



Scalable Software Services for Life Science

Josep Ll. Gelpí – WP7

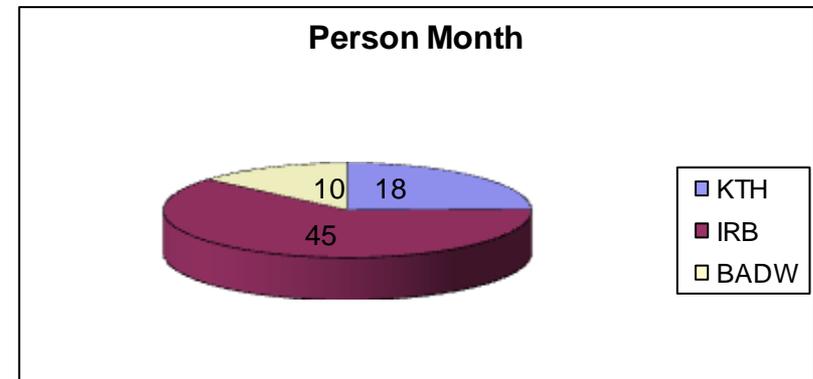
IRB - BSC



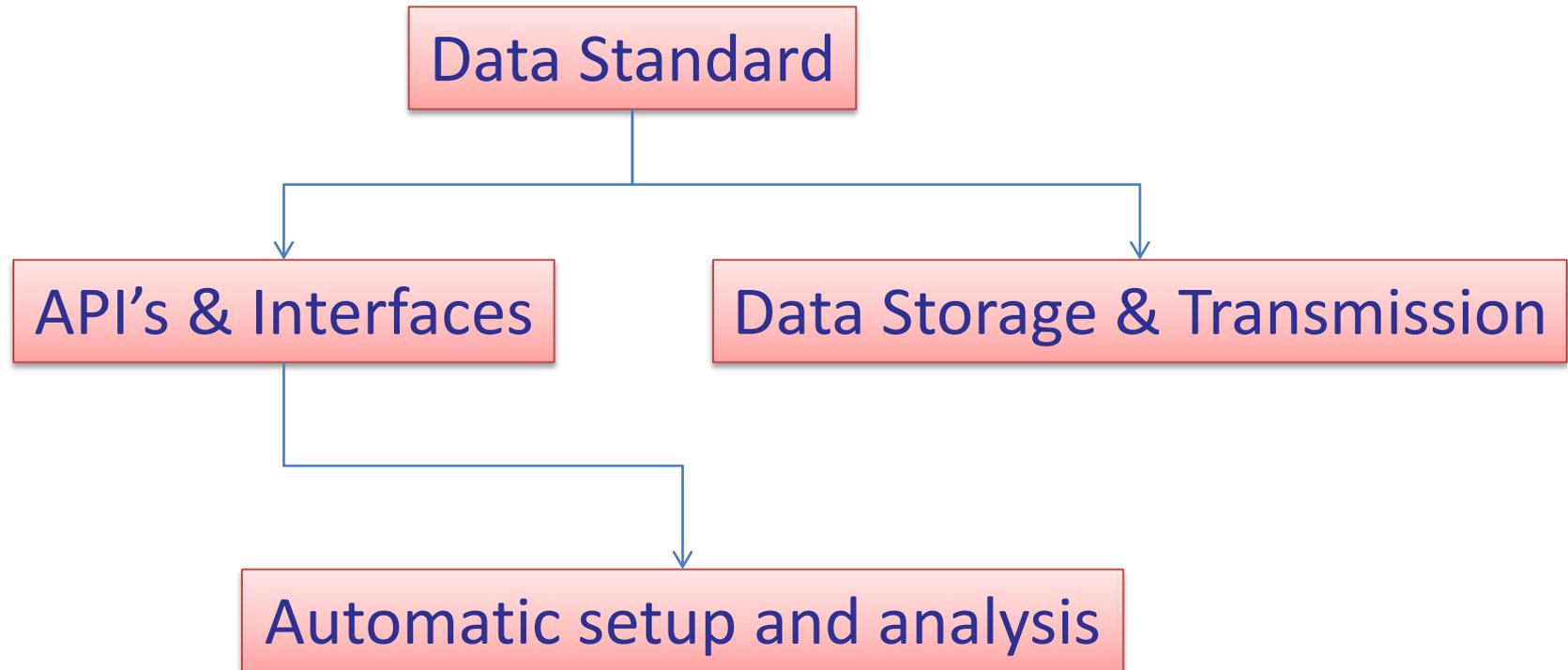
Participants

- Kungliga Tekniska Hogskolan (KTH)
- Leibniz-Rechenzentrum (LRZ)
- Institute of Research in Biomedicine (IRB)

1-38 month activity
Total 73 person month
3 partners involved



- 7.1. Analysis of requirements on data storage and exchange formats (LRZ)
- 7.2. File format standards for data and job description (KTH)
- 7.3. APIs for molecular modeling (IRB)
- 7.4. HT-set-up and analysis tool (IRB)
- 7.5. Simulation databases (IRB)
- 7.6. Management. (IRB)



BioXSD

MathML

RCSB **PDB**
PROTEIN DATA BANK

EDAM

CML

BioCatalogue ^{beta} 
"The Life Science Web Service Registry"

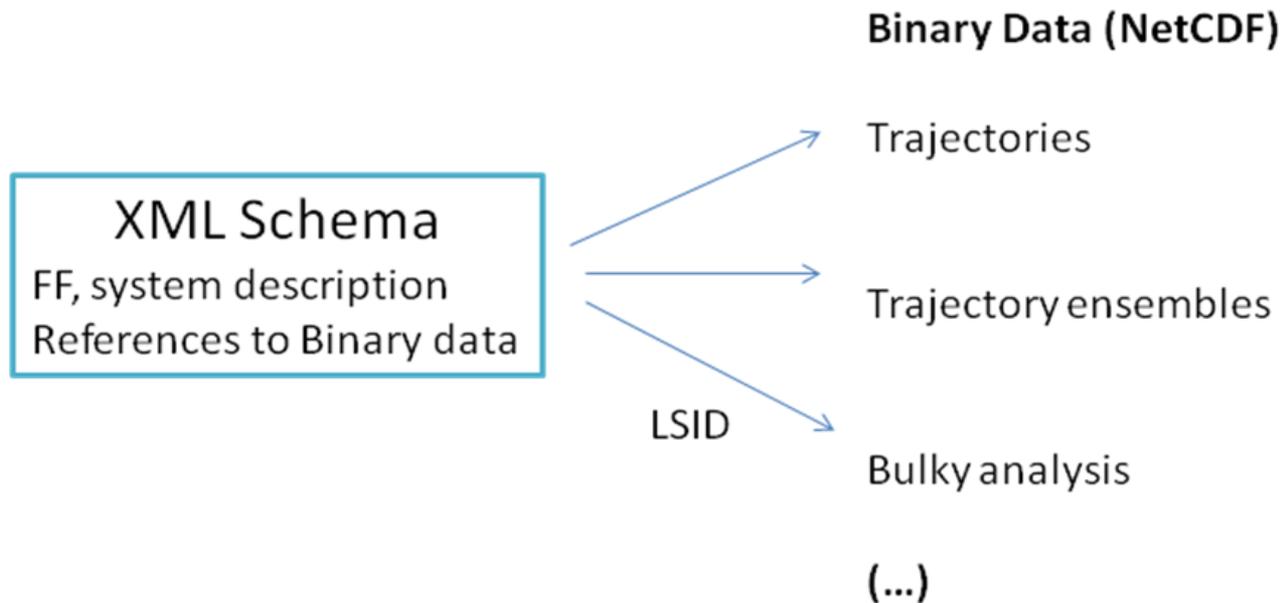
 ELIXIR

 OBO

moby 



- XML Based (FF, system and sim. Descriptions)
 - Data ontology required
- Use available standards for binaries (NetCDF, ...)
- Allow for distributed storage (LSID)



- FF description:
Equations used,
Unified symbols
for FF parameters

- MathML

- FF atom types
and parameters:
Parameter values,
incremental
datasets

```
<math title="k (x - x_o) ^2 ">  
  <mstyle>  
    <mi>k</mi>  
    <msup>  
      <mrow>  
        <mo>(</mo><mi>x</mi>  
        <mo>-</mo>  
        <msub>  
          <mi>x</mi><mi>o</mi>  
        </msub>  
        <mo>)</mo>  
      </mrow>  
      <mn>2</mn>  
    </msup>  
  </mstyle>  
</math>
```

- **Compound/Residue libraries:** Libraries as collection of compounds.

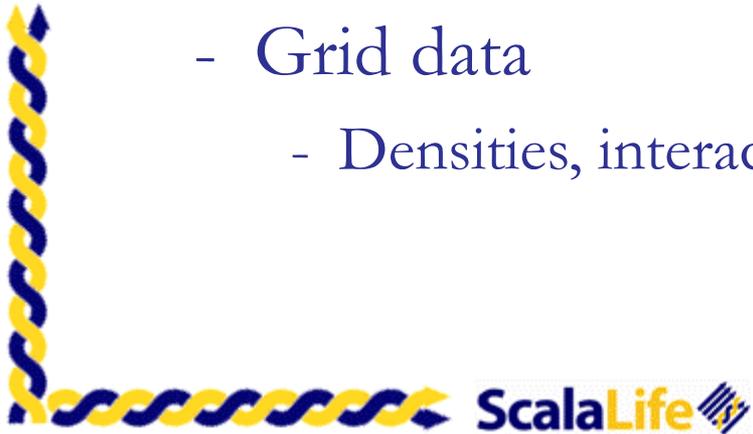
PDBML: the representation of archival macromolecular structure data in XML.
John Wesbrook, Nobutoshi Ito, Haruki Nakamura, Kim Henrick and Helen M. Berman,
Bioinformatics, 21(7), 988-992, 2005. <http://pdbml.pdb.org/>

- Covers small molecules, macromolecule residues
- Includes molecular geometry, standard and non-standard residues
- Compatible with PDB's XML

- Simulated system
 - Cross-references to biological DBs
 - BioXSD for biological data
 - Physical description using compound library
 - Differences from reference system
- Simulation details
 - Conditions, logs,...
 - PBC boxes
- Simulation results
 - Coordinates, trajectories, velocities, ensembles, ...
 - Restart, checkpoints
 - Binary / Compressed / Distributed



- Analysis data organized according to data structure
 - Single values
 - Global or averaged
 - 1D,
 - Residue or time based
 - 2D,
 - Contacts, Structure clustering, ...
 - Grid data
 - Densities, interaction potentials,...



- Distributed approach
 - Large, binary data could be kept at source
 - Perhaps, including analysis software modules

- LSID: Life Sciences Identifiers

- Unique
- Include a standard procedure to locate/retrieve data
- LSID:<Authority>:<Namespace>:<ObjectID>[:<Version>]

mmb.pcb.ub.es:ScalaLife:1I6F_GROMACS3.2_G43-SPC_0_TRAJ

- Central catalogue of data.
- Data providers should implement a simple retrieval

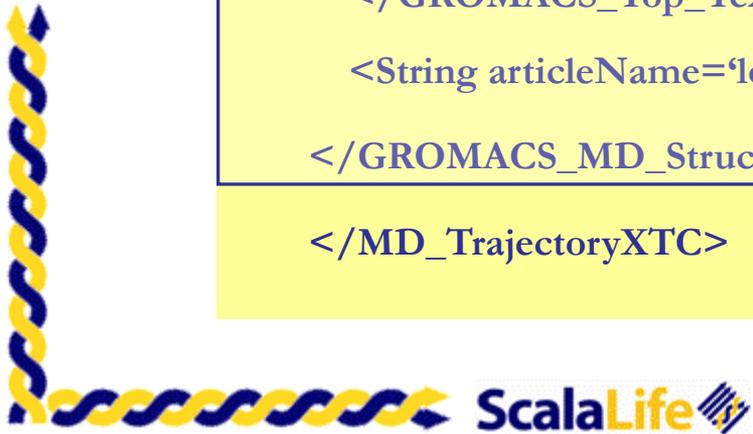
```
<MD_TrajectoryXTC id='2ki5' namespace='PDB'>
```

```
  <String articleName='coordinates'>  
    ...XTC Data...  
  </String>
```

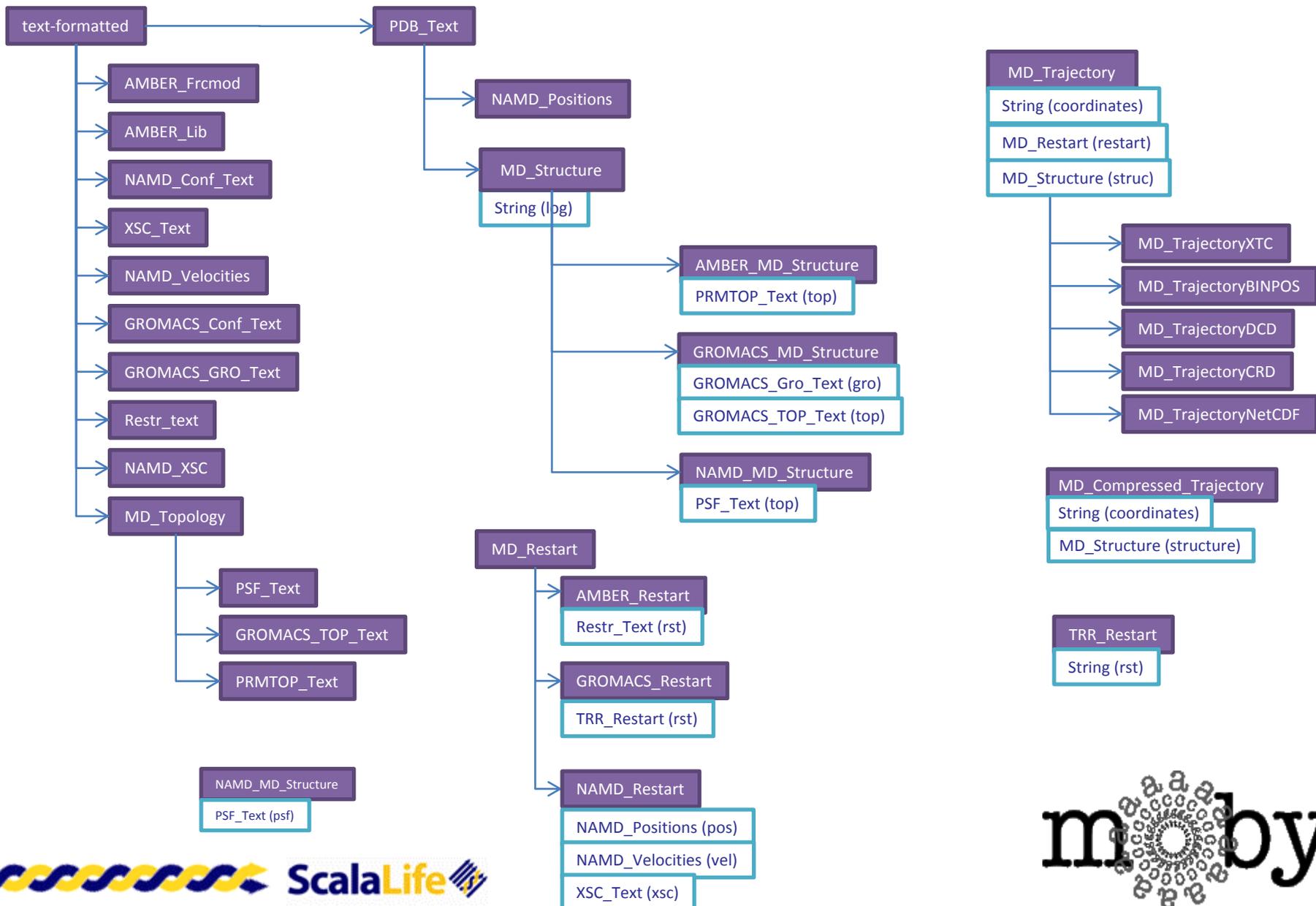
```
  <GROMACS_Restart articleName='restart'>  
    <TRR_Restart>...TRR data...</TRR_Restart>  
  </GROMACS_Restart>
```

```
  <GROMACS_MD_Structure articleName="struc">  
    <PDB_Text articlename="structure"> ---PDB Data... </PDB_Text>  
    <GROMACS_Gro_Text articleName="gro">  
      ...GRO data...  
    </GROMACS_Gro_Text>  
    <GROMACS_Top_Text articleName="top">  
      ...Top data---  
    </GROMACS_Top_Text>  
    <String articleName='log'>...log data...</String>  
  </GROMACS_MD_Structure>
```

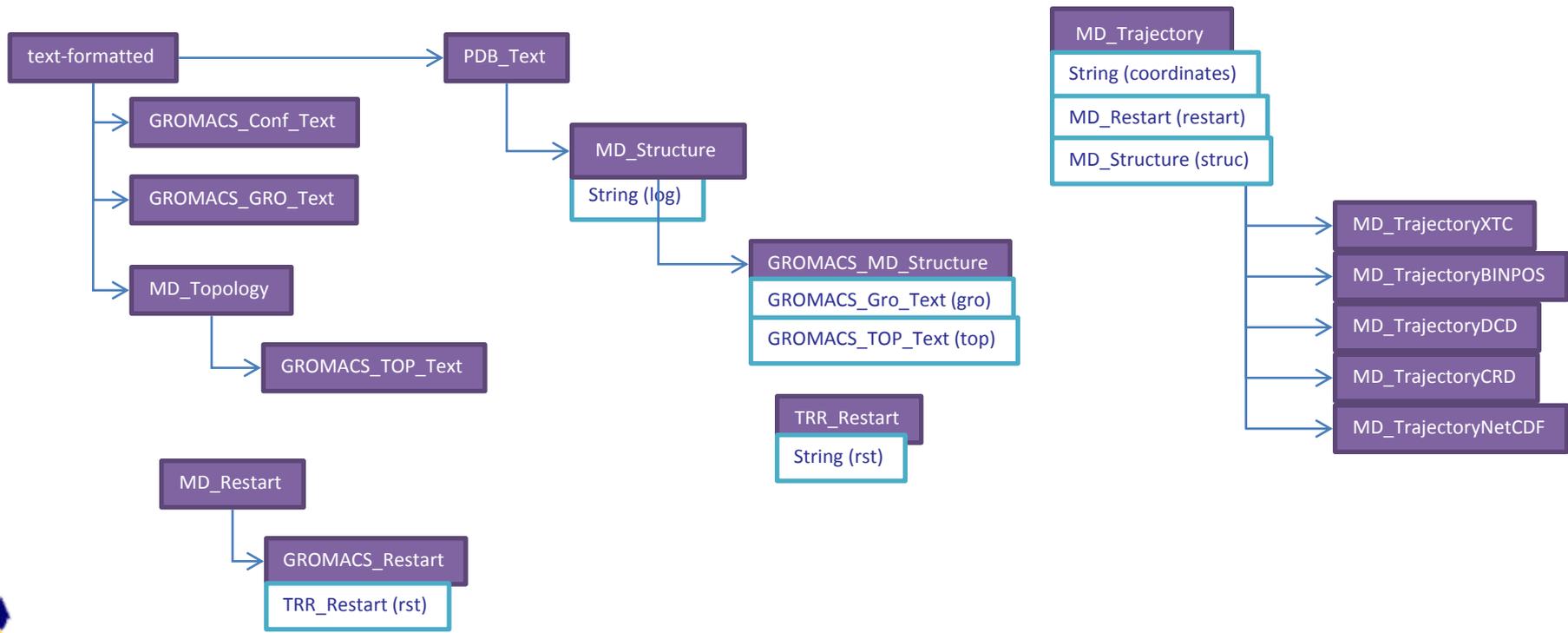
```
</MD_TrajectoryXTC>
```



MDMoby ontology



A zoom in GROMACS,...



- Validation protocol
 - 1st draft of XML schema agreed by WP7/Scalalife partners
 - Expert consultants community selected
 - Scalalife partners
 - MD / QM developers
 - Visualization / analysis developers
 - Key MD users (alpha users community?)
 - Feedback from consultants discussed and incorporated
 - Approval by Scalalife and release



- Open questions
 - Existing standards to use
 - Transferability of parameters/FF
 - GROMACS oriented?
 - Integration in BioXSD / MLPDB / OBO / ...
 - ...



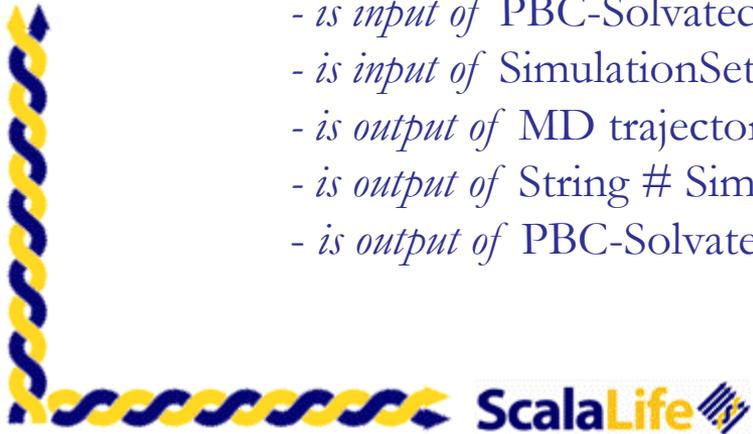
- Aim: Define API's adapted to standard datatypes, to be developed and optimized at the lower level.
- MD or analysis codes could be built on top the these universal MD API's
- Requires Operations Ontology
 - Highly modular and hierarchical set of operations with well defined input and output
- Building an ontology requires agreement !!

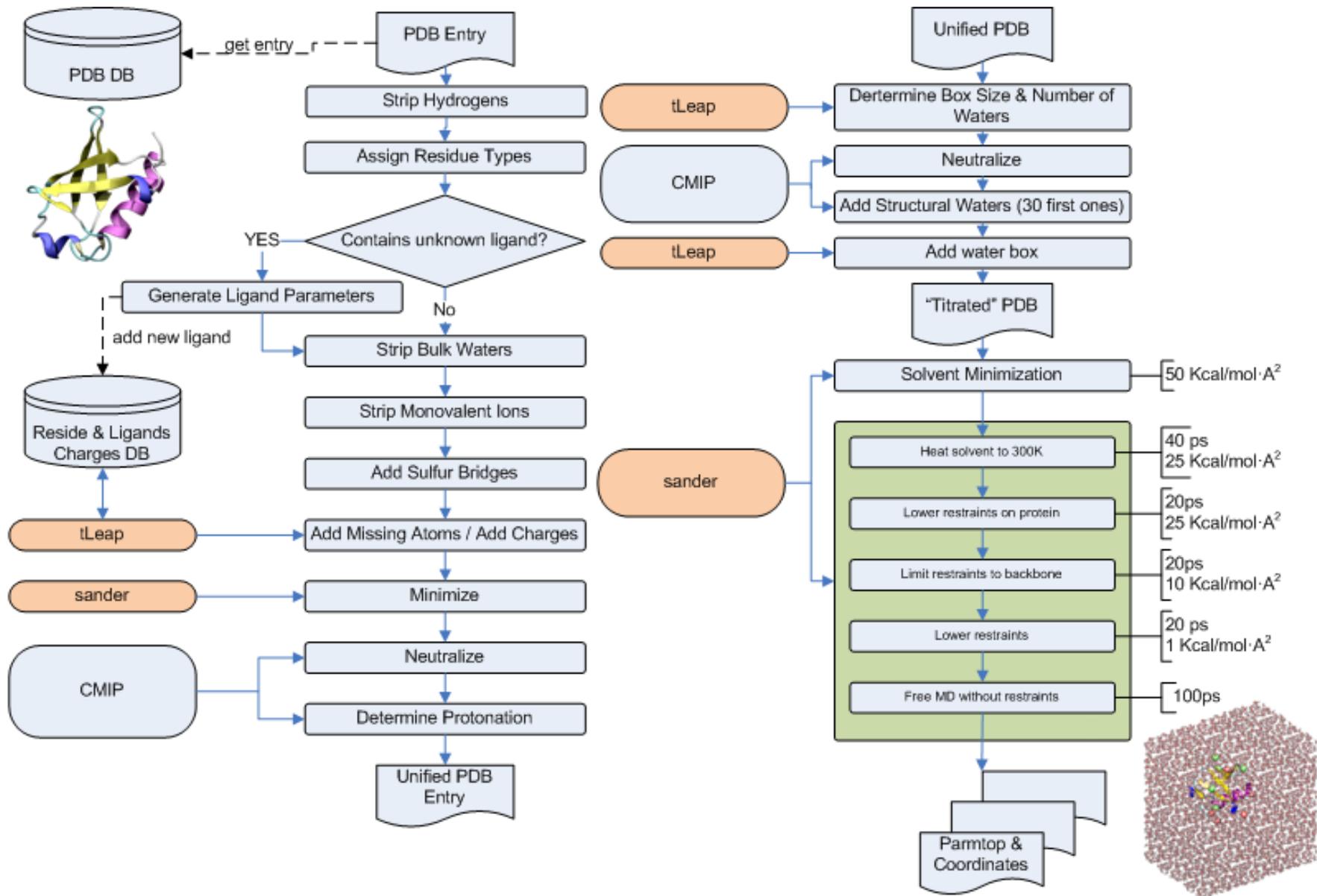
-System setup

- *is a* Structure checking
 - *is input of* PDB
 - *is output of* PDB #Corrected structure
 - *is a* chirality check
 - *is a* Thr chirality check
 - *is a* CA chirality check
- *is a* Residue Protonation
 - *is input of* PDB
 - *is output of* PDB #Protonated structure
 - *is a* CMIP based
 - *is a* ProtpKa based

-Simulation

- *is a* NPT Simulation
 - *is input of* PBC-Solvated system
 - *is input of* SimulationSettings
 - *is output of* MD trajectory
 - *is output of* String # Simulation log
 - *is output of* PBC-Solvated system # Restart data





- Automatic setup and analysis
 - Build on top of API's operations
 - Adapted to a wide selection of end users
 - Black-box workflows for non-expert
 - Ex: Full MD Setup using XXX FF including solvation
 - Ex. Stability after residue mutation
 - Detailed workflow control for experts
 - Runnable from
 - Web interfaces (MDWeb)
 - Web Services (MDMoby)
 - Command-line scripting (MDMoby – MobyLite API)




```

my $pdbStruct = new PDB_Text($pdbCode,"PDB");
$pdbStruct->content($pdbContent,1);

my $cleaned = cleanStructureFromPDBText ('structure' => $pdbStruct) ->
{'structure'};

my $pdb_h = getGROMACS_MD_StructureFromPDBText (
'structure' => $cleaned, 'forcefield' => $forcefield ) -> {'structure'};

my $min_h = optimizeStructureFromGROMACS_MD_Structure (
'structure' => $pdb_h, 'minimize' => 500, 'md_type' => 1) -> {'structure'};

my $min_h2 = optimizeStructureFromGROMACS_MD_Structure (
'structure' => $min_h, 'minimize' => 500, 'md_type' => 2) -> {'structure'};

my $solv = solvateStructureFromGROMACS_MD_Structure (
'structure' => $min_h2, 'ions' => 'true', 'boxsize' => 0.8, 'boxtype' =>
'octahedron', 'ionic_concentration' => 0.05) -> {'structure'};

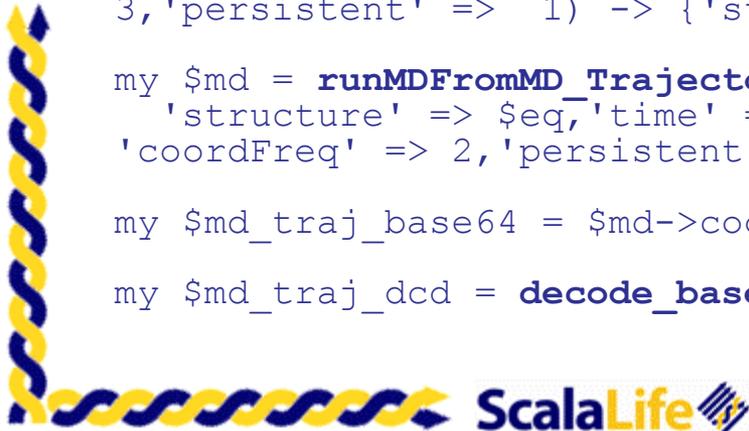
my $seq = runMDFromGROMACS_MD_Structure_async (
'structure' => $solv, 'time' => 1, 'timestep' => 0.001, 'md_type' =>
3, 'persistent' => 1) -> {'structure'};

my $md = runMDFromMD_TrajectoryXTC_async (
'structure' => $seq, 'time' => 5, 'md_type' => 4, 'timestep' => 0.004,
'coordFreq' => 2, 'persistent' => 1) -> {'structure'};

my $md_traj_base64 = $md->coordinates;

my $md_traj_dcd = decode_base64($md_traj_base64);

```



MDWeb

Molecular Dynamics on Web

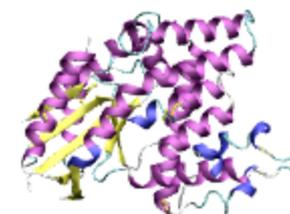
User: Josep Ll. Gelpi


[Home](#)
[Start new project](#)
[Close workspace](#)
[Help](#)

test 2ki5_a (MDWeb4d80f45acca7c)

Last modification on: 16/03/2011 18:33

Disk Usage: 1.9 MB



Stored structures

Click on structure title to deploy the toolbox.

- Base structure (196.4 kB)

Select the desired operation.

 Title: Comment:

List of Operations:

- List of Operations:
- Check for disulphide bonds
- Clean PDB
- Fix Side Chains
- Mutate residue
- Amber FULL MD Setup
- Amber MD Setup

test 2ki5_a (MDWeb4d80f45acca7c)

Last modification on: 16/03/2011 18:44

Disk Usage: 2.5 MB

Stored structures

Click on structure title to deploy the toolbox.

- Base structure (196.4 kB)
- Structure Topology for Gromacs_01 (

Select the desired operation.

Title: Comment:

Box Type:

Add Ions:

Terminado

Forcefield:



Structure

- Atoms
- Ligands
- Wireframe
- Cartoon
- Show hydrogen bonds

Cartoon color

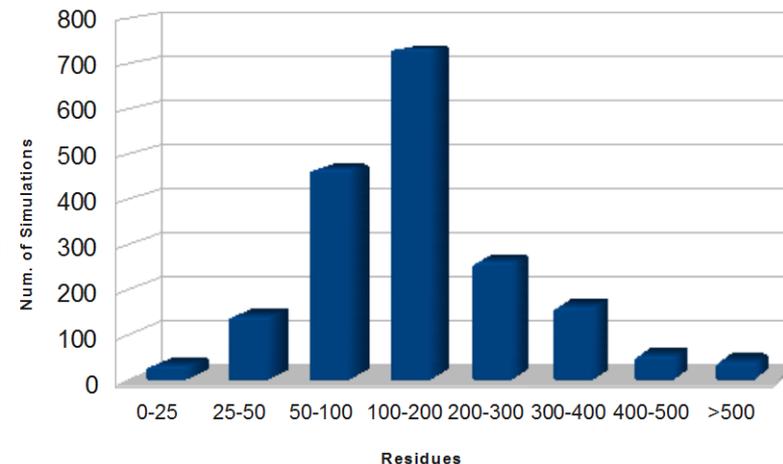
Structure Chain

Hide Hydrogens

Jmol script terminated

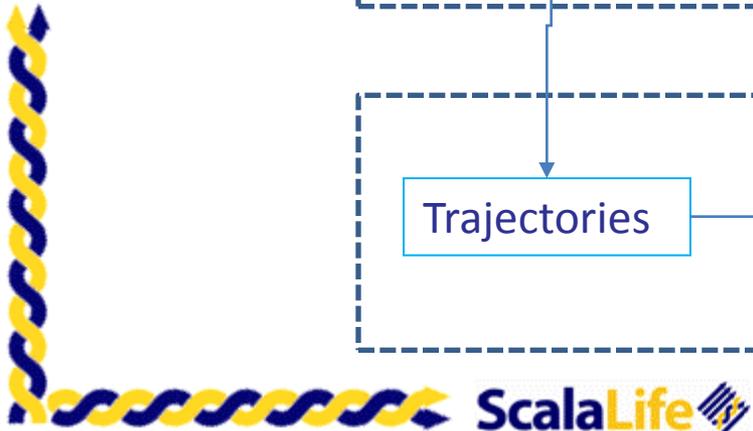
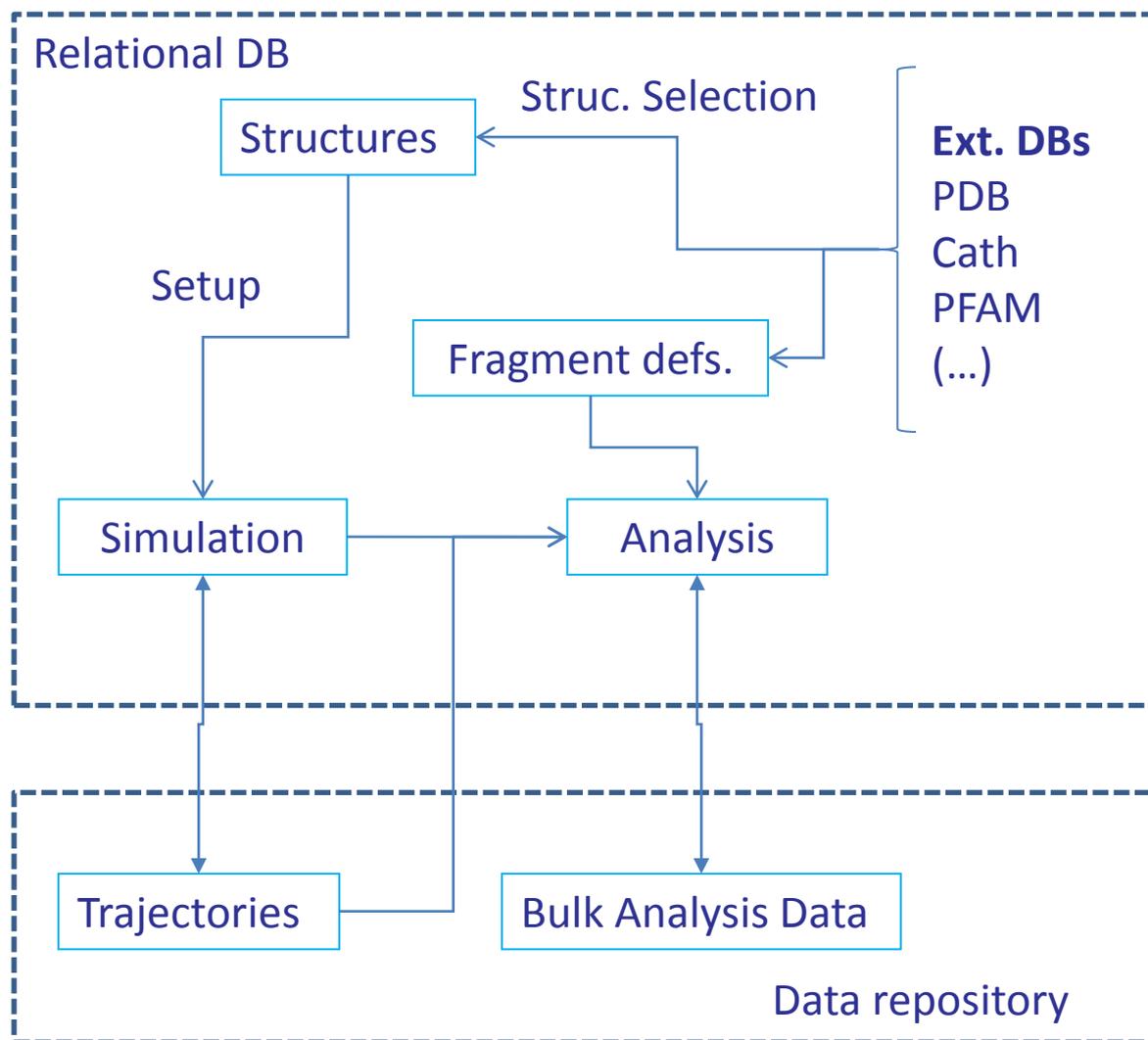
- Open questions
 - Programming library vs. Web Services library
 - Grid/Cloud
 - Access in HPC
 - Integration with bioinformatics (non structural)
 - Collection of “black-box” workflows

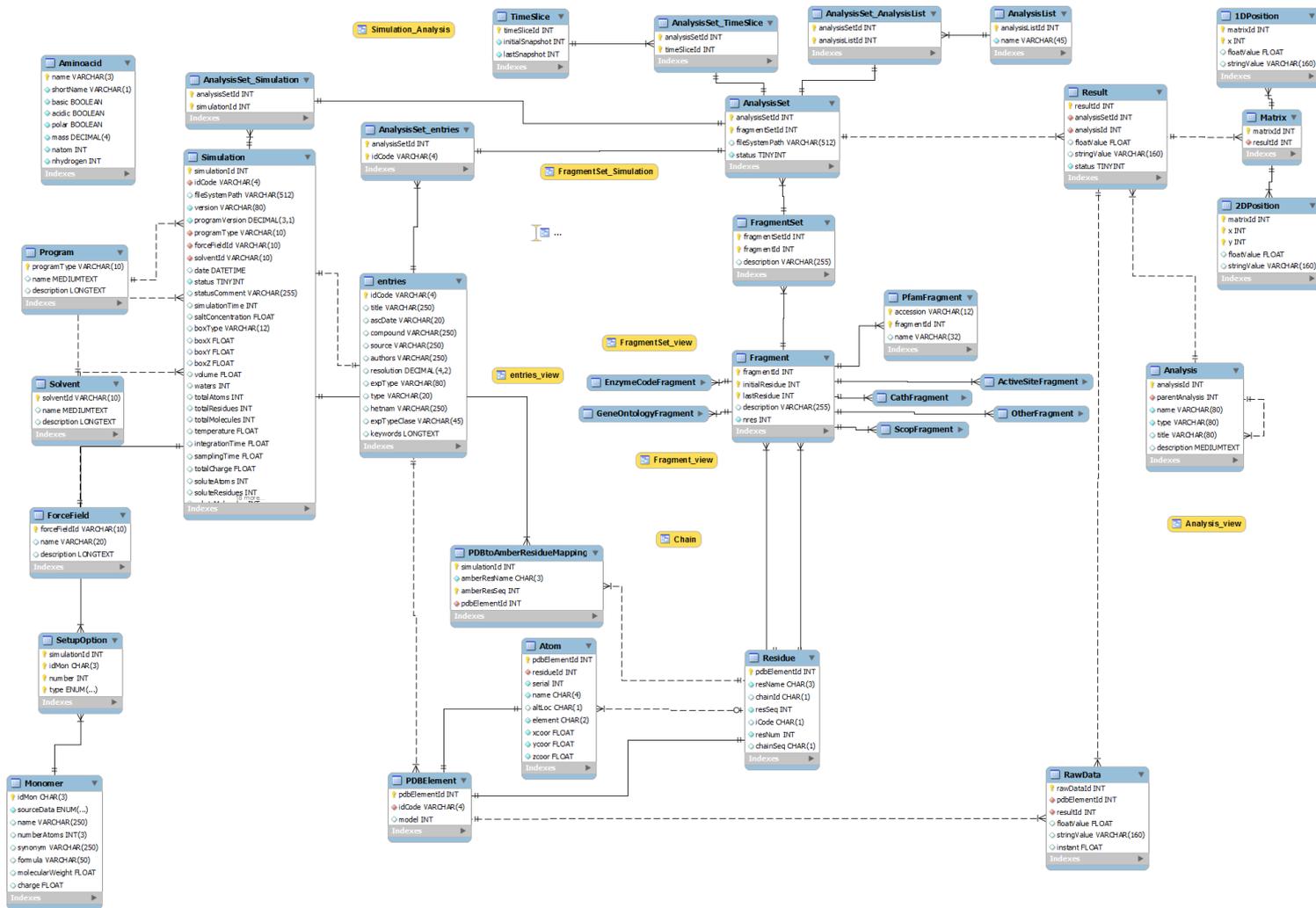
- Reference Example MoDEL project
 - 1875 simulations, 1595 systems
 - Disk usage:
 - Total: 18Tb
 - Raw Solvent trajectories: 14Tb
 - Dry trajectories: 1.7Tb
 - Analysis: 1.1Tb
 - Setup: 1.2Tb
 - MySQL DB: 23Gb



- ▶ Rueda M . et al. PNAS 2007, 104, 798-801
- ▶ Meyer T. et al. Structure 2010. 18, 1399-1409
- ▶ <http://mmb.pcb.ub.es/MoDEL>







model 2.3





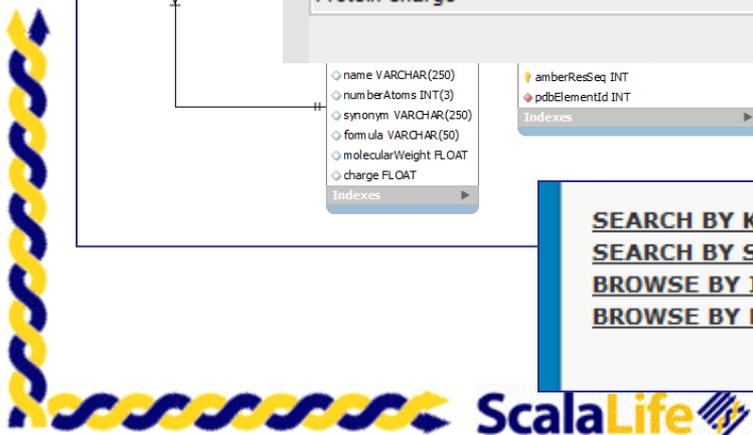
Simulation info for 1A3D - AMBER9.0 parm99 tip3P (P99 T3P) simulation

PDB code	1A3D		
Program Version	AMBER 9.0		
Force field	parm99 (standard amber forcefield)		
Salt Concentration	60.24 mM		
Total atoms	19574		
Total residues	119		
Total molecules	1		
Box	Octahedral (66.55 Å, 66.55 Å, 66.55 Å)		
Volume: 227000.0 A3	Solvent molecules: 5935		
Simulation time	84500 ps		
Integration time	0.0020 ps		
Sampling time	1.0 ps		
Temperature	300.0 K		
Shake	vdW cutt-off	PBC	Ewald
2	8.0 A	2	1
Protein charge	-3.0 e		

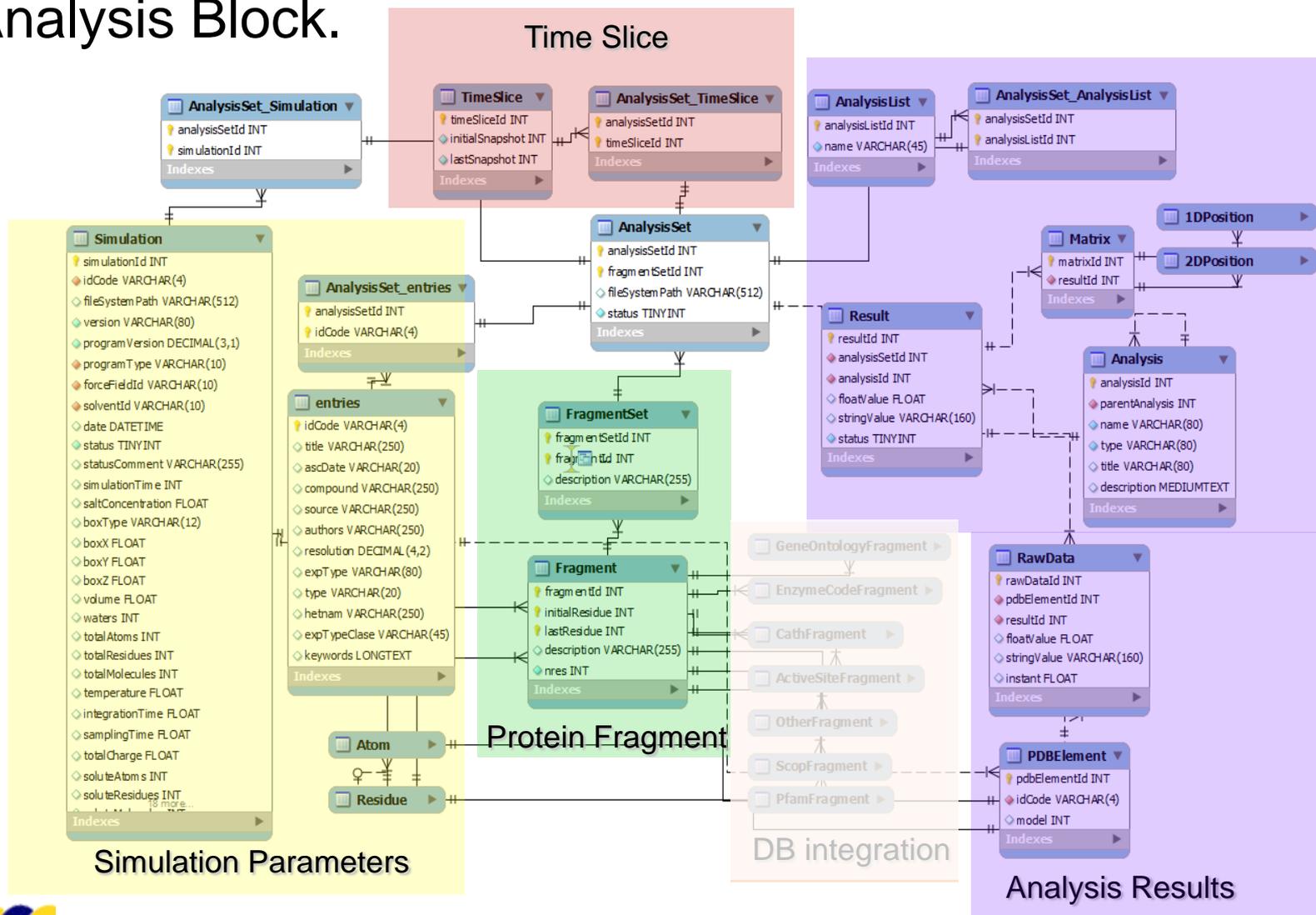
SEARCH BY KEYWORD
SEARCH BY SEQUENCE
BROWSE BY ID
BROWSE BY FOLD

structure

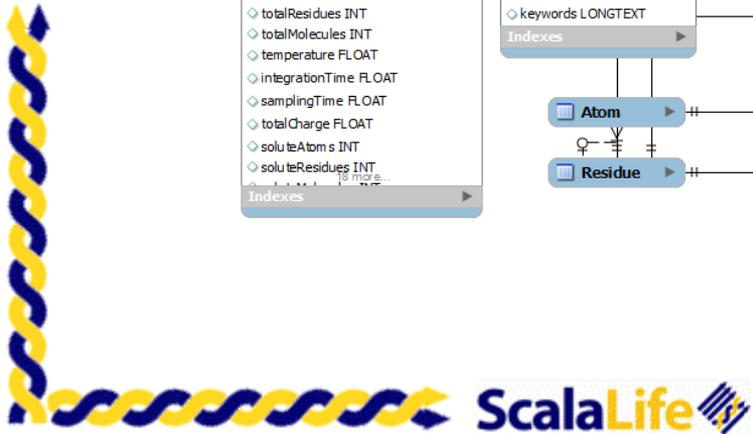
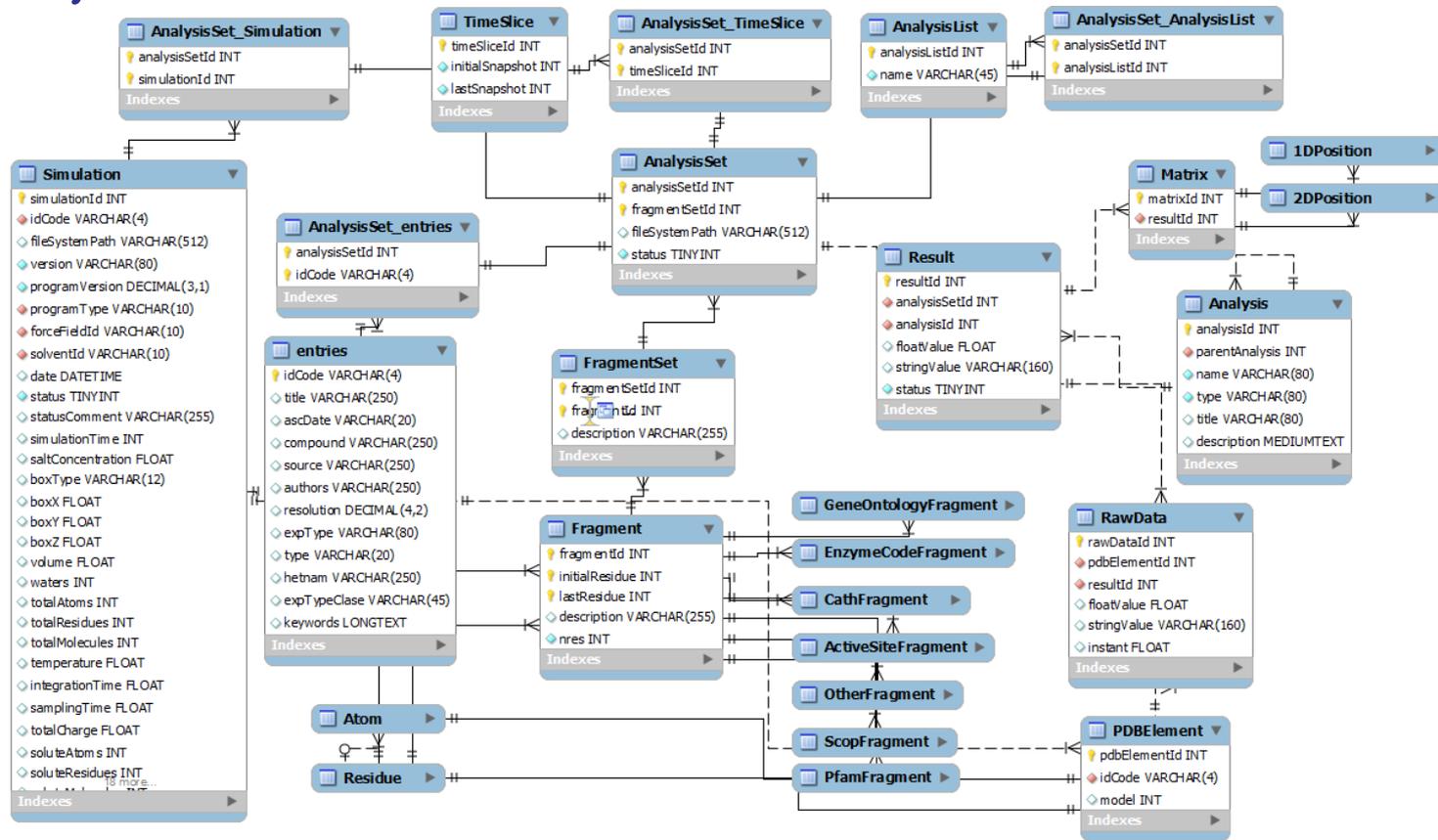
- PDB is cross-referenced to Uniprot, CATH, PFAM, etc.



Analysis Block.



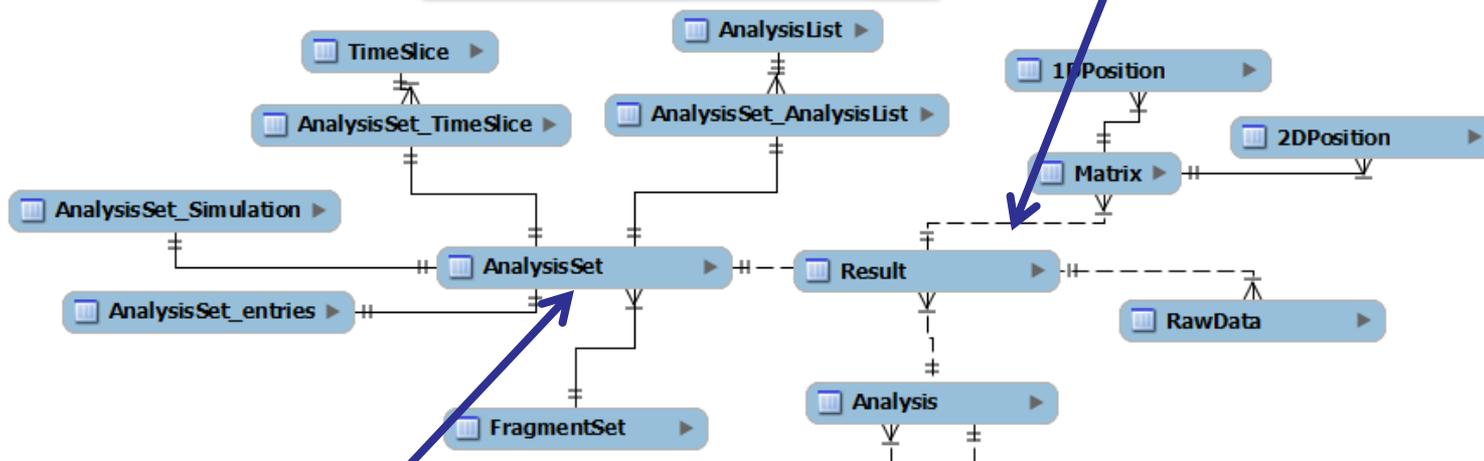
- Analysis Block



- Analysis Block. Analysis Selection

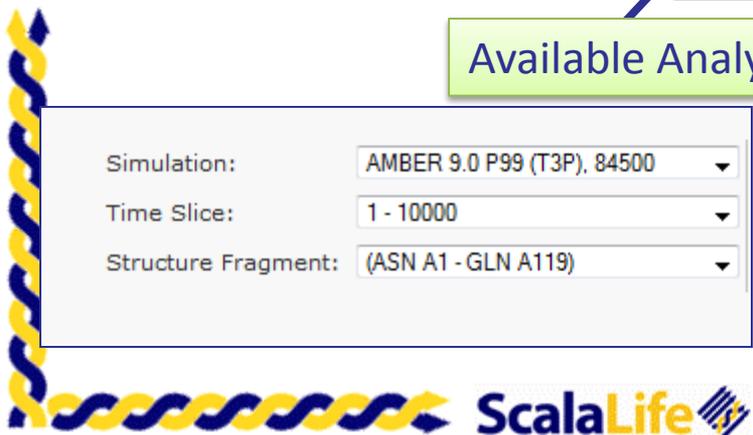
Predefined lists, possible personalization

Analysis Data. Common data structure for all analysis

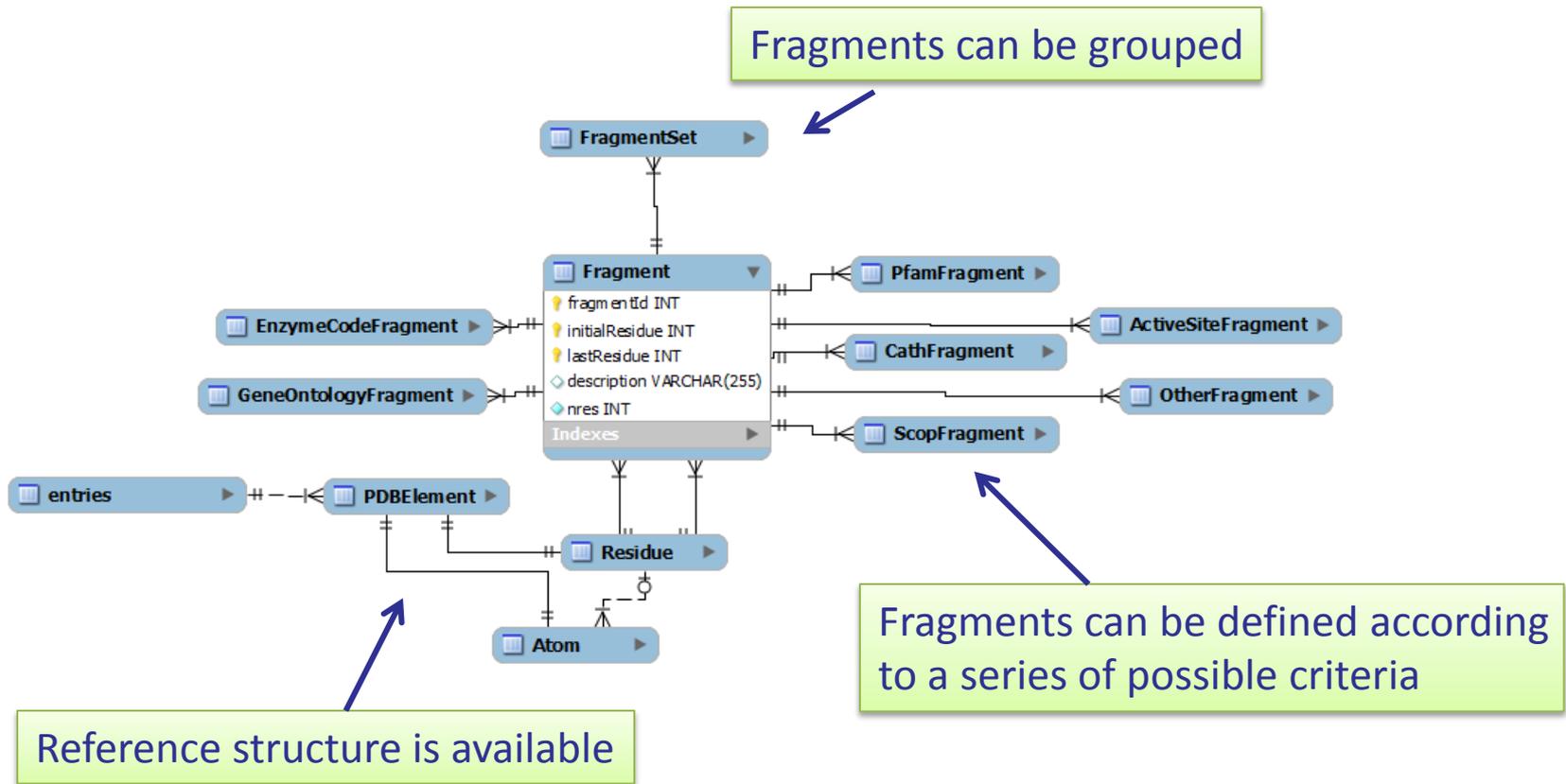


Available Analysis: Simulation + Time Slice + FragmentSet

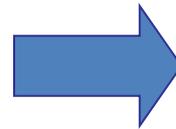
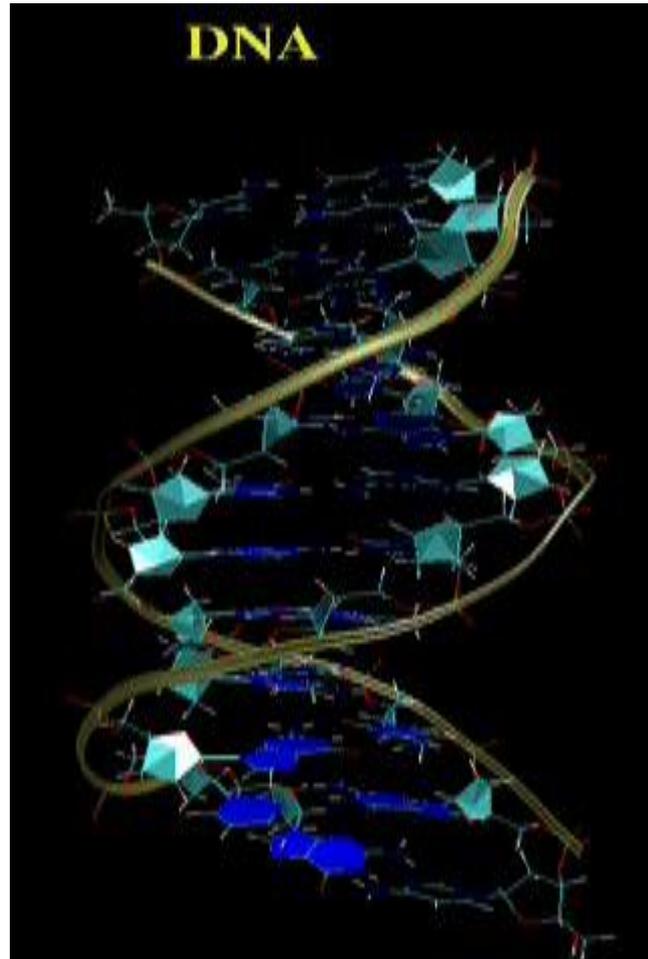
Simulation:	AMBER 9.0 P99 (T3P), 84500
Time Slice:	1 - 10000
Structure Fragment:	(ASN A1 - GLN A119)



- Analysis Block. Fragment Selection



High-compression tools (PCAZIP)



Diagonalization of
covariance matrix

Eigenvectors

Select essential
space

Project cart coord
Essential space

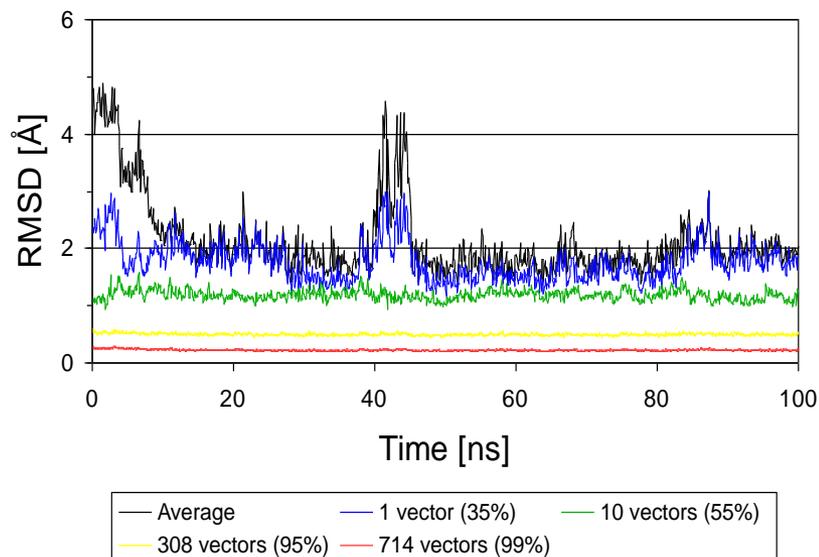
Store projections



Quality
threshold

PCAZIP data reduction

1idr RMSD



	95% cutoff	
Protein	RMSd	File size
1ark	0.36	8.5
1cei	0.36	7.8
1sr0	0.45	6.0
2gb1	0.29	10.0
3ci2	0.36	8.6
2icb	0.33	8.8
1idr	0.50	5.1

- Open questions
 - Distributed approach
 - Distributed analysis software
 - Long term storage vs recalculation
 - Compression of raw solvent trajectories
 - Integration with ELIXIR core databases

