Preparing MD-ready receptors for all therapeutic protein targets and beyond

Jesper Sørensen, Ph.D. Head of Biomodeling at OpenEye

CUP XXI, 2022



## The Goal

- Prepare the entire PDB
- Prepare AlphaFold2
- Organize the data

- PDB: 187,844 experiments
- AF2v2-human: 23,391 models
- Target and family classification:
  So many choices...

All data sources as of 3/4/2022

• Mine the data

• So many things to ask



### The first iteration – focus on pharma targets

MMDS	MMDS - Macromolecular Data Service				MMDS - Macromolecular Data Service				
Sessions	Projects Experiments		,7	Sessi	ons Projects	Experiments			
	Q CDK2 Expand Collapse			Protei	n/Ligand IDs:				
					Tags:				
	- RCSB	38762 sites		18003 n	eulte Page 1 of	1801 Next			
	<ul> <li>Enzymes</li> </ul>	28308 sites		103001	sound. I age I of				
	<ul> <li>Kinases (EC 2.7.x.x)</li> </ul>	9200 sites		ID	Experiment	Projects			
	<ul> <li>CMGC: Containing CDK, MAPK, GSK3, CLK families</li> </ul>	1920 sites		6745	5VUN	Neuronal NOS			
	<ul> <li>Cyclin-dependent kinase (CDK) family</li> </ul>	722 sites		6746	241.4	amina avidese containing 2			
	<ul> <li>CDK1 subfamily</li> </ul>	621 sites	ל ד	0740	JALA	anine oxidase copper containing 5			
	CDK2 (cyclin dependent kinase 2)	608 sites	- (	67/7	206/11	Estraran recentor & Job 2 I nuclear recentor coactivator 2			

• 2,673 Targets in the family tree







# Acknowledgements

- Zhe Mei
- Addison Smith
- Bob Tolbert
- David LeBard
- Mike Word
- Stan Wlodek
- Krisztina Boda
- Greg Warren



- Ant Nicholls
- Christopher Bayly
- Gaetano Calabro
- Hyesu Jang
- She Zhang
- Orion/Cloud Teams
- Customer Success Group



### Connecting data sources, compute, and analysis



## SPRUCE – OpenEye's Protein Preparation Tool



- Reads a variety of inputs as PDB/MMCIF files + MTZ Files
  - Experiments: X-ray, NMR, cryoEM
  - Models: Homology models, AI models, one-off models
- Forms the basis for a lot of our tools
  - OEDOCKING, SZMAP, SZYBKI, ZAP, STMD, NES, Enhanced Sampling





# Iridium Categorization

Active site:

Residues within 5A from the ligand Co-factors within 5A from the ligand

Density coverage:

	HT	MT	NT	
Ligand	> 0.90	> 0.5	< 0.5	
Active Site	> 0.95	> 0.5	< 0.5	

- HT → MT
  - DPI > 0.50
  - Alternate conformations in the ligand or active site residues
  - Any ligand heavy atom with an occupancy < 0.90</li>
  - Any active site heavy atom with occupancy < 0.5</li>
  - Packing residues: Interactions with ligand
  - Excipients: Interactions with ligand
  - Covalently bound ligand, using covalent definition from OEInteractionHints

### $* \rightarrow \text{NT}$

Rfree > 0.45 and Resolution < 3.5

Iridium criteria: Warren et al., Drug Disc. Today, 17, 1270 (2012).





Conformation of ligand from chain A is strongly affected by residue from chain B.

Original PDB coordinates Re-modeled in Merck kga dataset





• Protein residue perception improvements with capping group support





Fetch

Data

### • Fall 2021

• Improved Hydrogen placement for Histidine

Flag

issues

Improved Tautomer selection for ligands and co-factors

Gen

BU(s)

Split

system

- Added Tail modeling (same database approach)
- SiteHopper Toolkit released
- SiteHopper as a superposition method (joined Seq, Site Seq, DDM, SSE)

Prep

SC,

Loops

Caps

Prep

Place H

- Spring 2022 upcoming
  - Biounit extraction performance
  - Logging/messaging prep issues to customers (pre-filter and post prep)



### Tail Modeling

Super

pose



Eval

Do

work





## The Goal

- Prepare the entire PDB
- Prepare AlphaFold2
- Organize the data

- PDB: 187,844 experiments
- AF2v2-human: 23,391 models
- Target and family classification:
  - Guide to Pharmacology
  - ...

• Mine the data

• So many things to ask





### **Good Data Science Practices**

- Versioned workflows to access and prepare data
- Versioned datasets
- Consistent generation & update of datasets
- Reproducible science, reproducible validation
- Democratization of good science





### MMDS 1. Make RCSB PDB Collection

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

MMDS 2. Generate Guide to Pharmacology Datasets

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

MMDS 3. Target Reference Picker Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

MMDS 4. Family reference picker Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

MMDS 5. Structure Prep Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

#### MMDS 6. Add family data to MMDS

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

#### MMDS 7. Add context data to MMDS

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

#### MMDS 8. Add structures to MMDS

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

#### MMDS 9. Add receptors to MMDS

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description

#### MMDS 10. Add sequence alignments to MMDS

Package: OpenEye Biomodeler Floes v0.1.1b16 Show description



### The Floes

### and data organization





## **Updateable Source Collection**

- Orion storage of each PDB/MMCIF and MTZ file
- The workfloe makes the data updatable
  - New revisions of PDBs, Removal of deprecated codes
  - New versions of AlphaFold models
- Currently 764 GB of data
- EM Maps inclusion is planned



## Data organization

- No convention on "organized" n
- A target can be organized by
  - Class (e.g. enzymes, receptors, trai
  - Evolutionary relationship (i.e. sequ
  - Therapeutic area (e.g. cancer, neu
  - ... your scheme here ...



- We have some great public databases
  - PDB at RCSB, Alphafold and UniprotKB at EBI
  - Lots of cross-links and annotations per structure
  - They are indexed, but not "organized"



# Target definition

- Targets are defined in a hierarchy
- Each target has a reference structure
  - Defines the biological unit/form
  - Defines the frame of reference
- Each family (parent) can have a ref
  - If each child below can superpose
- Automated reference structure picking
  - Based on consensus of structures for the target

### MMDS - Macromolecular Data Service

Session

5	Projects	Experiments			
	Q CDK		▼ Expand	▲ Collapse	
	- RCSB				38762 sites
	- Enzy	mes			28308 sites
	<b>–</b> K	inases (EC 2.7.x.x)			9200 sites
	-	CMGC: Containing CDK, MAPK, GSK3, C	LK families		1920 sites
		+ CLK family			108 sites
		<ul> <li>Cyclin-dependent kinase (CDK) family</li> </ul>			722 sites
		CCRK subfamily			
		CDK10 subfamily			
		<ul> <li>CDK1 subfamily</li> </ul>			621 sites
		CDK1 (cyclin dependent kinase	1)		13 sites
		CDK2 (cyclin dependent kinase	2)		608 sites
		<ul> <li>CDK4 subfamily</li> </ul>			10 sites
		CDK6 (cyclin dependent kinase	6)		10 sites
		<ul> <li>CDK5 subfamily</li> </ul>			10 sites
		CDK5 (cyclin dependent kinase	5)		10 sites
		<ul> <li>CDK7 subfamily</li> </ul>			4 sites
		CDK7 (cyclin dependent kinase	7)		4 sites
		<ul> <li>CDK8 subfamily</li> </ul>			32 sites
		CDK8 (cyclin dependent kinase	8)		32 sites
		<ul> <li>CDK9 subfamily</li> </ul>			26 sites
		CDK9 (cyclin dependent kinase	9)		26 sites
		<ul> <li>CRK7 subfamily</li> </ul>			16 sites
		CDK12 (cyclin dependent kinase	∋ 12)		14 sites
		CDK13 (cyclin dependent kinase	ə 13)		2 sites
		PITSLRE subfamily			
_		<ul> <li>TAIRE subfamily</li> </ul>			3 sites
)		CDK16 (cyclin dependent kinase	e 16)		( 3 sites )



### Problems in first iteration...

Initiato

Chain

Chain Peptid Chain

Peptic Chain Chain

Chain Chain Chain Chain

- Targets not in GtoP
  - Viral targets\*
  - Bacterial targets
  - e.g. Haemoglobin

are processing				
r methionine <sup>i</sup>		Removed; by host 🧳 By similarity		
(PRO_0000261261)	2 - 1447	Gag-Pol polyprotein	🏦 Add 🔧 BLAST	1446
(PRO_0000042285)	2 - 132	Matrix protein p17 🔗 By similarity	📾 Add 🔧 BLAST	131
(PRO_0000042286)	133 - 363	Capsid protein p24 🕜 By similarity	🏦 Add 🔧 BLAST	231
e <sup>i</sup> (PRO_0000042287)	364 - 377	Spacer peptide 1 🕜 By similarity	📾 Add 🔧 BLAST	14
(PRO_0000042288)	378 - 432	Nucleocapsid protein p7 🔗 By similarity	📾 Add 🔧 BLAST	55
e <sup>i</sup> (PRO_0000246710)	433 - 440	Transframe peptide 🕜 Sequence analysis		8
(PRO_0000042289)	441 - 500	p6-pol 🔗 Sequence analysis	📾 Add 🔧 BLAST	60
(PRO_0000038647)	501 - 599	Protease 🕜 By similarity	🏦 Add 🔧 BLAST	99
(PRO_0000042290)	600 - 1159	Reverse transcriptase/ribonuclease H 🔗 By similarity	📾 Add 🔧 BLAST	560
(PRO_0000042291)	600 - 1039	p51 RT 💊 By similarity	📾 Add 🔧 BLAST	440
(PRO_0000042292)	1040 - 1159	p15	📾 Add 🔧 BLAST	120
(PRO_0000042293)	1160 -	Integrase 🛛 By similarity	🗎 Add 🔧 BLAST	288

Gag Polyprotein: P04585

- Problems mapping using UniprotKB IDs
  - Viral targets\*

	En En	ntry 🗢	Entry name 🗢		Protein names 🖨 🛛 🔊	Gene names 🔷	Organism 🗢	Length 🗘
C	] P0	)4585	POL_HV1H2	☆	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype B (isolate HXB2) (HIV-1)	1,435
	P0	03369	POL_HV1A2	<b>☆</b>	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype B (isolate ARV2/SF2) (HIV-1)	1,437
	<b>P0</b>	05961	POL_HV1MN	<mark>☆</mark>	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype B (isolate MN) (HIV-1)	1,441
	Q9	9QSR3	POL_HV1VI	<b>☆</b>	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype F1 (isolate VI850) (HIV-1)	1,430
	Q9	9Q720	POL_HV1V9	<b>☆</b>	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype H (isolate VI991) (HIV-1)	1,436
	Q7	75002	POL_HV1ET	<b>☆</b>	Gag-Pol polyprotein	gag-pol	Human immunodeficiency virus type 1 group M subtype C (isolate ETH2220) (HIV-1)	1,439
	08	νοσου	DU H1/103	<b>.</b>	Gan-Pol nolyprotein	dad-bol	Human immunodeficiency virus type 1 group M subtype E1 (isolate Q3RD020) (HTV-1)	1 430



## **UniProt Features with Sequence Splitting**

Sequence splitting occurs when a protein is split into non-overlapping features. This is easily seen in the HIV target Gag Polyprotein (P04591).

### UniProtKB - P04591 (GAG\_HV1H2)

Feature key	Position(s)	Description	Actions	Graphical view	Length	
Initiator methionine <sup>i</sup>		Removed; by host 🕜 By similarity				
Chain <sup>1</sup> (PRO_0000261216)	2 - 500	Gag polyprotein	🏦 Add 🔧 BLAST		499	
Chain <sup>1</sup> (PRO_0000038593)	2 - 132	Matrix protein p17 🛷 By similarity	🏦 Add 🔧 BLAST		131	
Chain <sup>1</sup> (PRO_0000038594)	133 - 363	Capsid protein p24 🔗 By similarity	🏦 Add 🔧 BLAST		231	
Peptide <sup>1</sup> (PRO_0000038595)	364 - 377	Spacer peptide 1 🕜 By similarity	🏦 Add 🔧 BLAST		14	
Chain <sup>1</sup> (PRO_000038596)	378 - 432	Nucleocapsid protein p7 🕜 By similarity	🏦 Add 🔧 BLAST	-	Uncategorized	
Peptide <sup>1</sup> (PRO_0000038597)	433 - 448	Spacer peptide 2 🕜 By similarity	🖮 Add 🔧 BLAST		<ul> <li>Gag polyprotein</li> </ul>	
Chain <sup>1</sup> (PRO_000038598)	449 - 500	p6-gag 🔗 By similarity	🛱 Add 🔧 BLAST		Capsid protein p27, alternate cleaved 2	
					Phosphorylated protein pp24	
					Capsid protein p24	
UniprotKB "targe	et" is brok	ken into			Gag protein	
smaller feature-t	argets th	at			Matrix protein p10 Matrix protein p15	
represent the druggable targets				Matrix protein p17		
in HIV.					Protease p15	
				Updates the fam	ily and	
				target tree in MN	ADS <b>ODENEUE</b>	

### **Complex Target Features and Inconsistencies**

Combining both sequence splitting and sequence hierarchy can give much greater protein detail. Inconsistencies in UniProt submissions inappropriately convolutes.

The Gag-Pol polyprotein and Myeloperoxidase targets are complicated examples where sequence splitting occurs inside sequence hierarchies of variable depths.



Description	Actions	Graphical view
Ubiquitin-related	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🗃 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	
Ubiquitin	🏦 Add 🔧 BLAST	

In the case of Ubiquitin, variations in UniProt submission produce splits that do not actually contribute to uniqueness, and generate redundancies that complicate its utility.



## Problems in first iteration...

- Targets not in GtoP
  - Viral targets\*
  - Bacterial targets
  - e.g. Haemoglobin
- Problems mapping using UniprotKB IDs
  - Viral targets\*
  - UniprotKB IDs map to multiple structures, e.g. macromolecular assemblies





## Problems in first iteration...

- Targets not in GtoP
  - Viral targets\*
  - Bacterial targets
  - e.g. Haemoglobin
- Problems mapping using UniprotKB IDs
  - Viral targets\*
  - UniprotKB IDs map to multiple structures, e.g. macromolecular assemblies
- Target reference issues
  - No liganded structures in a target which "site" is of interest? Big problem for AlphaFold models



### **Current Status**

- Prepare the entire PDB
- Prepare AlphaFold2
- Organize the data

• Mine the data

- 105,779 experiments were prepared
  ~200,000+ design units
- 1,325 models were prepared
  1,325 design units
- Organized in MMDS
  - Guide to pharmacology tree
  - Uncategorized tree
  - 12,728 (4x improvement)
- Let's turn to this next...









### Structure Prep at Scale



### Summary

- Closing in on goal of preparing entire PDB + AlphaFold using Orion
- Organizing the prepared data in MMDS
- Solved some interesting problems with linking databases of experiments to a definition of "target"
- Data mining and learning from this data has just begun



# Thank You

## The End



For more information, please contact:

sales@eyesopen.com info@eyesopen.com

www.openeye.inc

+1-505-473-7385

