# Biologics Modelling
# Proteins & Peptides

Annette Höglund, Jianxin Duan

Helsingfors

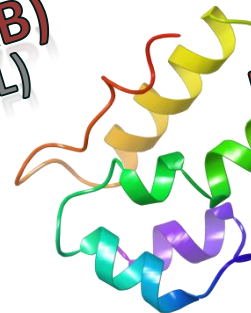November 12, 2015

SCHRÖDINGER.

# BioLuminate Features

- Protein-protein docking
- Antibody structure prediction from sequence
- Antibody humanization
- Fast homology model generation
- Accurate long loop predictions
- Residue scanning
- Affinity Maturation
- Cysteine scanning
- Crosslink design
- Peptide QSAR
- Aggregation hot spot ID
- Free energy perturbation

SCHRÖDINGER.

# Using BioLuminate to Go From Sequence to Model of Antibody/Antigen Complex

- **Starting point:**
  - Sequence of antibody
    - FAB13B5
  - Crystal structure of known antigen (unbound)
    - HIV-1 Capsid Protein (P24) (Dimerization Domain)
    - 1A43

- **Can we use computational methods to predict structure of antibody/antigen complex?**
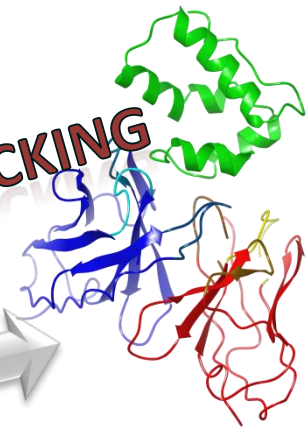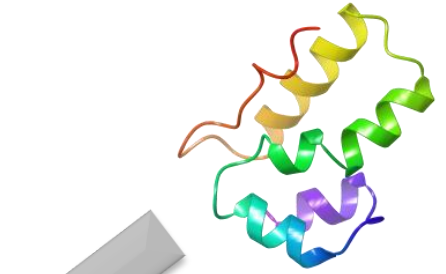  - Epitope ID



HOMOLOGY MODEL (AB)
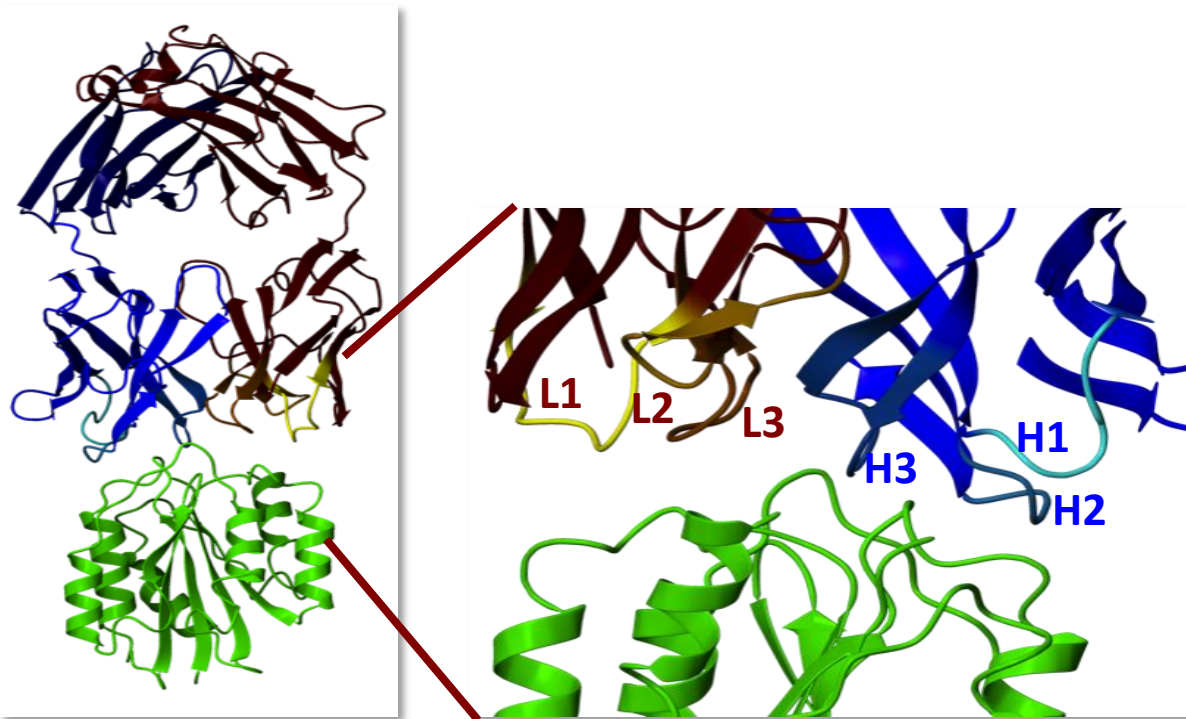(ANY HOMOLOGY MODEL)

X-RAY

PROTEIN PROTEIN DOCKING

SCHRÖDINGER®

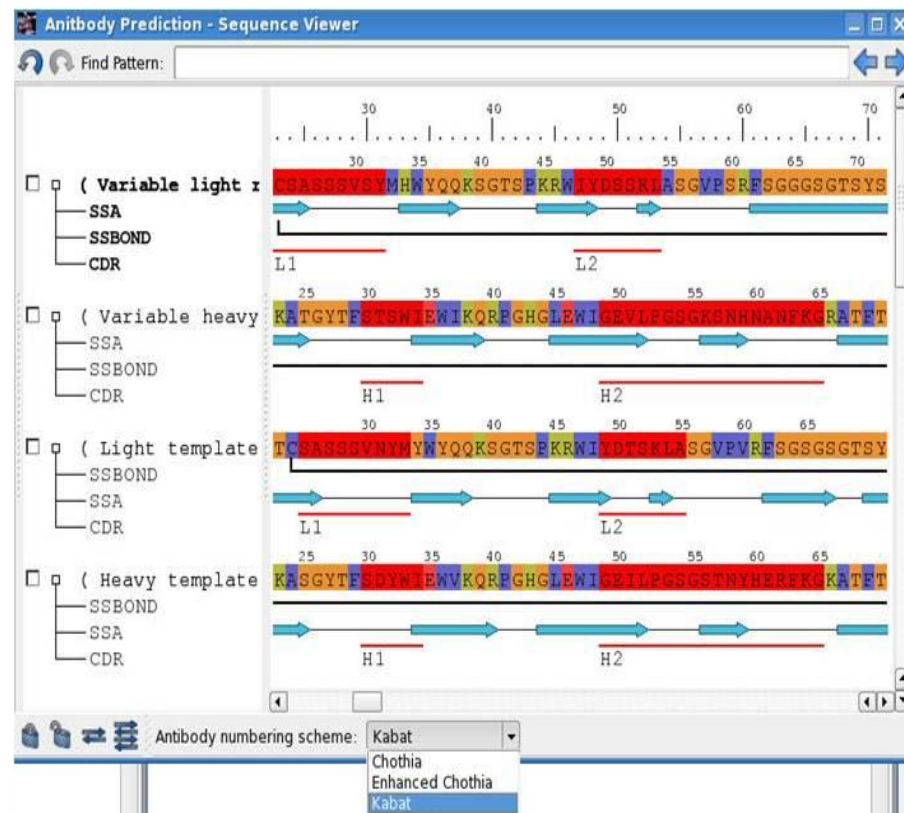# Antibody Modeling Using BioLuminate...

## Workspace:



Recognizes & colors chains and CDRs

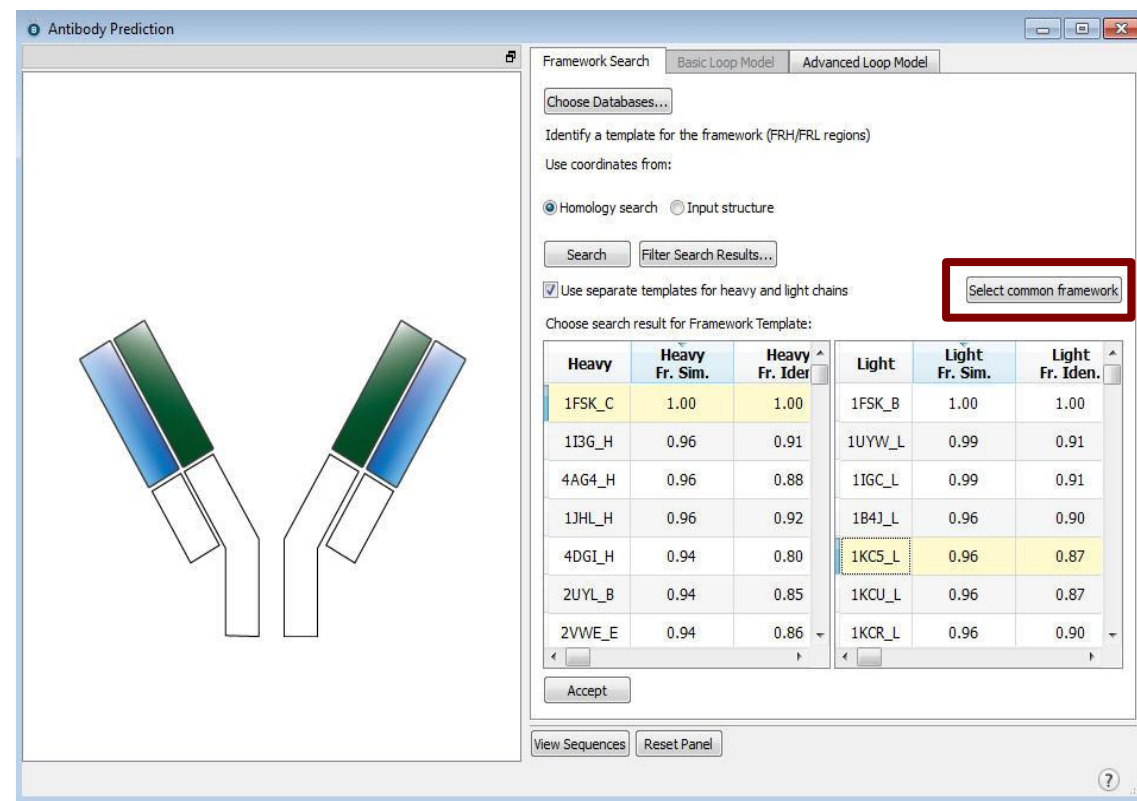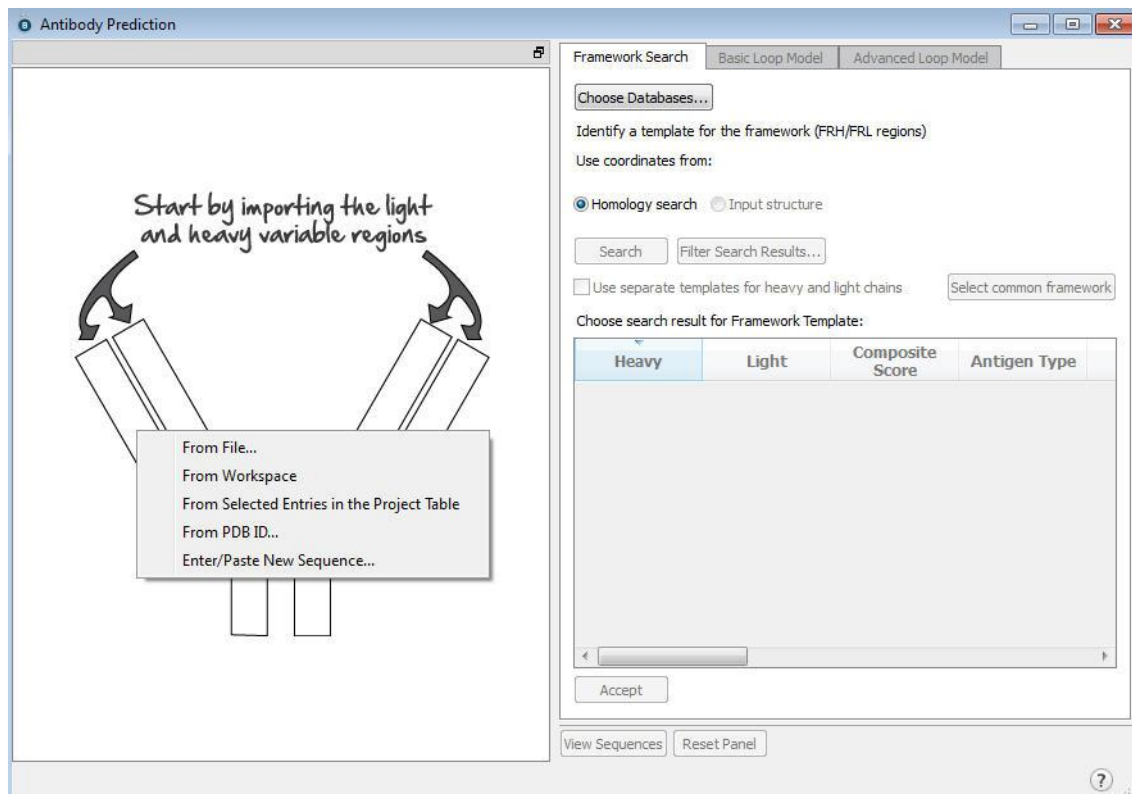Antibody Aware Environment

## Sequence Viewer:



- Detects and labels CDR loops
- Multiple numbering schemes

SCHRÖDINGER.

# Antibody Modeling – CDR Prediction: Input & Framework



- Antibody-specific workflow
- Search public or in-house structures for templates

- Framework selection:
  - Separate control over L/H chain framework templates
  - Control over framework used to align chains

SCHRÖDINGER.

# H3 Antibody Loop Prediction Remains Difficult Using Homology Methods

| Program | \<L1\> | \<L2\> | \<L3\> | \<H1\> | \<H2\> | \<H3\> | # Best | # Worst |
|---------|--------|--------|--------|--------|--------|--------|--------|---------|
| Accelrys | 1.2 | 0.7 | 1.4 | 1.1 | 1.6 | 3.0 | 1 | 3 |
| CCG/MOE | 0.7 | 0.5 | 1.5 | 1.3 | 1.1 | 3.6 | 1 | 1 |
| PIGS server | 1.0 | 0.4 | 1.4 | 1.1 | 0.8 | 3.2 | 3 | 0 |
| Rosetta | 1.0 | 0.5 | 1.6 | 1.7 | 1.1 | 3.3 | 0 | 2 |
| **BioLuminate** | **1.0** | **0.5** | **1.2** | **1.1** | **1.1** | **2.2** | **3** | **0** |

Red - Best average RMSD for loop

Gray - Worst average RMSD for loop

Use *de novo* approach to predict H3

SCHRÖDINGER.

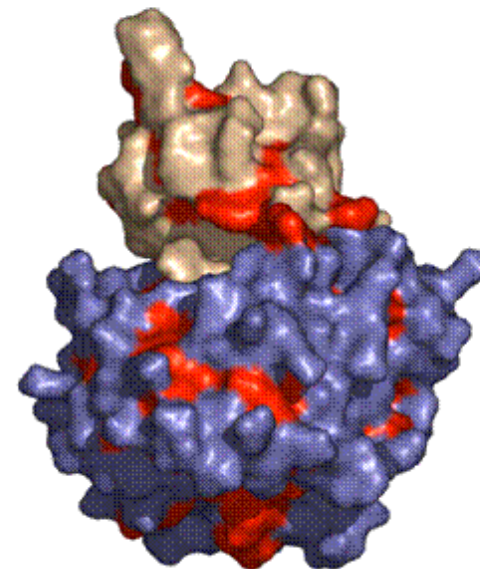# Protein-Protein Docking: How do Two Proteins Best Fit Together?

- Licensed from Vajda group at Boston University
  - Kozakov et al. (2006) *Proteins: Struct, Funct, Bioinf* **65** 392-406
- #1 server in most recent CAPRI competition
  - Competitive with human groups

| #  | Human groups: | Automatic Servers: |
|----|---------------|---------------------|
| 1  | Sandor Vajda | CLUSPRO |
| 2  | Martin Zacharias | HADDOCK |
| 3  | Xiaoqin Zou | GRAMM-X |
| 4  | Haim Wolfson, Miriam Eisenstein | SKE-DOCK |
| 5  | Huan-Xiang Zhou, Zhiping Weng | PatchDock, FireDock, FiberDock |
| 6  | Alexandre Bonvin | TOP-DOWN |
| 7  | Juan Fernandez-Recio | |
| 8  | Jeffrey Gray | |

CAPRI rankings
(Nir London, Rosetta Design
Group, 2010)

Piper/Cluspro:
- **#1 group**
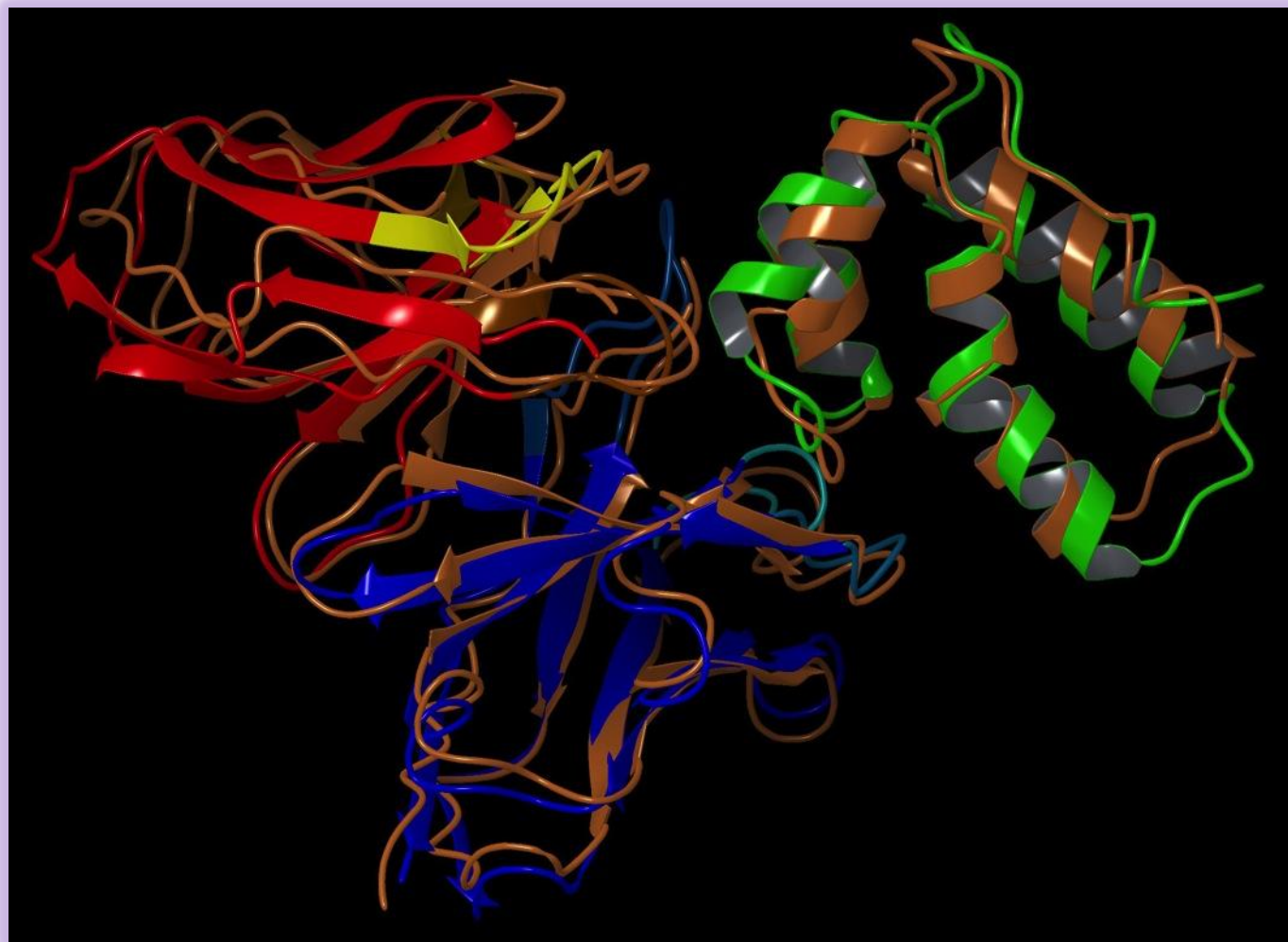- **#1 server**

SCHRÖDINGER.

# Antibody Prediction Using BioLuminate



Predicted CDR region FAB13B5 versus experiment (1E6J, light green)

SCHRÖDINGER.

# Antibody/Antigen Complex: Predicted Versus Experiment



Modeled FAB13B5 CDR docked with crystal structure of unbound antigen P24 (orange) versus x-ray complex 1E6J. 3rd ranked complex shown.

SCHRÖDINGER.

# Antibody Modeling – Humanization

## CDR-Grafting (Framework Replacement)



- Easy to use
- Automatically IDs clashing residues for back mutation

## Homology-based suggestions



Compare to human sequences

- Degree of variability
- 3D information

# Aggregation Prediction

- Aggregation can be viewed as recognition by a large sphere
  - Roll large probe sphere
  - Detect patches of exposed hydrophobic residues

- Reference: SAP (Spatial Aggregation Propensity)
  - Validated through collaboration with Novartis
  - Chennamsetty et al. (2009) *PNAS* 106 11937



Normal surface uses water probe

Much larger effective "protein" probe

Color surface red to reflect hydrophobicity of contributing residues
Red hydrophobic "hot spots" are likely aggregation regions

SCHRÖDINGER.

# Aggregation Surface Analysis

# Protein Stability Thermodynamics



Folding/unfolding equilibrium

**But what does the unfolded state really look like?**

SCHRÖDINGER

# Schematic Thermodynamic Cycle



- Simulate the (non-physical) protein side chain transformation

- Standard approach in FEP

- All non-physical terms cancel in the final result

$$\Delta\Delta G_{stability} = \Delta G_1 - \Delta G_2$$
$$= \Delta G_A - \Delta G_B$$

# Setting up a Cycle

Unfolded state is modeled by capped peptide



$\mathrm{D}G_{protein}$

$\mathrm{D}G_{unfold1}$

$\mathrm{D}G_{unfold2}$

$\mathrm{D}G_{peptide}$

$$\mathrm{DD}G = \mathrm{D}G_{unfold1} - \mathrm{D}G_{unfold2} = \mathrm{D}G_{protein} - \mathrm{D}G_{peptide}$$

SCHRÖDINGER.

# Current Results, Part of the Fold-X Test Set

| System | PDB ID | # Mutations | R²-value | MUE | RMSE | ΔΔG Sign correct |
|---|---|---|---|---|---|---|
| T4-Lysozyme | 2LZM | 66 | 0.67 | 1.2 | 1.6 | 92% |
| Human Lysozyme | 1REX | 45 | 0.66 | 1.3 | 1.8 | 80% |
| Peptostrept. Magn. Prot. L | 1HZ6 | 44 | 0.59 | 1.1 | 1.3 | 89% |
| B1 IG binding protein G | 1PGA | 24 | 0.37 | 1.1 | 1.4 | 79% |
| Fibronectin II domain | 1TEN | 32 | 0.60 | 1.4 | 1.7 | 91% |
| FK506 BP | 1FKB | 27 | 0.4 | 1.6 | 4.9* | 85% |
| All | | 238 | 0.55 | 1.2 | 1.7 | 87% |

**Note: No charge changes in this set!**

Errors in kcal/mol

*Result strongly influenced by some outliers

SCHRÖDINGER.

# Correlation Plot



**Slope = 1.3**
**Offset = -0.2**

Blue:
Outliers of more than 3x MUE

# Comparison to Other Tools

- FEP+ performs well, but comparable to other tools
- For FEP+, no parameterisation was necessary, so results are more transferable

| Software | R$^2$-value achieved* | Stabilizing/destabilizing % correct | MUE [kcal/mol] |
|---|---|---|---|
| CC/PBSA | 0.31 | 79% | 1.0 |
| EGAD | 0.35 | 71% | 1.0 |
| FoldX | 0.25 | 70% | 1.3 |
| Hunter | 0.20 | 69% | 1.1 |
| I-Mutant2.0 | 0.29 | 78% | 1.1 |
| Rosetta | 0.07 | 73% | 1.7 |
| FEP+ (smaller data set!) | 0.55 | 87% | 1.2 |

* Calculated from R-values given in Tab I of Potapov, 2009, Prot. Eng. Des. Sel., 22, 553

SCHRÖDINGER.

# The development of peptide therapeutics

**Lead Identification**

- Literature search
- Library screening

**Hit to Lead**

- Affinity and selectivity optimization
- Property Optimization

**Lead Optimization**

- Half-life extension

SCHRÖDINGER.

# Computational tools can accelerate each step

**Lead Identification**

- Peptide docking

**Hit to Lead**

- Residue Scanning
- Affinity Maturation
- Peptide QSAR/QSAM

**Lead Optimization**

- Peptide QSAR/QSAM

SCHRÖDINGER.
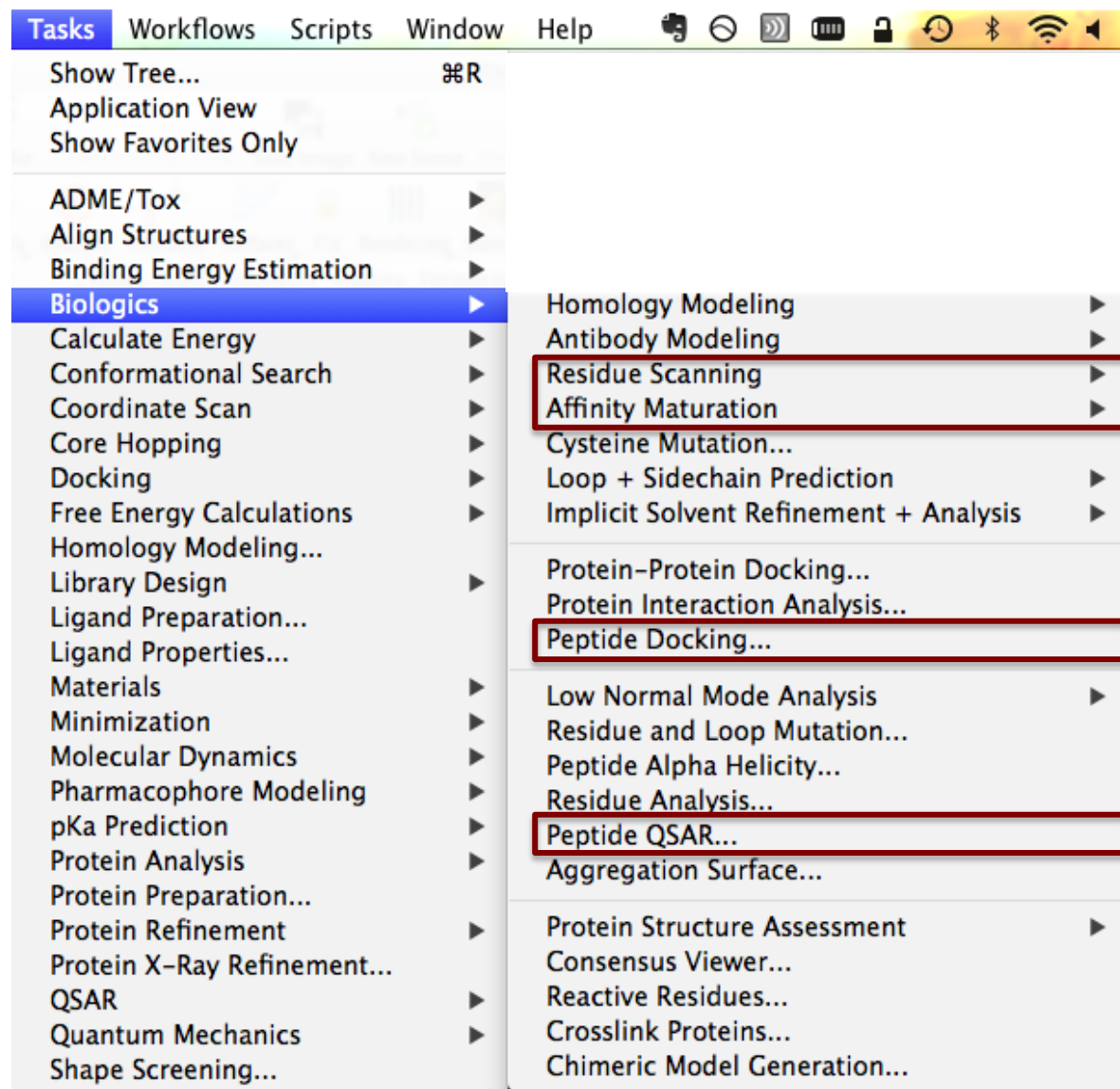
# Peptide Modeling with the Biologics Suite

# Lead Identification

Peptide Docking

SCHRÖDINGER.

# Polypeptide docking and Glide

- Several brute force sampling methods performs well in polypeptide docking but require hundreds or thousands of CPU hours per polypeptide docking
  - Rosetta FlexPepDock *ab-initio* (Raveh et al., *PloS one* **2011**, *6*, e18934)
  - DynaDock (Antes, *Proteins*, **2010**, *78*, 1084)
  - HADDOCK (Trellet et al., *PloS one* **2013**, *8*, e58769)
- Small molecule docking programs such as Glide are comparatively fast and accurate for docking of small molecules

**Question:  Can Glide SP dock 4-11 residue polypeptides well?**

Raveh B et al. *PloS one.* **2011**; *6*: e18934.

SCHRÖDINGER.

# Performance tested on a dataset from Raveh et al.

| PDB ID Holo/Apo | Sequence | Atoms | Rotatable Bonds | Residues | Secondary Structure |
|---|---|---|---|---|---|
| 1AWR/2ALF | **HAGPIA** | 80 | 19 | 6 | C |
| 1ER8 | **HPFHLLVY** | 145 | 35 | 8 | C |
| 1N7F/1N7E | **AVTRTYSC** | 124 | 39 | 8 | b+C |
| 1NLN | **QVQSLKRRRCF** | 198 | 62 | 11 | b+C |
| 1NVR/2QHN | **ASVSA** | 61 | 18 | 5 | b+C |
| 1QKZ | A**NGGASGQVK** | 124 | 40 | 10 | C |
| 1RXZ/1RWZ | KS**TQATLERWF** | 193 | 58 | 11 | b+C |
| 1SSH/1OOT | GP**PPAMPARP**T | 154 | 36 | 11 | C |
| 1TW6 | **AVPI** | 62 | 12 | 4 | C |
| 1W9E/1RJ6 | **NEFYF** | 92 | 25 | 5 | b+C |
| 1Z9O | **EDEFYDALS** | 137 | 44 | 9 | C |
| 2C3I/2J2I | **KRRRHPSG** | 147 | 44 | 8 | C |
| 2FGR/2FGQ | **DNWQNGTS** | 116 | 37 | 8 | C |
| 2FNT | **RQVNFLG** | 120 | 34 | 7 | C |
| 2J6F | **PPKPRPRR** | 152 | 38 | 8 | C |
| 2O9V/2O9S | V**PPPVPPPPS** | 144 | 25 | 10 | C |
| 2P1K | **SATSAKATQTD** | 148 | 50 | 11 | b+C |
| 2VJ0/1B9K | **FEDNF**VP | 106 | 27 | 7 | C |
| 3D1E | **GQLGLF** | 91 | 25 | 6 | C |

Raveh B et al. *PloS One.* **2011**; *6*: e18934.

SCHRÖDINGER.

# Regular Glide Performance is Poor

- Metric of success: iRMSD of any of top 10 poses < 2.0 Å

  iRMSD: RMSD of peptide backbone atoms within 8 Å of protein

- Only **21%** of systems have an accurate pose (iRMSD ≤ 2.0 Å) within top 10 ranked poses by Glide SP (as compared to **63%** with Rosetta FlexPepDock)

- α-helical polypeptides not considered (ConfGen does not generate such conformations)

SCHRÖDINGER.

# Optimized SP-PEP parameters improve results

| Parameters | ConfGen | Rough Scoring | Refinement | Minimization | Final Pose |
|---|---|---|---|---|---|
| **Glide Default** | 17 | 10 | 9 | 5 | 4 |

24 experiments

| | | | | | |
|---|---|---|---|---|---|
| **SP-PEP** | 17 | 11 | 10 | 8 | 7 |

Number of systems with at least one iRMSD < 2.0 Å pose

**SP-PEP Parameters:** 10 conformers generated using ConfGen, dock each conformer using Glide SP

Tubert-Brohman I et al. *Chem. Inf. Model.* 2013; 53(7): 1689-99

SCHRÖDINGER.

# Classification of Complexes Based on Accuracy

| PDB | Highest ranking of accurate pose | PDB | Highest ranking of accurate pose | PDB | Highest ranking of accurate pose | PDB | Highest ranking of accurate pose |
|-----|-----|-----|-----|-----|-----|-----|-----|
| 1N7F | 2 | 1AWR | 1 | 2J6F | 321 | 1QKZ | 92 |
| 1NLN | 1 | 1ER8 | 5 | 2O9V | 13 | 1RXZ | - |
| 1NVR | 1 | 1SSH | 1 | | | 1Z9O | - |
| 1TW6 | 1 | 1W9E | 18 | | | 2C3I | - |
| 2FNT | 1 | 1P1K | 9 | | | 2FGR | - |
| 3D1E | 1 | 1VJ0 | 14 | | | | |

| **Easy** | **Medium** | **Hard** | **Very Hard** |
|---|---|---|---|

SCHRÖDINGER.

# Conclusions

| | Standard SP | SP-PEP | SP-PEP + MMGBSA | Rosetta FlexPepDock |
|---|---|---|---|---|
| % cases where top 10 iRMSD < 2Å | 21% | 41% | 58% | 63% (but 100x slower) |

- ConfGen performed well finding <2Å RMSD pose in 100% of cases

- α-helical peptides cannot be docked with Glide

- Performance best on short, extended, non-ionizable peptides

- More work is needed to achieve small-molecule like accuracy

SCHRÖDINGER.

# Hit to Lead

Residue Scanning, Affinity Maturation

SCHRÖDINGER.

# Residue Scanning: Overview



Original Sequence

Analogs

Measure Δ

Affinity
Stability
pKa
Etc.

- Used to determine what effect specific amino acid positions have on properties such as binding, stability, etc.

- Tells us what mutations may be beneficial, and what may be harmful

- Can be a very laborious and difficult task to do in the lab.
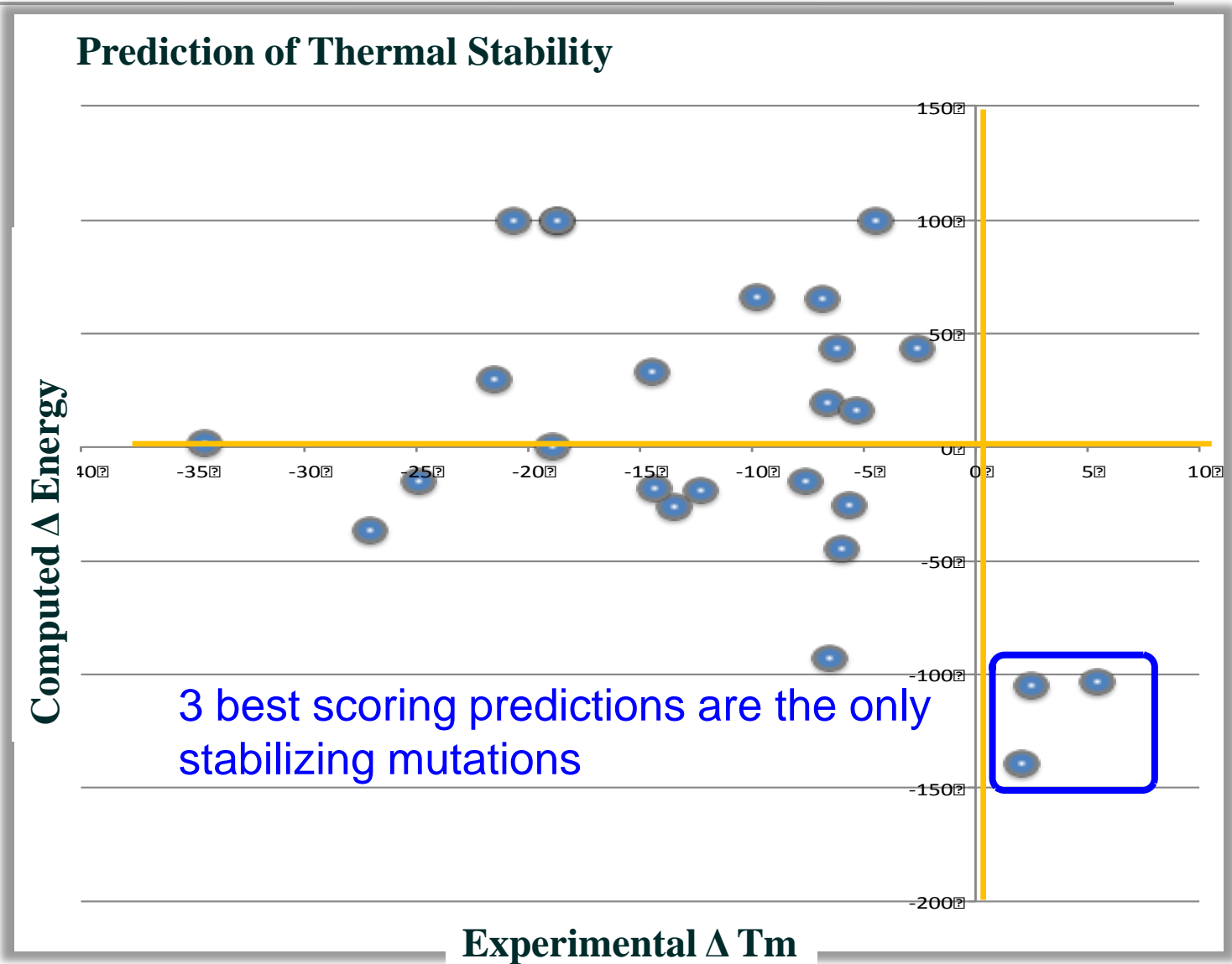
SCHRÖDINGER.

# Residue Scanning in BioLuminate



- Select any protein residues to be mutated

- Run time ~30sec/mutation

- See how properties change

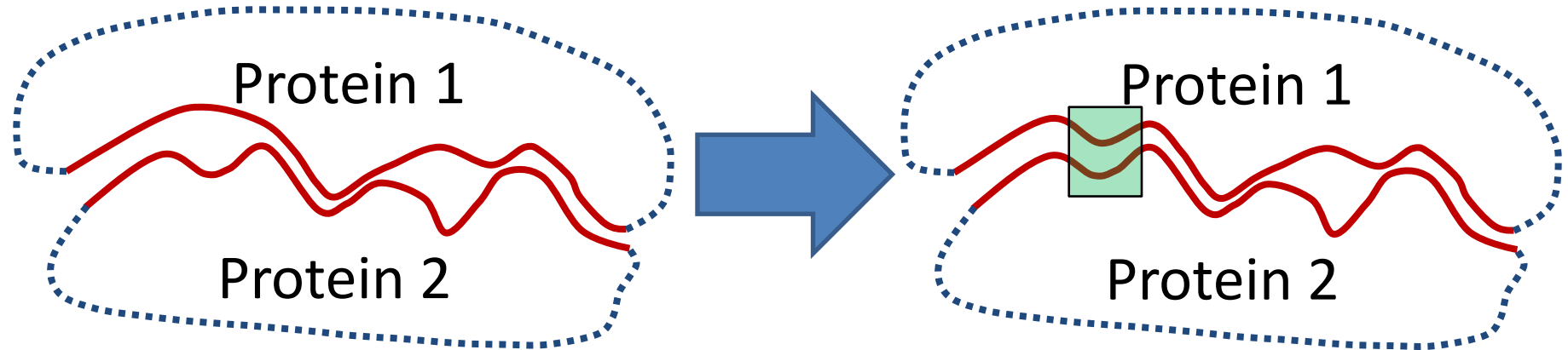  – Affinity, stability, hydrophobicity, SASA, etc.

# Prospective Example: Thermal Stability of SH3 Domain Mutants

- 2 mutation locations
  - Glu107
  - Ser124

- 25 mutations made and tested experimentally

- Only 3 mutations lead to increased thermal stability
  - E107D
  - S124K
  - S124R

- Residue scanning IDs these 3 mutations



**Prediction of Thermal Stability**

Computed Δ Energy

Experimental Δ Tm

3 best scoring predictions are the only stabilizing mutations
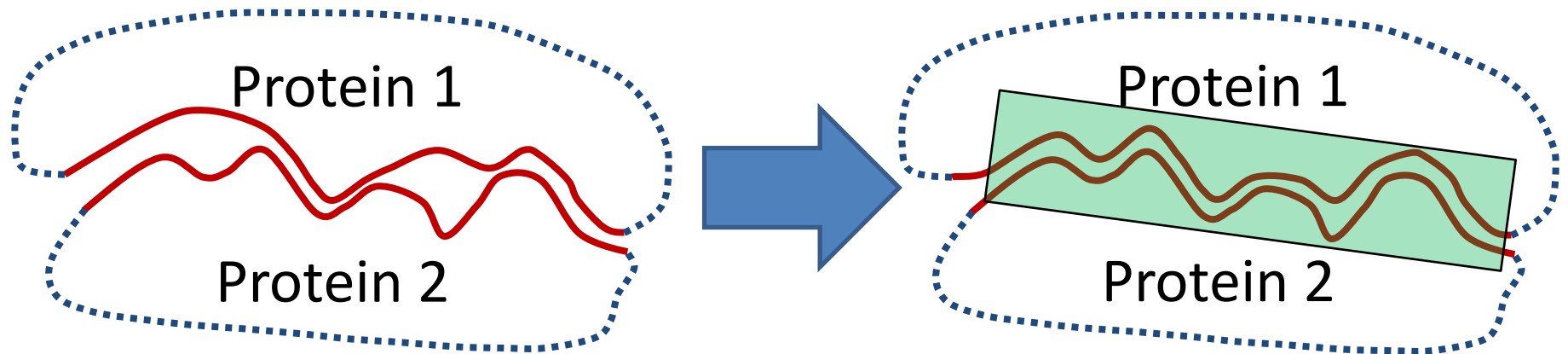
SCHRÖDINGER.

# From Residue Scanning to Affinity Maturation

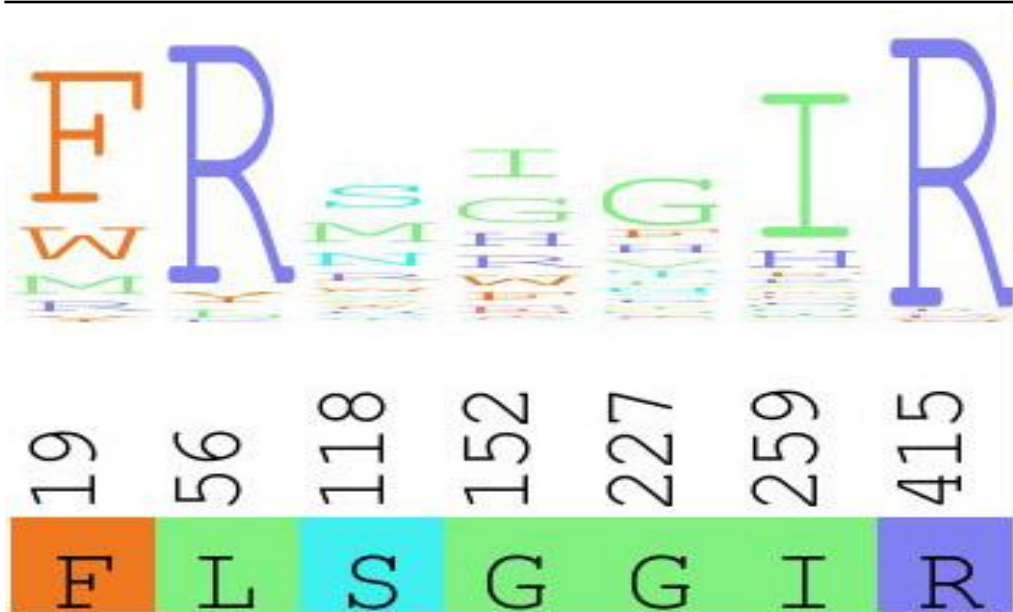Residue Scanning (single mutations):



Affinity Maturation/Protein Design (multiple simultaneous mutations):



SCHRÖDINGER.

# Affinity Maturation in BioLuminate

- Search multiple residue pos simultaneously for changes

- Use to suggest new sequences, or to influence random library design

# Lead Optimization

Peptide QSAR

SCHRÖDINGER.

# What is QSAM modeling?

- In traditional QSAR modeling, structural features of biomolecules are used to develop models for activity
  - i.e. Activity = $f$ (molecular structure)

- QSAM stands for **Q**uantitative **S**equence **A**ctivity **M**odeling:
  - As compared to small molecule QSAR approaches, QSAM models **sequence** information directly using **sequence descriptors**
  - i.e. Activity = $f$ (peptide sequence)

SCHRÖDINGER.

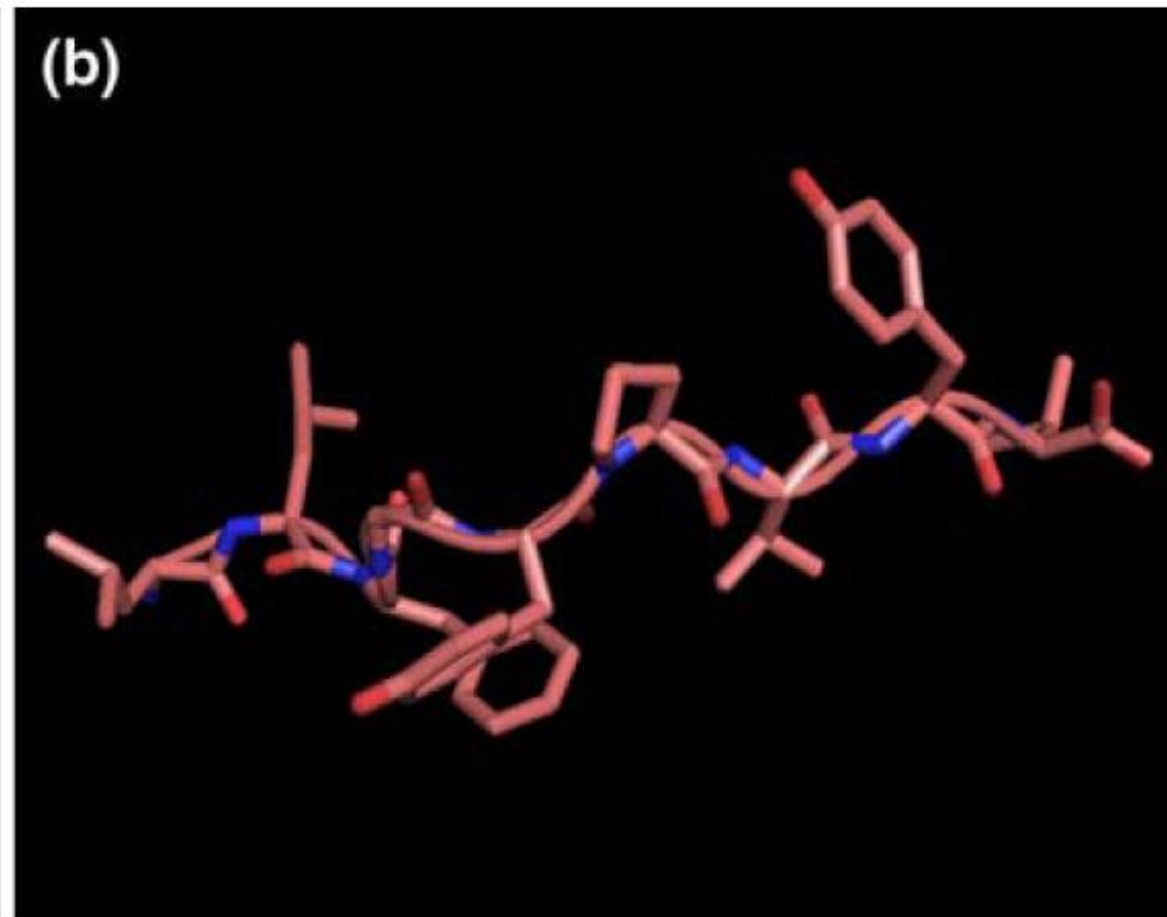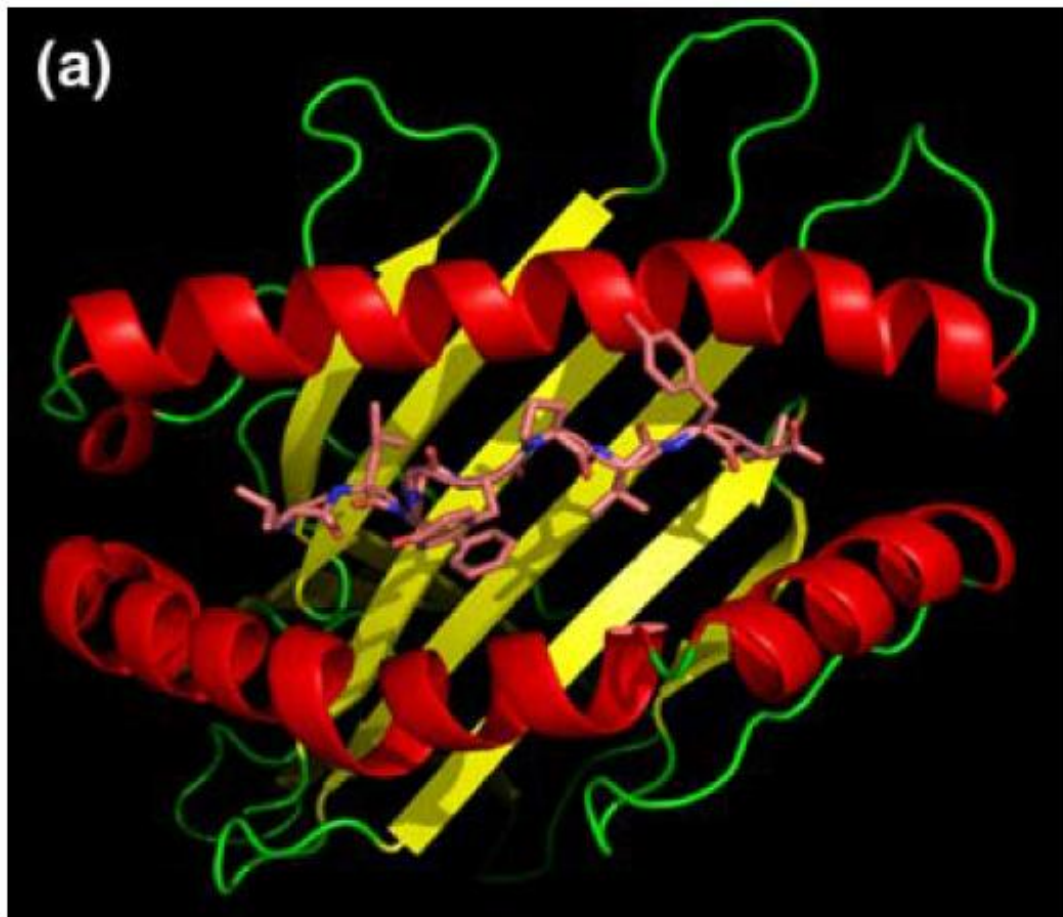# Sequence descriptors are similar to molecular descriptors

- They are based on physicochemical properties of the individual amino acids that comprise the sequence
  - i.e. size, shape, charge, etc
- Three Examples:
  - Zvalue: derived from principle components analysis (PCA) of 29 physicochemical properties of the 20 natural AAs
    - Hellberg et al. *J Med Chem*. 1987; 30: 1126-1135.
  - EZvalue: derived from principle components analysis (PCA) of 29 physicochemical properties for 87 AAs (natural and modified)
    - Sandberg et al. J *Med Chem*. 1998; 41: 2481-2491
  - DPPS: 10 score vectors derived from PCA of 109 properties of the 20 natural Aas
    - Properties include 23 electronic properties, 37 steric properties, 54 hydrophobic properties and 5 hydrogen bond properties
    - Tian et al. *Amino Acids*. 2009; 36: 535-554

SCHRÖDINGER®

# QSAM Modelling: Pros and Cons

- Pros:
  - **Very quick calculation**
    - There is no need for any sort of 3D-structure
      - And certainly no requirement for alignment/docking
    - Can be used to filter through large lists of sequences very rapidly
- Cons:
  - **Immediate interpretation is difficult**
    - The underlying amino-acid descriptors do have physical interpretability
      - Theoretically it is possible to understand what residues are required at each position*
  - **Success depends on having descriptors for each amino acid present**
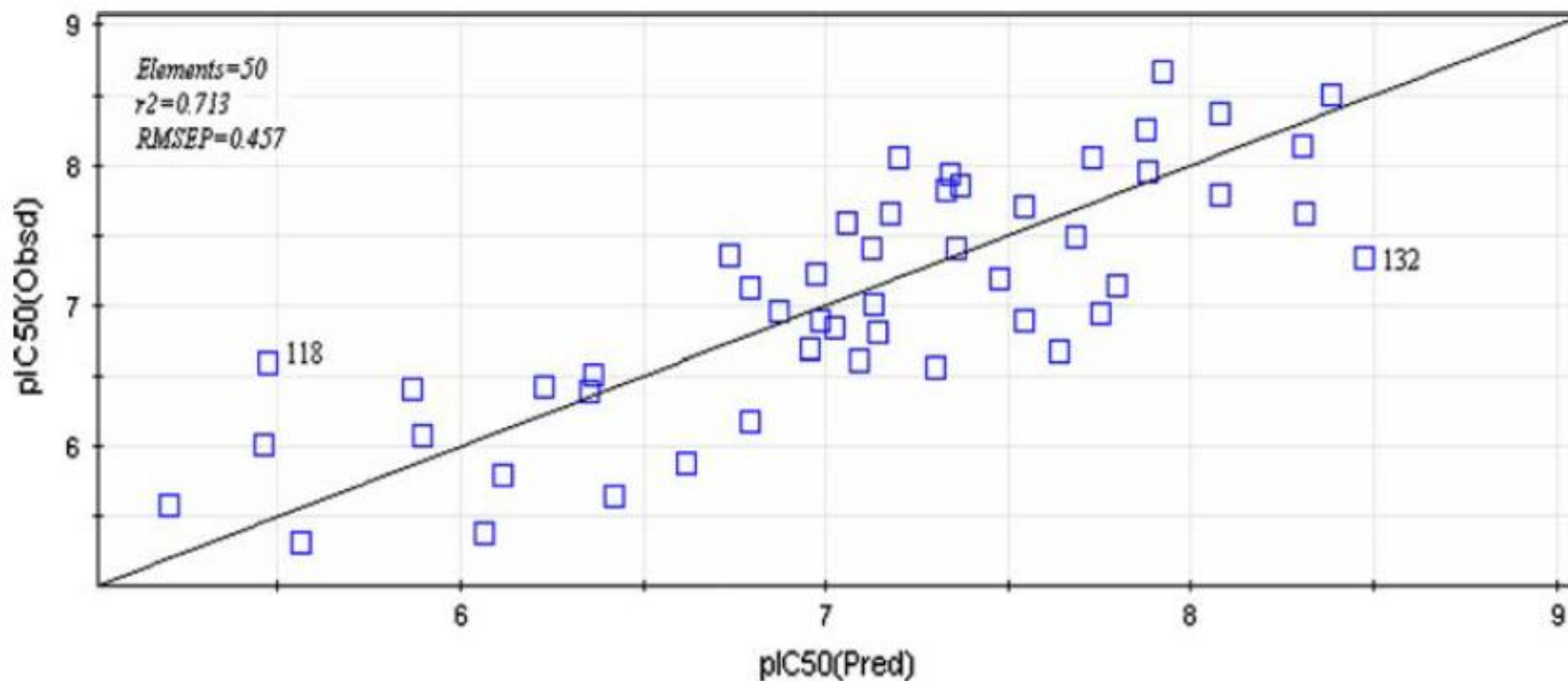    - Handling un-natural amino-acids can be difficult

\* This is as much a limitation of the underlying Canvas PLS implementation as it is of QSAM and the Bioluminate panel. More advanced PLS tools (e.g. Umetric's SIMCA package) would enable a more detailed analysis to be performed.

SCHRÖDINGER.

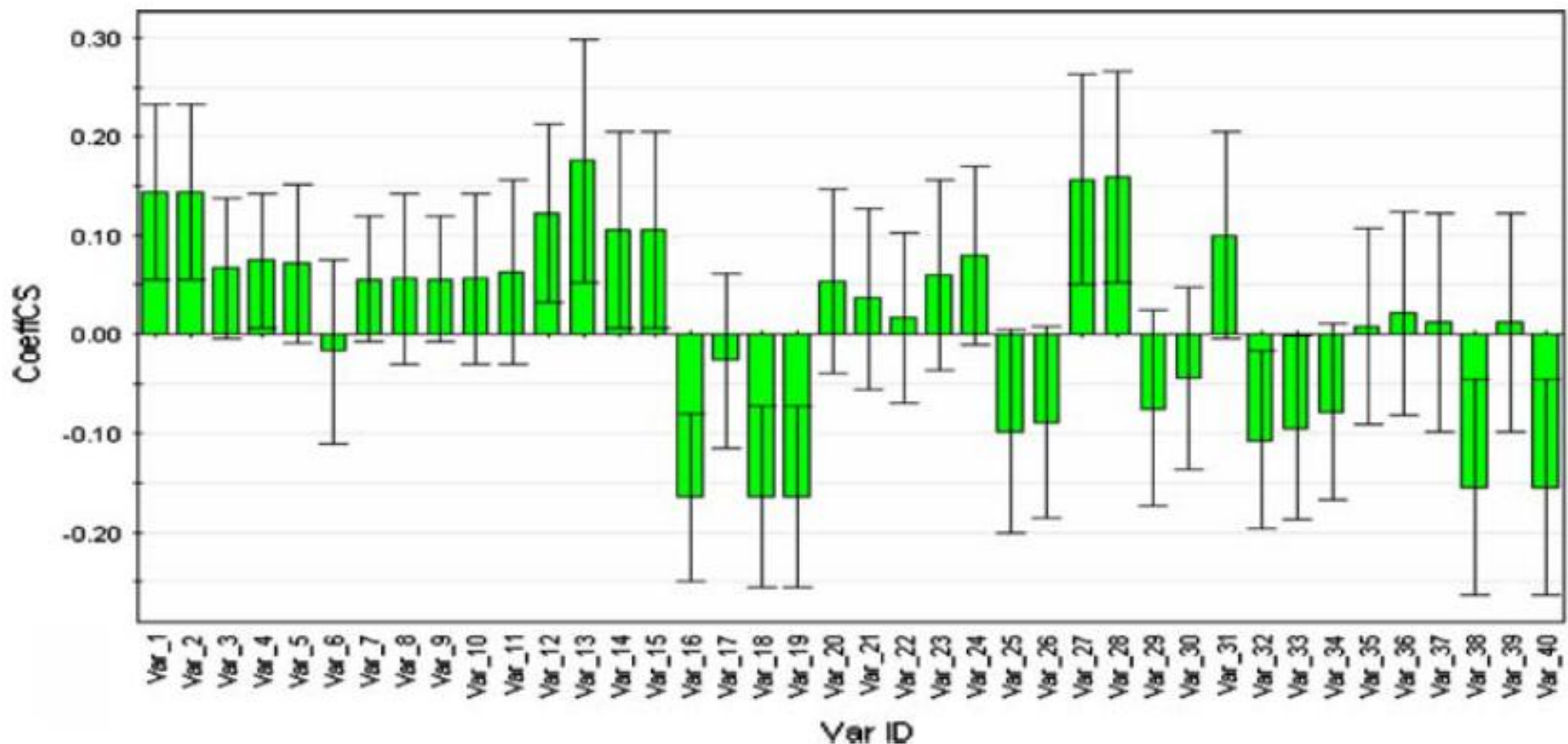# Example: Modeling antigenic peptide binding to MHC



Tian et al. *Amino Acids*. 2009; 36: 535-554

SCHRÖDINGER.

# The model performs very well



- The model was derived using the DPPS descriptor on a dataset of 152 sequences
- Partial least squares (PLS) regression was used to generate the model

Tian et al. *Amino Acids*. 2009; 36: 535-554

SCHRÖDINGER.

# But physical interpretation of the model is tricky ...



- Standardized coefficients of 40 selected variables from the model.
- Each variable corresponds to a peptide sequence position.

Tian et al. *Amino Acids*. 2009; 36: 535-554

SCHRÖDINGER.

# Acknowledgement

SCHRÖDINGER.