

# Mining PDB Subpockets to Rebuild Ligand Binding Conformations

FELIX

Number of PDB

structures:

**EXAMPLE OF SUCCESFUL** 

LIGAND

RECONSTRUCTION

MED-SuMo in combination with FcLigand

software were able to successfully rebuilt vinyl

sulfone inhibitor, a ligand of cathepsin K (ID:

1MEM) by hybridization of retrieved 3D-

fragments including a substructure constraint

to the inhibitor. Four fragments which are

associated with different proteins in the PDB

(cathepsin K, falcipain-3, cruzain protein,

cathepsin S) were necessary for this

reconstruction. Top panel: the four fragments

hybridization. Lower panel: The hybrid

molecule (in colour) and vinyl sulfone inhibitor

(in grey) in the binding site of cathepsin K.

onto their targets

before

📴 🖫 ▼ 💢 🖫 | Carbons ▼ 📴 ▼ | >— >— >> 🐿 🔕 | 🚳 🕉

📂 🖫 ▼ 💢 🖳 📳 | Carbons ▼ 📴 ▼ | >— >— >> 🐿 🔕 | 🚳 🕉

Jean-Yves Trosset<sup>1</sup>, Maud Vieillevoye<sup>2</sup>, Stewart Adcock<sup>2</sup>, François Delfaud<sup>2,3</sup>

<sup>1</sup>BIRL Sup'Biotech, 66 rue Guy Môquet, 94800 Villejuif, France <sup>2</sup>Felix Concordia SARL, 400 av de Roumanille, Bat 7, 06906 Sophia Antipolis, France <sup>3</sup>Medit SA, 2 rue du Belvédère, 91120 Palaiseau, France

> Contact: fdelfaud@felixc.eu http://www.felixc.eu/



#### **ABSTRACT**

We constituted a database of subpocket-fragment complexes (Fragmentor) through deconvolution of each PDB (Protein Data Bank) ligand in all possible fragments that match one of the chemical molecule contained in the Pubchem database. After application of a Matriochka filter, we obtained a set of 28.482 2D-fragments and 398.236 3D-fragments (PDB conformations). Subpockets were defined as protein surfaces located at 4,5 Å around fragments.

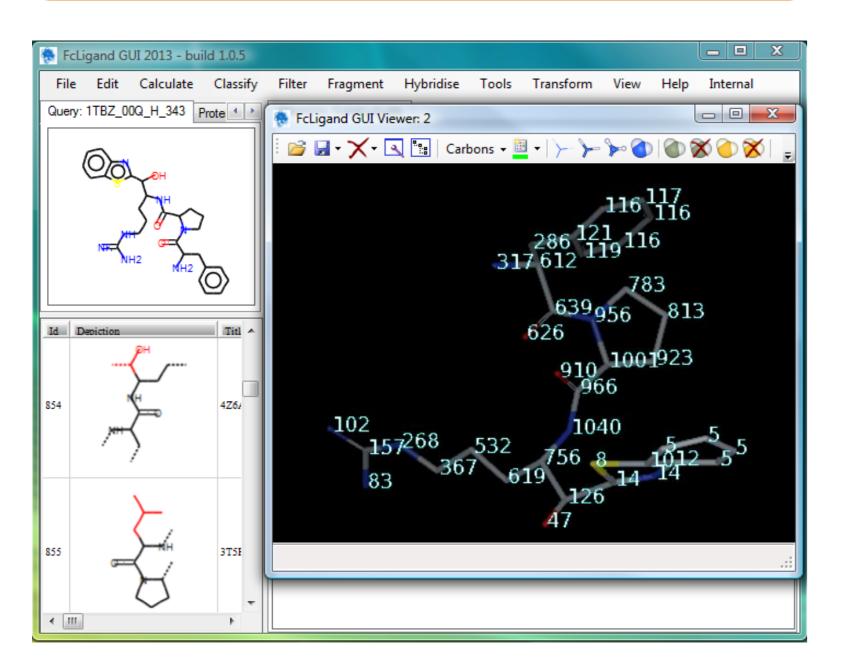
The goal of this work was to determine if there is enough information in the PDB to successfully rebuilt the binding mode of ligands, starting from the target protein structure. To test this hypothesis, we selected ligands in unique PDB entry to build a test set of 2292 protein-ligand complexes and tried to recover the position and conformation of the ligands using our Fragmentor database. We were able to predict at least 80% of the 3D-structure from 1091 ligands (48%). In conclusion, this study highlights the quality of the information contained in the PDB and supports the use of its structural information for docking tools or fragment-based drug design.

### **WORKFLOW**

1 Compare the binding site of the target protein with Fragmentor database

2 Retrieve subpockets and associated 3D-fragments that superpose onto the binding site of the target

3 Calculate how many times the atoms of the reference ligand are recovered at a deviation less than 1Å with a common substructure of collected 3D-fragments



Pairwise comparisons of subpockets were **MEDP-SiteClassifier** out with the carried software, an evolution of MED-SuMo software. The left column shows the overlapping chemical moieties (in black) between selected 3Dfragments and the reference ligand (on top); a MCS (Maximum Common Substructure), coupled a 3D deviation detection, turns non overlapping atoms in red. A javascript in FcLigand viewer illustrates the number of times the atoms of the reference ligand were recovered with overlapping 3D-chemical fragments.

### of 2D 3000

5000 **[101-10'000] [11-100]** 2000 **[2-10]** 1000 **[1]** 11 12 13 14 15 16

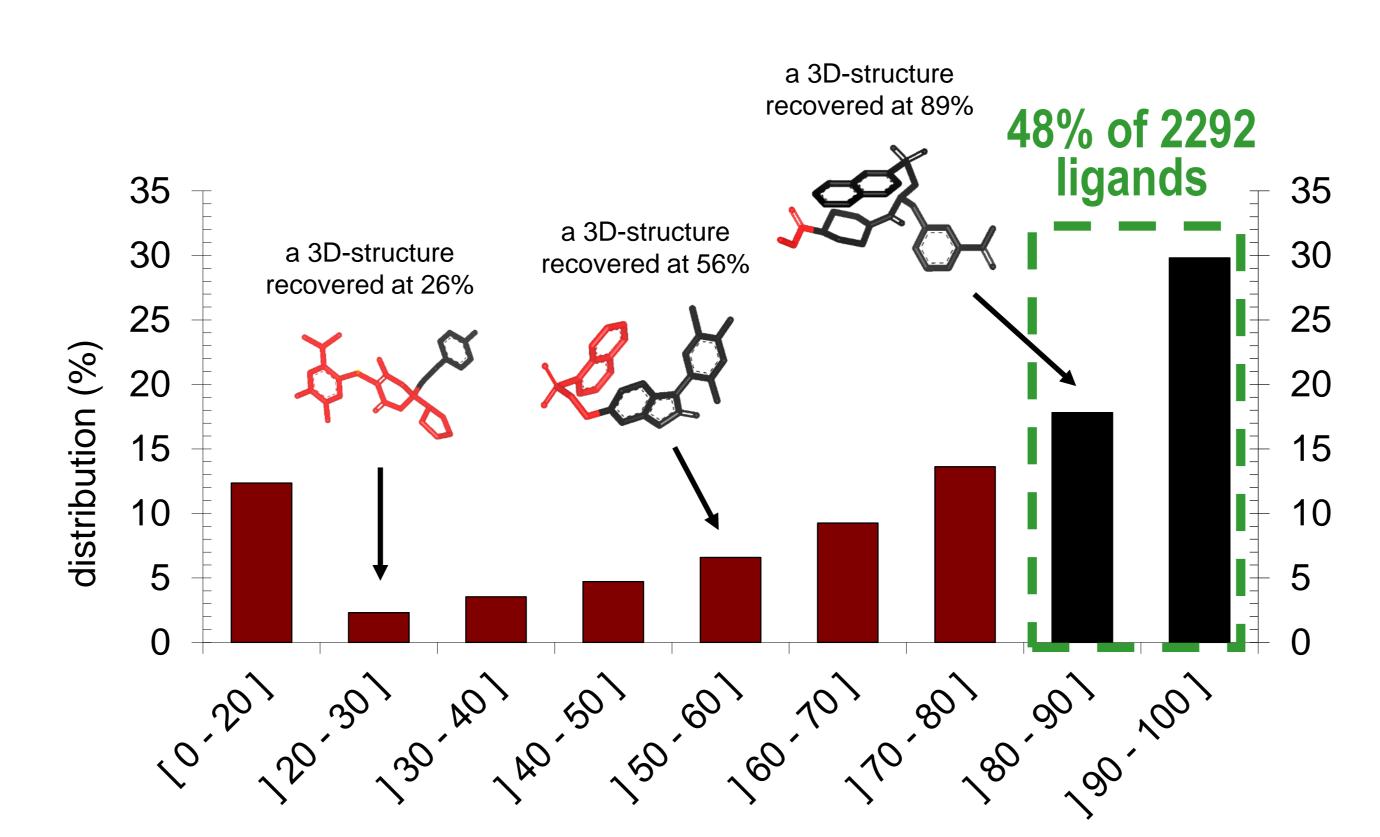
fragment size (heavy atom count)

FRAGMENTOR: A DATABASE CONTAINING

398.236 SUBPOCKET-FRAGMENT COMPLEXES

We constituted a database of subpocket-fragment complexes derived from PDB data. 40% of the 28.482 2D-fragments are present only in one PDB file (grey columns) whereas the other 60% can be retrieved from multiple PDB files (blue columns).

### PERCENTAGE OF 3D-STRUCTURE RECOVERY **OF 2292 LIGANDS:** 48% CAN BE ALMOST RECONSTRUCTED



% of 3D-ligand structure recovery by overlapping fragments

The analysis of 2292 protein-ligand complexes reveals that the PDB contains sufficient structural information to rebuilt 48% of ligands with at least 80% of their 3D-structure. Three ligands are depicted to illustrate the graph (in red: atoms which were not predicted).

#### INFLUENCE OF PFAM AFFILIATION ON LIGAND RECOVERY

0,48

The horizontal axis lists the 15 most frequent Pfam families present in the PDB. The grey curve (right axis) indicates the total number of structures that belong to the considered Pfam family in the PDB. The complexes for which the ligand was well predicted are represented by red columns

**PFAM DESCRIPTION COLUMNS COLUMNS** mmunoglobulin C1-set domain 0,43 mmunoglobulin V-set domain 0,55 0,57 Protein kinase domain PF00089 1,00 Proteasome subunit 1,00 0,00 0,50 Nitric oxide synthase, oxygenase domain 0,00 Photosynthetic reaction centre protein 0,10 Eukaryotic-type carbonic anhydrase 0,56 0,44 Protein tyrosine kinase 0,14 C-type lysozyme/alpha-lactalbumin PF00062 0,00 0,66 Eukaryotic aspartyl protease 1,00 0,00 Ras family 0,00 1,00 PF00959 Phage lysozyme 0,54 WEIGHTED MEAN

TOTAL MEAN (2292 complexes)

It seems that the quantity of structural data available for the Pfam family of the target does not significantly influence the percentage of recovery of the ligand (success rate: 54% for ligands stemming from the 15 most frequent families versus

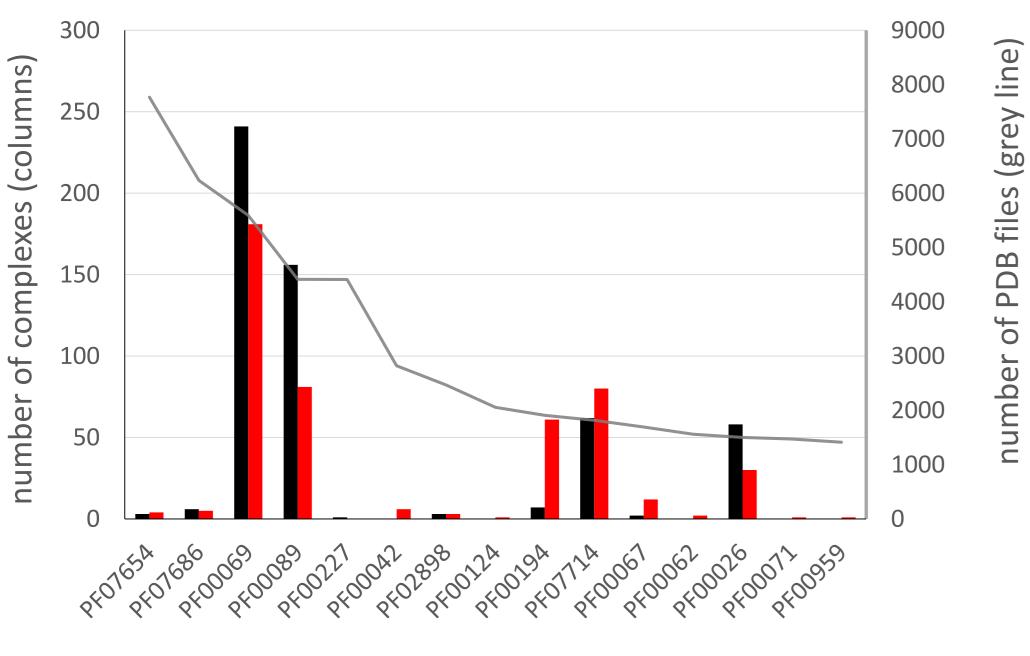
48% for the whole dataset)

whereas the complexes for which the ligand was predicted with less than 80% of its 3Dstructure are represented by black columns (left axis).

## CONCLUSION

superposed

- We generated of a rich subpocket-fragment database from PDB data which was used for the present proof of concept
- 48% of ligands could be recovered by 80% of their 3D-structure using a protocol based on **MEDP-SiteClassifier** software
- The percentage of ligand recovery does not depend on the pfam family of the target
- We were able to successfully reconstruct a vinyl sulfone inhibitor in FcLigand software after hybridization of fragments collected by MED-SuMo.



- Complexes from best recovered 3D-structures (48%) Complexes from less well recovered 3D-structures (52%)
- Entire PDB