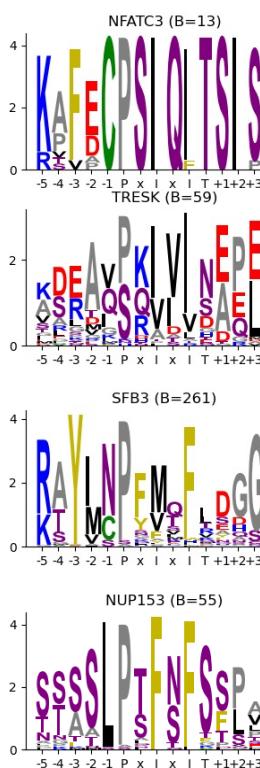
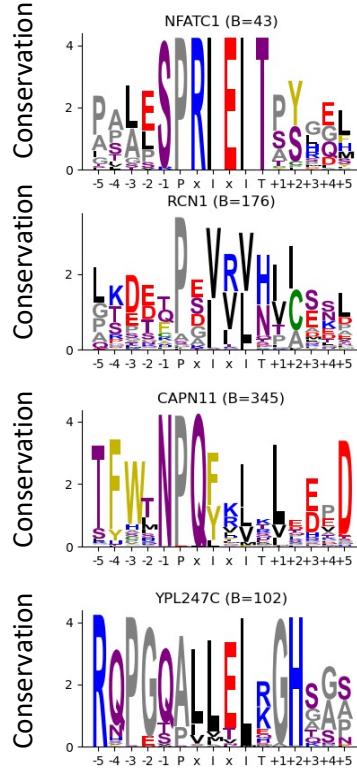
MirrorTree method
(Pazos and Valencia 1994)

Find pairing that maximizes key fingerprints of interacting proteins:

- Interacting partners sequences tend to mutate at similar rates
- Binding sites tend to coevolve

Marmier et al PLOS CB 2019)

Calcineurin-binding fragments are highly diverse



Step 2: Sequence Generative Modeling (compositional Restricted Boltmann Machines)

Graphical model constituted by two coupled sets of random variables

$$P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp [-E(\mathbf{v}, \mathbf{h})]$$

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^N g_i(v_i) + \sum_{\mu=1}^M U_\mu(h_\mu) - \sum_{i,\mu} w_{i\mu}(v_i)h_\mu$$

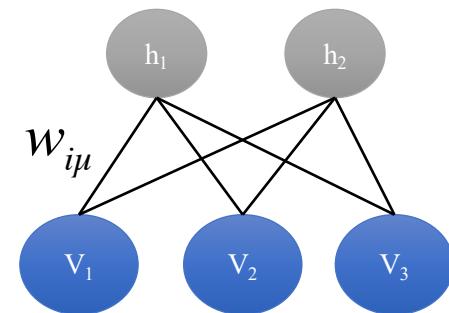
The marginal defines a probability distribution over the data space

$$P(v) = \int \prod_\mu dh_\mu P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp \left[\sum_i g_i(v_i) + \sum_\mu \Gamma_\mu \left(\sum_i w_{i\mu}(v_i) \right) \right]$$

*Trainable, non-quadratic function
(generalizes over pairwise models)*

*Sparse weight matrix
(Interpretability, compositionality)*

Hidden (Representation) layer



Visible (Data) layer

Ackley Sejnowski Hinton 1985
Smolensky 1986

The conditional distribution defines a representation of data

$$\langle \mathbf{h}_\mu | \mathbf{v} \rangle = \Gamma'_\mu \left[\sum_i w_{i\mu}(v_i) \right]$$

Tubiana, Monasson PRL 2017

Tubiana, Cocco, Monasson eLife 2019

Tubiana, Cocco, Monasson Neur. Comp. 2019

Training algorithm for RBM

Data set: $\{\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^B\}$

Want to maximize log-likelihood: $\mathcal{L} = \sum_b \log P(\mathbf{v}^b | \{w_{i\mu}\}, \{g_i\}, \{\mathcal{U}_\mu\})$

Stochastic gradient ascent:

parameters Θ

$$\log P(\mathbf{v}) = -E_{\text{eff}}(\mathbf{v}) - \log Z$$

$$\frac{\partial \mathcal{L}}{\partial \Theta_a} = \left\langle \frac{\partial E_{\text{eff}}(\mathbf{v}|\Theta)}{\partial \Theta_a} \right\rangle_{\mathbf{v} \sim \text{RBM}} - \left\langle \frac{\partial E_{\text{eff}}(\mathbf{v}|\Theta)}{\partial \Theta_a} \right\rangle_{\mathbf{v} \sim \text{Data}}$$

Moment
Matching
Equations

Requires MCMC
sampling

Computed directly
from data

Learning algorithms : Boltzmann Machine Learning (Ackley Hinton Sejnowski 1985),
PCD (Tieleman Hinton 2008)

Sampling from RBM

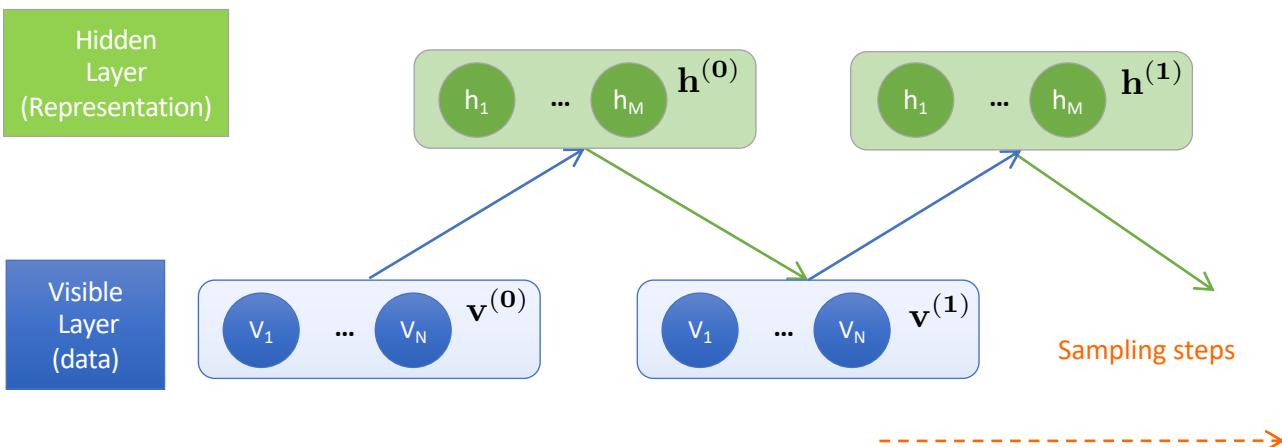
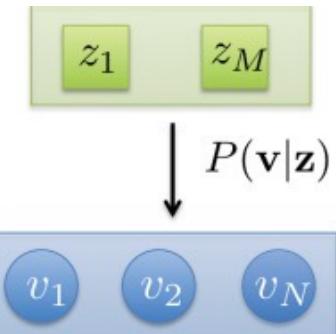
- Compute Hidden units Inputs
- Sample each hidden unit independently
- Compute the visible layer inputs
- Sample each visible unit independently

$$\begin{aligned} I_\mu &= \sum_i w_{i\mu} v_i \\ P(h_\mu | I_\mu) &\propto \exp [-\mathcal{U}_\mu(h_\mu) + h_\mu I_\mu] \\ I_i &= \sum_\mu w_{i\mu} h_\mu \\ P(v_i | I_i) &\propto \exp [(g_i + I_i) v_i] \end{aligned} \quad \left. \right\}$$

Directed Latent variables model

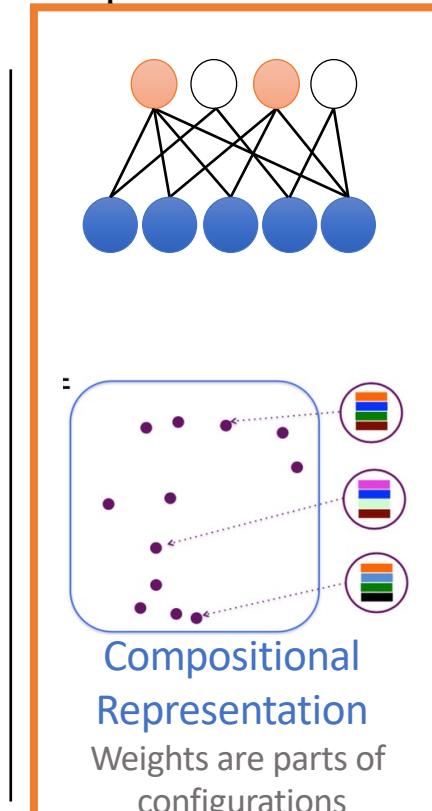
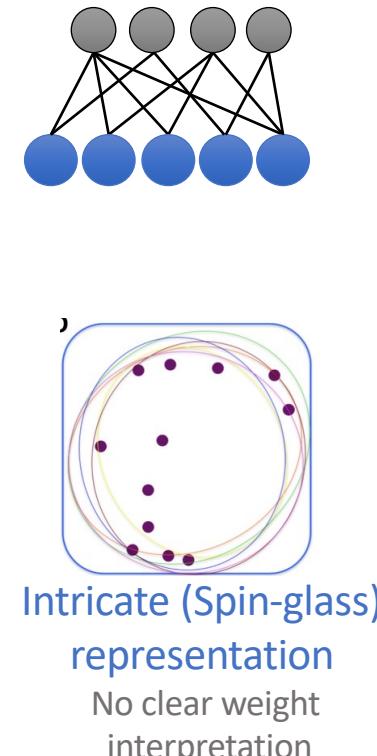
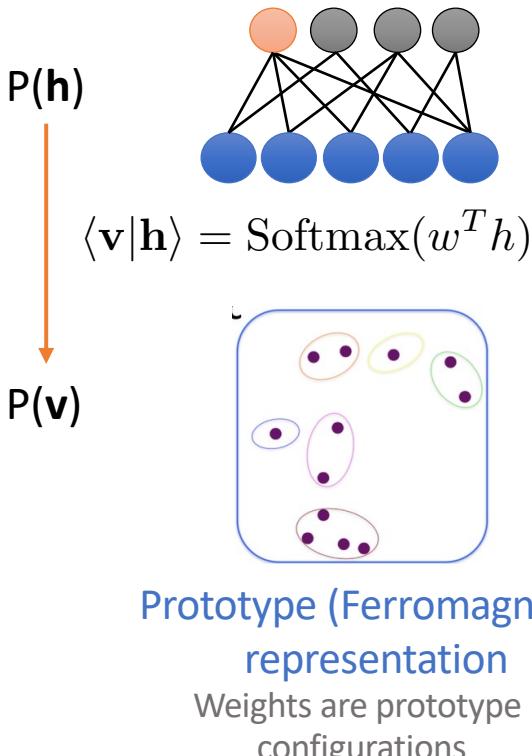
- PCA, ICA
- Sparse dictionaries
- Variational Autoencoders

$$P(\mathbf{z}) = \prod_\mu P_\mu(z_\mu)$$

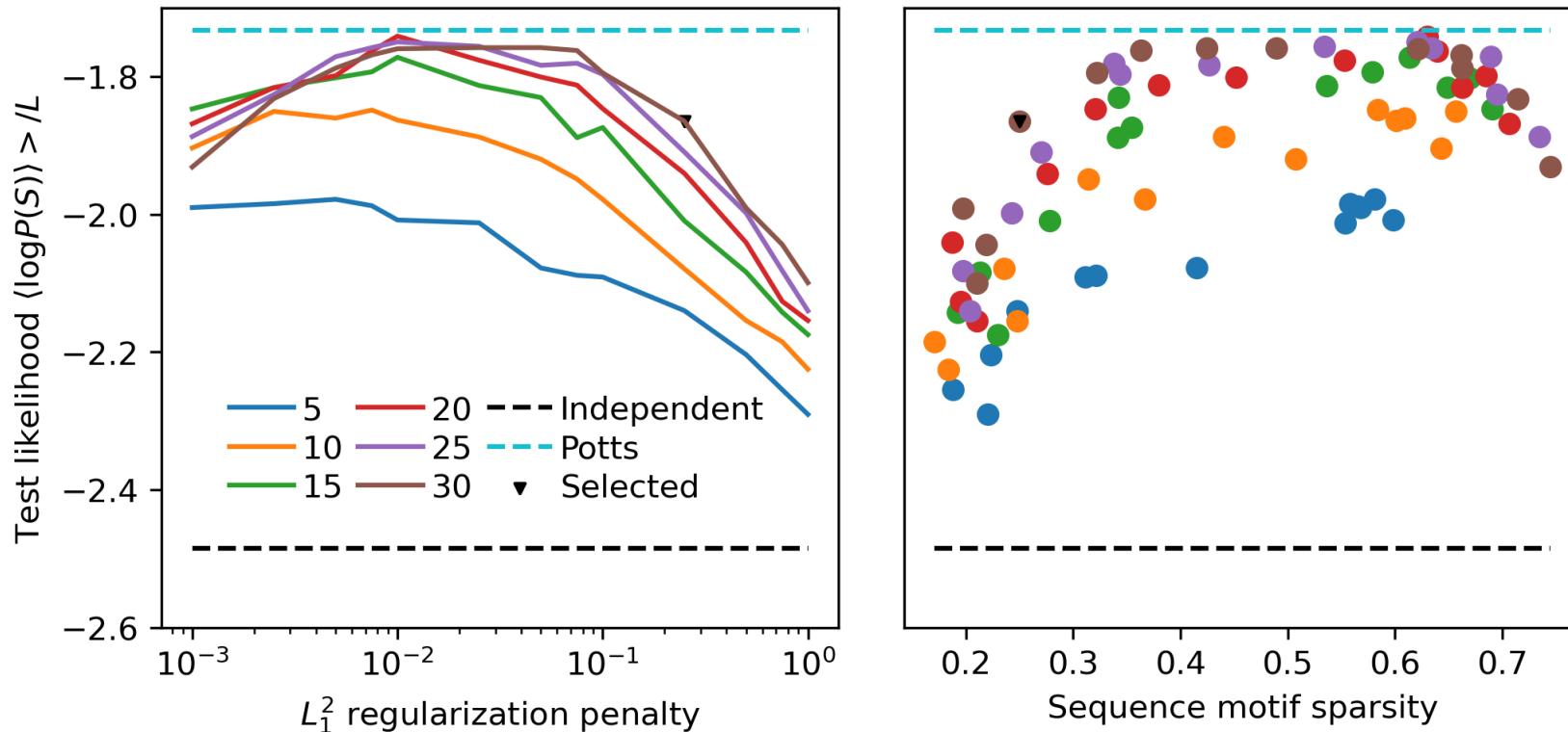


Compositional Restricted Boltzmann Machines

For latent variable generative models, hidden unit distribution guides **weight interpretation & extrapolation regime**. For RBMs, it is unspecified.



Learning cRBM: the interpretability-accuracy trade-off



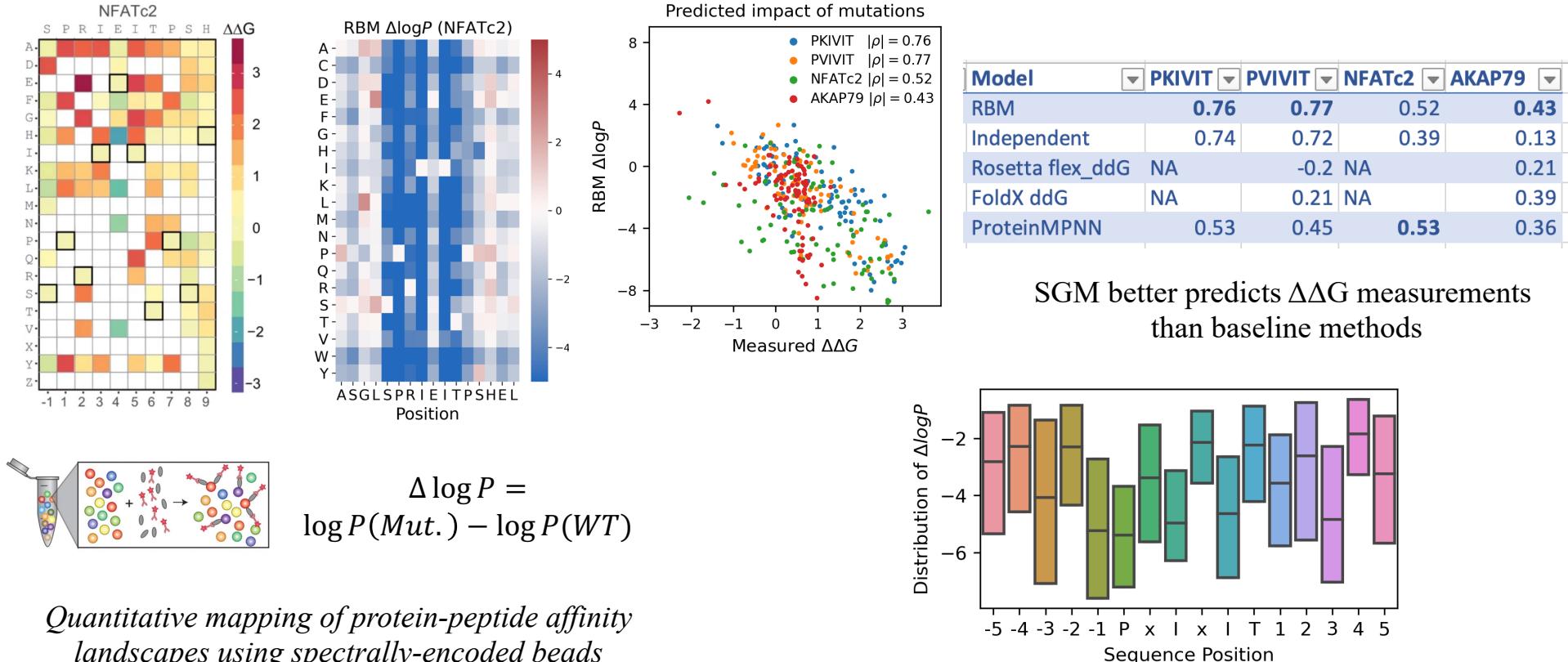
Practical considerations for learning cRBMs

<https://github.com/jertubiana/PGM>

- Objective function: Maximum likelihood + regularization penalties
- Sampling algorithm: MCMC, PCD.
- Optimizer: RMSprop (adaptive learning rates, improves convergence rates).
- Hidden unit potential: dReLU (adaptive non-linearity, for fitting non-gaussian distributions)
- Parameterization: Batch normalization (improves hessian conditioning, promotes homogeneity).
- Regularization: L_2 on fields, L_1^2 on weights (promotes sparsity+homogeneity).
- Partition function estimation: Annealed Importance Sampling.

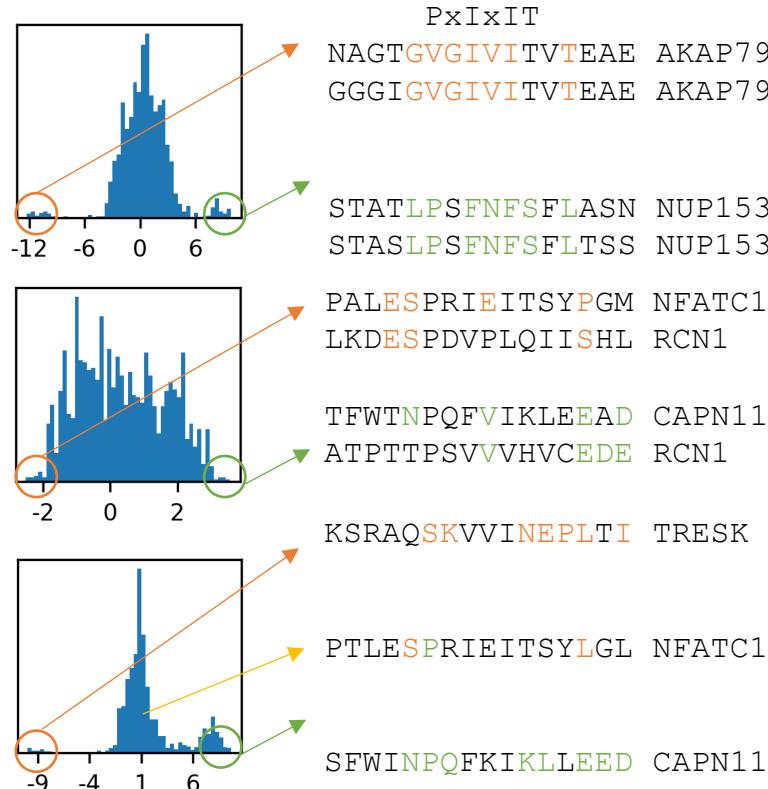
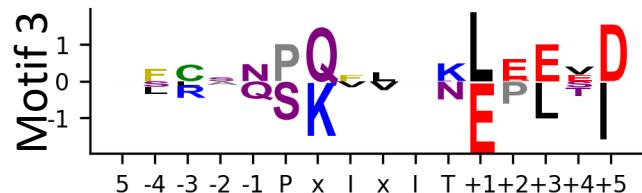
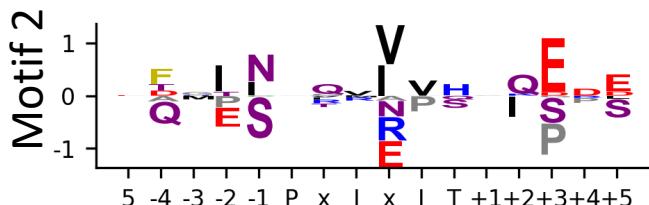
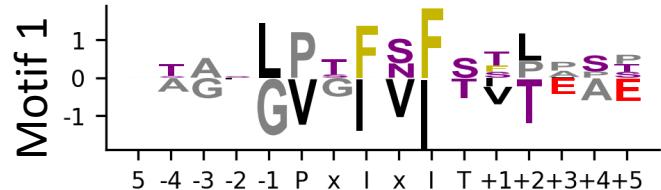


SGM predicts binding affinity changes upon mutation

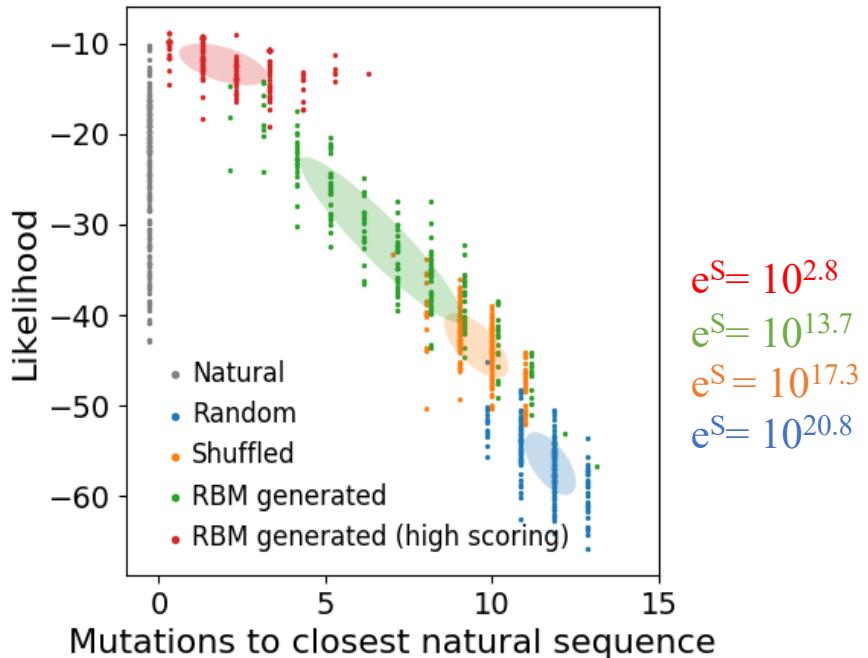


SGM learns sequence motifs shared between evolutionary-unrelated binders

$$P(S) = \frac{1}{Z} \exp\left[\sum_{i=1}^N g_i(s_i) + \sum_{\mu=1}^M \Gamma_\mu\left(\sum_{i=1}^N w_{i\mu}(s_i)\right)\right]$$



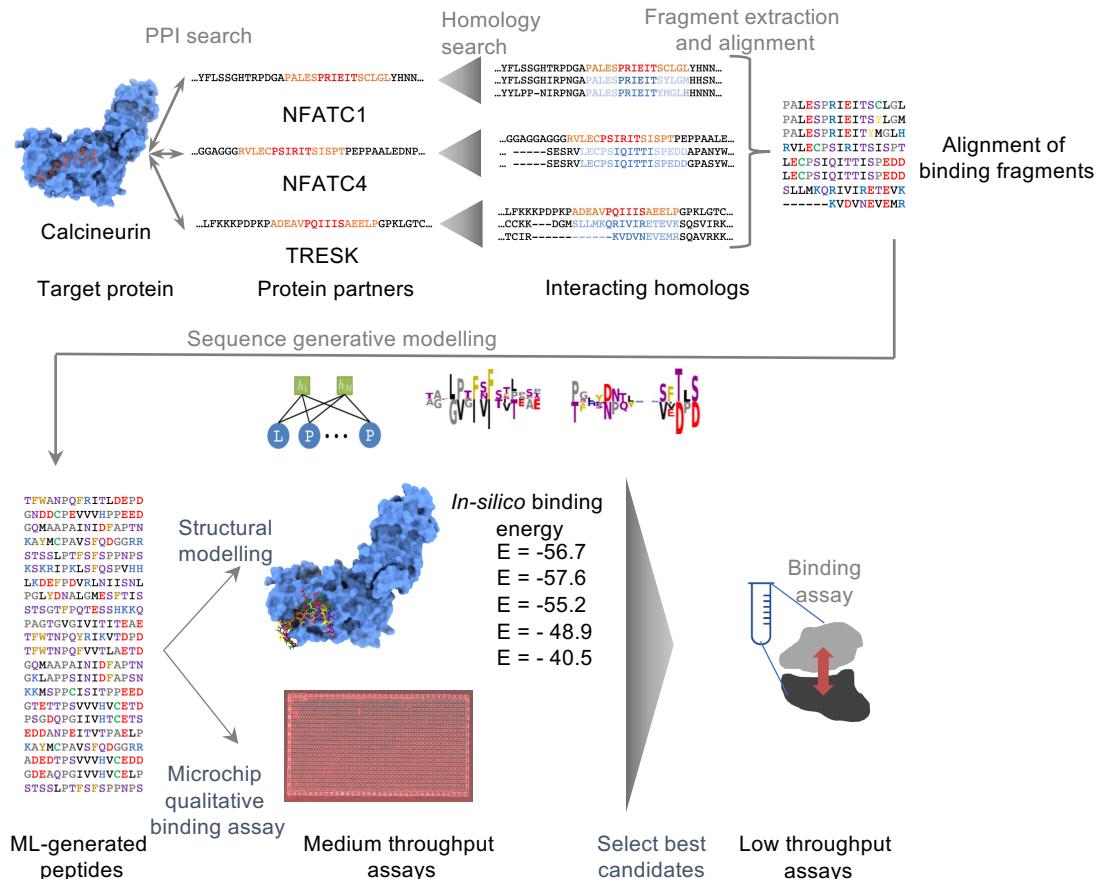
Peptide Library generation



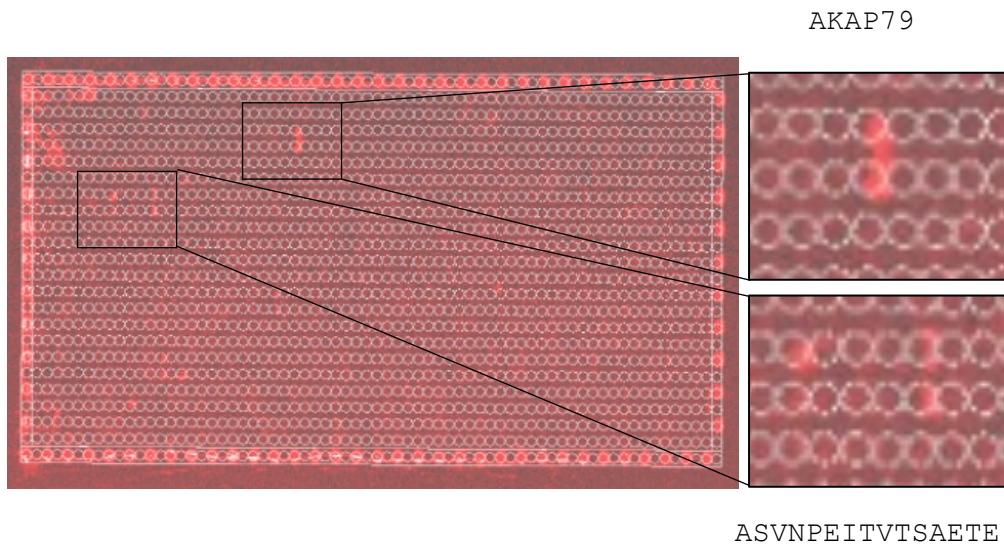
Novelty/Diversity-Quality trade-off

Source	Role	#Num sequences
Random peptides	Negative control	36
Literature designs	Positive control	2
Natural peptides	Positive (?)	75
Independent (PSSM)	Baseline design	72
cRBM, $\beta = 1$	Design	180
cRBM, $\beta = 2$	Design	361

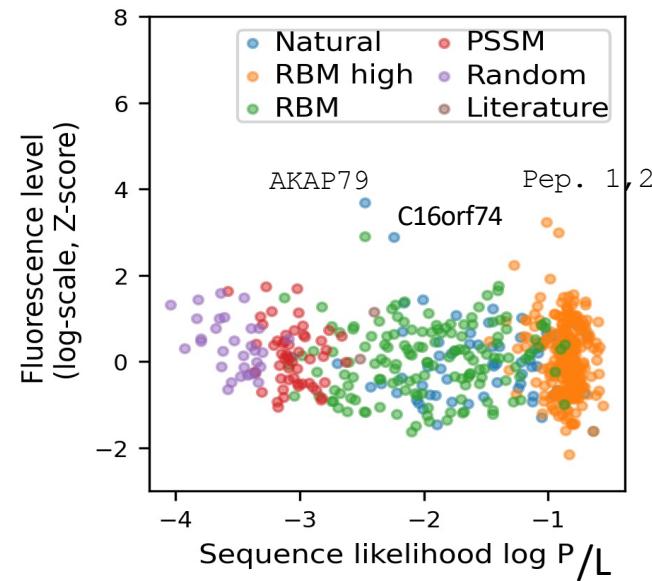
An integrative PPI inhibitor peptide design protocol



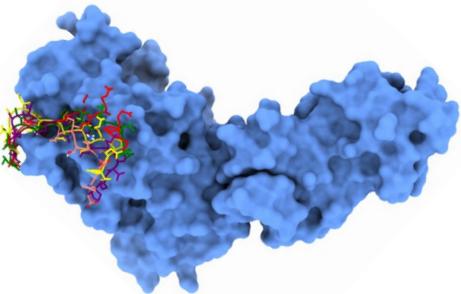
Step 3: Library filtering by microarray screening experiment



PepPerChip Microarray screening

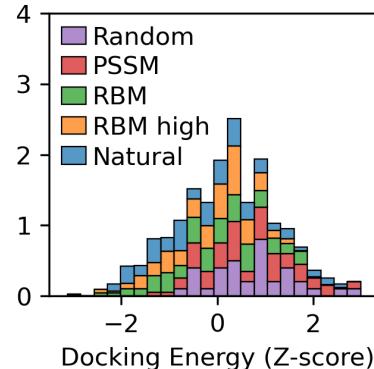


Step 3': Library filtering by molecular docking

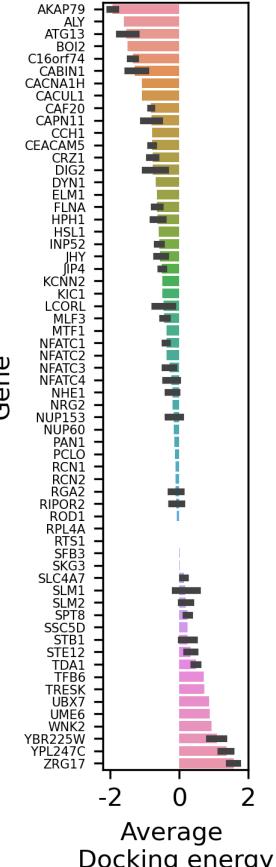
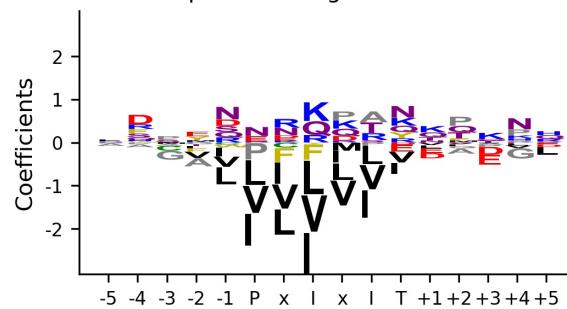


In-silico docking energy score:

- Ensemble of five bound crystal structures as templates
- Threading with Modeller
- Flexible refinement and scoring using PepCrawler (average of minimum energies).



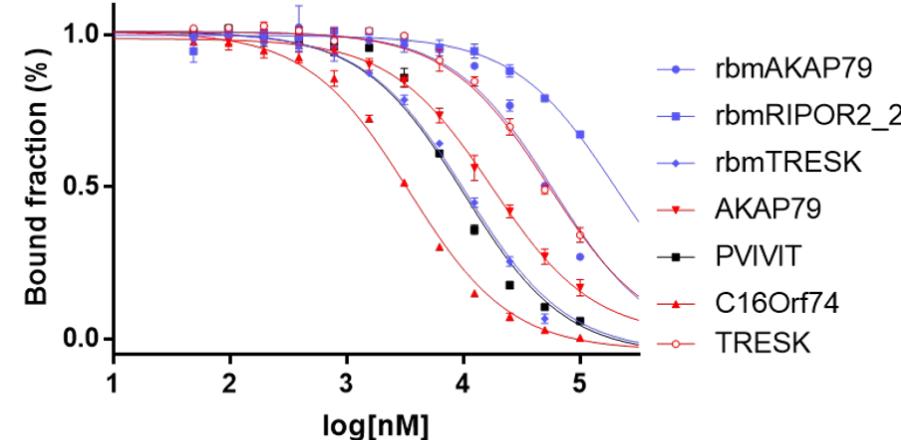
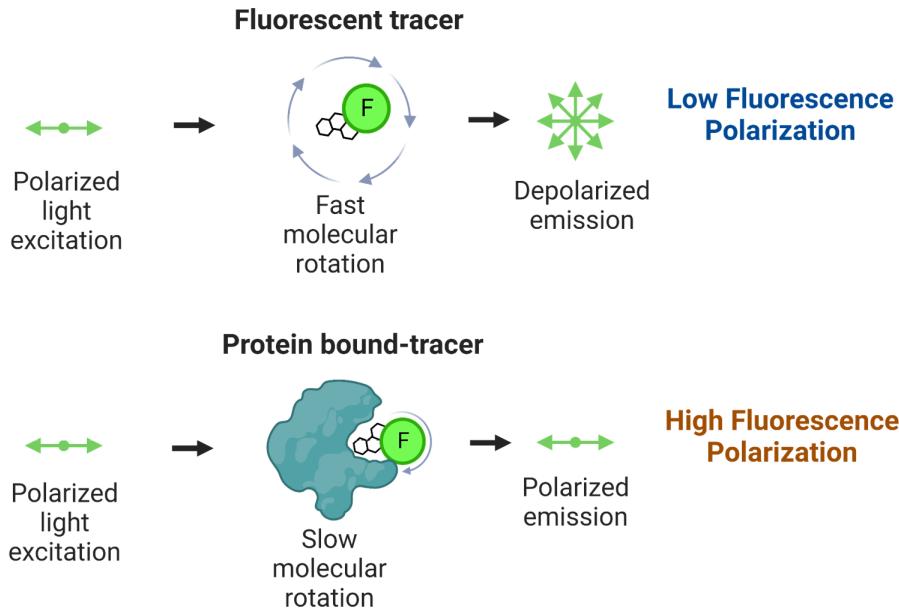
Equivalent single-site model



PepCrawler: a fast RRT-based algorithm for high-resolution refinement and binding affinity estimation of peptide inhibitors

Donsky and Wolfson Bioinformatics 2011

Step 4: Experimental validation by FP binding assay



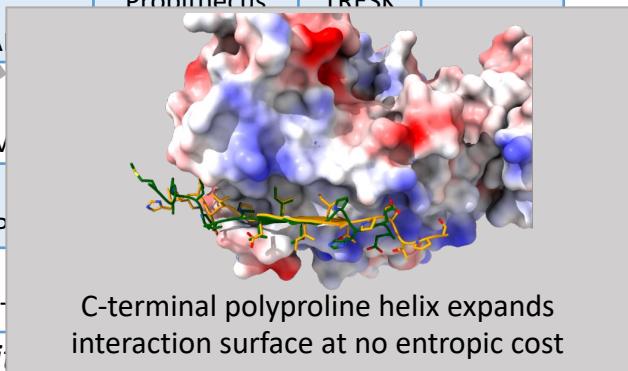
Credit: <https://bpsbioscience.com/product-types/biochemical-assay-kits-by-format-type/fluorescence-polarization>

Step 4: Experimental validation by FP binding assay

Name	Sequence	IC50 (uM)	Source	Closest natural sequence	Motif recombination:
C16Orf74	KHLDVPDIITPPPTP	1.17	Natural	KHLDVPDIITPPPTP	ADEAIPEIVISKPEEP [Design] ADEAI PQIVIDAGADE [TRESK, 50uM]
PVIVIT	MAGPHPVIVITGPHEE	10.2	Positive control	/	SPSNP PEIVIS SREDN [KCNN3] TFTWTNPQFKIYL PEED [CAPN11]
rbmTRESK	ADEAIPEIVISKPEEP	14	Designed (low T)	ADEAVPQIIASAELP	Lack PxIxIT SLIM: AAGAG VGIVIT VTEAE ADGAG VGIVIT VTEAE
AKAP79	KRMEPIAIITDTEIS	17.5	Natural	KRMEPIAIITDTEIS	
TRESK	ADEAVPQIIASAELP	54	Natural	ADEAVPQIIASAELP	Homo sapiens
rbmAKAP79	AAGAGVGIVITVTEAE	57	Designed (low T)	AAGAG /GIVITVTEAE	Pelecanus crispus
rbmTRESK_2	ADEAIPEITITSAELP	60	Designed (low T)	ADEAIQPITITAELP	AKAP79
rbmAKAP79_2	ADGAGVGIVITVTEAE	69	Designed (low T)	ADGAGVGIVITVTEAE	Prosimis
rbmRIPOR2	ASVSNPEITVTSATE	79	Designed	QSQSNPEITVTP	TRESK
rbmRIPOR2_2	HVSSSPRITITPTQHR	200	Designed (low T)	HVSSSPDITATPTQHR	

7/10 designs
3/4 natural compete with binding of substrate to Cn

Tubiana*, Adriana-Lifshitz

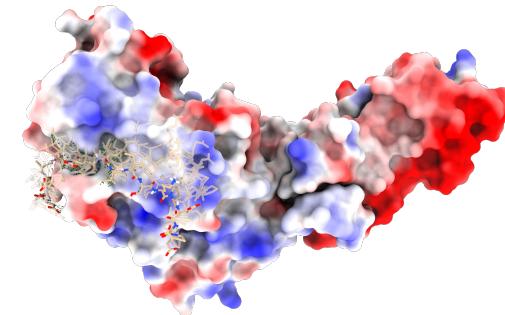


13



Summary and future directions

- Peptides are attractive candidates for PPI inhibitor design, but design is challenging.
- We proposed and validated an integrative design protocol based on a Sequence Generative Model trained from native binders of the target protein.
- The SGM captures key sequence patterns important for binding, and recombines them to generate novel and diverse peptide binders.
- Flexible molecular docking efficiently complements SGMs by differentiating between weak and strong native binders.



Next steps for Calcineurin:

- Cellular assays
- Display experiments to further optimize binding affinity, multivalent constructs
- Pharmacophore-based drug discovery / HTS via competition



Available online at www.sciencedirect.com
ScienceDirect

Current Opinion in
Structural Biology

Machine learning for evolutionary-based and physics-inspired protein design: Current and future synergies
Cyril Malbranque^{1,2}, David Bikard², Simona Cocco¹, Rémi Monasson¹ and Jérôme Tubiana³



Acknowledgements

Tel Aviv University

Haim Wolfson

Mark Rozanov

Naama Hurwitz

Michael Nissan

Yoav Lotem

Maayan Gal

Lucia Lifshits

Daniel Bar

TAU CS system team

Sonia Lichtenzveig Sela



The Hebrew University
of Jerusalem

Dina Schneidman-Duhovny

Merav Breitbart

Lirane Bitton

Matan Halfon

Shon Cohen

Tomer Cohen

Edan Patt

Ecole Normale Supérieure

Rémi Monasson

Simona Cocco

Icahn School of Medicine

Mount Sinai NY

Yi Shi

Yufei Xiang

Zhe Sang

University of Pittsburgh

Kong Chen

Li Fan

