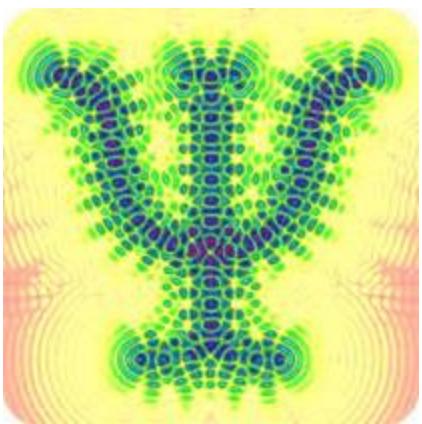


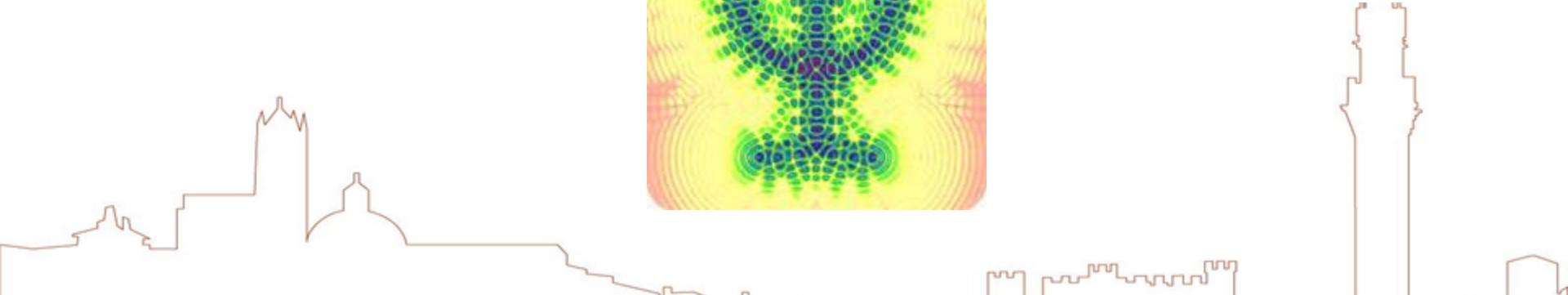
Understanding color tuning rules in the Arch family with a fully computational QM/MM approach



RESMOL GROUP



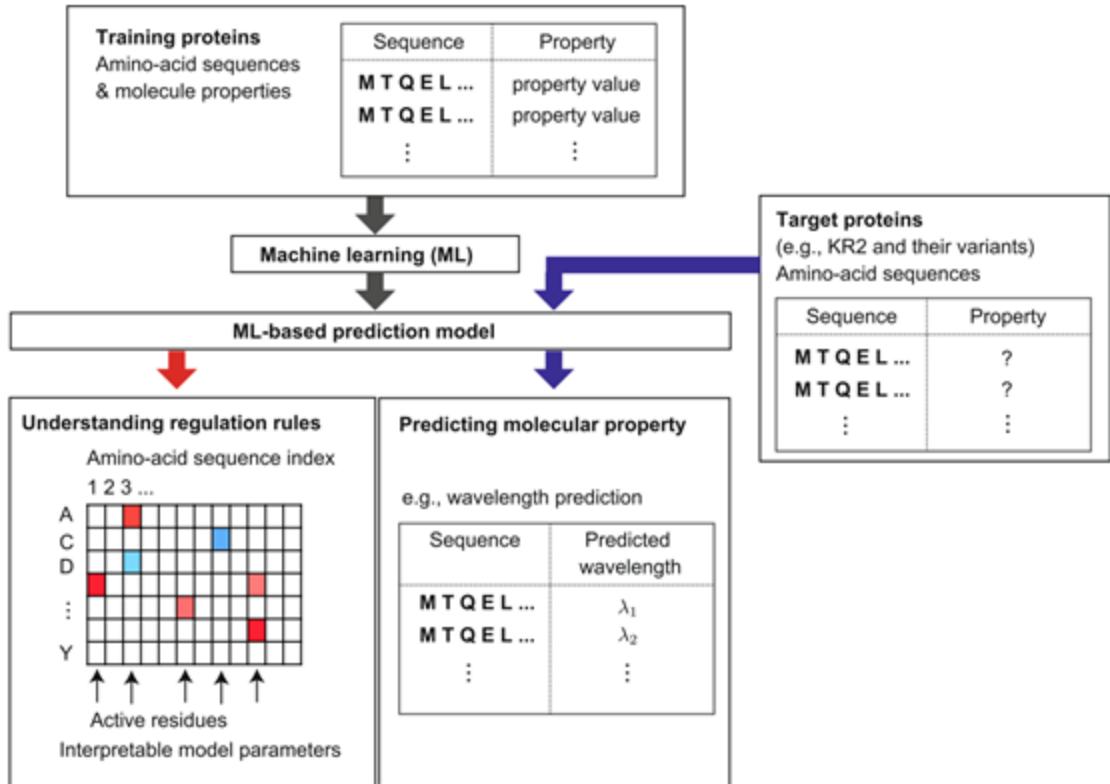
UNIVERSITÀ
DI SIENA 1240



- Data-driven determination of Arch3 variants excitation energies from gene sequences
- Self-consistent computational approach
- Understanding color tuning determinants in the Arch family
- Identifying Arch red-shifted candidates for optogenetics

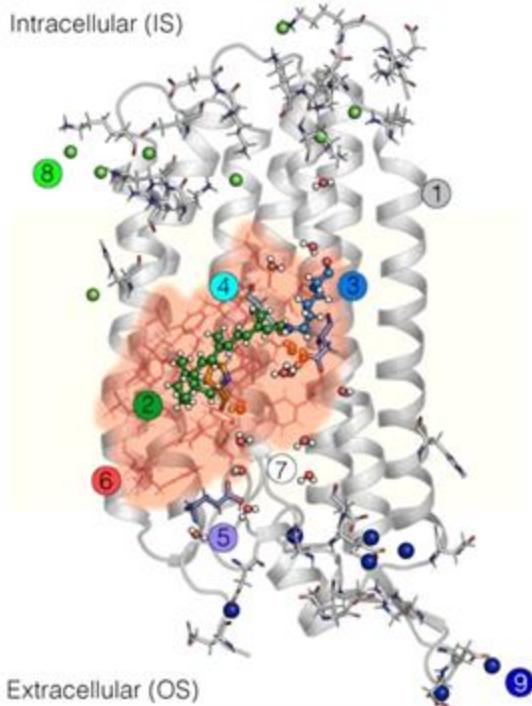
HPC project

- Data-driven determination of excitation energies from gene se

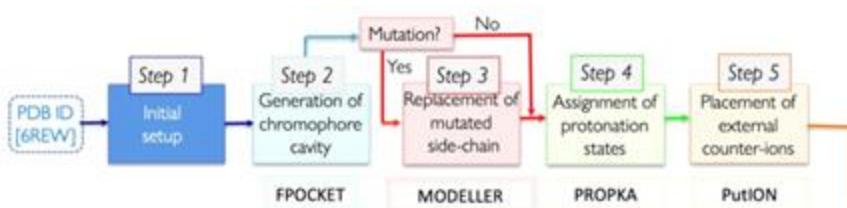


- Self-consistent computational approach

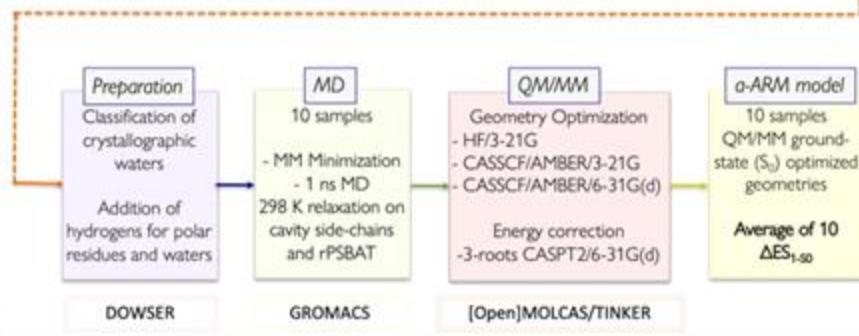
a α -ARM QM/MM Model



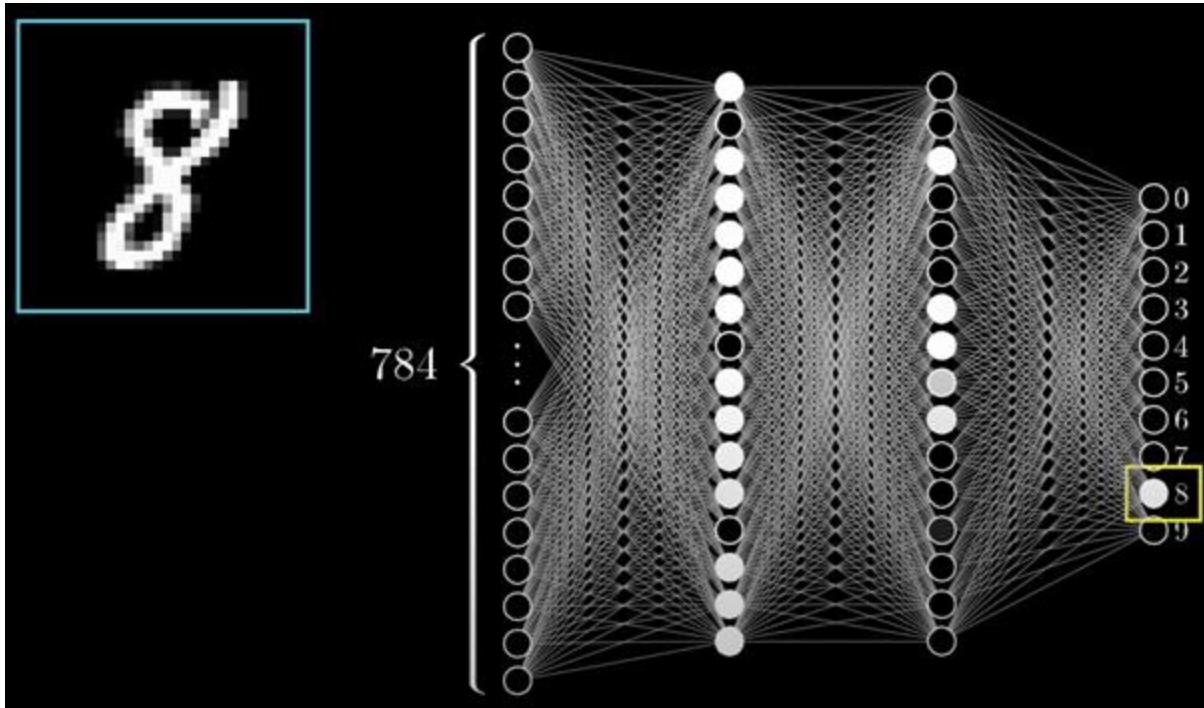
b Phase I: α -ARM Input File Generator (~ 5 min without user manipulation)



c Phase II: α -ARM QM/MM Model Generator (~ 24 h without user manipulation)

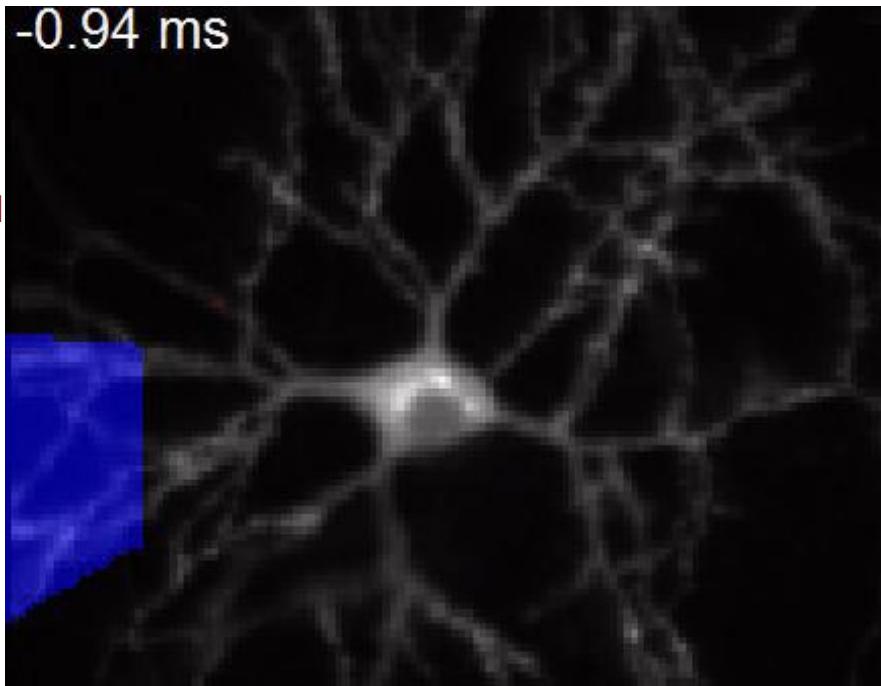
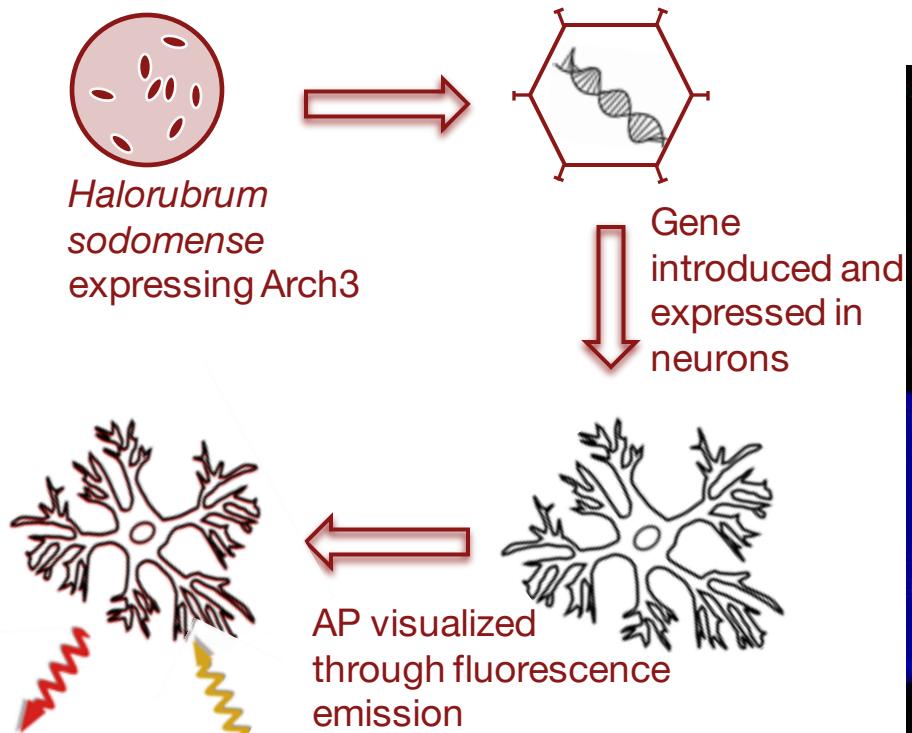


- Understanding of color tuning determinants in the Arch family

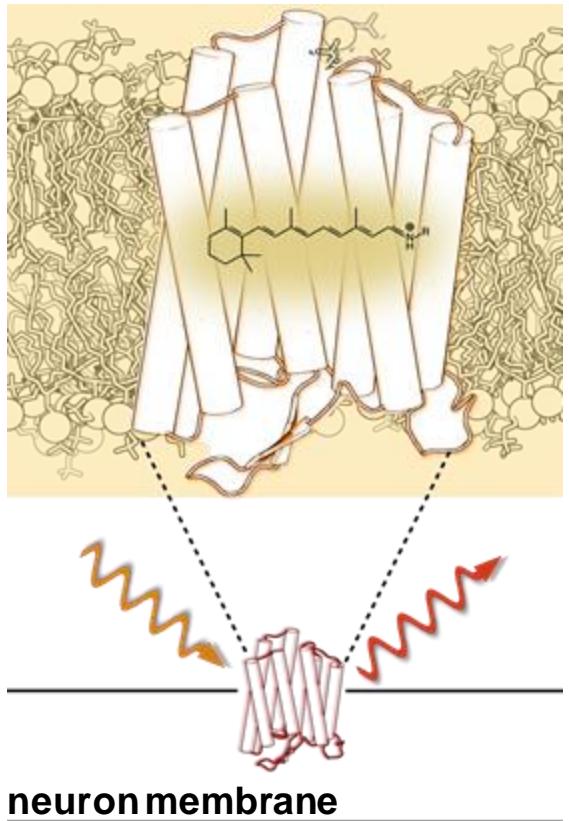


HPC project

- Identifying Arch3 red-shifted candidate for optogenetics



Hochbaum, D et al. *Nature methods* 11, 825, 2014.



Sensors

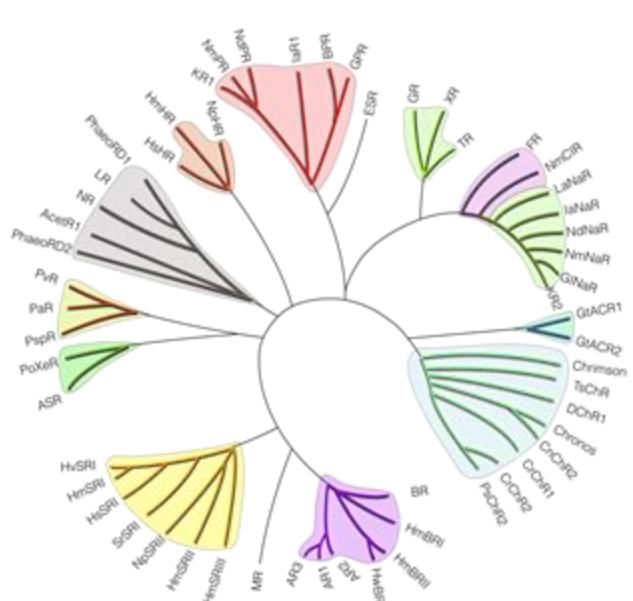
- Spectral orthogonality with actuators
- Deep tissue penetration
- Low phototoxicity
- Long fluorescence lifetime
- High voltage sensitivity

Rhodopsin family

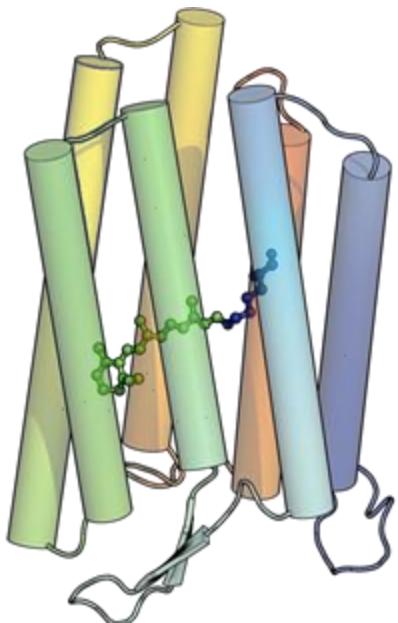


UNIVERSITÀ
DI SIENA
1240

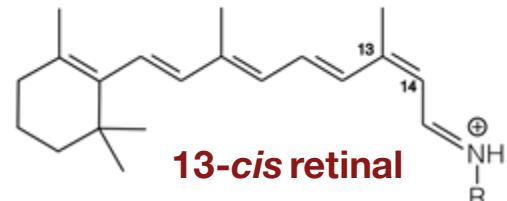
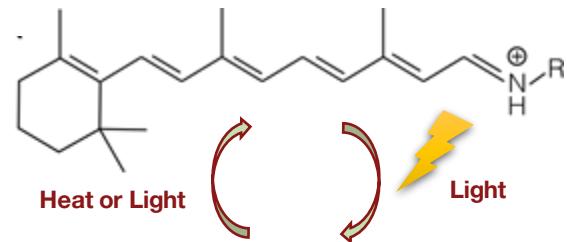
- Huge phylogenetic variability



- Opsin family
 - 7 transmembrane helices
 - Retinal chromophore



- Light-induced photoisomerization



Karasuvama, M et al. *Scientific reports* 8:15580, 2018.

Grdinaru, V. (CalTech) McIsaac, R. S et al. Directed evolution of a far-red fluorescent rhodopsin. *Proceedings of the National Academy of Sciences* **2014**, *111*, 13034-13039.

Arnold, F. H. (CalTech) McIsaac, R. S. et al. Recent advances in engineering microbial rhodopsins for optogenetics. *Curr. Opin. Struct. Biol.* **2015**, *33*, 8-15.

Cohen, A. E. (Harvard) Hochbaum, D. R. et al. All-optical electrophysiology in mammalian neurons using engineered microbial rhodopsins. *Nat. Methods* **2014**, *11*, 825-825.

Boyden, E. S. (MIT) Piatkevich, K. D. et al. A robotic multidimensional directed evolution approach applied to fluorescent voltage reporters. *Nat. Chem. Biol.* **2018**, *14*, 352-360.

Protein	$\Delta E_{S1-S0}^{a,Exp}$ (kcal mol ⁻¹)	$\lambda_{max}^{a,Exp}$ (nm)	$\Delta E_{S1-S0}^{f,Exp}$ (kcal mol ⁻¹)	$\lambda_{max}^{f,Exp}$ (nm)	FQY	Ref.
Arch2	51.5	555	NR ^a	NR ^a	NR	22
Arch3	51.4	556	41.6 ^b	687 ^b	$1-9 \times 10^{-4}$	23
QuasAr1	49.3	580	40.0, 38.6	715, 740	$6.5 \times 10^{-3}, 8.0 \times 10^{-3}$	24, 25
Archon2	48.8	586	38.9	735	1.1×10^{-2}	26
QuasAr2	48.5	590	40.0	715	4.0×10^{-3}	24
Arch7	46.4	616	39.3	727	1.2×10^{-2}	23
Arch5	46.0	622	39.1	731	8.7×10^{-3}	23

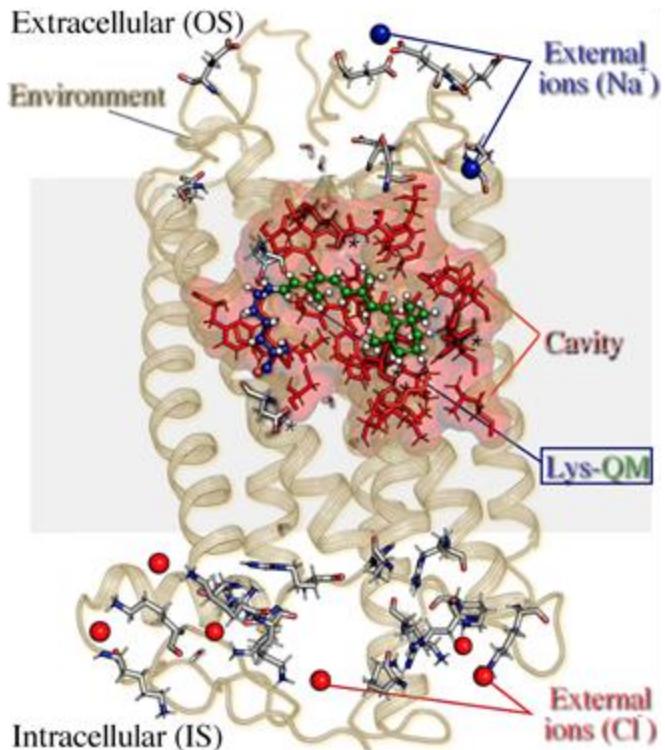
^a Not Reported

^b The λ_{max}^f reported is likely coming from the Q-intermediate rather than from the dark-adapted state (see ref. ²⁷)

The a-ARM protocol

Phase I: Model Structure setup

AUTOMATED setup



Phase II: QM/MM calculation setup

AUTOMATED setup

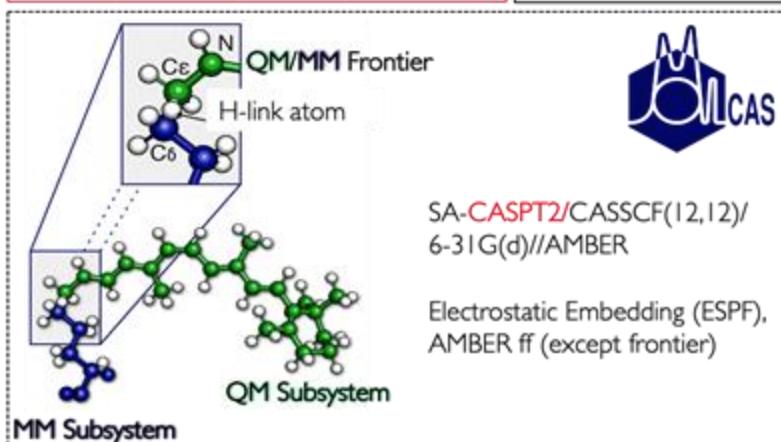
MM Subsystem

Protein environment (*aa*, ions and waters)
Fixed at crystallographic/homology geometry

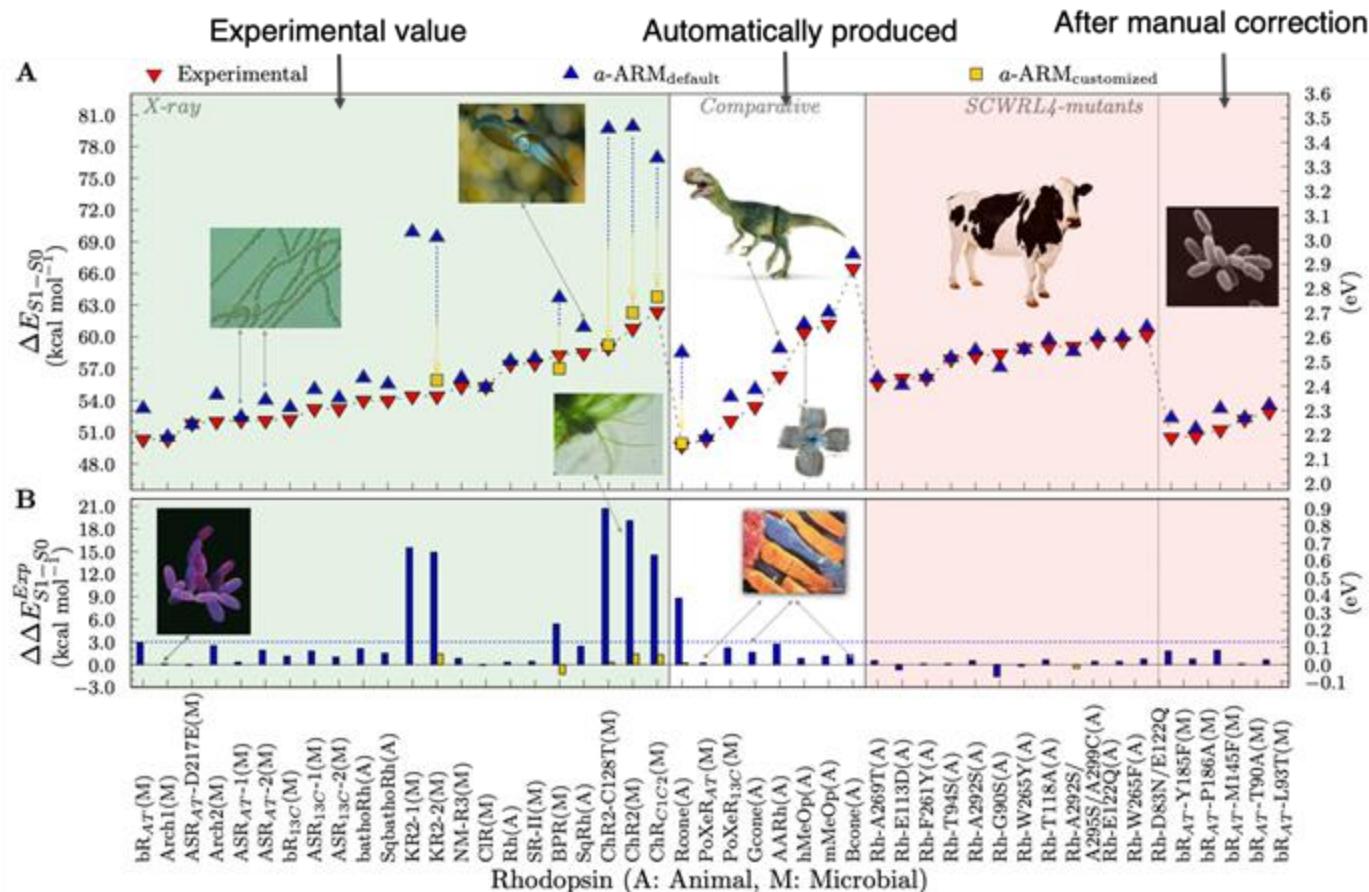


MM Subsystem

Chromophore cavity (*aa* and waters)
MM optimization(AMBER ff)



The a-ARM protocol

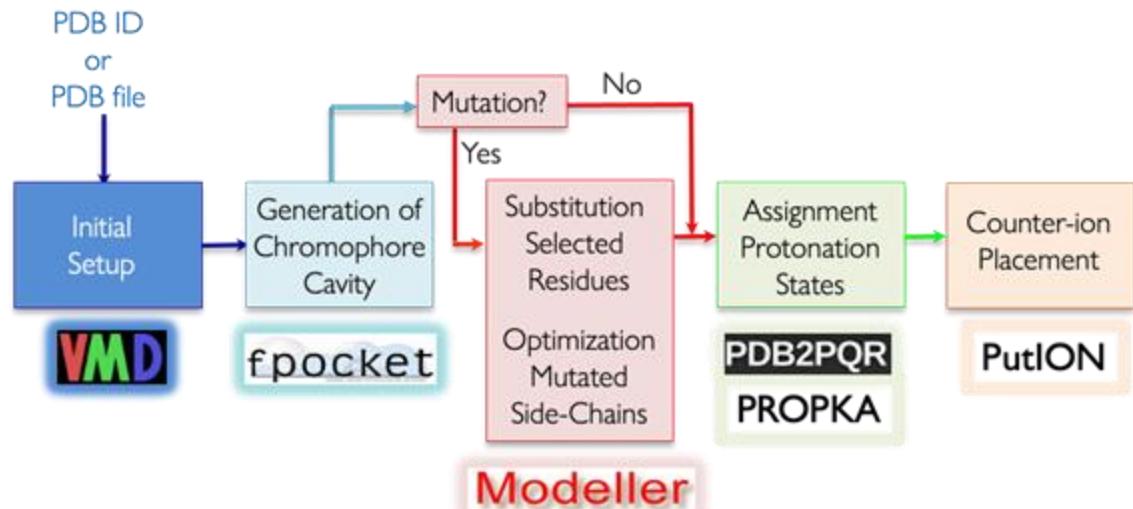
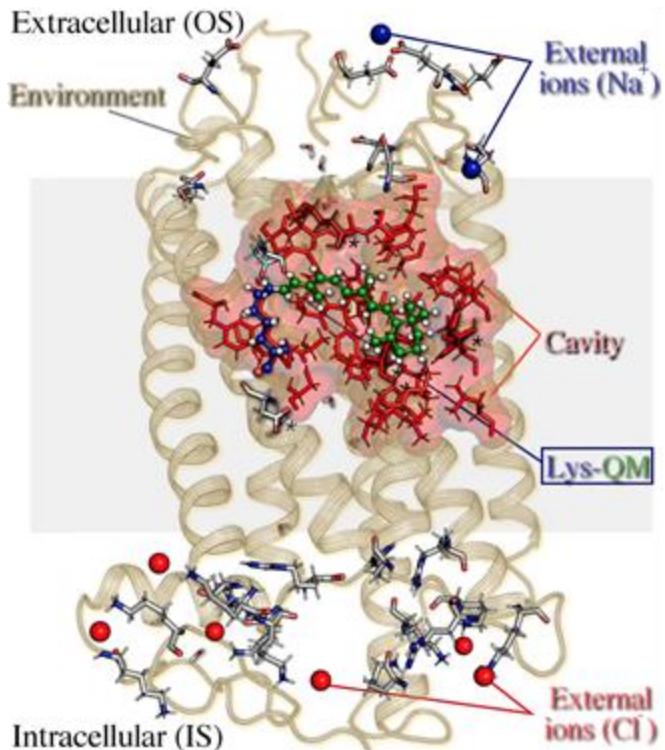


**MAE for a-ARM
is 0.9 kcal/mol**

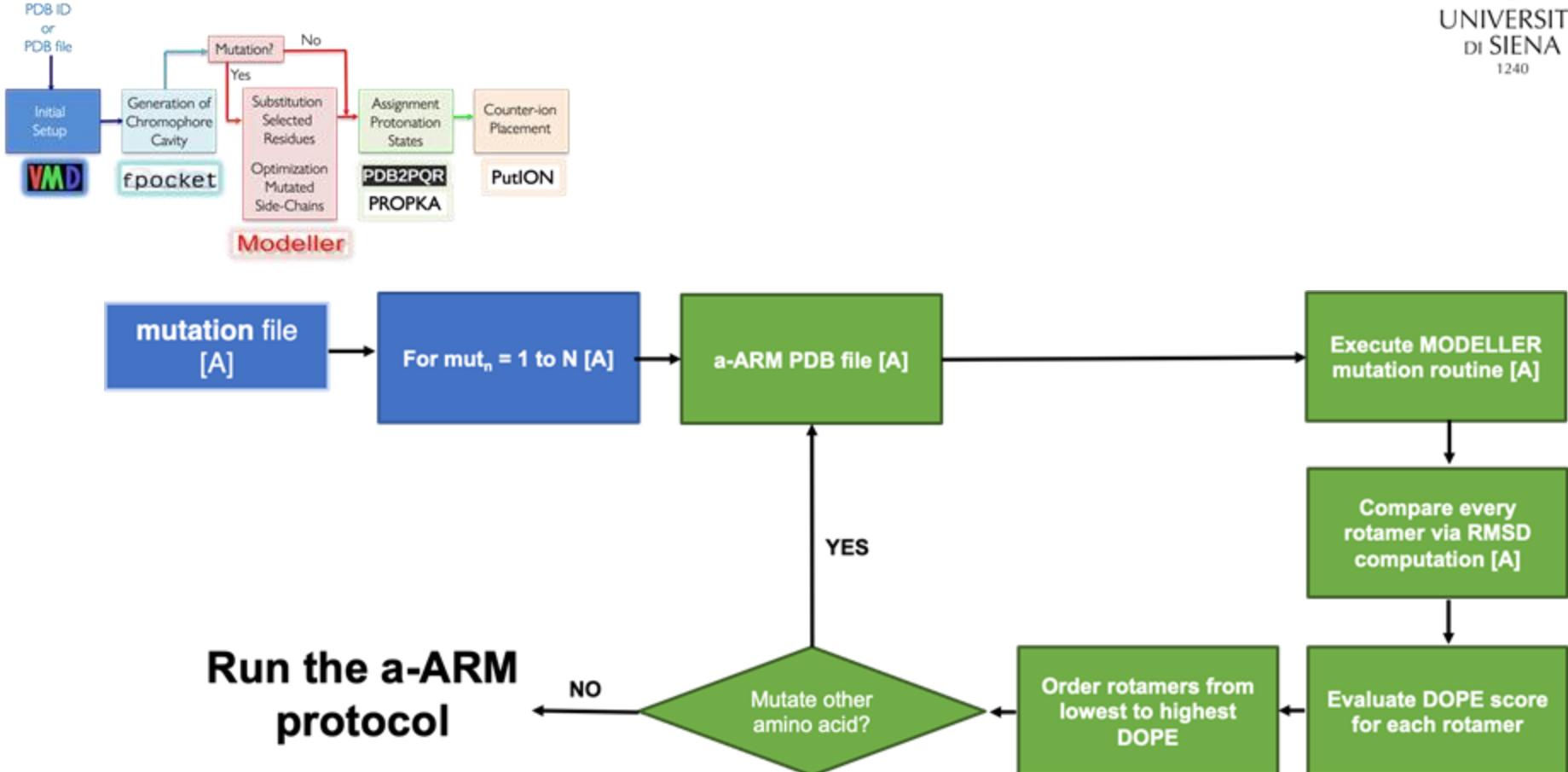
Pedraza-González, L.; De Vico, L.; Marín, M. D. C.; Fanelli, F.; Olivucci, M.. *J. Chem. Theory Comput.* 2019, 15, 3134-3152.

The a-ARM protocol

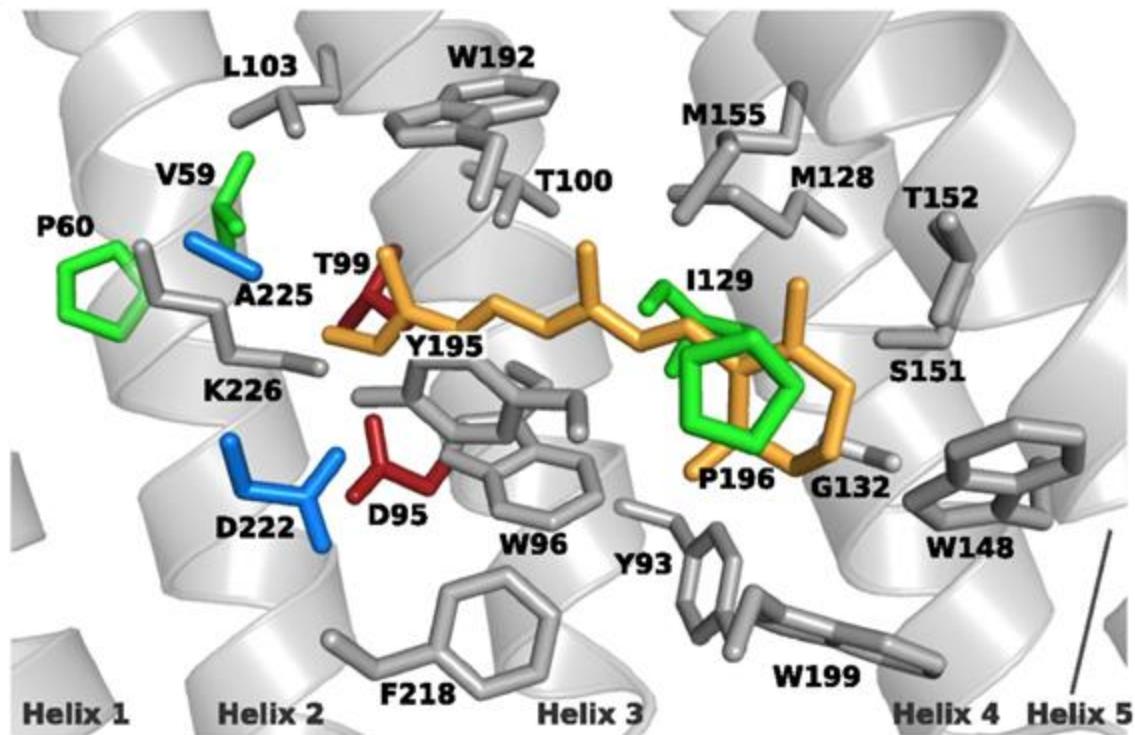
Phase I: Model Structure setup AUTOMATED setup



Design of the mutation strategy



Design of the mutation strategy



Positions V59, P60, D95, T99, P196, D222, A225 were found to be important in mutagenesis experiments

Arnold, F. H. (CalTech)

McIsaac, R. S. et al. Recent advances in engineering microbial rhodopsins for optogenetics. *Curr. Opin. Struct. Biol.* 2015, 33, 8-15.

Cohen, A. E. (Harvard)

Hochbaum, D. R. et al. All-optical electrophysiology in mammalian neurons using engineered microbial rhodopsins. *Nat. Methods* 2014, 11, 825-825.

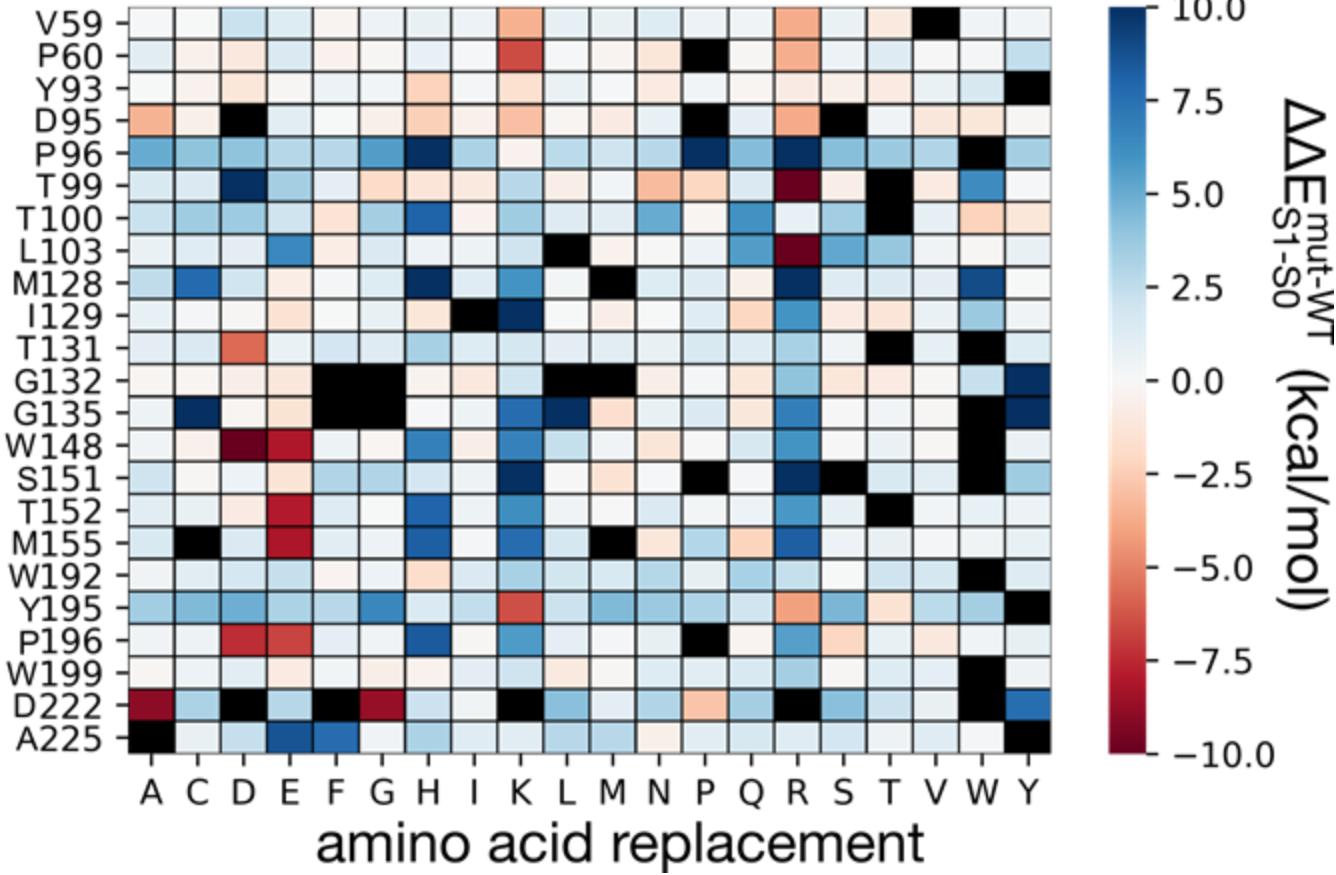
Preliminary results - Single Mutants



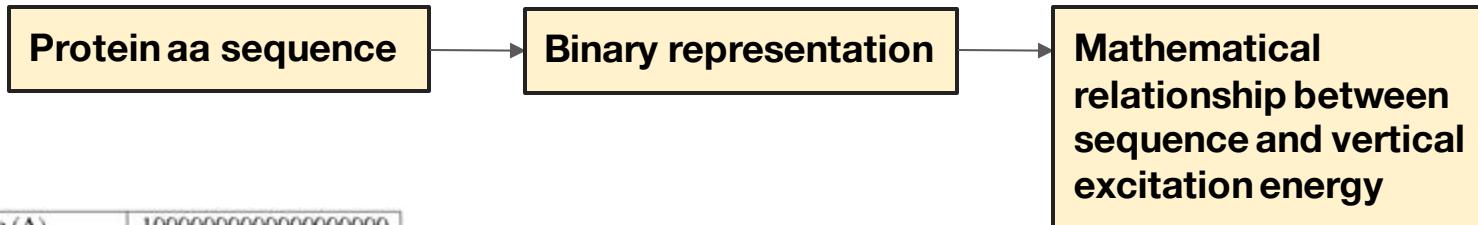
UNIVERSITÀ
DI SIENA
1240

Biggest shifts are displayed for substitutions to charged aa

cavity position



Sequences representation



Alanine (A)	1000000000000000000000000
Cysteine (C)	0100000000000000000000000
Aspartate (D)	0010000000000000000000000
Glutamate (E)	0001000000000000000000000
Phenylalanine (F)	0000100000000000000000000
Glycine (G)	0000010000000000000000000
Histidine (H)	0000001000000000000000000
Isoleucine (I)	0000000100000000000000000
Lysine (K)	0000000010000000000000000
Leucine (L)	0000000001000000000000000
Methionine (M)	0000000000100000000000000
Asparagine (N)	0000000000010000000000000
Proline (P)	0000000000001000000000000
Glutamine (Q)	0000000000000010000000000
Arginine (R)	0000000000000001000000000
Serine (S)	0000000000000000100000000
Threonine (T)	0000000000000000010000000
Valine (V)	0000000000000000001000000
Tryptophan (W)	0000000000000000000010000
Tyrosine (Y)	0000000000000000000000100

$$f(\mathbf{x}) = \beta_0 + \beta_{1,1} x_{1,1} + \beta_{1,2} x_{1,2} + \dots + \beta_{1,N} x_{1,N} \text{ Amino-acid A}$$
$$+ \beta_{2,1} x_{2,1} + \beta_{2,2} x_{2,2} + \dots + \beta_{2,N} x_{2,N} \text{ Amino-acid C}$$
$$\vdots \qquad \vdots \qquad \vdots$$
$$+ \beta_{M,1} x_{M,1} + \beta_{M,2} x_{M,2} + \dots + \beta_{M,N} x_{M,N} \text{ Amino-acid Y}$$

Residue 1 Residue 2 Residue N

Sequences representation

M = 20 (binary representation)

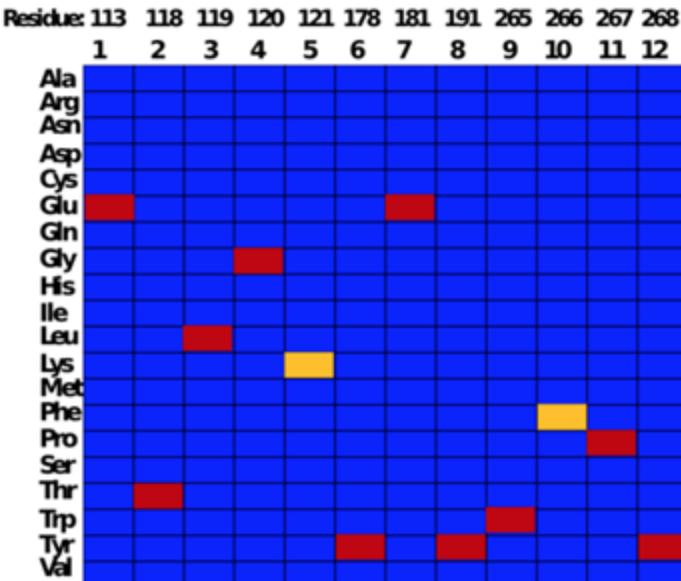
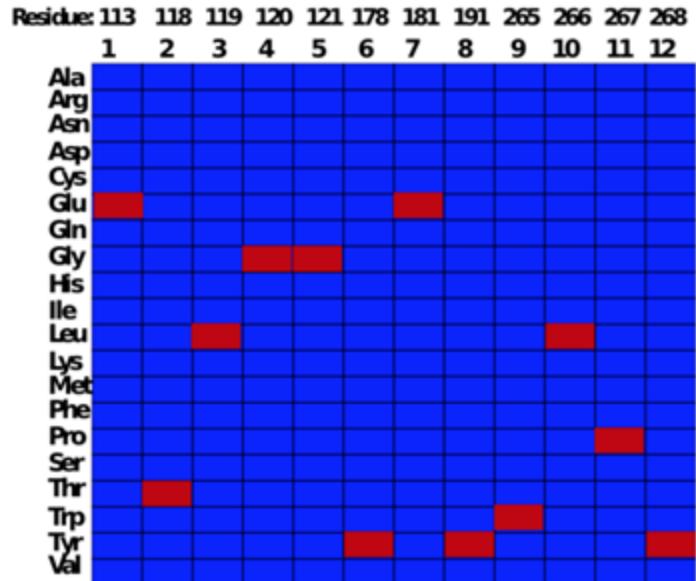
N = sequence length

K = number of variants

Alanine (A)	10000000000000000000
Cysteine (C)	01000000000000000000
Aspartate (D)	00100000000000000000
Glutamate (E)	00010000000000000000
Phenylalanine (F)	00001000000000000000
Glycine (G)	00000100000000000000
Histidine (H)	00000010000000000000
Isoleucine (I)	00000001000000000000
Lysine (K)	00000000100000000000
Leucine (L)	00000000010000000000
Methionine (M)	00000000001000000000
Asparagines (N)	00000000000100000000
Proline (P)	00000000000010000000
Glutamine (Q)	00000000000001000000
Arginine (R)	00000000000000100000
Serine (S)	00000000000000001000
Threonine (T)	00000000000000001000
Valine (V)	00000000000000000100
Tryptophan (W)	000000000000000000010
Tyrosine (Y)	000000000000000000001

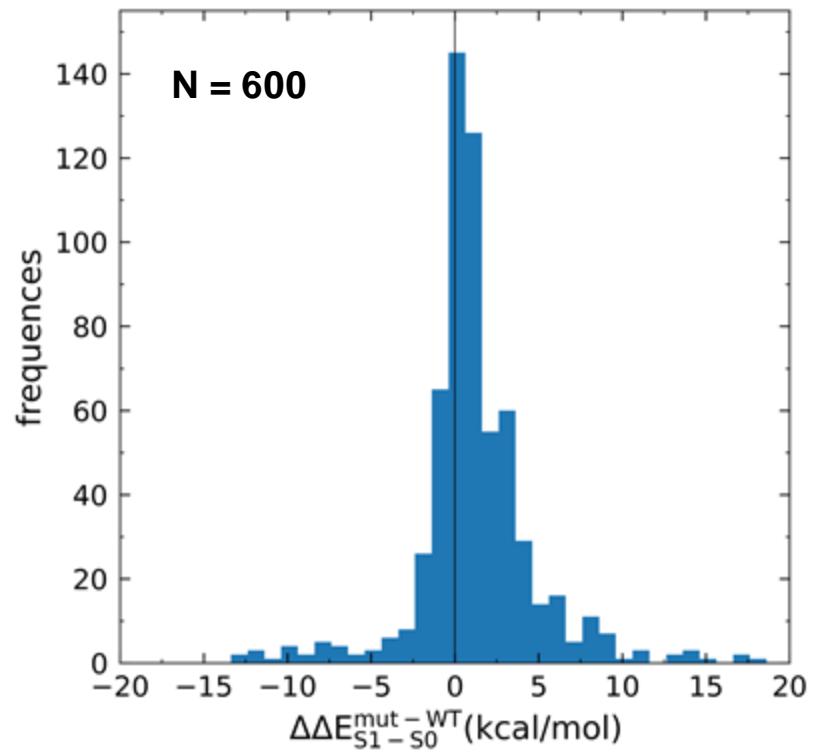
To represent K
rhodopsins I need a
matrix with dimensions:
K, MxN

Sequences representation

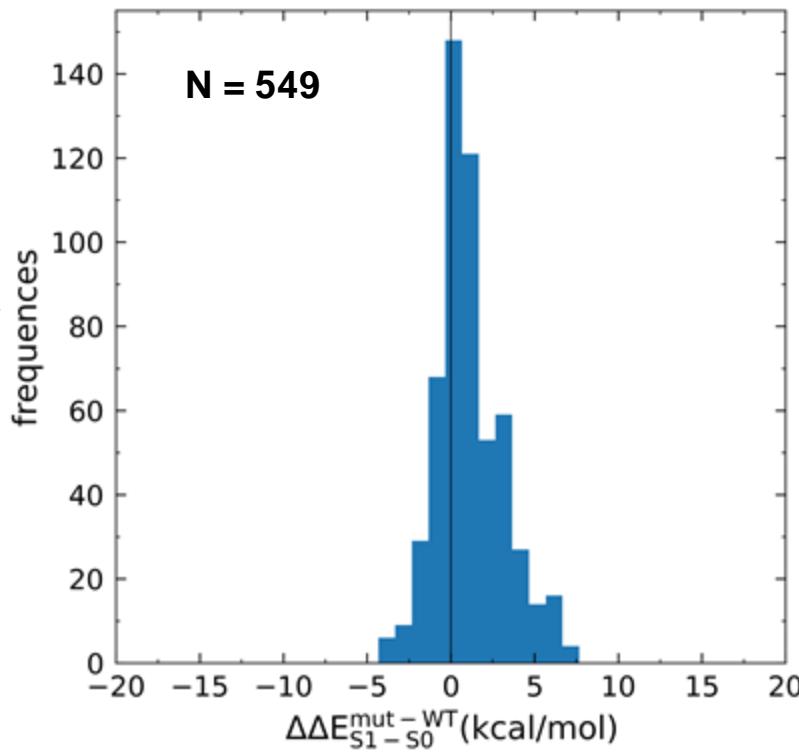


Preprocessing of input data - towards ML

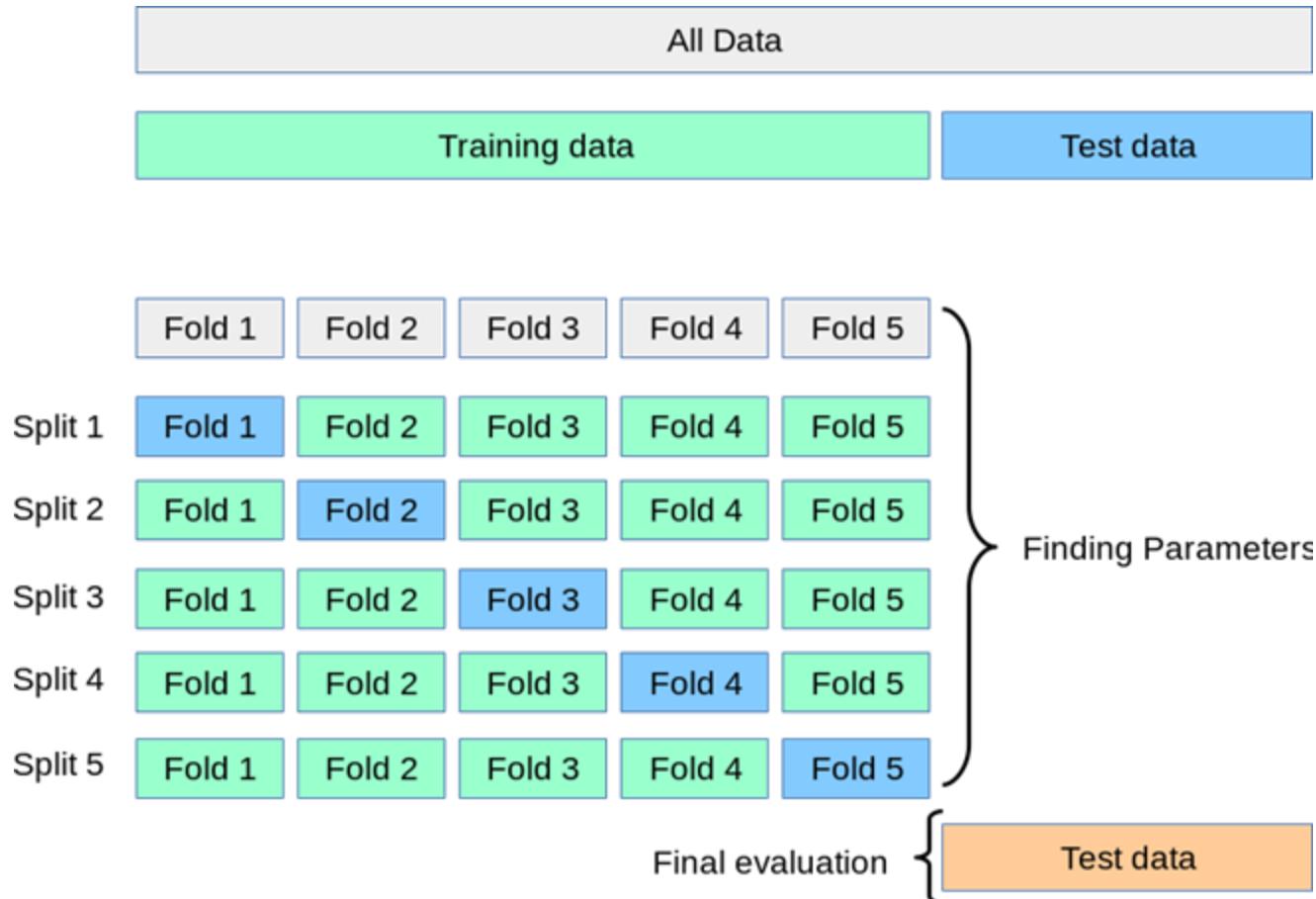
Are outliers actually outliers?



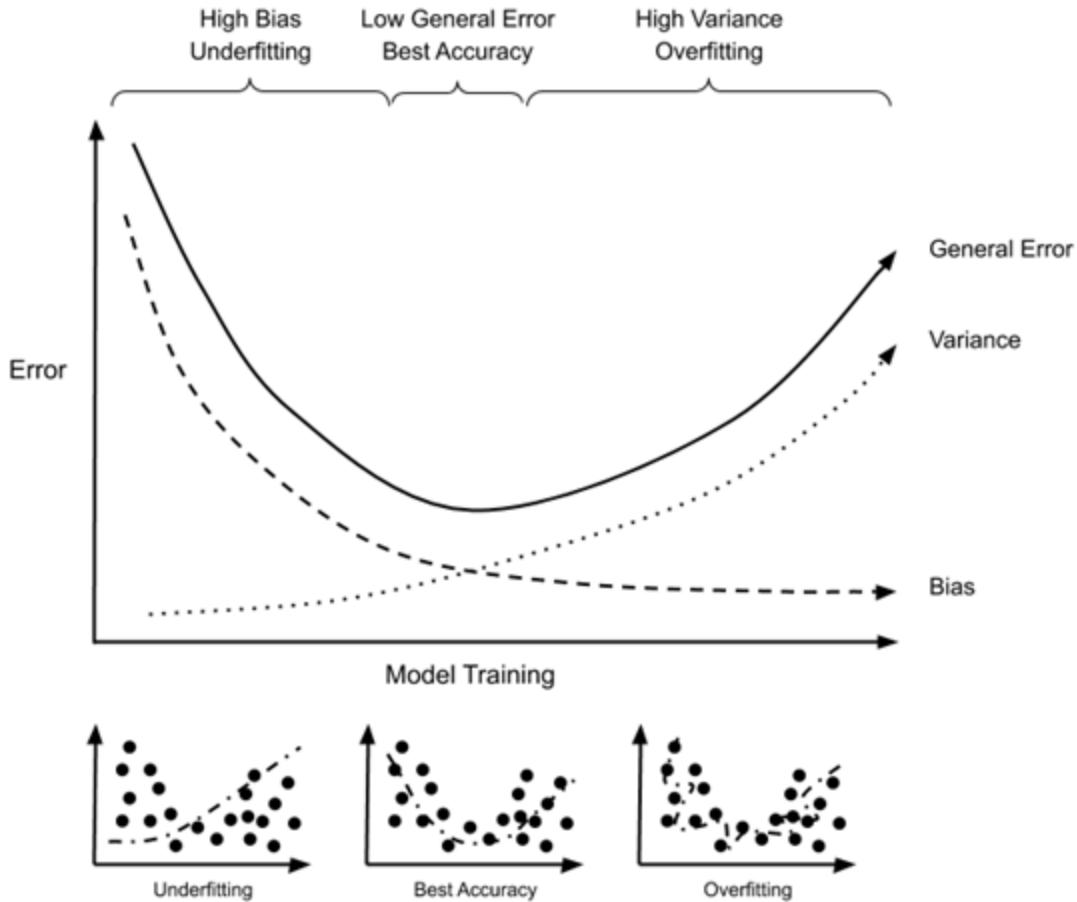
IQR range rule



Preprocessing of input data - towards ML



Preprocessing of input data - towards ML

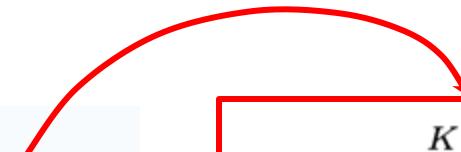


M = 20 (binary representation of a single AA)

N = sequence length

K = number of variants

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^M \sum_{j=1}^N \beta_{i,j} x_{i,j}$$


$$\min_{\beta, \beta_0} \sum_{k=1}^K (\Delta E_{S1-S0}^{(k)} - \beta_0 - \sum_{i=1}^M \sum_{j=1}^N \beta_{i,j} x_{i,j}^{(k)})^2$$

or , in matrix notation:

$$\mathbf{y} = \boldsymbol{\beta}_0 + \mathbf{X}\boldsymbol{\beta} \longrightarrow \|\mathbf{e}\|^2 = \| \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \|^2 \rightarrow \min \| \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \|^2 \rightarrow \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Linear Function

Define as least square problem

Solve analytically or via GD

Linear approaches - Ridge Regression

M = 20 (binary representation)

N = sequence length

K = number of variants

$$\min_{\beta, \beta_0} \sum_{k=1}^K (\Delta E_{S1-S0}^{(k)} - \beta_0 - \sum_{i=1}^M \sum_{j=1}^N \beta_{i,j} x_{i,j}^{(k)})^2 + \lambda \sum_{i=1}^M \sum_{j=1}^N \beta_{i,j}^2$$

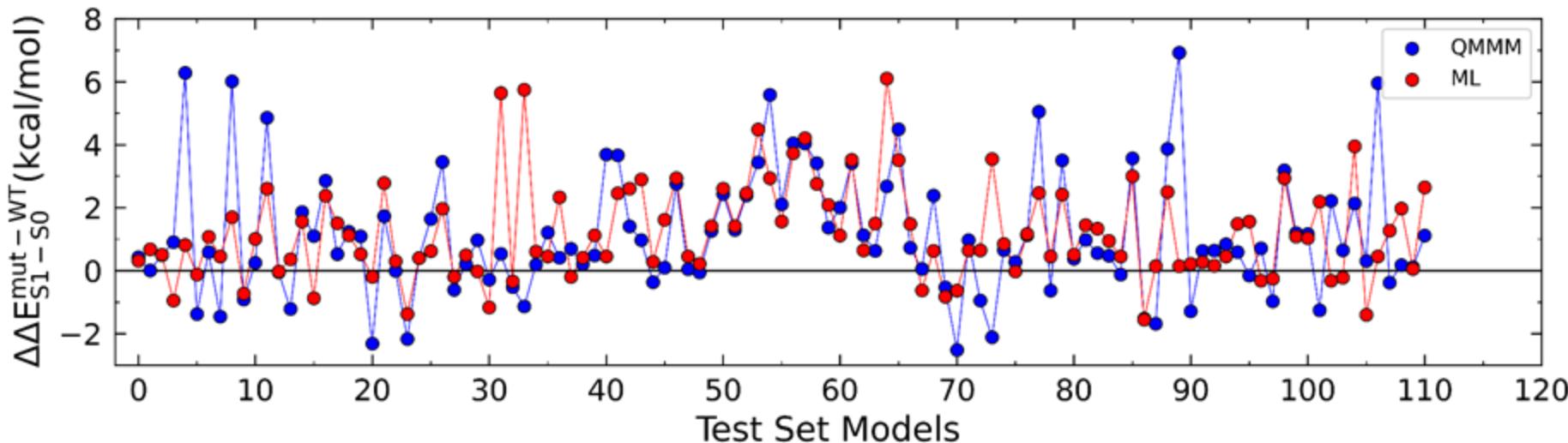
or , in matrix notation:

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + I\lambda)^{-1} \mathbf{X}^T \mathbf{y}$$

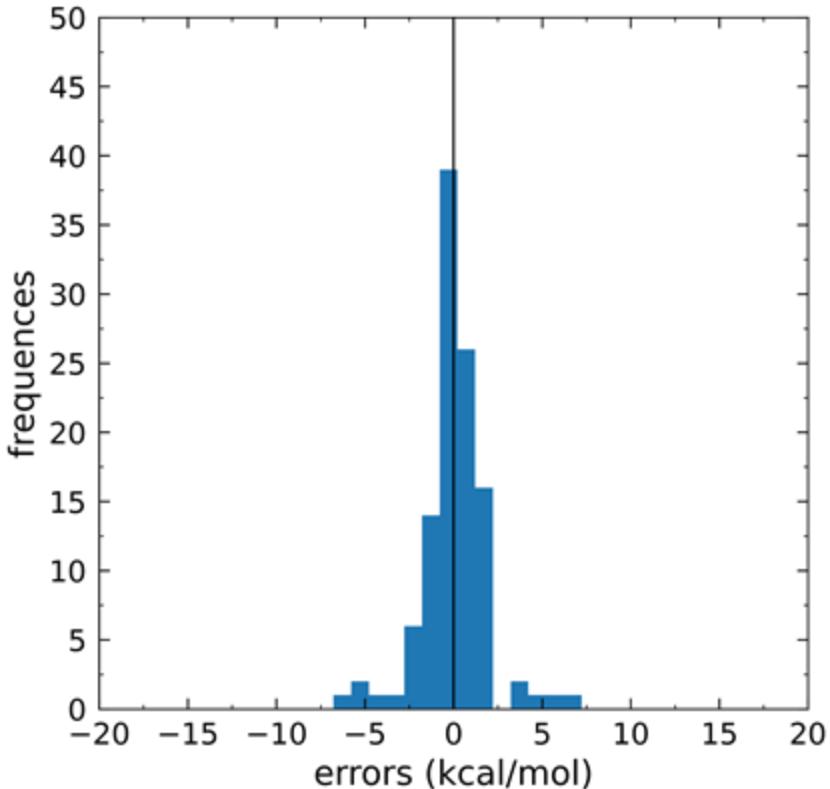
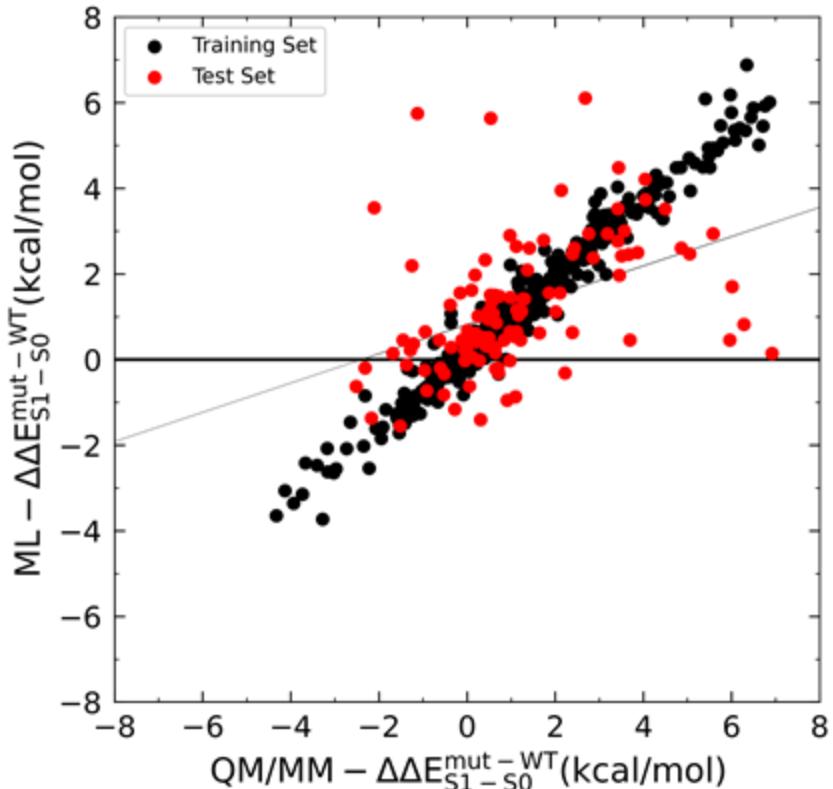
Linear approaches - Ridge Regression

Ridge Regression

MAE $\rightarrow 1.10 \text{ kcal/mol}$



Linear approaches - Ridge Regression



Linear approaches - Lasso Regression

M = 20 (binary representation)

N = sequence length

K = number of variants

$$\min_{\beta, \beta_0} \sum_{k=1}^K (\Delta E_{S1-S0}^{(k)} - \beta_0 - \sum_{i=1}^M \sum_{j=1}^N \beta_{i,j} x_{i,j}^{(k)})^2 + \lambda \sum_{i=1}^M \sum_{j=1}^N |\beta_{i,j}|$$

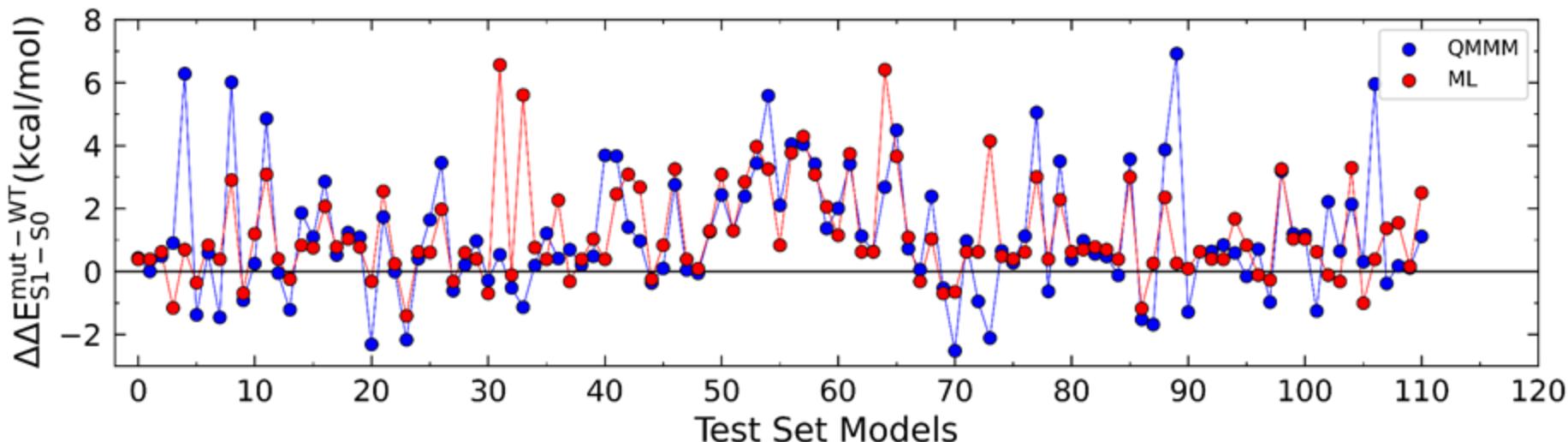
Linear approaches - Lasso Regression

Lasso Regression

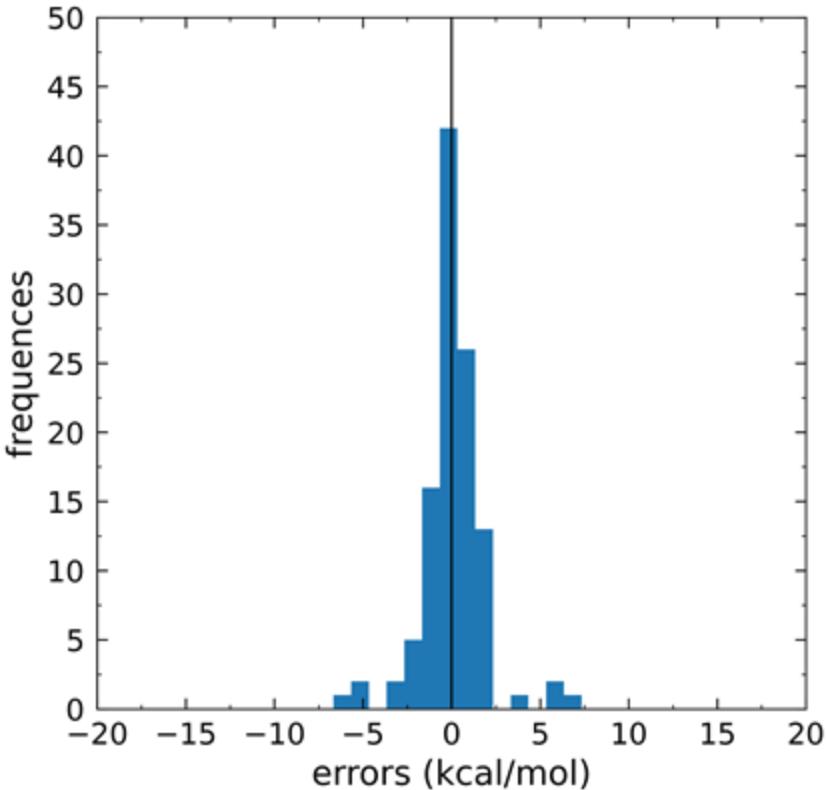
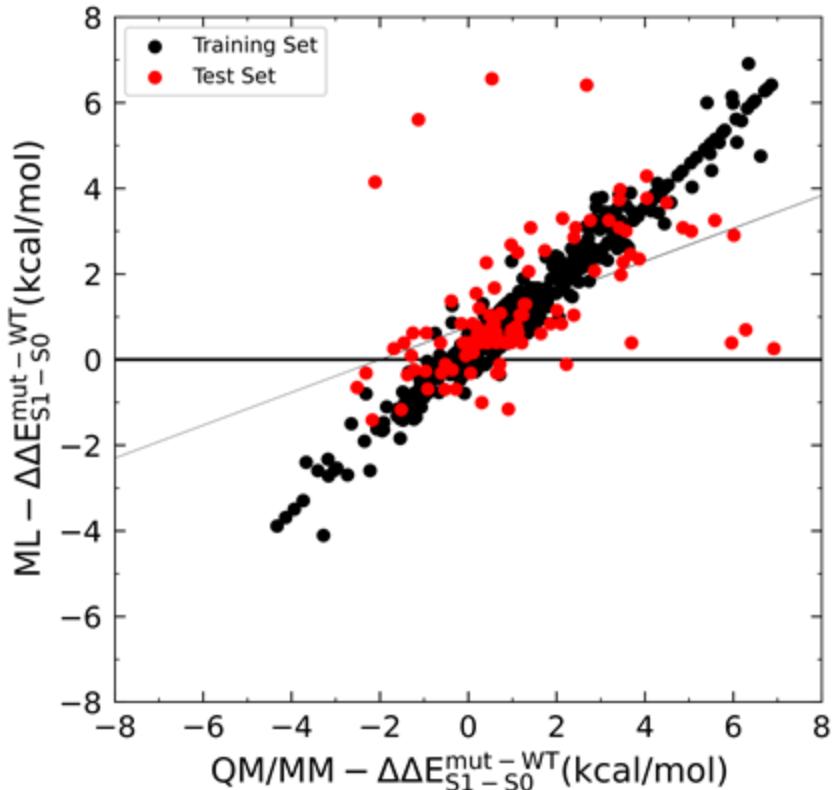
MAE -> 1.05 kcal/mol

Ridge Regression

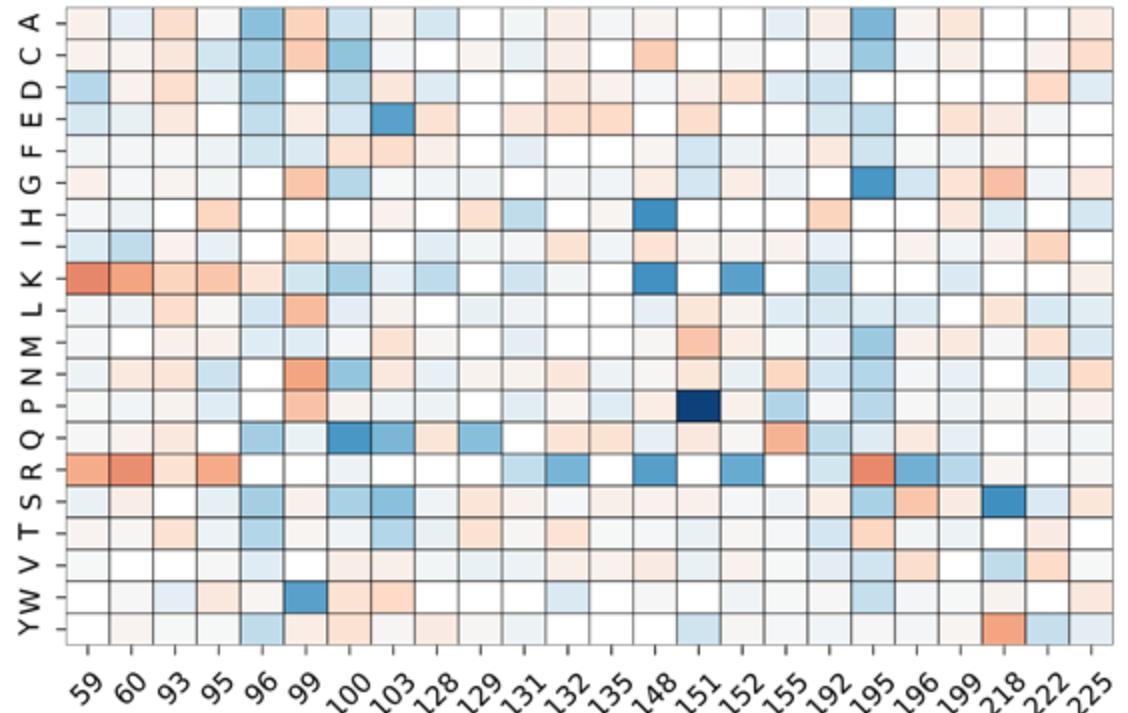
MAE -> 1.10 kcal/mol



Linear approaches - Lasso Regression

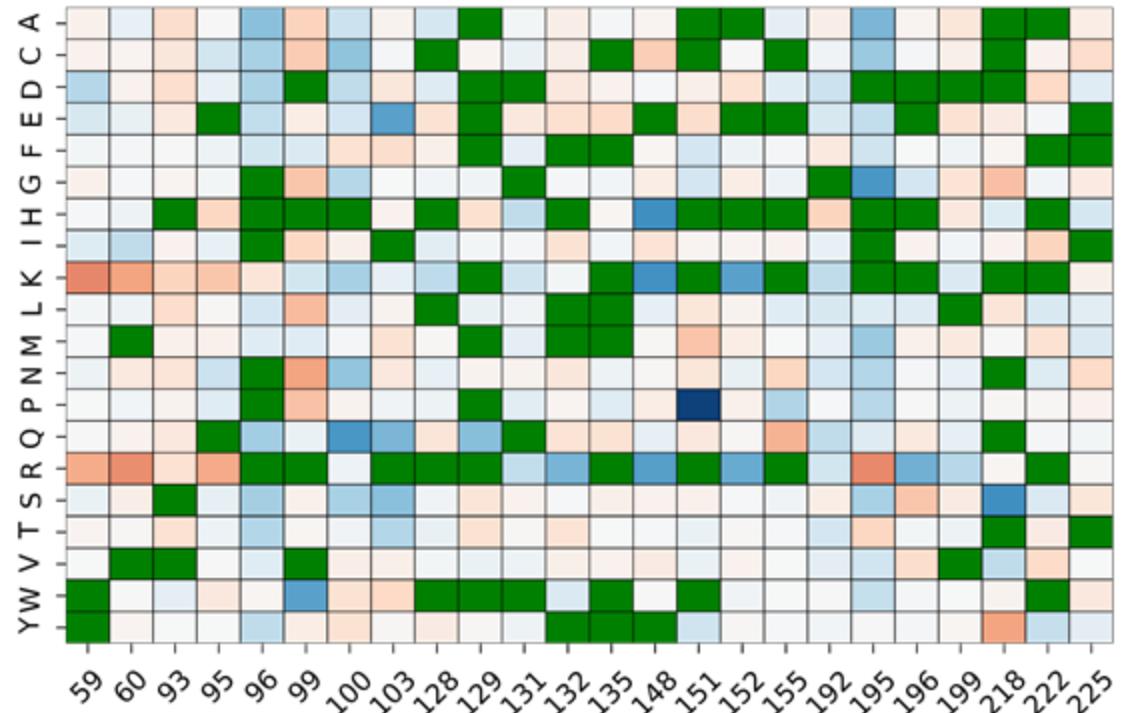


Lasso Regression - coefficients representation



$$\sum_{i=1}^M \sum_{j=1}^N x_{i,j} \beta_{i,j}$$

Lasso Regression - coefficients representation



$$\sum_{i=1}^M \sum_{j=1}^N x_{i,j} \beta_{i,j}$$

Non-Linear approaches - Random Forest

RF Regression

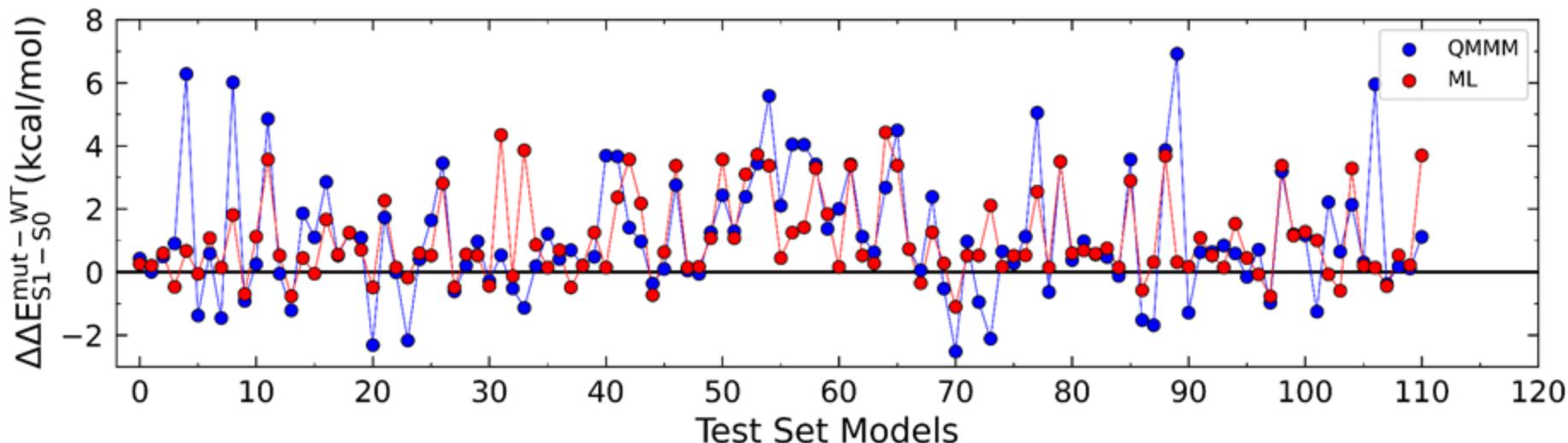
MAE -> 1.00 kcal/mol

Lasso Regression

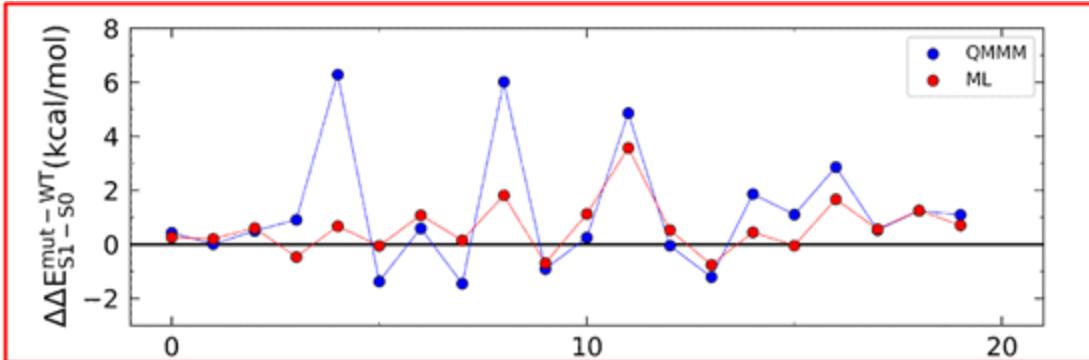
MAE -> 1.05 kcal/mol

Ridge Regression

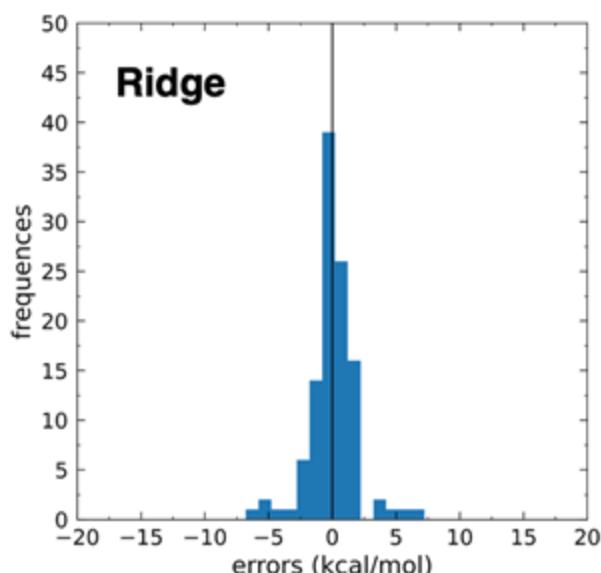
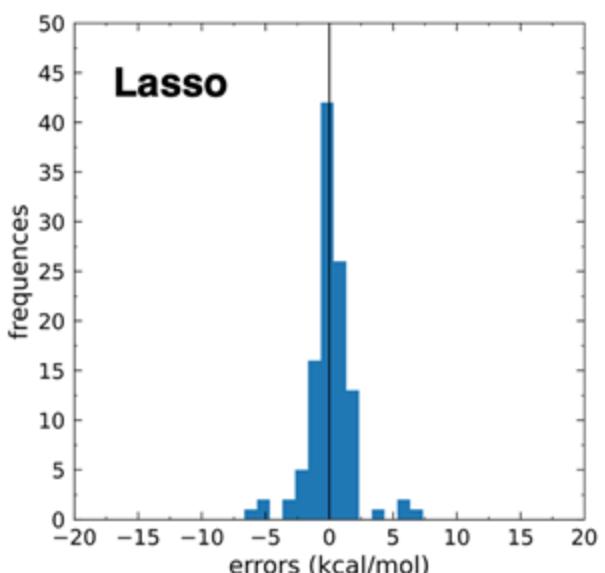
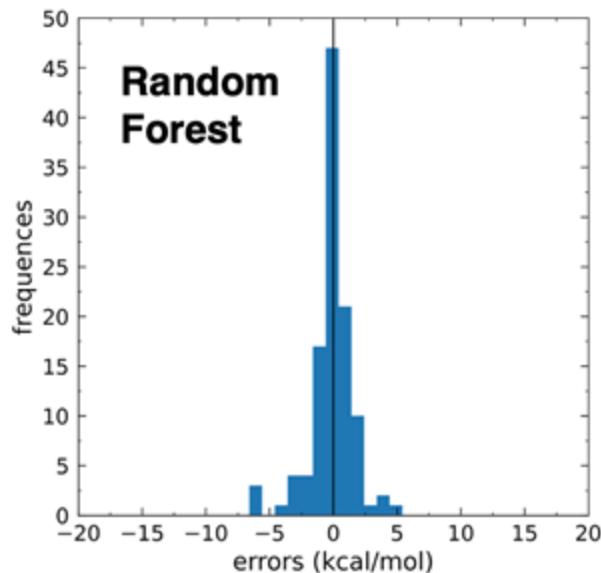
MAE -> 1.10 kcal/mol



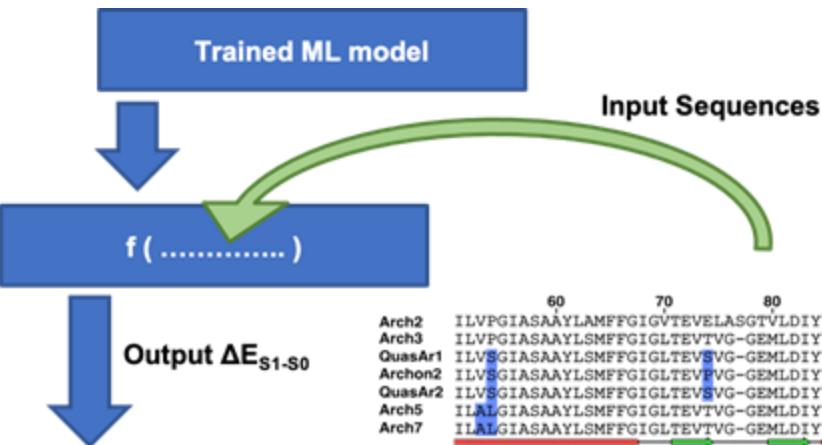
Non-Linear approaches - Random Forest



Non-Linear approaches - Random Forest



Meaningfulness of the models



60 70 80 90 100

Arch2	ILVPGIASAAYLAMFFGIGGVTEVELASGTVLIDIIYARYADWLFTTPLLLLL
Arch3	ILVPGIASAAYLSMFFGIGLTTEVTVG-GEMLDIYYARYADWLFTTPLLLLL
QuasAr1	ILV G IASAAYLSMFFGIGLT E V S VG-GEMLDIYYARYA W LFTTPLLLLL
Archon2	ILV G IASAAYLSMFFGIGLT E V P VG-GEMLDIYYARYA W LFTTPLLLLL
QuasAr2	ILV G IASAAYLSMFFGIGLT E V S VG-GEMLDIYYARYA Q WLFTTPLLLLL
Arch5	ILALGIASAA Y LSMFFGIGLT E V T VG-GEMLDIYYARYA E WLFTTPLLLLL
Arch7	ILALGIASAA Y LSMFFGIGLT E V T VG-GEMLDIYYARYA E WLFTTPLLLLL

110 120 130 140 150

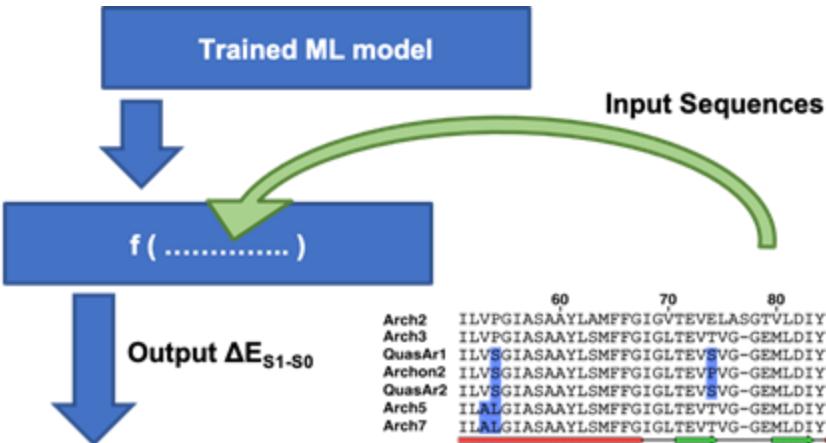
Arch2	DLALLAKVDRVTIGTLIGVDALIMVTGLIGALS K PLARYTWWLFSTIAF
Arch3	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM
QuasAr1	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM
Archon2	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM
QuasAr2	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM
Arch5	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM
Arch7	DLALLAKVDRVTIGTLVGVDALIMVTGLIGALSHTAIARYSWWL F STICM

Protein	Exp ΔE_{S1-S0} (kcal/mol)	Exp shift (kcal/mol)	ML shift (kcal/mol)
Arch3	51.4	0	0
D95E/T99C/V59A	46.0	-5.4	-2.9
D95E/T99C/P60L	45.8	-5.6	-2.7
D95E/T99C/P196S	45.5	-5.9	-4.3
Arch5	46.0	-5.4	-5.1
Arch7	46.4	-5.0	-3.1

Meaningfulness of the models

ML predicts shifts in the same direction as Exp.

Good QM/MM models?



	60	70	80	90	100
Arch2	ILVPGIASAAYLAMFFGIGTVTEVELASGTVLIDIIYARYADWLFTTPLLLL				
Arch3	ILVPGIASAAYLSMFFGIGLTETVVG-GEMLDIYYARYADWLFTTPLLLL				
QuasAr1	ILVPGIASAAYLSMFFGIGLTETVVG-GEMLDIYYARYA WLF FTTPLLLL				
Archon2	ILVPGIASAAYLSMFFGIGLTETVVG-GEMLDIYYARYA WLF FTTPLLLL				
QuasAr2	ILVPGIASAAYLSMFFGIGLTETVVG-GEMLDIYYARYA WLF FTTPLLLL				
Arch5	ILALGIASAAYL SMFFGIGLTETVVG-GEMLDIYYARYAWLFFTTP LLLL				
Arch7	ILALGIASAAYL SMFFGIGLTETVVG-GEMLDIYYARYAWLFFTP LLLL				

	110	120	130	140	150
Arch2	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS KTPLARYT WLF STIAF				
Arch3	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				
QuasAr1	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				
Archon2	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				
QuasAr2	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				
Arch5	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				
Arch7	D LALLAKVDRVTIGTLIGVDALIMVTGLIGALS HATAIRYS WLF STICM				

Protein	Exp ΔE_{S1-S0} (kcal/mol)	Exp shift (kcal/mol)	ML shift (kcal/mol)
Arch3	51.4	0	0
D95E/T99C/V59A	46.0	-5.4	-2.9
D95E/T99C/P60L	45.8	-5.6	-2.7
D95E/T99C/P196S	45.5	-5.9	-4.3
Arch5	46.0	-5.4	-5.1
Arch7	46.4	-5.0	-3.1

- Statistical Analysis of the excitation energies
- Increase the training set by generating more double mutants
- Define the problem as a classification task
- Use non linear-models including NN
- Outliers (?)