

# COMPUTATIONAL DESIGN OF PROTEINS AND ENZYMES

VAITEA OPUU

*Laboratoire de Biologie Structurale de la Cellule,  
École Polytechnique*

October 29, 2020

## *Jury Members*

Anne-Claude Camproux	Université Paris-Diderot	Rapporteur
Marc Delarue	Institut Pasteur	Rapporteur
Martin Weigt	Sorbonne Université	Examinateur
Sebastian Will	École Polytechnique	Examinateur
Thomas Simonson	École Polytechnique	Directeur de thèse

## WHOLE PROTEIN DESIGN

- Complete redesign of a PDZ domain involved in protein-protein interactions

## ENZYME DESIGN

- Design of MetRS binding site to increase activity for unnatural amino acids

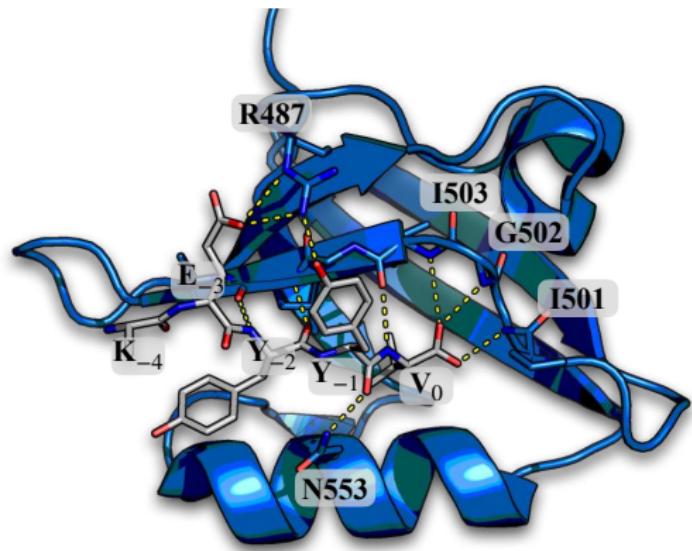
## EXPLORE THE BENEFITS/LIMITS OF THE PHYSICS-BASED APPROACH

- Molecular mechanics model
- Boltzmann sampling

## DESIGN PDZ PAIRS WITH OVERLAPPING GENES

- Algorithm development
- PDZ application

# PDZ DOMAINS



## PDZ STRUCTURE

- ~ 90 Amino acids
- 6  $\beta$  sheets
- 2  $\alpha$  helices ( $\alpha_1$  and  $\alpha_2$ )

## BIOLOGICAL CONTEXT

- Protein-Protein recognition
- Proteins activity modulation

## DESIGN GOAL

Search for sequences that adopt this fold

⇒ Solve inverse folding problem

# MAIN INGREDIENTS

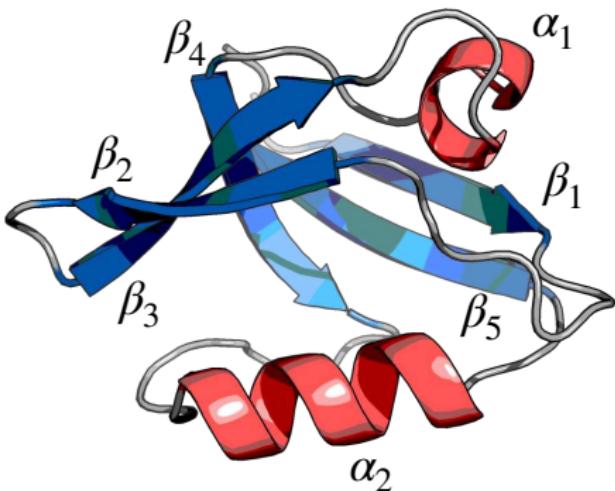
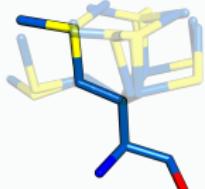
## MOLECULAR MECHANICS + CONTINUUM ELECTROSTATICS

$$E = E^{MM} + E^{GB} + E^{\text{NONPOLAR}}$$

## EMPIRICAL UNFOLDED STATE

$$\sum_{aa} E_{aa}^{uf}(\text{type}_{aa})$$

## DISCRETE CONFORMATIONS



- Experimental geometry for the backbone
- Fixed backbone during the exploration
- Discrete rotamers

# NUMERICAL PROCEDURE

## DESIGN SOFTWARE:

Proteus (Mignon, Druart, Michael, Opuu, Polydorides, Villa *et al*, JPCA, 2021)

## ENERGY FUNCTION:

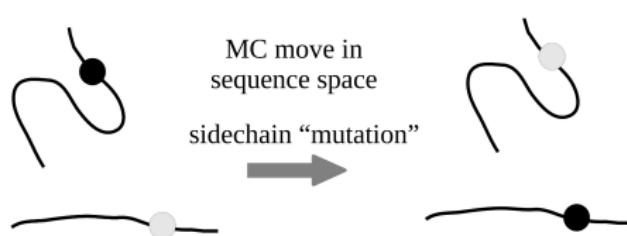
- Amber ff99SB
- pre-calculation in a lookup table
- Unfolded state parametrized empirically

## SEQUENCE SPACE:

- Proline & glycine not allowed to mutate
- 13 positions involved in peptide binding are not allowed to mutate
- All 61 others mutate freely:  $10^{76}$  possible sequences

## SAMPLING METHOD:

- Replica exchange MC
- $10^9$  Monte Carlo steps per replica



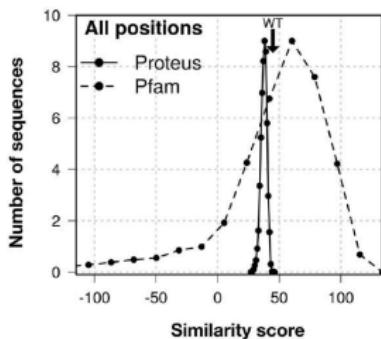
# DESIGNED SEQUENCES RESEMBLE NATURAL SEQUENCES

$\approx 10^5$  DESIGNED SEQUENCES

	10	20	30	40	50	60	70	80
Proteus_0/1-84	R T M I K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_1/1-84	R T M V K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M D G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_2/1-84	L T M I K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_3/1-84	R T M I K T L T K K E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_4/1-84	R T M I K O L K K T E E E P M G I	T L R Q D E E E K I K I D R I M D G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_5/1-84	R T M V K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G L T T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_6/1-84	R T M V K T L T K K E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_7/1-84	R T M I K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_8/1-84	L T M V K T L K K T E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G L M T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						
Proteus_9/1-84	R T M I K T L T K R E E E P M G I	T L R Q D E E E K I K I D R I M E G G E V D R Q G M L Q I G T Q I L O V N G T K T K D L L V E M L Q R L L K D S K G E I Q V L I L P D						

⋮  
⋮

37% MEAN IDENTITY



## CORE POSITIONS: PROTEUS VS PFAM



# SELECTING SEQUENCES FOR EXPERIMENTAL TESTING

Top 2 000 sequences



Proteus energy<1.5 kcal/mol

Solubility

Isoelectric point

1268 sequences

Negative design

Superfamily score  
above average

692 sequences

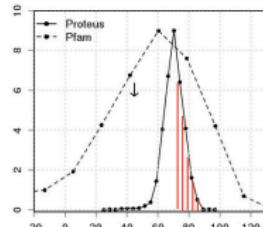
Ad-hoc

Similarity

215 sequences

Net charge, ionic mutations

15 sequences

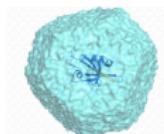


Cluster

6 sequences

Molecular dynamics

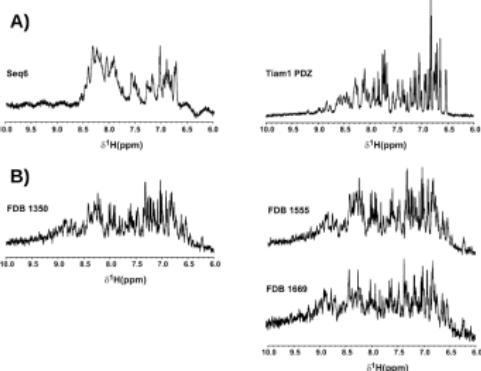
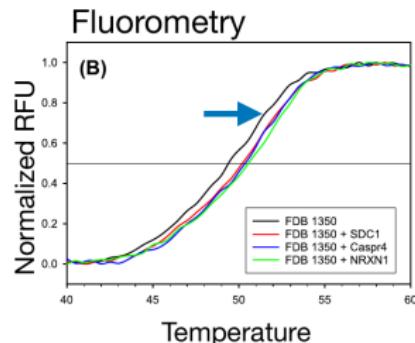
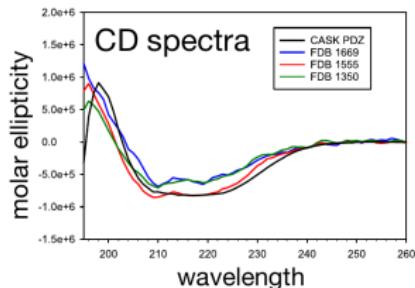
**3 for experiments**



**52~53 mutations out of ~90 positions in tested sequences!**

# EXPERIMENTS CONFIRM ALL 3 DESIGNS

(E. Fuentes, U. Iowa)



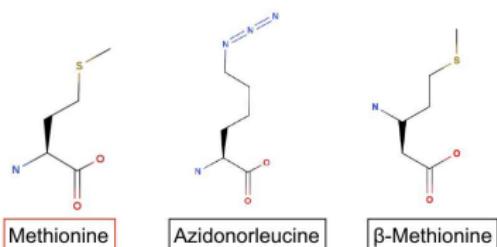
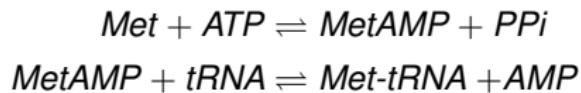
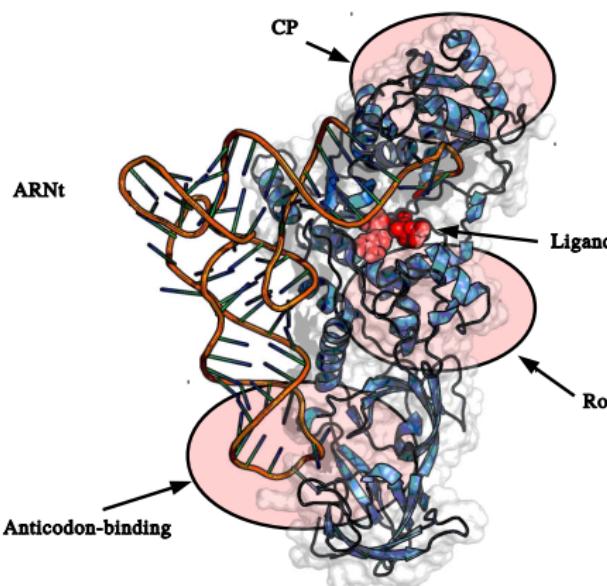
## RESULTS

- PDZ-like secondary structure content
- 1D-NMR spectra typical of folded protein
- Thermal denaturation upshifted by peptide binding

⇒ Designed sequences adopt PDZ fold

**First successful protein redesign with a physics based energy function**  
Opuu *et al* (2020) Sci. Rep.

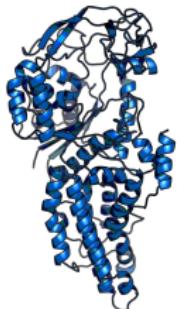
# ENGINEERING METRS FOR LIGAND:SUBSTRATE BINDING AND CATALYTIC POWER



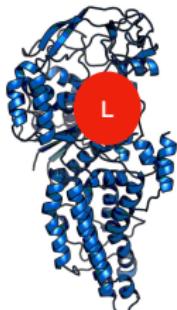
**UAAs**

# DESIGNING FOR AFFINITY OR CATALYSIS IS VERY DIFFICULT

APO



Affinity



## CALCULATION STRATEGIES

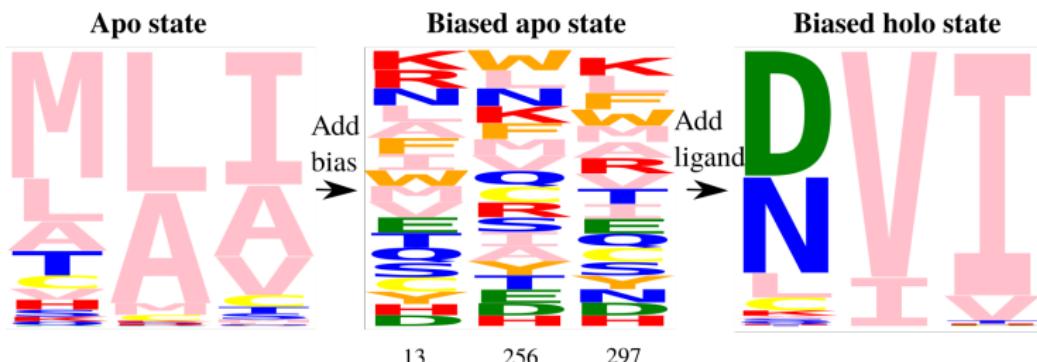
- Large combinatorial space
- Simultaneous optimization directions: stability, affinity, catalysis
- Positive design for the bound state and negative design for the unbound state
- Existing methods are heuristic or very expensive
  - optimize the bound state energy (Rosetta)
  - exhaustive enumeration of states (Osprey)

## PHYSICAL DESCRIPTION

- Transition state is difficult to characterize
- Electrostatics is an important component of the affinity

# A RIGOROUS METHOD TO DESIGN FOR AFFINITY: ADAPTIVE LANDSCAPE FLATTENING

Villa, Panel, Chen, Simonson 2018; Bhattacherjee Walling 2013

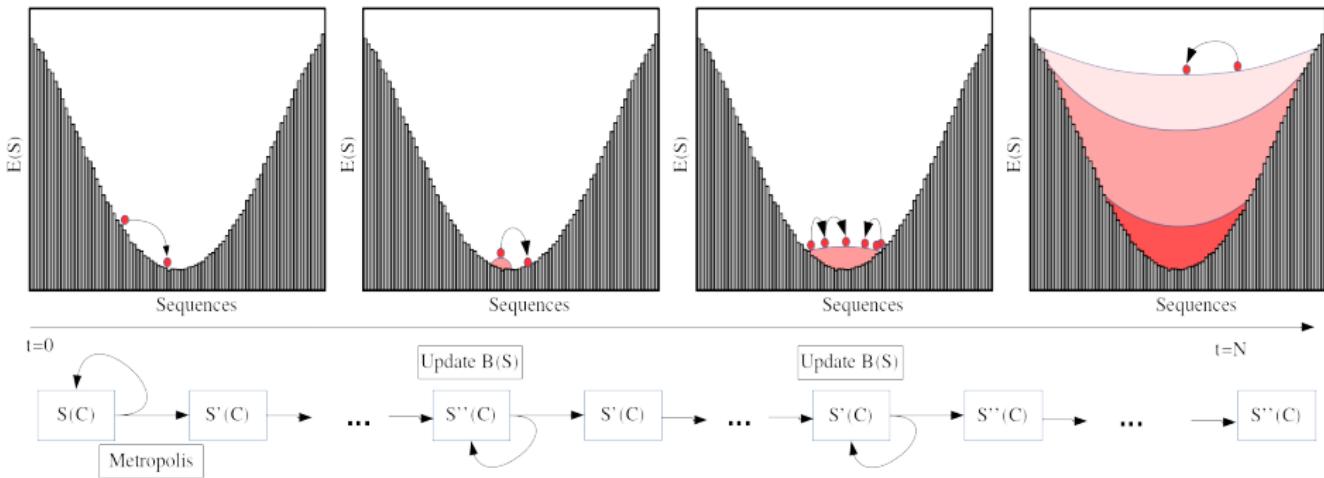


SIMULATE HOLO STATE WITH THE SAME BIAS:

The bias **subtracts out** the apo free energy

Sequences populated according to their binding free energy

# THE ADAPTIVE LANDSCAPE METHOD APPLIED TO SEQUENCES



- Update scheme borrowed from WellTempered Metadynamic

$$\text{inc}(B(S_i(t); t)) = e_0 \times \exp\left(-\frac{E_i^B(S_i(t); t)}{E_0}\right)$$

$$\text{inc}(B(S_i(t), S_j(t); t)) = e_0 \times \exp\left(-\frac{E_i^B(S_i(t), S_j(t); t)}{E_0}\right)$$

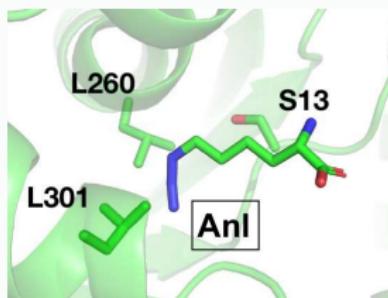
# FIRST TEST: REDESIGN METRS FOR AZIDONORLEUCINE BINDING

## AZIDONORLEUCINE (AnL)



- Cell localization
- Cell identification
- Protein labeling

## EXPERIMENTAL DATA



(Tanirkulu *et al*, 2009)

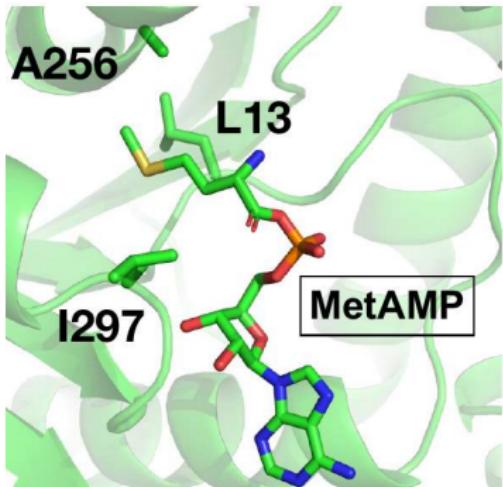
- 21 experimental variants
- 3 variable positions
- 2744 possible variants considered

## Main results:

5 of 6 most active variants are among the top 100 predictions

2 most active variants are among the top 10 specificity predictions

## 2<sup>ND</sup> TEST: REDESIGN METRS FOR MET BINDING



### SYSTEM DATA

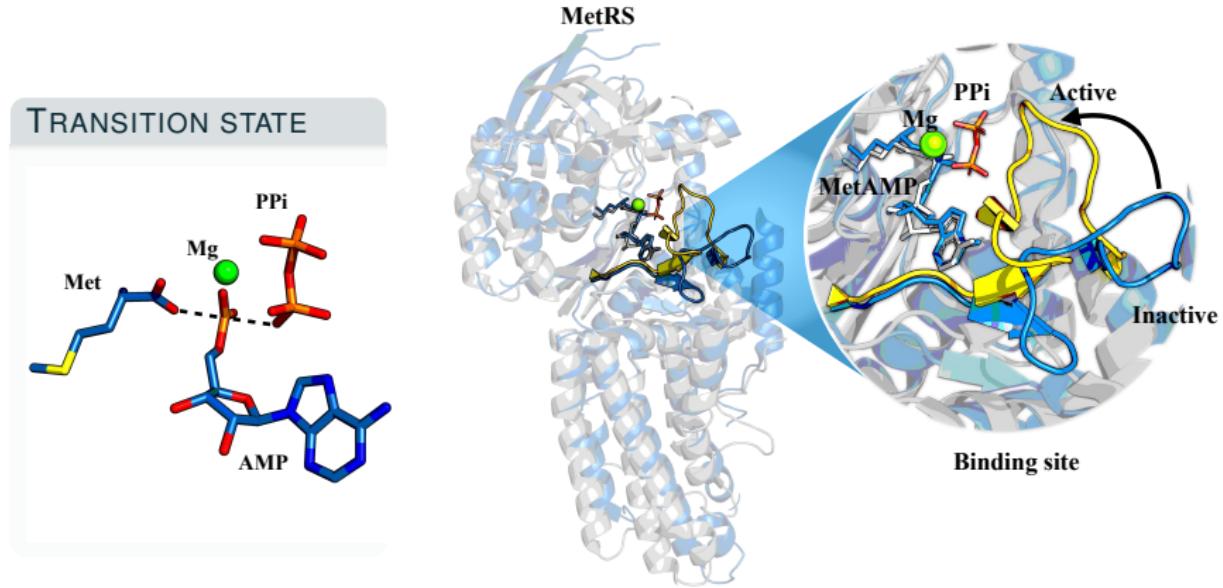
- 3 variable positions
- KMSKS loop in its "inactive" conformation

- 528 variants sampled during simulations out of 5832 allowed
- 20 experimental variants tested, based on the variants predicted
- Among top 40 predictions, 5 are active
- Computed affinities in good agreement with experiment:

0.9 kcal/mol rms error, 0.75 correlation

Opuu *et al*, Plos Comp. Bio. 2020

# ADAPTIVE LANDSCAPE FLATTENING ALLOWS US TO DESIGN FOR CATALYSIS



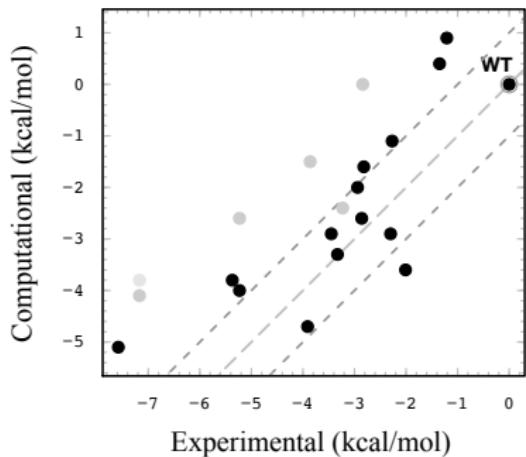
- Mg and ATP are **rigid**
- "APO" state: MetRS:ATP
- HOLO state: MetRS:[Met:ATP]<sup>‡</sup>

**Variants can now be selected by their catalytic efficiency  $k_{cat}/K_M$**

## DESIGN FOR CATALYTIC EFFICIENCY

Catalytic efficiency =  $k_{cat}/K_M$

$$\Delta G^\ddagger \stackrel{def}{=} kT \times \ln \left( \frac{k_{cat}}{K_M} \right)$$

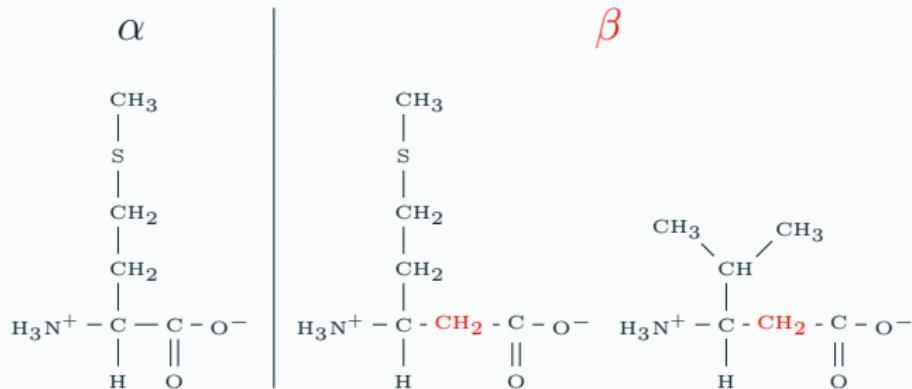


- 20 experimental values (above)
- 13 among top predictions
- Good agreement for  $k_{cat}/K_M$
- 0.8 correlation
- 1.1 kcal/mol rms error

**Structural models + physical description capture the important effects**

# DESIGN METRS FOR $\beta$ AMINO ACID ACTIVITY

## $\beta$ AMINO ACIDS



## $\beta$ POLYPEPTIDE

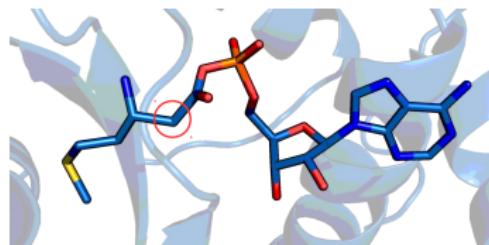
- Allow new backbone geometries
- Reduce peptide degradation

## CHALLENGES

- Structural models
- Force field parameters
- Which positions to mutate

# 1<sup>ST</sup> ROUND: DESIGN 3 POSITIONS FOR BETA-METAMP AFFINITY

## $\beta$ -METAMP



## SYSTEM DATA

- ligand:  $\beta$ -MetAMP
- 3 mutable positions: 13, 256, 297
- APO state: MetRS
- HOLO state: MetRS: $\beta$ -MetAMP

- 206 variants sampled by simulations
- **20** variants experimentally tested, 11 from MC (blind tests)
- Fair agreement for  $k_{cat}/K_M$
- **5 out of 11 have a weak but measurable activity**



**3 variants have decreased  $\alpha$ -Met selectivity (factors 2-8)**



**Activity lower than WT  $\beta$ -Met activity**

## 2<sup>ND</sup> ROUND: DESIGN 3 POSITIONS FOR CATALYTIC EFFICIENCY

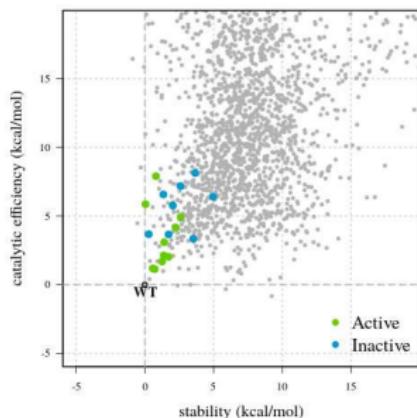
### SYSTEM DATA

- Same positions can mutate
- Apo state: MetRS:ATP
- Holo state: MetRS: $\beta$ -Met $^{\ddagger}$

$$\Delta G^{\ddagger}(\beta\text{-Met})$$

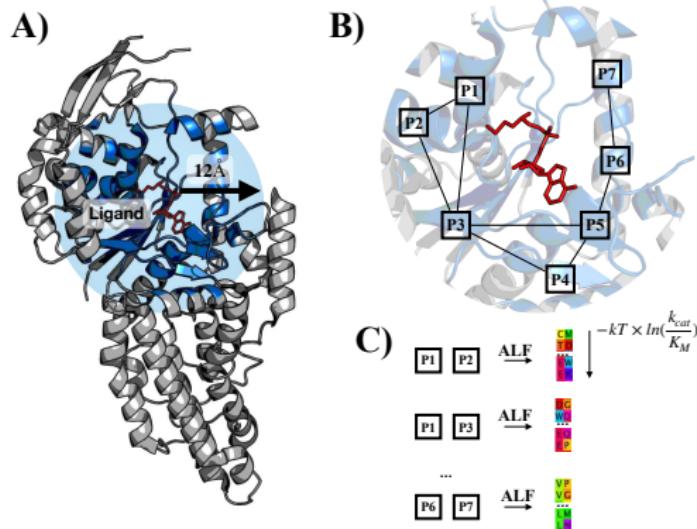
variants	predicted	experiments
CAC	4.5	2.2
CAI	1.0	1.6
CAV	2.4	3.6
LAC	3.2	2.0
LAI	0.0	0.0
LAT	2.3	3.1
LAV	1.4	3.1
MAC	8.1	2.5
MAV	6.1	3.1
SAC	5.3	2.9
SAI	2.0	1.9

- Fair agreement for  $k_{cat}/K_M$
- Still no variants with improved activity

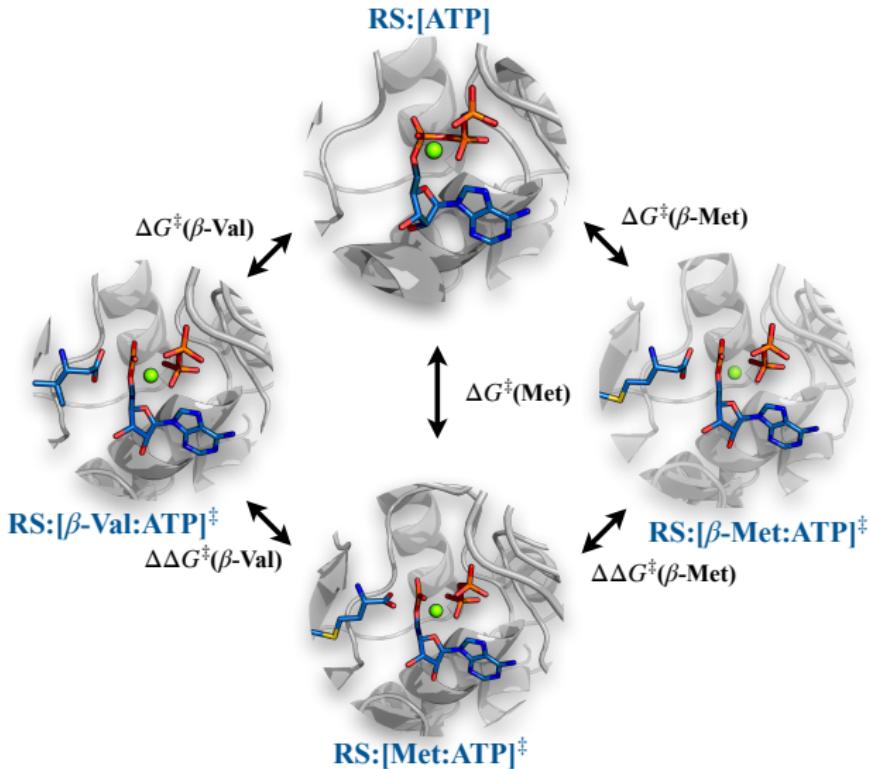


### 3<sup>RD</sup> ROUND: SCAN THE WHOLE ACTIVE SITE

- Select all the positions in the active site (19 positions)
- Design pairs of positions
- Compute a "pair-activity" score
- Select positions with high scores, to form a few quartets
- Redesign quartets



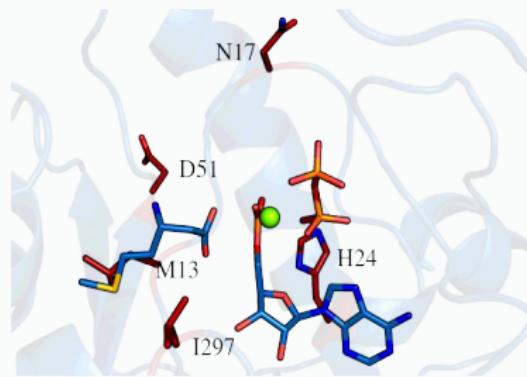
## SYSTEMS CONSIDERED



**Design four positions at once**

## QUARTET SELECTIONS

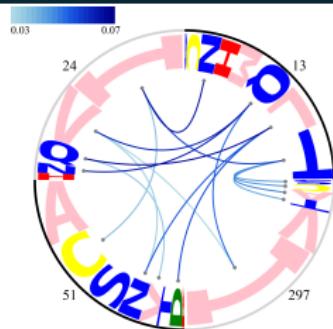
- Positions are in the first and second layer of the binding site
- 3 quartets of 4 positions considered



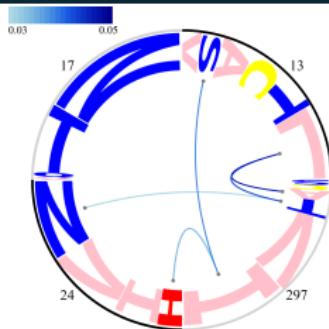
## SELECTED VARIANTS

- 2970 + 3168 + 720 variants sampled by MC
- 89 + 81 + 37 variants satisfy various thresholds (stability, catalytic efficiency and selectivity)

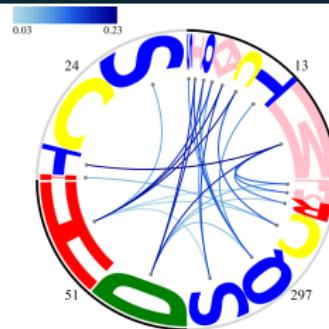
## $\beta$ -MET DESIGNS



Quartet 1



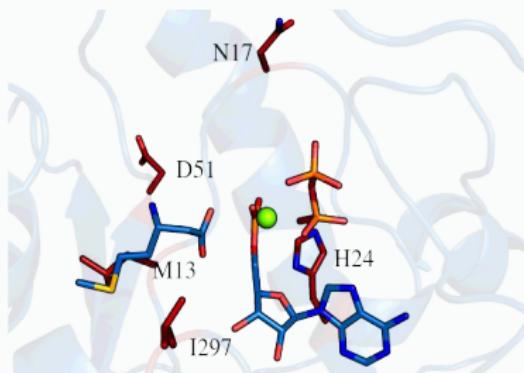
Quartet 2



Quartet 3

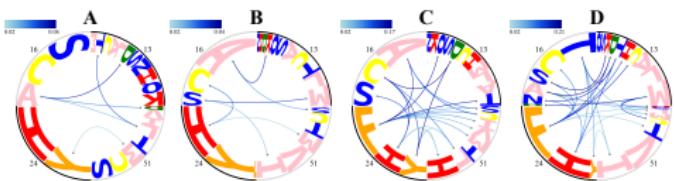
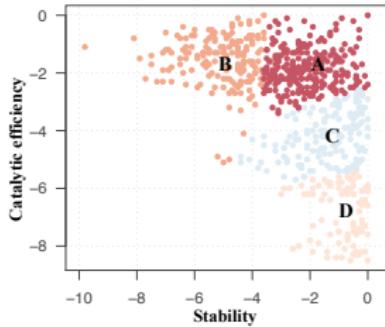
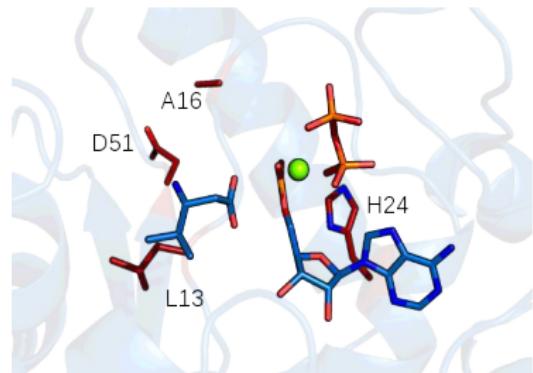
## DESIGNS

- 2 positions are proposed to be mutated into smaller and more polar types
- ASP 51 is proposed to be mutated into HIS
- 2 positions are maintained in their native types



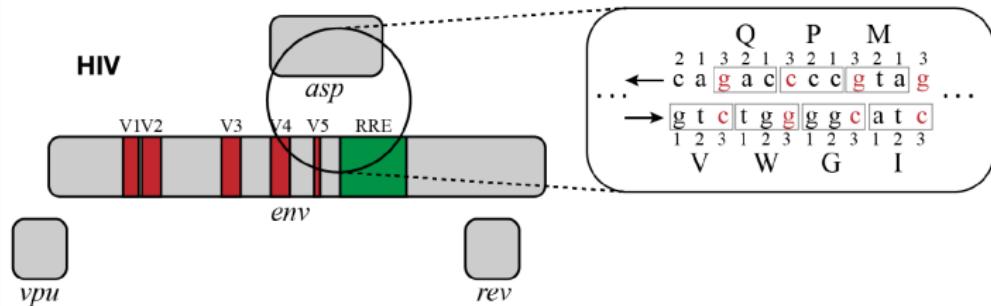
## SELECTION OF POSITIONS

- 1 quartet of 4 positions considered
- 661 variants predicted
- H24F is predicted as a good choice for  $\beta$ -Val (part of HIGH)



# DESIGN OF OVERLAPPING GENES

HIV (CASSAN *et al* 2016)



## VERY RARE OUTSIDE VIRUSES

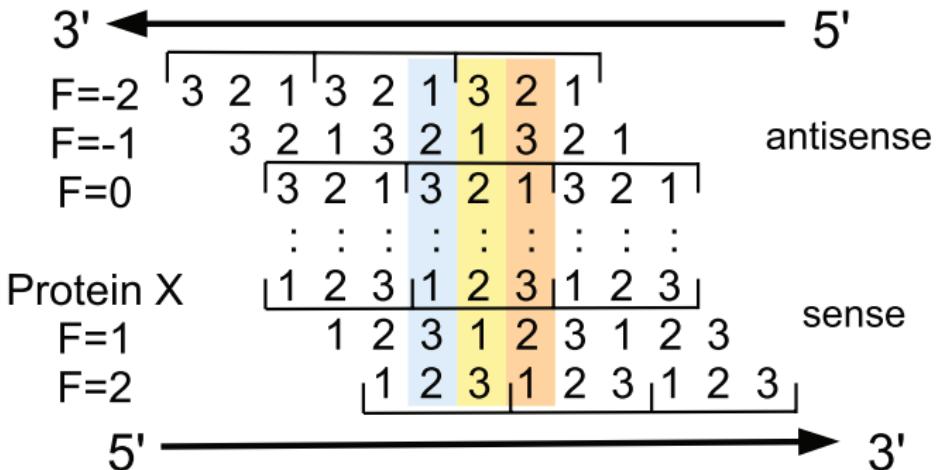
- Several in SARS-COV 2
- Cancer related genes: INK4a and ARF
- Autophagy related genes

...

## WHY DESIGN OVERLAPPING GENES

- Limit genetic drift
- bio-confinement of GMOs
- Smaller genomes
- Genetic regulations
- Viral evolution

## DESIGN OVERLAPPING GENES: THE PROBLEM STATEMENT



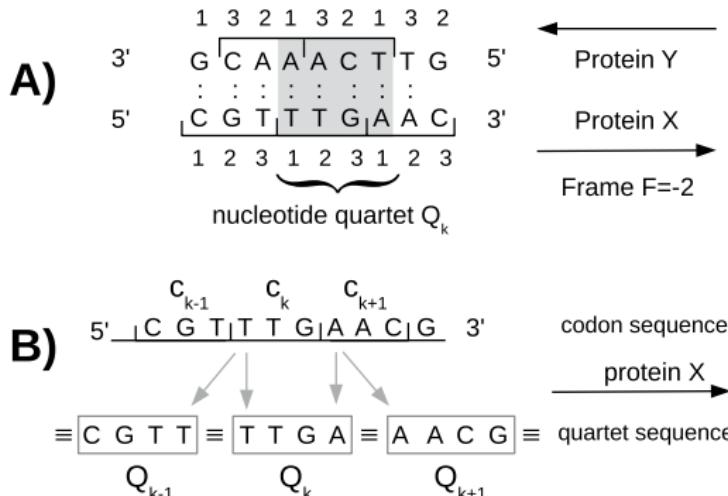
Codon overlap + base pairing  $\implies$  Strong constraints on sequences

For arbitrary pair of proteins, **overlapping coding rarely exists!**

**Based on two arbitrary sequences, find the closest homologs that have an overlapping coding scheme.**

## EXACT SOLUTION BASED ON LINKED QUARTETS

Opuu *et al*, Sci. Rep, 2017



### QUARTET DEFINITION

- Redefine the DNA sequence as linked list of quartets
- Similarity based score for each quartet

Reduction to the problem of finding a best path by **dynamic programming**

# PDZ APPLICATION

## DESIGNS

- 5 PDZ with known structures
- All 5 phases, different offsets
- 1715 designs of overlapping PDZ pairs

## CHARACTERIZATION

- Similarities to Pfam PDZ database
- Superfamily recognition
- Pi, Charge
- Disorder measures (Iupred)
- Cavity detection (McVol)

## VALIDATION

- Molecular dynamics in explicit solvent 0.5-3  $\mu$ s
- Experimental tests (underway)

X . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

X . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

X . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

X . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
?

X R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

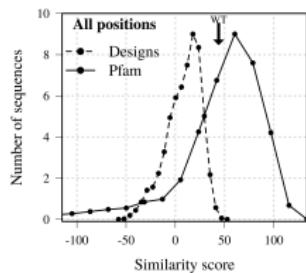
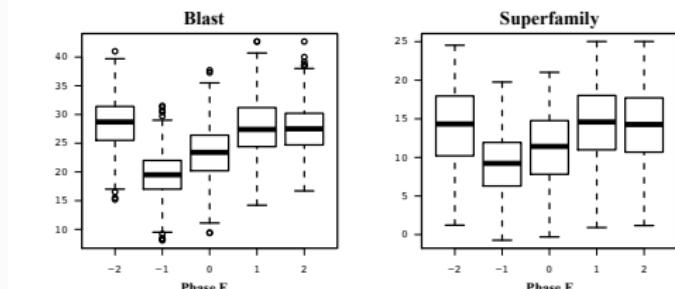
X R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

X R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
Y . . . R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V A R I M H G G M I H R . . .  
→

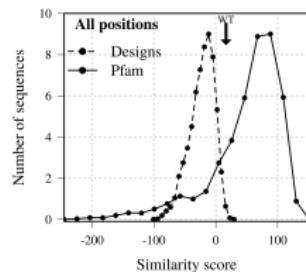
# SEQUENCE CHARACTERIZATION

## RESULTS

- Over-represented types:  
LEU, ARG, SER
- Under-represented types:  
ALA, GLY, VAL, ASP
- Phase F = -2 is the best
- Phase F = {-1, 0} are the poorest



Cask

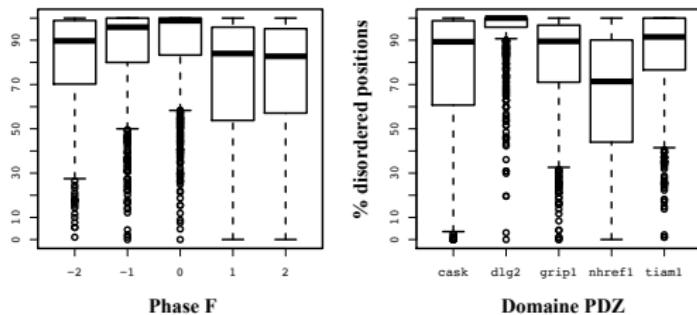


Nhrf1

# STRUCTURAL CHARACTERIZATION

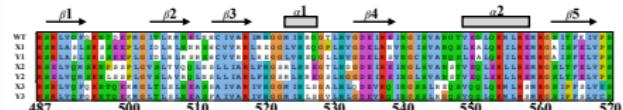
## CAVITY AND DISORDER PROPENSITIES

- 3D structure reconstructed with the experimental template structures
  - 715/1715 (41.7 %) pairs have no cavities
- Disorder propensities are computed using the sequence only
  - With the threshold of 80% of non-disorder positions: 65% pairs are non-disordered



## MD VALIDATION

### CHOSEN SEQUENCES

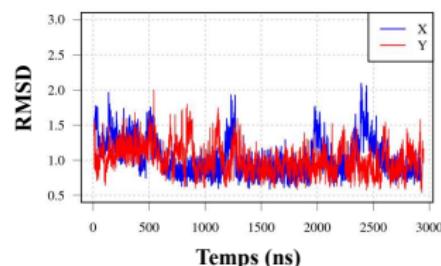
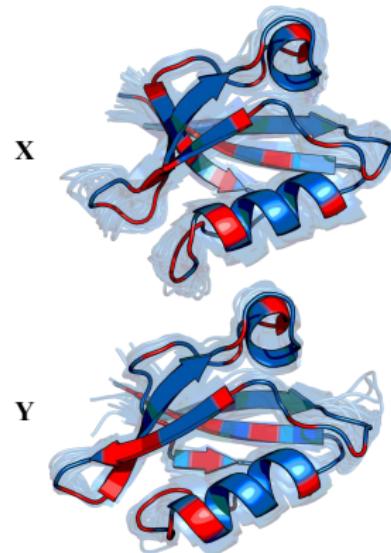


- 3 pairs of Cask sequences selected
- Phase is F=-2

### MD STABILITY ASSESSMENT

- 1 pair is stable for 3  $\mu$  s (27-28 mutations)
- 1 pair stable for 0.5  $\mu$  s (34-36 mutations)
- 1 pair showed some instabilities (31-32 mutations)

Experiments underway (G. Travé)



## OVERALL CONCLUSIONS

### COMPLETE REDESIGN OF A PDZ DOMAIN

- 1st successful redesign of a PDZ domain with a physics-based approach

### ENZYME REDESIGN FOR AFFINITY AND CATALYTIC EFFICIENCY

- Affinity based designs can discover WT ligand binders
- First example of selection explicitly for TS binding

### ENZYME REDESIGN FOR THE ACTIVATION OF UAAs

- Discovery of active variants for  $\beta$ -Met
- 3 variants with decreased  $\alpha$ -Met selectivity
- Systematic scan of active site (experiments underway)

### OVERLAPPING CODING DESIGNS

- Sequence based design algorithm combined with post-filtering
- 1 pair of overlapping PDZ designs very stable in MD

## COMPUTATIONAL GROUP, EXPERIMENTAL GROUPS

### COMPUTATIONAL GROUP

Thomas Gaillard      Francesco Villa  
David Mignon      Alexandrine Daniel  
Nicolas Panel      Thomas Simonson

### METRS COLLABORATORS

Giuliano Nigro      Yves Mechulam  
Christine Lazennec-Schurdevin      Emmanuelle Schmitt

### PHYSICS-BASED PDZ REDESIGN

Young Joo Sun  
Titus Hou  
Ernesto J. Fuentes

### OVERLAPPING PDZ REDESIGN

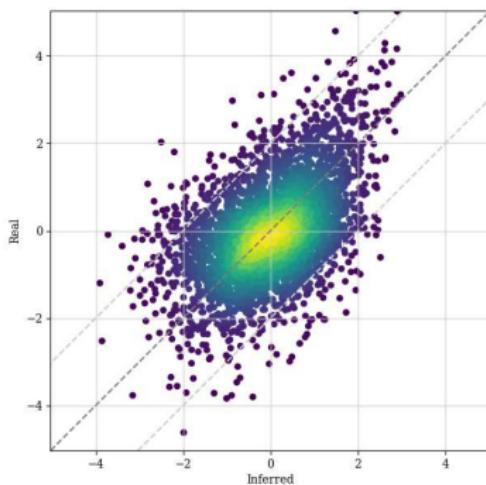
Gilles Travé

## MSA BOOTSTRAP + METROPOLIS TEST

- Sequences are modelled with the Bias only:  
 $E(S) = \sum_i B(S_i) + \sum_{i < j} B(S_{ij})$
- Sequences are sampled randomly from the MSA

## TOY SYSTEM TESTINGS

- 5 positions:  $21^5 (10^{6.6})$  sequences
- $5 \times 21 + 10 \times 441$  (4515) parameters
- Correlation = 0.5, mean error  $\sim 2$  kcal/mol



# DESIGN OF 6 OVERLAPPING PDZ SEQUENCES

3' ←

X S P V I K F T I S G R M E R L M K Q L Q E V T Q N A V S I  
X' S P S I D F I C S V R S S G K P N E A L E F S N S C T N Y  
Y S P V I K F T I S G R M E R L M K Q L Q E V T Q N A V S I  
Y' L P L Y R F P Y L F G E F E G Q P K R C A R V L Q F M H Q I  
Z S P V I K F T I S G R M E R L M K Q L Q E V T Q N A V S I  
Z' P P S I S F V L F G R V G R P T K Q L S S R T P V H T T D  
T T C C C C T T A T A G T T T G T T T T G G G A C T T G A G G G G A C C C C A A A A G A C C T G A G C T T G T C A A C C T T G T A C A C A C A T A  
A G G G G G G A G A T A C G A A A T A C A G A A M A C U C T G A A C T C C C C T T G G G G T T T T G C A A G C T G A A C A G T T G G A A C A T G T G T G T A T  
U' R G G D I E N T R N P R T P L G V F C K L E R V G T C V V  
U R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V  
V' E G E I S K I Q E T L E L P L G F S A S S N E L E H V L Y  
V R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V  
W' R G R Y V Q R K Y K K P S N S P W G F L Q A R T S W N M C C I  
W R S R L V Q F Q K N T D E P M G I T L K M N E L N H C I V  
5' →

X G N I E R I E D G V H L T G Q R H I M G G H M I R A V I C  
X' G G I D R P P D G V R L P G R H A P L G G W L R P I C C  
Y G N I E R I E D G V H L T G Q R H I M G G H M I R A V I C  
Y' G R Y R A P P R G C A L S G Q R T P S S G G L S Y P P Y L V  
Z G N I E R I E D G V H L T G Q R H I M G G H M I R A V I C  
Z' G S I E R P T E W V N C P V G Q T H T S V G G S P V V C  
A G G G G T A T A G A G G G C C C C A G G G G T G C T T G G G G A C C C C T G G G G G G T C T T A T G C C C C T A T G T G T G  
T C C C C C G A T A T C G C G G G G G G T C T C C A C A G G A C C C C T T G G G G A C C C C C C A G A G G A A T A C G G G G G A T A C A A C  
U' S P D I S R G V S H T Q G T P C V W E T P P E N T G D T T  
U A R I M H G G M I H R Q G T L H V G D E I R E I N G I S V  
V' P P I S R G G G S P T R K G P L A C G R P P Q R I R G I Q H  
V A R I M H G G M I H R Q G T L H V G D E I R E I N G I S V  
W' P R Y T L A G G L P H A R D D P L R V G D P P R E Y G G Y N T  
W A R I M H G G M I H R Q G T L H V G D E I R E I N G I S V

X H N L E N M K L T I G M P E D T N K Q F Q V L R S R  
X' M N W S T R A Q L F G W P S T S P K K Y K R Y R G R  
Y H N L E N M K L T I G M P E D T N K Q F Q V L R S R  
Y' H E L E N S S A S F G L P F D L T E Q V K S I E G E  
Z H N L E N M K L T I G M P E D T N K Q F Q V L R S R  
Z' T G V R E L K C F V G L P L R P N R T S E I D G G R  
G T A C A G G T T G A G C A A G T C G A A G C T C T T T G G G T T C C C C T C A G C T C C C A A A G A C A T G A A G C T A T A G A G G G G A G C  
C A T G T C C A A C T G T T G A G G T C G A G A A A C C C A A G G G G A A G T C G A G G G T T T T G T A C T T T G A T A T C C C C C T C G  
U' H V P T R S S L Q K T P R G S R G F L V L S I S P P  
U A N Q T V E Q L Q K M L R E M R G S I T T F K I V P S  
V' M F Q L V R A C R K P Q G E V E G F L Y F R Y L P L  
V' A N Q T V E Q L Q K M L R E M R G S I T T F K I V P S  
W' C S N S F E L A E N P K G K S R V S C T F D I S P S  
W A N Q T V E Q L Q K M L R E M R G S I T T F K I V P S

- All 6 designs are based on Cask PDZ
- DNA sequence formulated as linked quintuplets
- 29.5% of identity
- 103 similarity scores