

Pre-calculated protein structure alignments at the RCSB PDB site

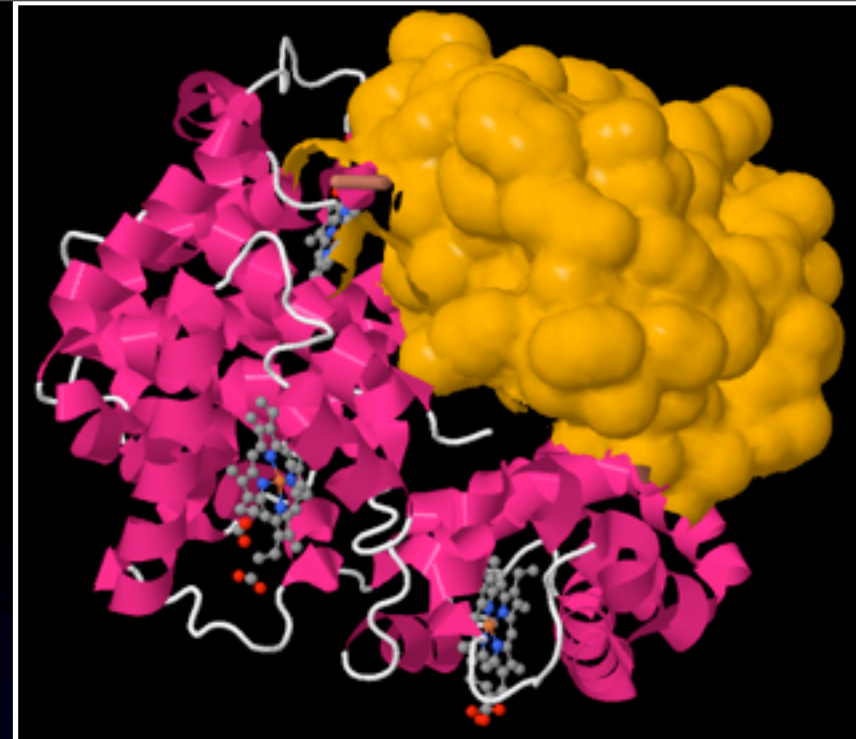
Andreas Prlić






RCSB Protein Data Bank

- Archive of experimentally determined 3D structures of proteins, nucleic acids, complex assemblies
- One of the largest scientific resources in life sciences

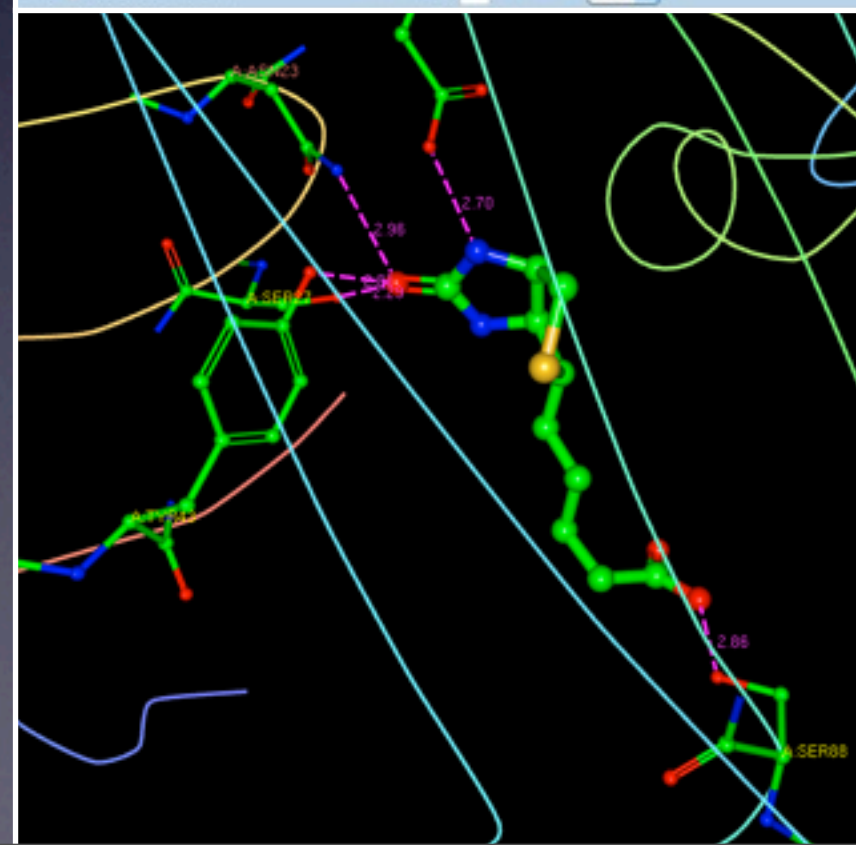
www.pdb.org



Click up/down. Click again to reverse order. Download options:   
Click on column headers to filter the data set. ?

Res Range	Cath Description	SCOP ID [Res Range]	Scop Domain	Scop Fold	PFAM Acc [Res Range]
[1-400]	• Glycosidases	• 21816 [496-581]	• Cyclomaltodextrin	• Immunoglobulin-like	• PF00128 [52-33]
[401-495]	• Golgi alpha-mannosidase	• 22485 [582-686]	• Cyclodextrin glycosyl hydrolase	• Prealbumin-like	• PF02806 [413-413]
[496-582]	• Immunoglobulins	• 27723 [407-495]	• Cyclodextrin glycosyl hydrolase	• Glycosyl hydrolase	• PF01833 [499-500]
[583-685]	• Immunoglobulins	• 28710 [1-406]	• Cyclodextrin glycosyl hydrolase	• TIM beta/alpha-barrel	• PF00686 [586-685]
[2-139]	• Dynein light chain	• 40873	• Profilin (actin-binding)	• Profilin-like	• PF00235 [1-13]
[2-139]	• Dynein light chain	• 40874	• Profilin (actin-binding)	• Profilin-like	• PF00235 [1-13]
[2-10]	• 5' to 3' exonuclease	• 138278 [241-341]	• DinB homolog (DB)	• Lesion bypass DNA polymerase	• PF00817 [12-14]
[78-166]		• 138279 [2-240]	• DinB homolog (DB)	• DNA/RNA polymerase	• PF11798 [173-217]
[167-233]					• PF11799 [217-233]
[244-336]					

Results Customize Columns Page 1 of 1 30

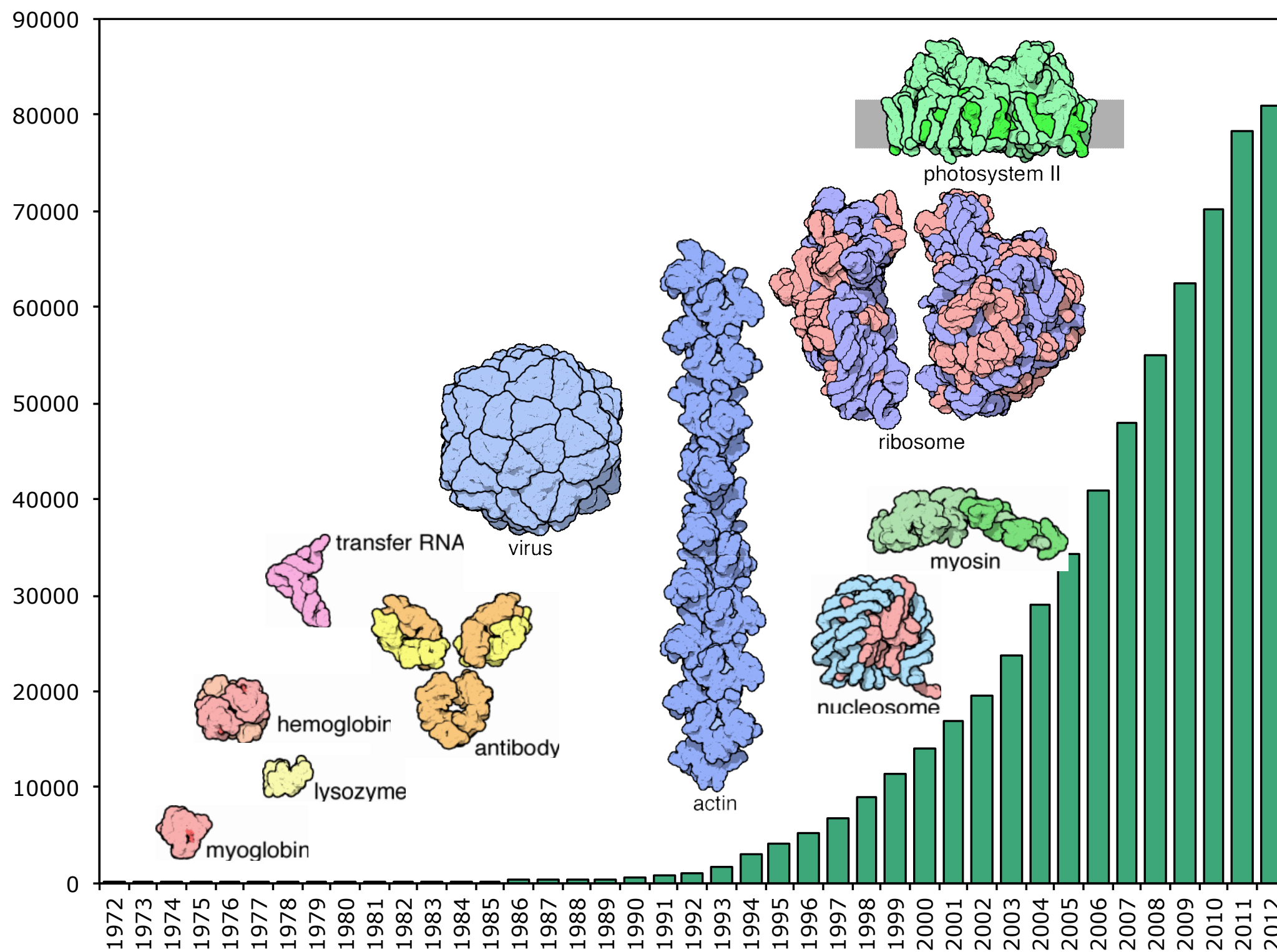


Jmol

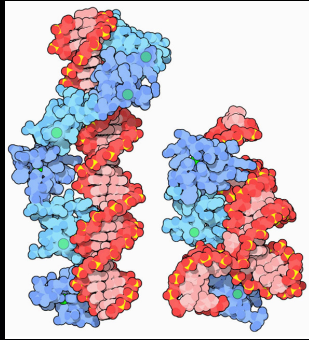
Custom
report

Ligand
Explorer

Number of released entries

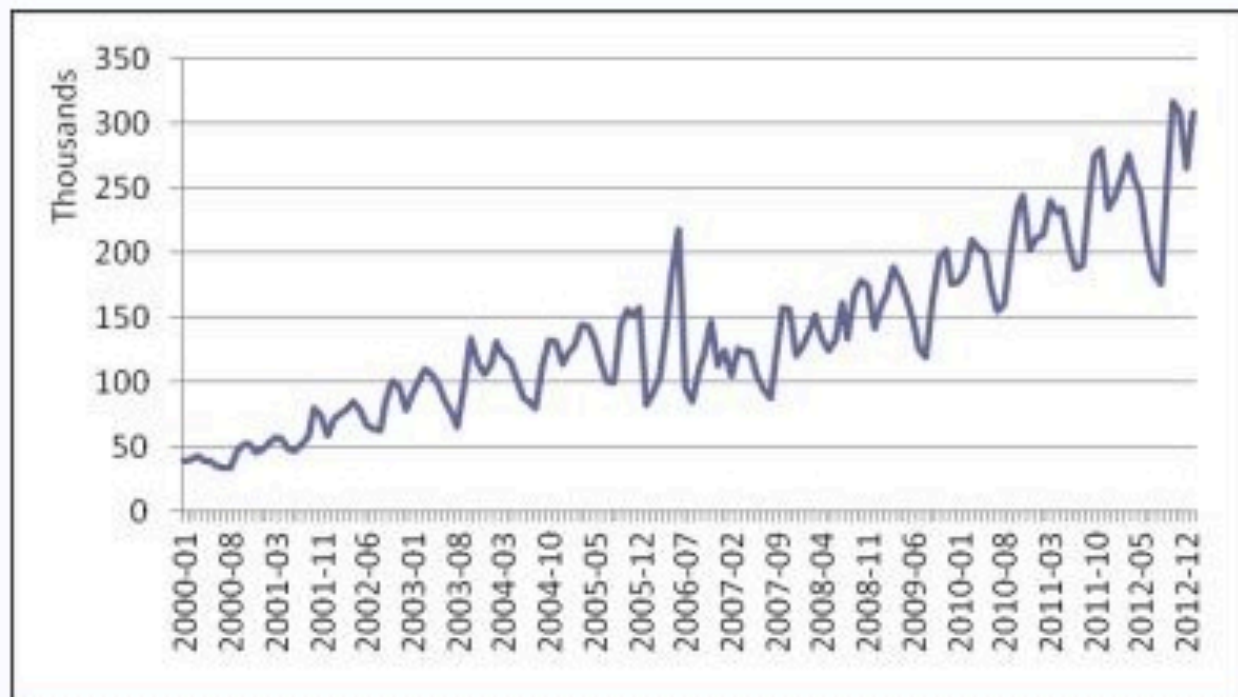


Year

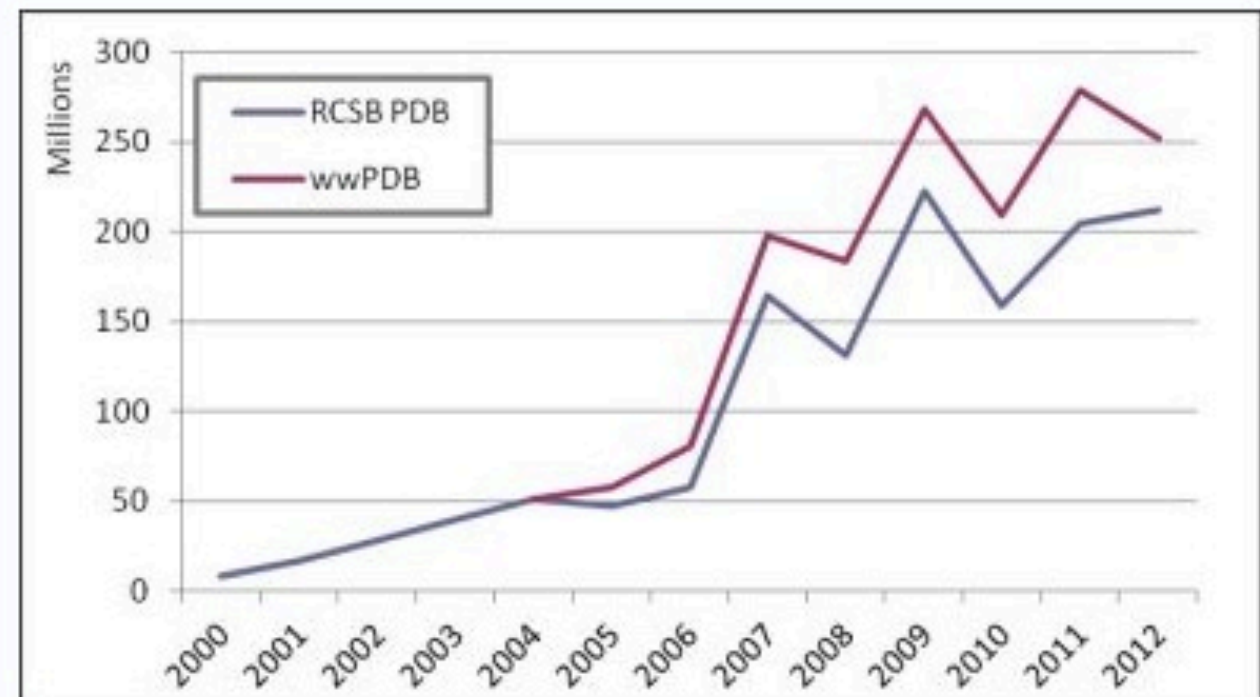


Growth in PDB Usage Over Time

RCSB PDB Website Unique Users by Month



FTP Downloads by Year



- More than 300,000 unique visitors per month
- Up to 300 concurrent users
- ~10 structures are downloaded per second 7/24/365
- Increasingly popular web services traffic

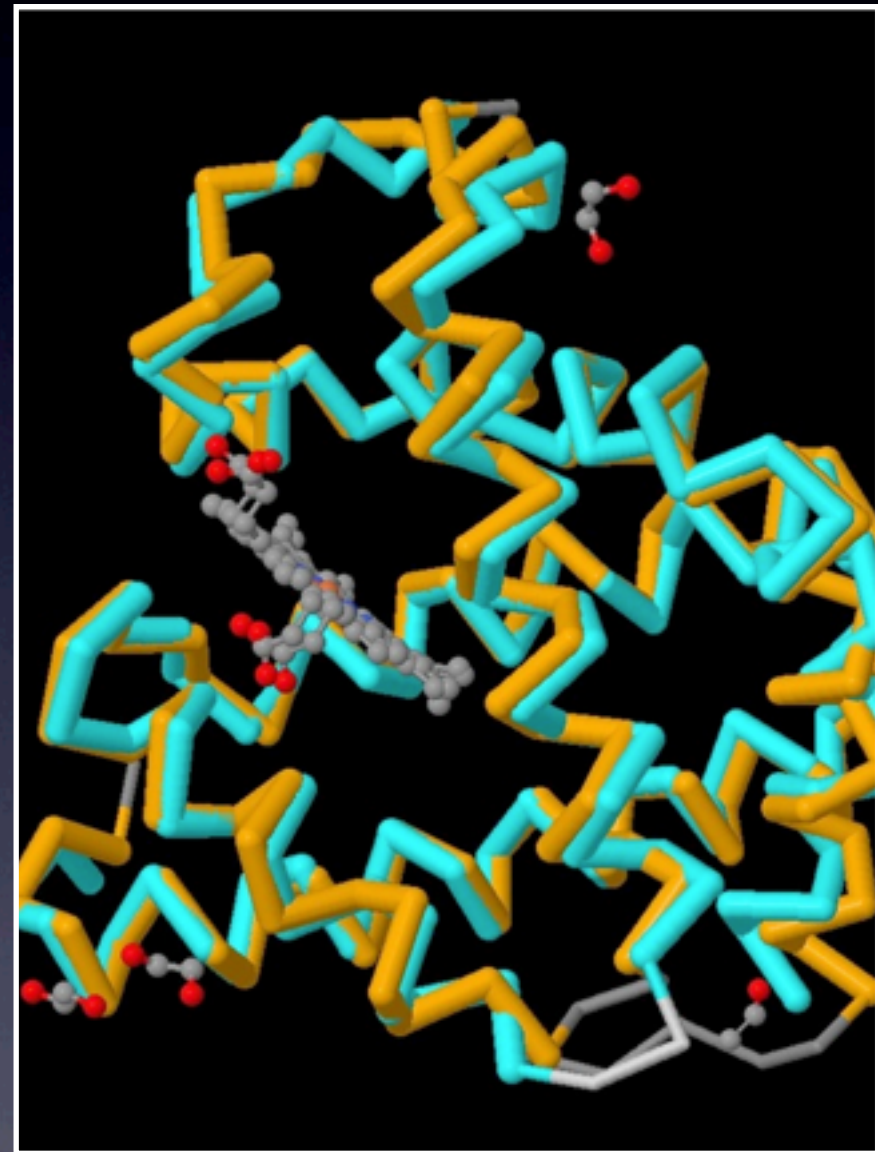
Technical Challenges



- How to provide efficient site
- Distribute data to user
- Roll out weekly updates to multiple data centers
- Availability
- How to provide interactive services
- Growing scale and complexity of structures

Our algorithms

- CPU demands
- Disk IO
- E.g. pre-calculated 3D protein structure clustering
- Parallelization, use of OSG



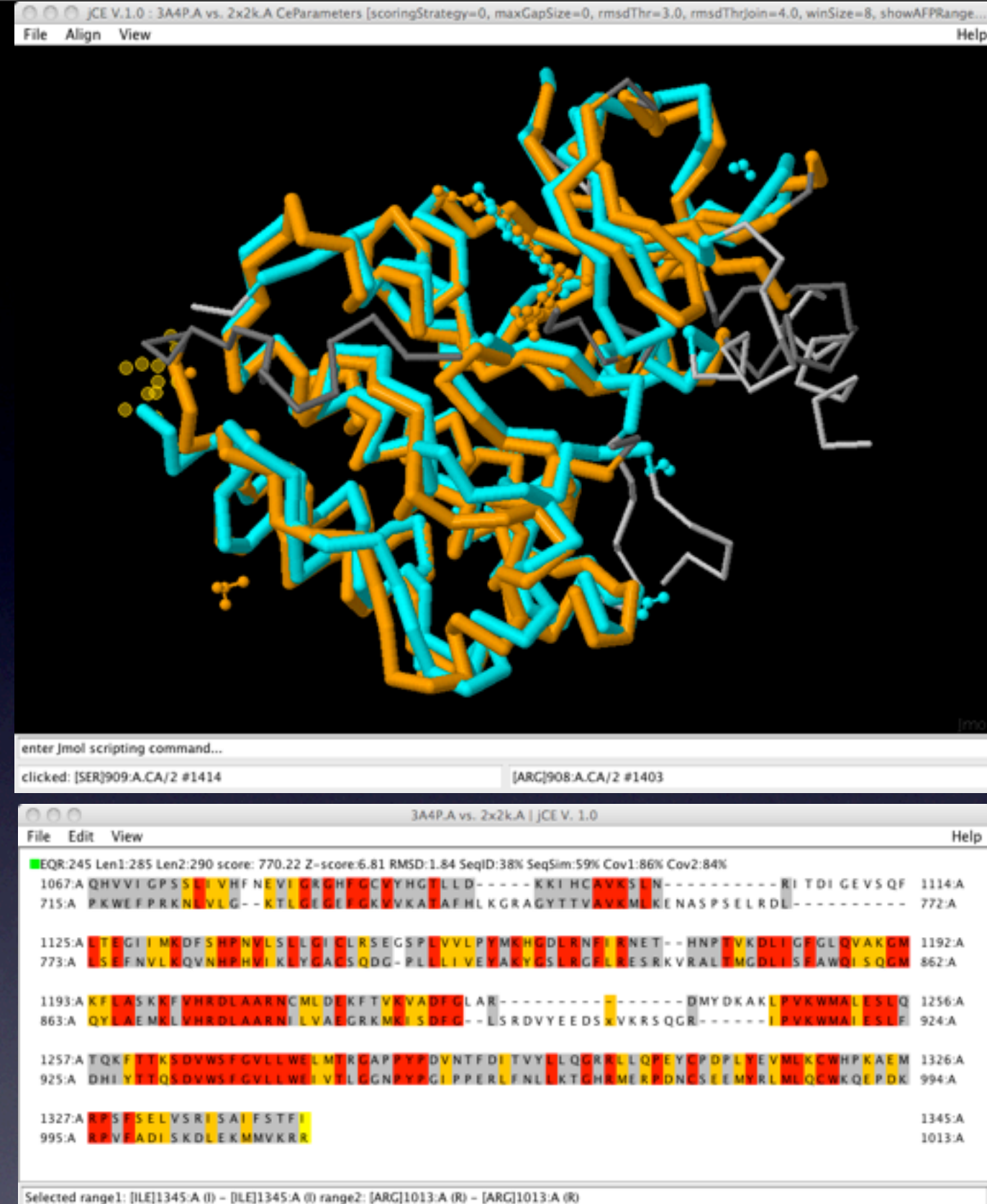
FATCAT / CE

Pairwise comparisons

Custom alignments - using Java
Web Start (stand alone
application)

LGPL v2 - BioJava


Web site - as a service



*c-MET kinase domain 3A4P, RET Tyrosine
Kinase 2X2K
Superimposition of ligand/inhibitor*

All vs All

- Using FATCAT-rigid for All vs.All comparison...
- Well, not true all vs. all, but using representative chains...
- Step 1: Cluster sequences
- Step 2: Calculate all vs. all for ~21,000 representative chains



Initial calculation of
1 billion alignments
(~160k CPU hours)
on OSG



Incremental weekly updates
(~1 million alignments)
<1000 CPU hours

Java Clients can
run anywhere



Open
Science
Grid



PDB

Custom Job
Management



Sends out instructions
to clients

Writes results
to disk



•
•
•

Simple XML protocol



“Hello, give me something to do”



“Run these alignments”



“Here are the results”

- Efficient clients
- < 10 MB install
- Fetch data on demand from public ftp server
- Can be used by PDB users for running custom pairwise alignments and database searches (using Java Web Start)

Systematic Structural Alignment

Objective: Find novel relationships

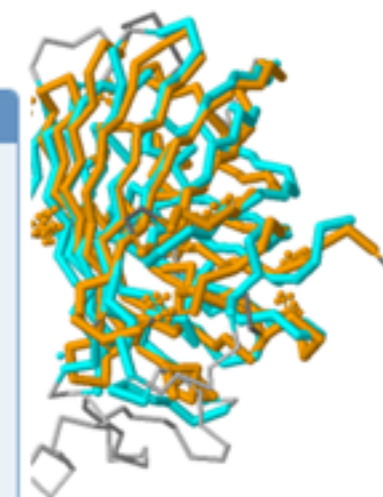
Rank	Result	Chain 2	Title	P-value	Score	Rmsd	Len1	Len2	%ID	%Cov1	%Cov2
1	view	2G2S.B	Green fluorescent p	0.0	478.36	0.93	226	165	96	73	100
2	view	2JAD.A	YELLOW FLUORESCI	0.0	665.32	1.01	226	346	96	100	65
3	view	3E5T.A	Red fluorescent pro	0.0	525.00	1.87	226	228	20	97	96
4	view	3EVP.A	Circular-permutate	0.0	407.39	0.35	226	223	99	61	62
5	view	3GB3.A	KillerRed	0.0	598.80	1.26	226	229	24	98	97
6	view	2G6Y.D	green fluorescent p	7.77E-16	489.59	2.22	226	214	18	93	98
7	view	3EVU.A	Myosin light chain k	2.89E-15	407.23	0.52	226	397	99	62	35
8	view	2A50.D	GFP-like non-fluores	3.06E-12	365.21	2.00	226	167	17	70	95
9	view	2G2S.A	Green fluorescent p	7.95E-10	167.91	0.22	226	64	0	27	97
10	view	1GL4.A	NIDOGEN-1	3.57E-7	295.62	3.01	226	273	9	94	78

Example: Green Fluorescent Protein

- Nidogen-1: similar 11-stranded beta-barrel and internal helices
- 3 Å RMSD, only 9% sequence identity
- Nidogen-1: component of basement membrane, no chromophore
- GFP and NID-1 may share common ancestor

Structure Alignment Results

Alignment Details:	Query: (colored orange/dark grey) GREEN FLUORESCENT PROTEIN	Subject: (colored cyan/light grey) NIDOGEN-1
P-value: 3.57e-07	PDB ID: 2WUR	PDB ID: 1GL4
Score: 295.62	Chain ID: A	Chain ID: A
RMSD: 3.01	Length: 226	Length: 273
%Id: 8.8%	Similarity: 94%	Similarity: 78%



3:A	KGEELFT- GVVPI LVELDGDVN- - - - - GHKFS- VSGEGEGDATY GKLT LKFI CT TGKLPVPWP TL VTTL	64:A
392:A	GROCVAEG SPQRV NGKVKGR I E V GSSOV PVVFE NTDLHSYVVMNHGR SYTA I STI PETV GYSLLPL APIG	461:A
68:A	- - - VQCF SRYP DHMK RHDF FKSAMPE GYVOERTI FF KDD- GNYKTRAEVKFEG- - DTLVNRI ELKGI DFK	131:A
462:A	GI I GWMFAVEQDGFK- - NGFSITGG- EFTROAEVTE L GH P GKLV LKQ QFSGID EHGHLTI STELEGRVP-	527:A
132:A	EDGNI LGHKLEYN YNSH NVYI MADK QKNGI KVNFKTRHNI E- - - - - DGSVQLADHY QQNTPI GDGPV	193:A
528:A	- - - - QI PYGASVHI EPYTELYHYSS- - SVITSSSTREYTVME PDODGAAPSHTHI YQWRTI TFQEC AHD	591:A
194:A	- - - - L LPDNHYL STOSALSKDPNEKRDH MVLL E FVTAAGIT	230:A
592:A	DARPALPSTQQLSVDSVEVLYN- KEERIL RYALSNSI GPVR	631:A

Bottlenecks

- Bottleneck no longer is CPU
- Running up to 1000 parallel jobs
- New Bottleneck: Disk IO

OSG and my environment

- Problem for many students: how to submit batch jobs
- Needed: Low entry level barrier
- Host a local UCSD - OSG tutorial for new users?
- Outreach/User support/ Helpdesk/Marketing is important
- I am on the OSG due to the Engage - outreach team
- Only later discovered the UCSD OSG group
- We have some old hardware which we would like to hook up

We are hiring

- Senior Java/Web developer

Acknowledgments

- Phil Bourne
- Peter Rose
- Chunxiao Bi
- Spencer Bliven
- Wolfgang Bluhm
- Cole Christie
- Dimitris Dimitropoulos
- Alex Gramos
- Gregory Quinn
- Chris Bizon
- Adam Godzik

