

# MD ANALYSIS

**A Python package for  
the rapid analysis of  
molecular dynamics  
simulations**

*SciPy 2016, Austin, TX*  
July 13, 2016

**Oliver Beckstein**  
Arizona State University

Richard J Gowers<sup>1</sup>, Max Linke<sup>2</sup>, Jonathan Barnoud<sup>3</sup>, Tyler J. E. Reddy<sup>4</sup>, Manuel Melo<sup>3</sup>, Sean. S. Seyler<sup>5</sup>, Jan Domanski<sup>4,6</sup>, David L. Dotson<sup>5</sup>, Sébastien Buchoux<sup>7</sup>, Ian Kenney<sup>5</sup>

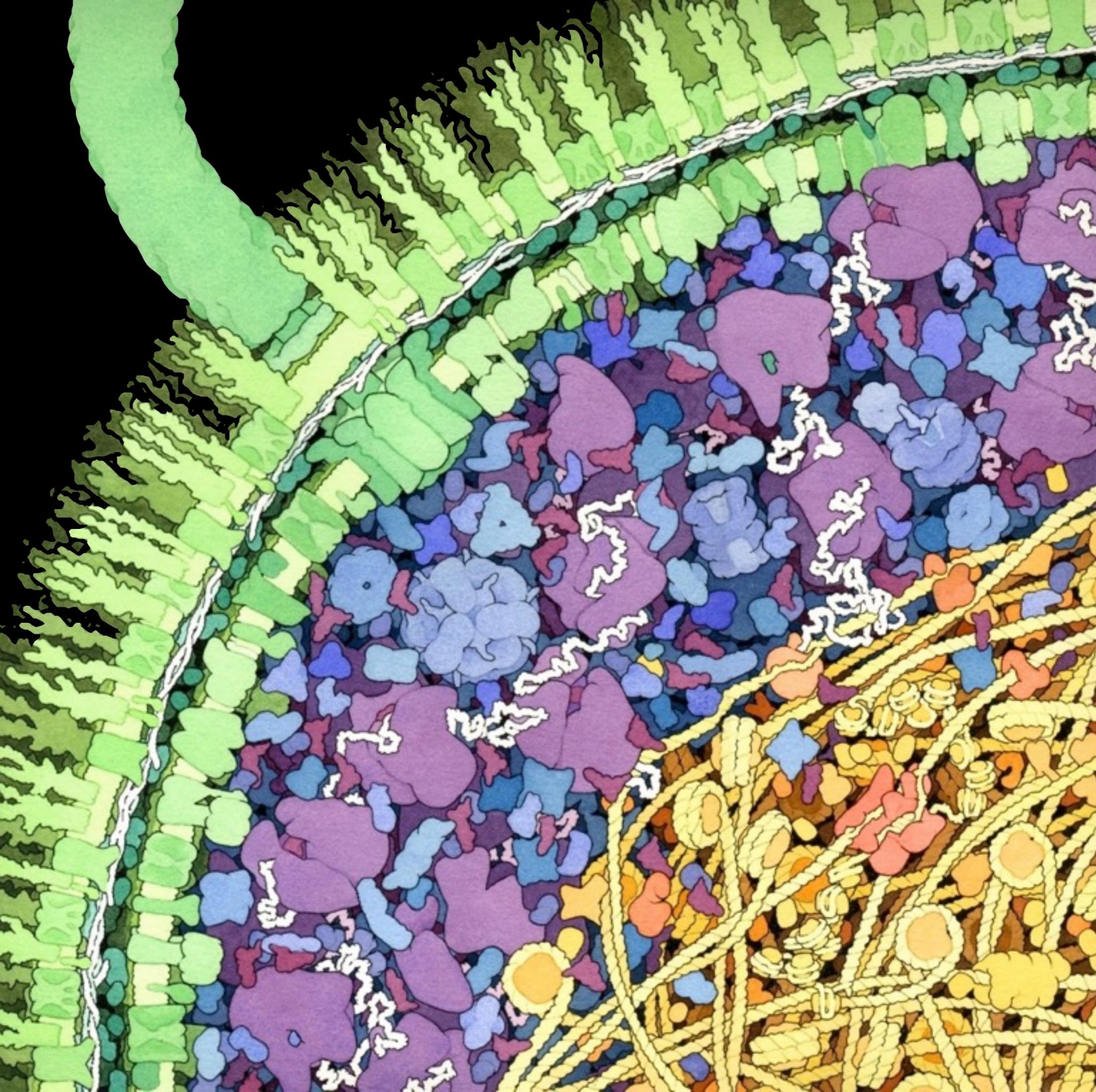
- <sup>1</sup>University of Edinburgh, UK,
- <sup>2</sup>Max Planck Institut für Biophysik, Frankfurt, Germany,
- <sup>3</sup>University of Groningen, The Netherlands,
- <sup>4</sup>University of Oxford, UK,
- <sup>5</sup>Arizona State University,
- <sup>6</sup>NIH NIDDK, Bethesda, MD,
- <sup>7</sup>Université de Picardie Jules Verne, France

LIFE

cells

molecules

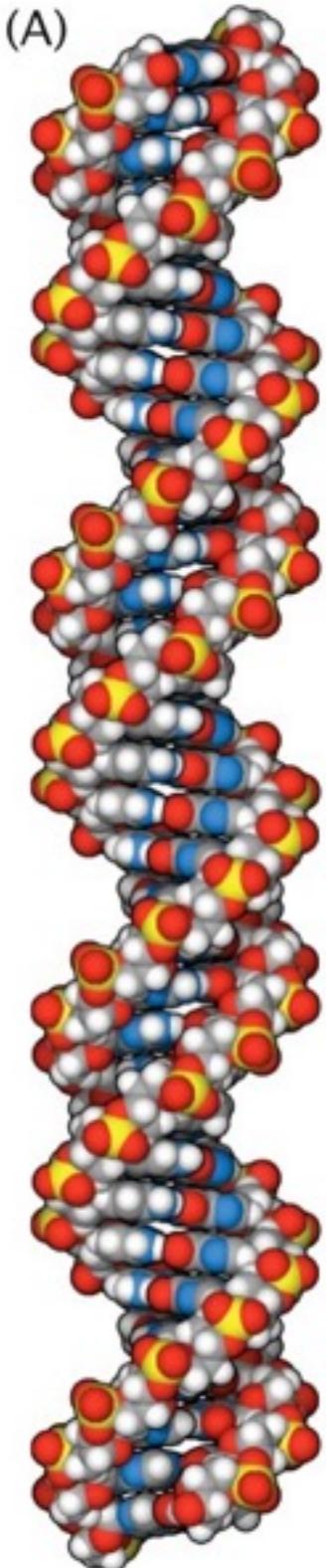
atoms



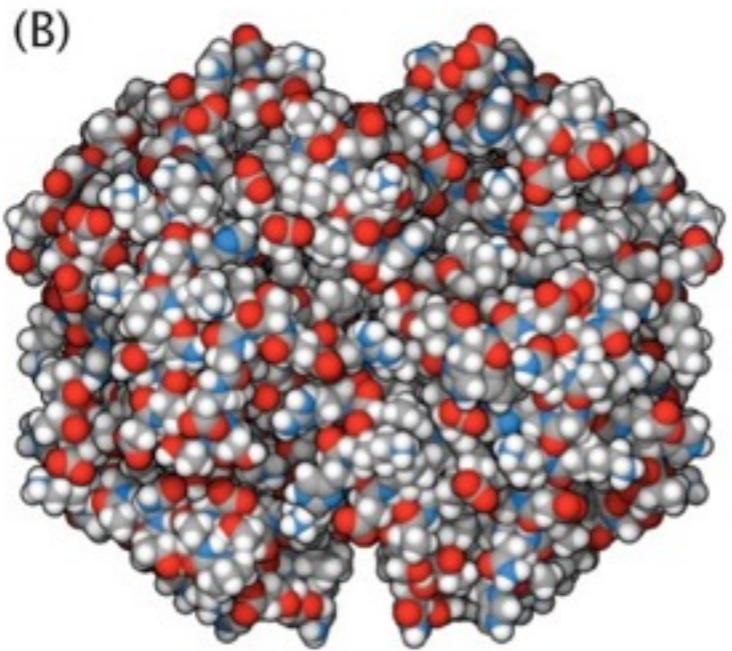
*Escherichia coli* (© 1999 David S. Goodsell, the Scripps Research Institute)

nucleic acids  
(DNA, RNA)

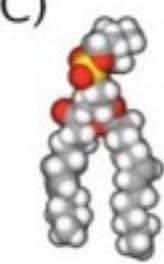
# Biomolecules



proteins



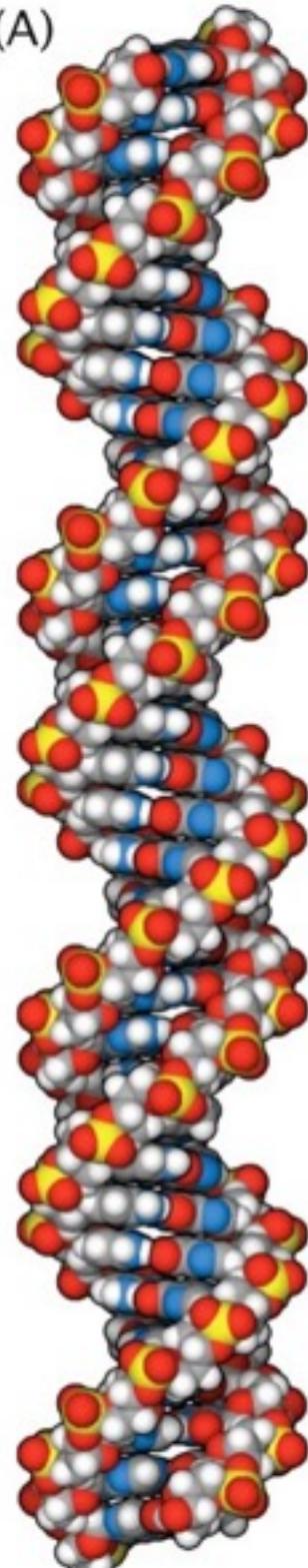
(C)



lipids

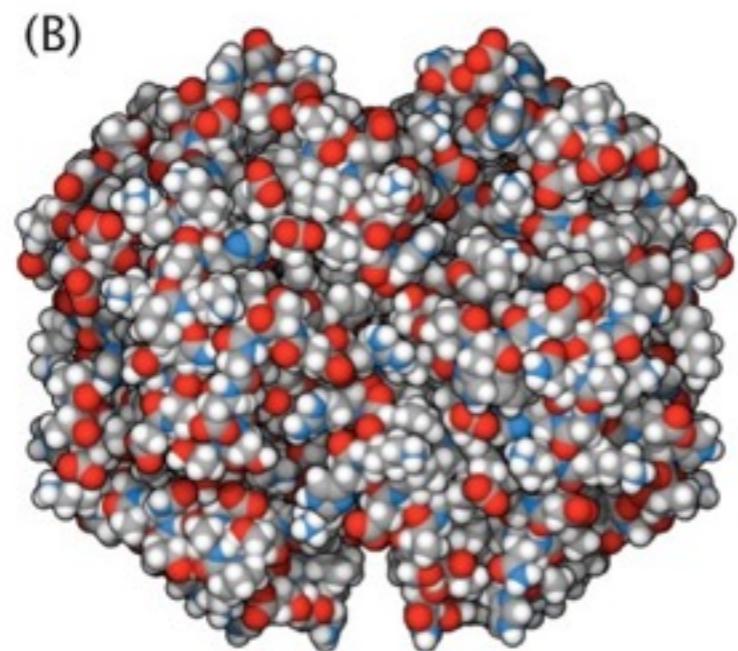
2 nm

nucleic acids  
(DNA, RNA)

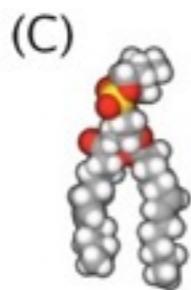


# Biomolecules

proteins

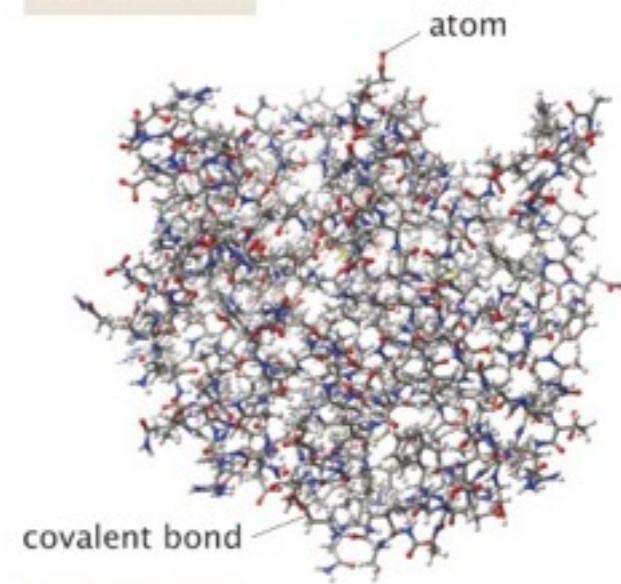


protein  
representations

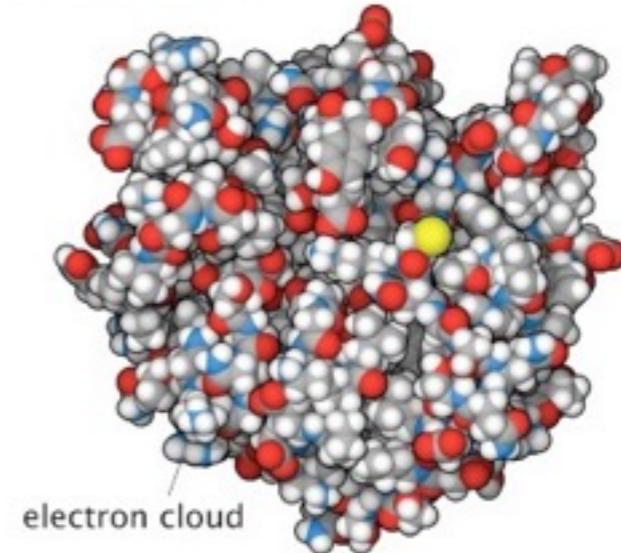


lipids

ball and stick



space-filling



ribbon

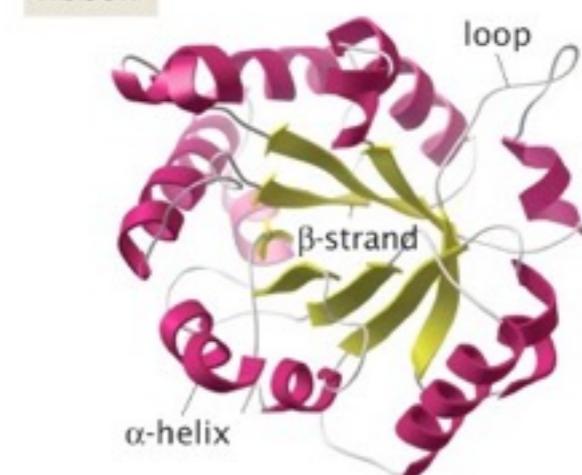
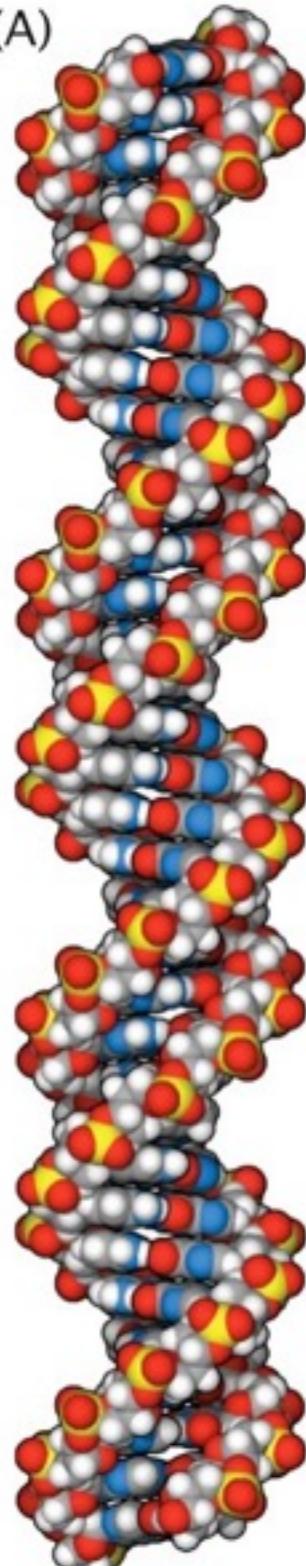


Figure 1.1 Physical Biology of the Cell, 2ed. (© Garland Science 2013)

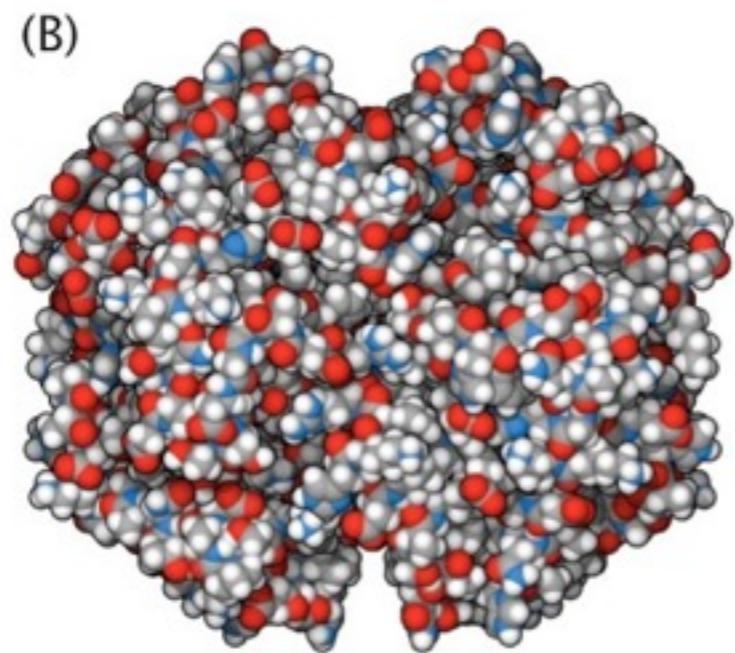
Figure 2.32 Physical Biology of the Cell, 2ed.

nucleic acids  
(DNA, RNA)

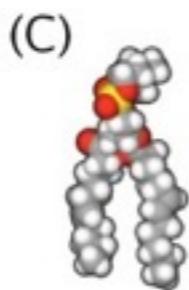


# Biomolecules

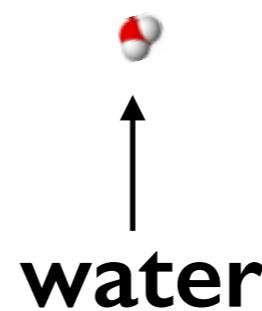
proteins



protein  
representations

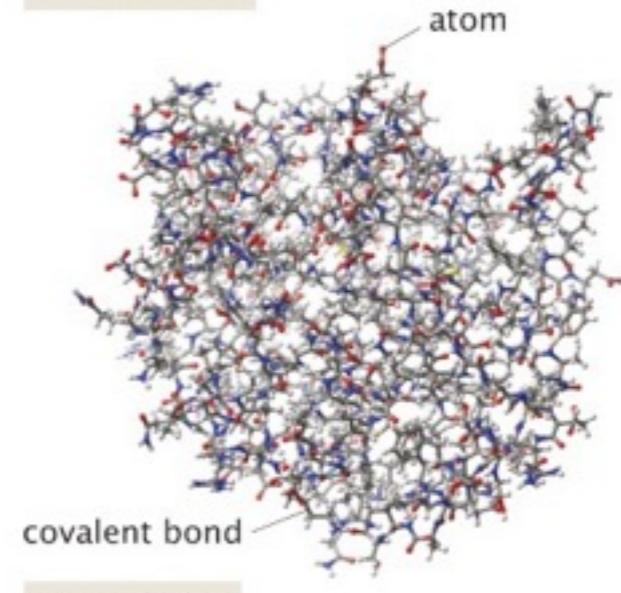


lipids

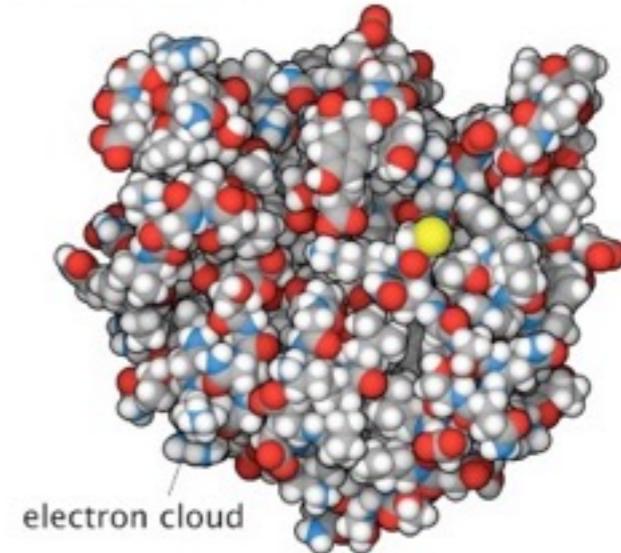


water

ball and stick



space-filling



ribbon

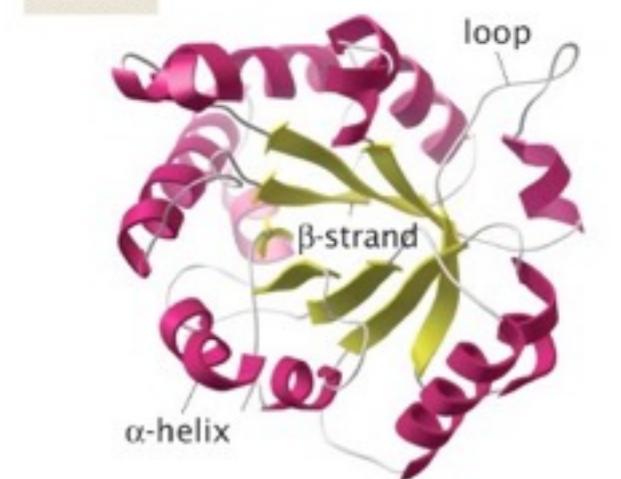
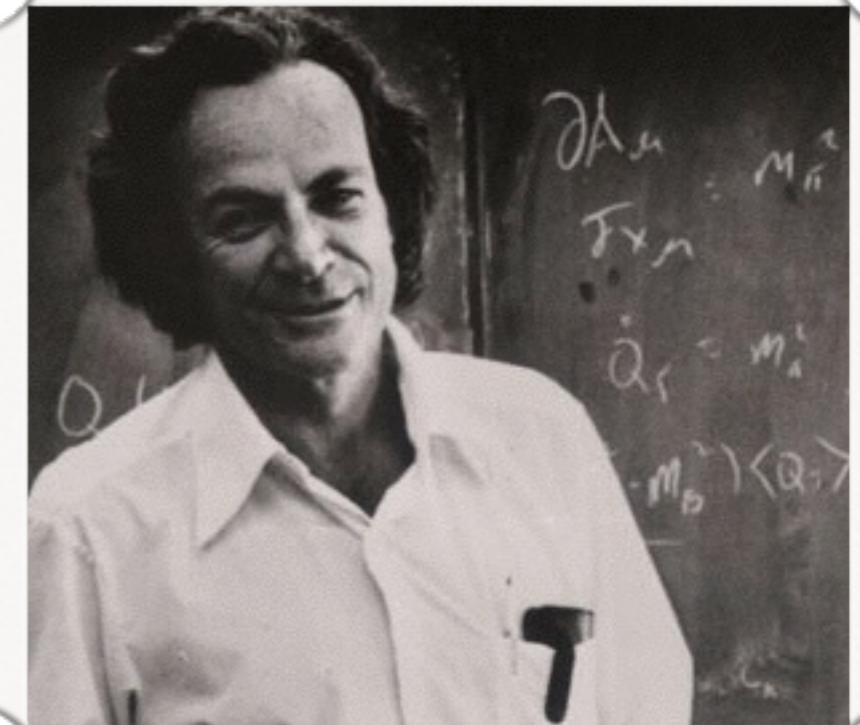


Figure 2.32 Physical Biology of the Cell, 2ed.

*“Everything that living things do can be understood  
in terms of the jiggling and wiggling of atoms.”*—  
Richard Feynman\*



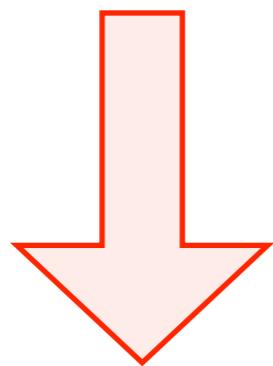
*Richard P. Feynman*

Wikimedia Commons

\* R.P. Feynman. *The Feynman Lectures on Physics*. 1963

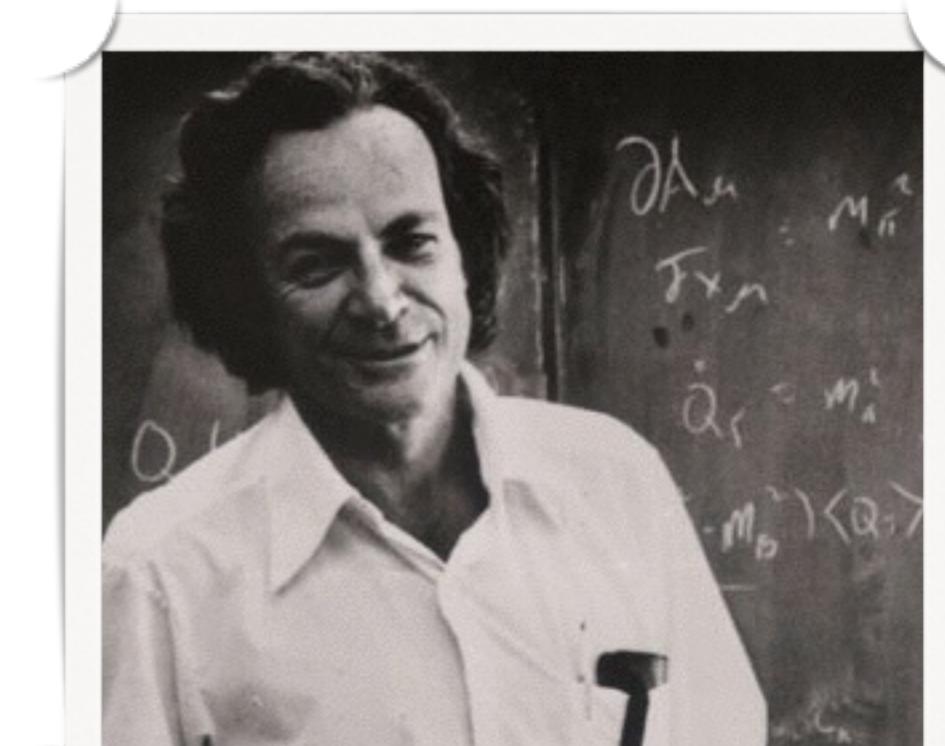
*“Everything that living things do can be understood  
in terms of the jiggling and wiggling of atoms.”* —

Richard Feynman\*



$$(\mathbf{r}_1(t), \dots, \mathbf{r}_N(t))$$

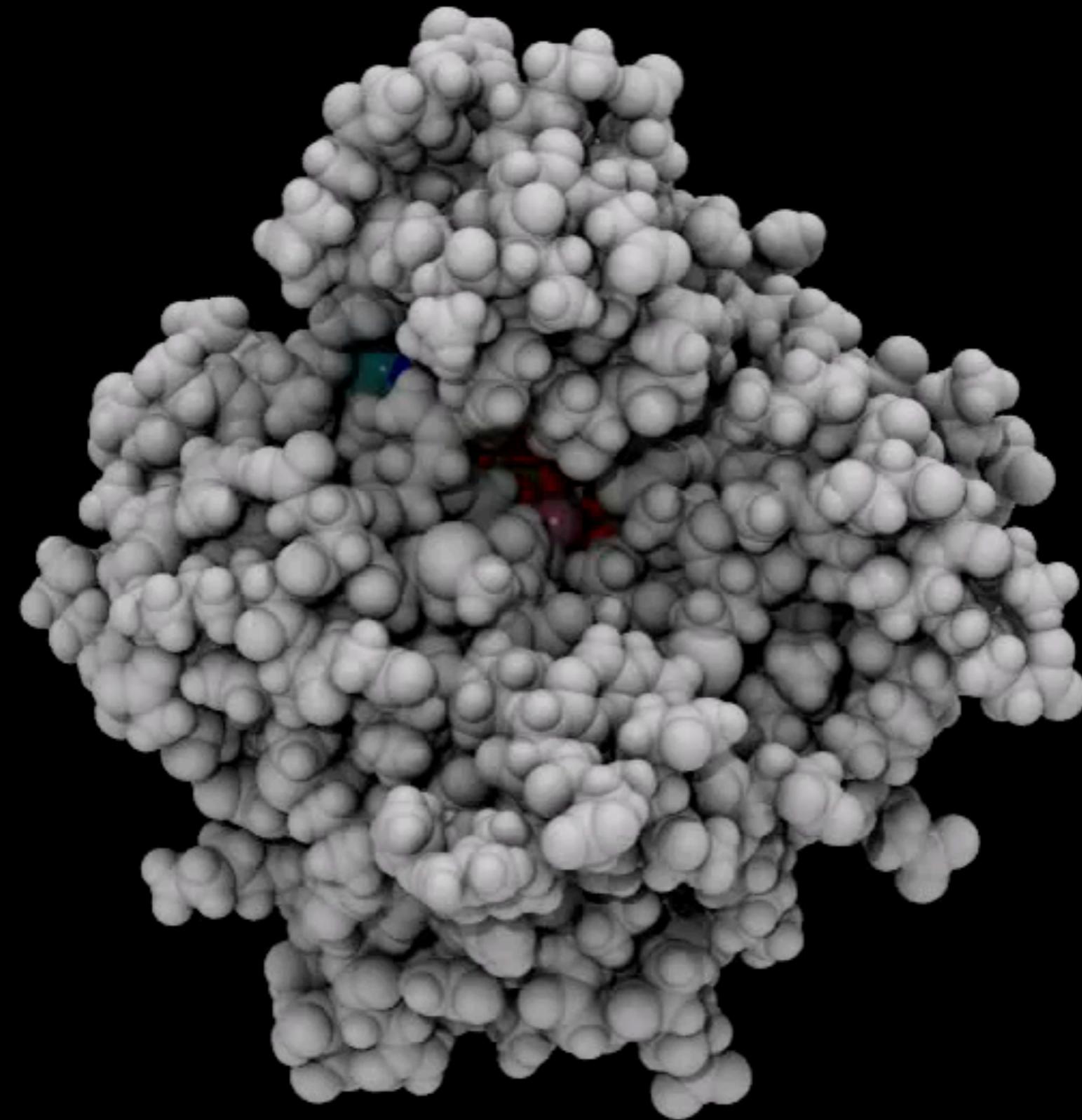
positions of all  $N$  atoms over time



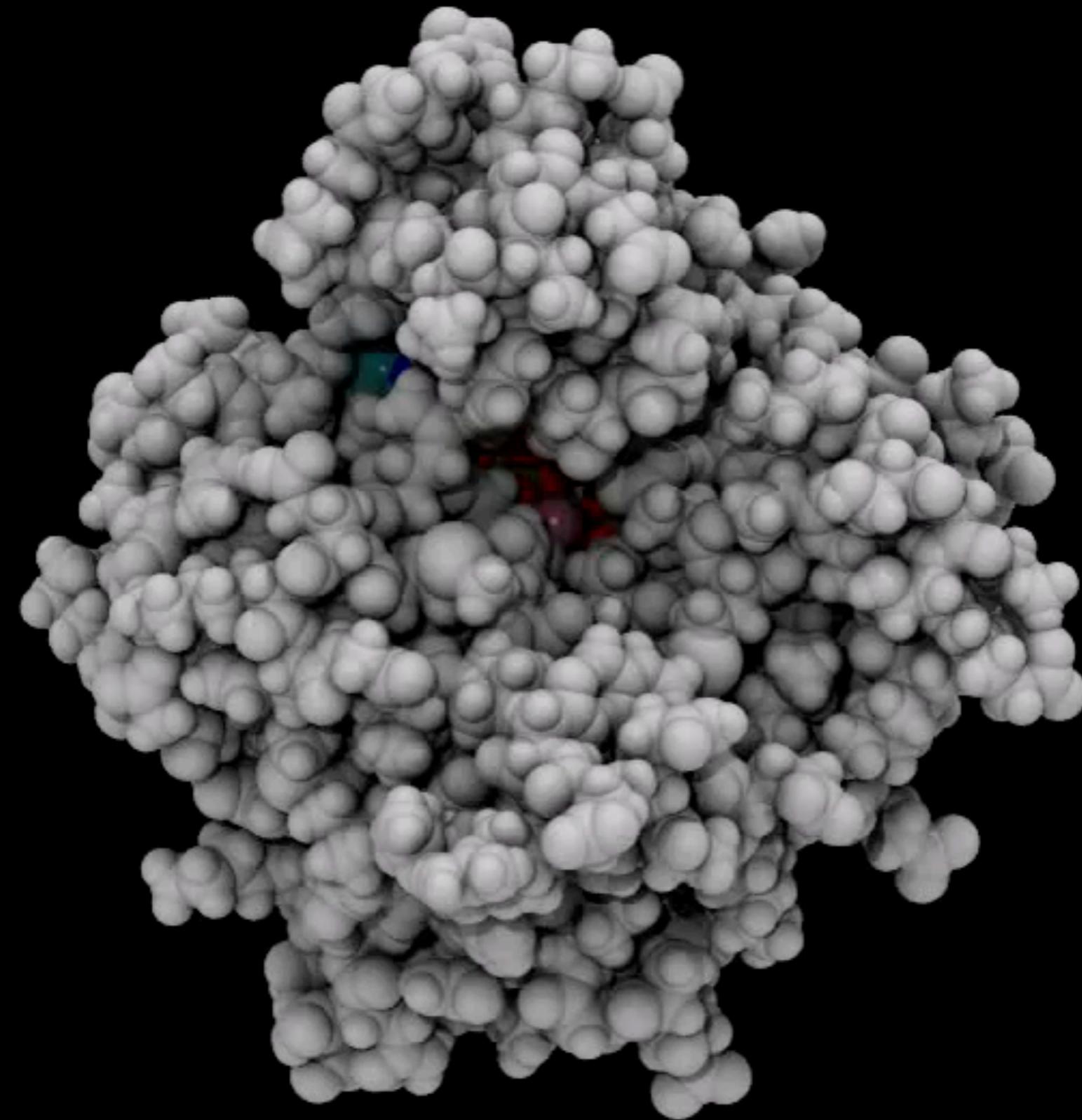
*Richard P. Feynman*

Wikimedia Commons

# Molecular Dynamics (MD) Simulations



# Molecular Dynamics (MD) Simulations



# **Molecular Dynamics (MD) Simulations**

**(classical)**

# Molecular Dynamics (MD) Simulations

(classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

← energy function = “force field”

# Molecular Dynamics (MD) Simulations

(classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

← energy function = “force field”

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

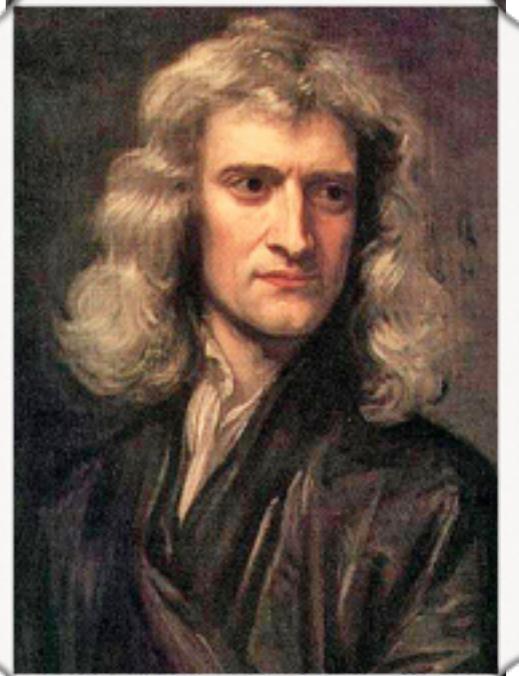
# Molecular Dynamics (MD) Simulations

(classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

← energy function = “force field”

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$



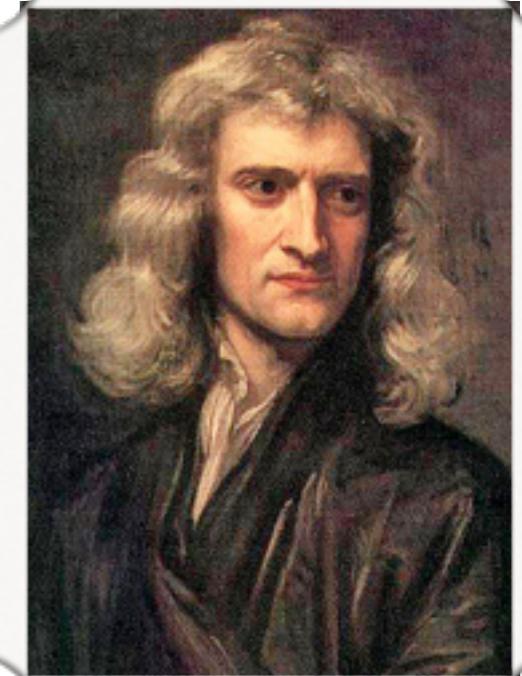
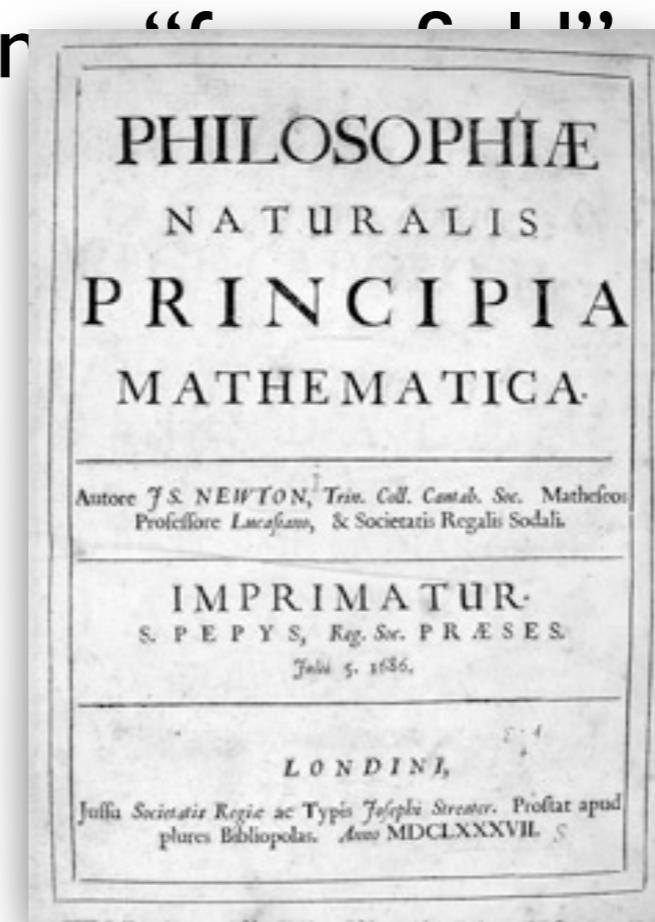
*Isaac Newton*

# Molecular Dynamics (MD) Simulations (classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$



J. S. Newton

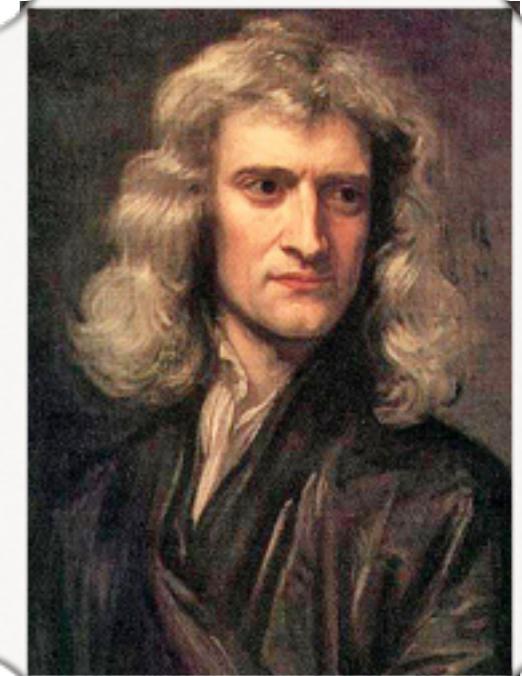
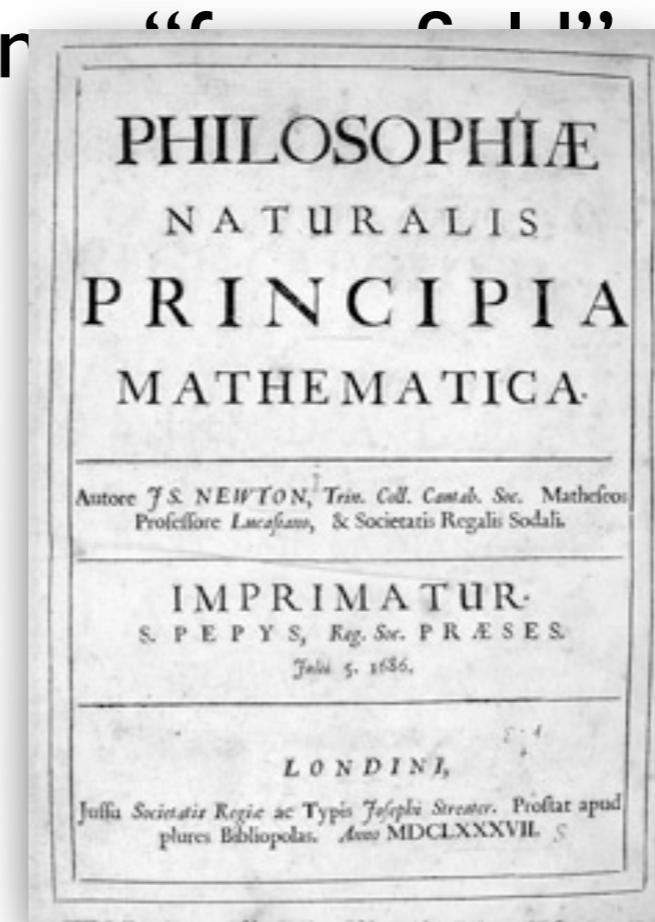
# Molecular Dynamics (MD) Simulations (classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

Newton's  
2nd law



J. S. Newton

# Molecular Dynamics (MD) Simulations (classical)

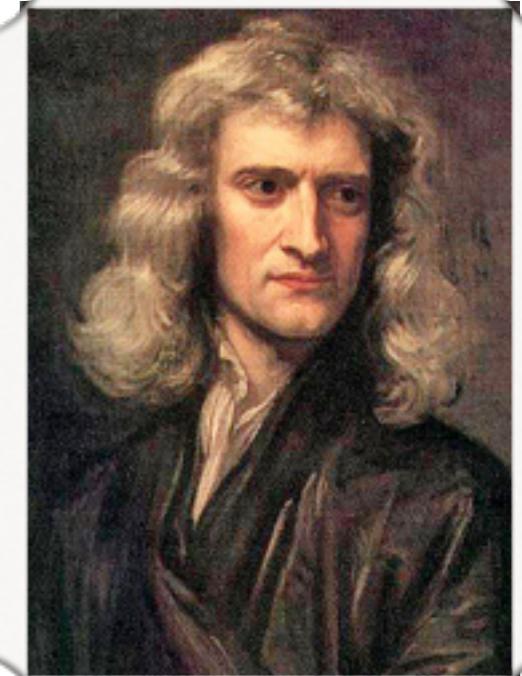
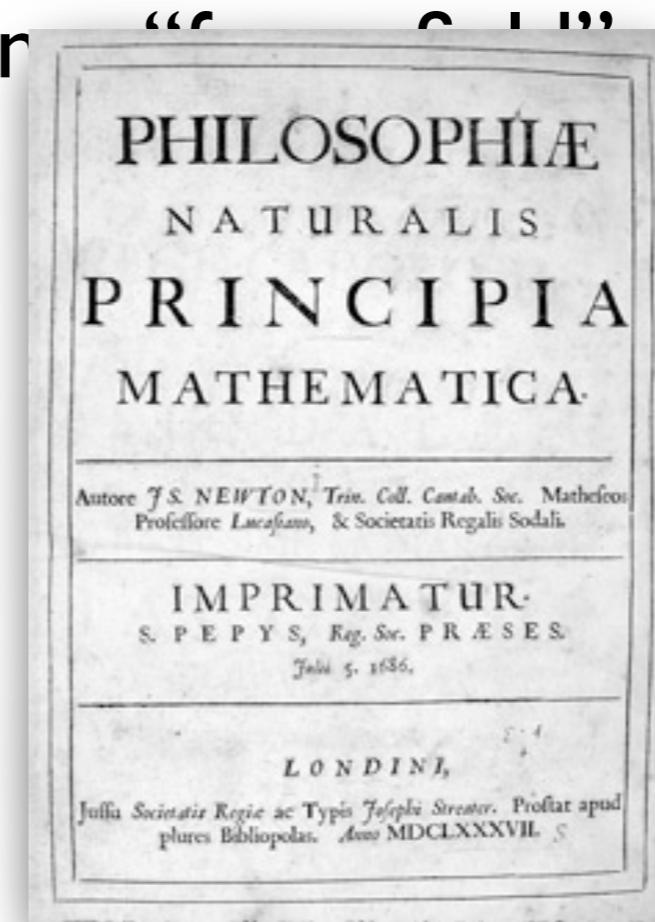
$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

Newton's  
2nd law

$$\frac{d^2 \mathbf{r}_i}{dt^2} = \frac{\mathbf{F}_i}{m_i}$$



J.S. Newton

# Molecular Dynamics (MD) Simulations (classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots)$$

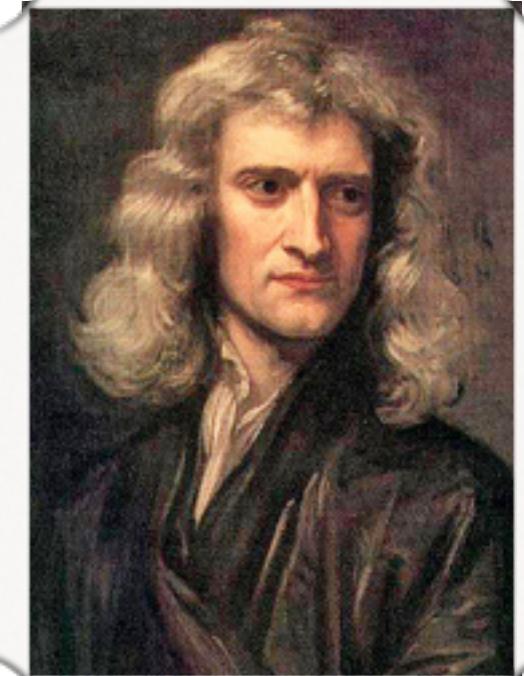
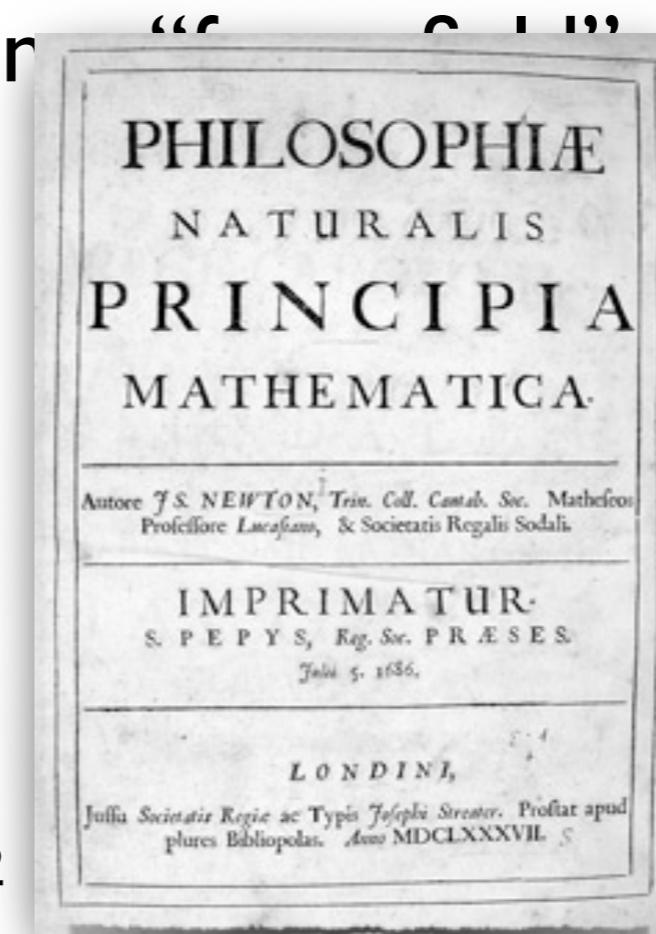
energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

Newton's  
2nd law

$$\frac{d^2 \mathbf{r}_i}{dt^2} = \frac{\mathbf{F}_i}{m_i}$$

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i}{m_i} \Delta t^2$$



J.S. Newton

# Molecular Dynamics (MD) Simulations (classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots, \mathbf{r}_N)$$

energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

$$\mathbf{F}_i = m_i \mathbf{a}_i$$

Newton's  
2nd law

$$\frac{d^2 \mathbf{r}_i}{dt^2} = \frac{\mathbf{F}_i}{m_i}$$

integrator

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i}{m_i} \Delta t^2$$



# Molecular Dynamics (MD) Simulations (classical)

$$U(\mathbf{r}_1, \dots, \mathbf{r}_N) = U_{\text{bonded}}(\mathbf{r}_1, \dots) + U_{\text{non-bonded}}(\mathbf{r}_1, \dots, \mathbf{r}_N)$$

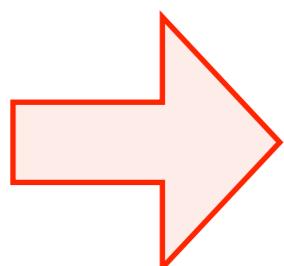
energy function

$$\mathbf{F}_i = -\frac{\partial}{\partial \mathbf{r}_i} U(\mathbf{r}_1, \dots, \mathbf{r}_i, \dots, \mathbf{r}_N)$$

$$\mathbf{F}_i = m_i \mathbf{a}_i \quad \xleftarrow{\text{Newton's 2nd law}}$$

$$\frac{d^2 \mathbf{r}_i}{dt^2} = \frac{\mathbf{F}_i}{m_i} \quad \xleftarrow{\text{integrator}}$$

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i}{m_i} \Delta t^2$$



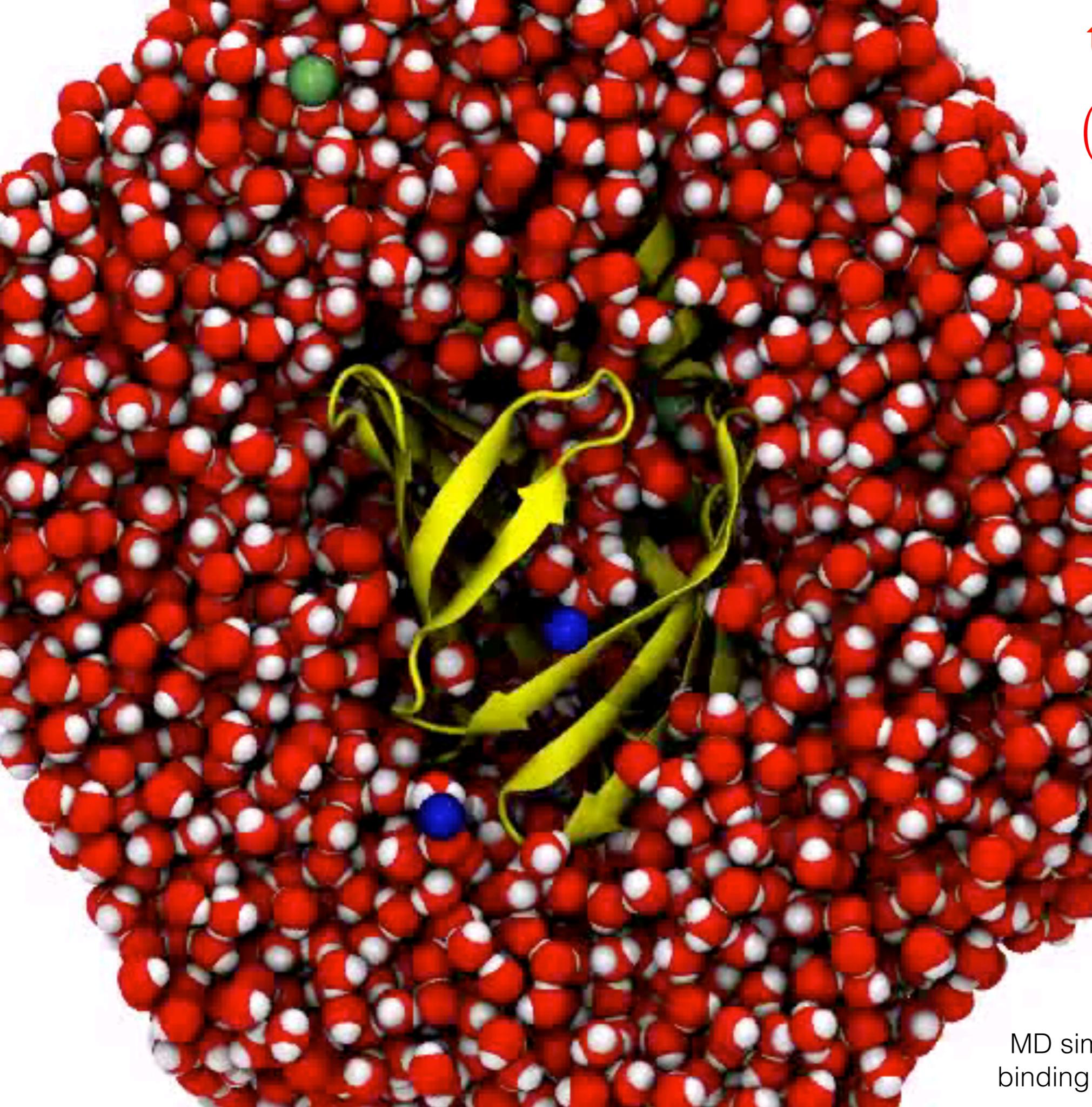
- $\mathbf{r}_1(0), \dots, \mathbf{r}_N(0)$
- $\mathbf{r}_1(\Delta t), \dots, \mathbf{r}_N(\Delta t)$
- $\mathbf{r}_1(2\Delta t), \dots, \mathbf{r}_N(2\Delta t)$
- $\mathbf{r}_1(3\Delta t), \dots, \mathbf{r}_N(3\Delta t)$

⋮

trajectory

$(\mathbf{r}_1(t), \dots, \mathbf{r}_N(t))$

$0 \leq t \leq \tau$



**trajectory**

$$(\mathbf{r}_1(t), \dots, \mathbf{r}_N(t))$$

$$0 \leq t \leq \tau$$

1 protein 

12 ions 

3432 waters



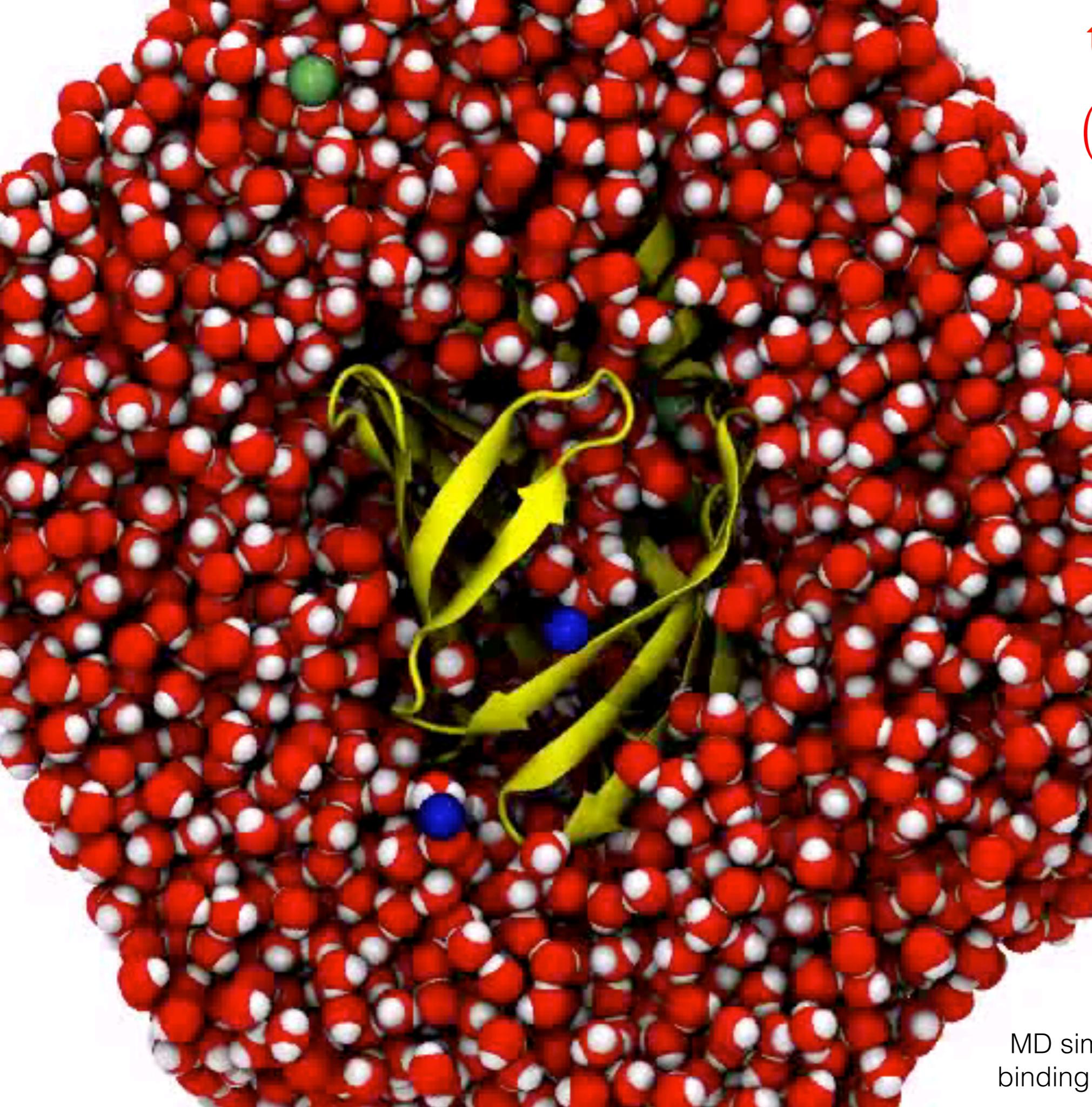
2113 atoms

12 atoms

10296 atoms

**12421 atoms**

MD simulation of intestinal fatty acid binding protein. Rendered with VMD.



**trajectory**

$$(\mathbf{r}_1(t), \dots, \mathbf{r}_N(t))$$

$$0 \leq t \leq \tau$$

1 protein 

12 ions 

3432 waters



2113 atoms

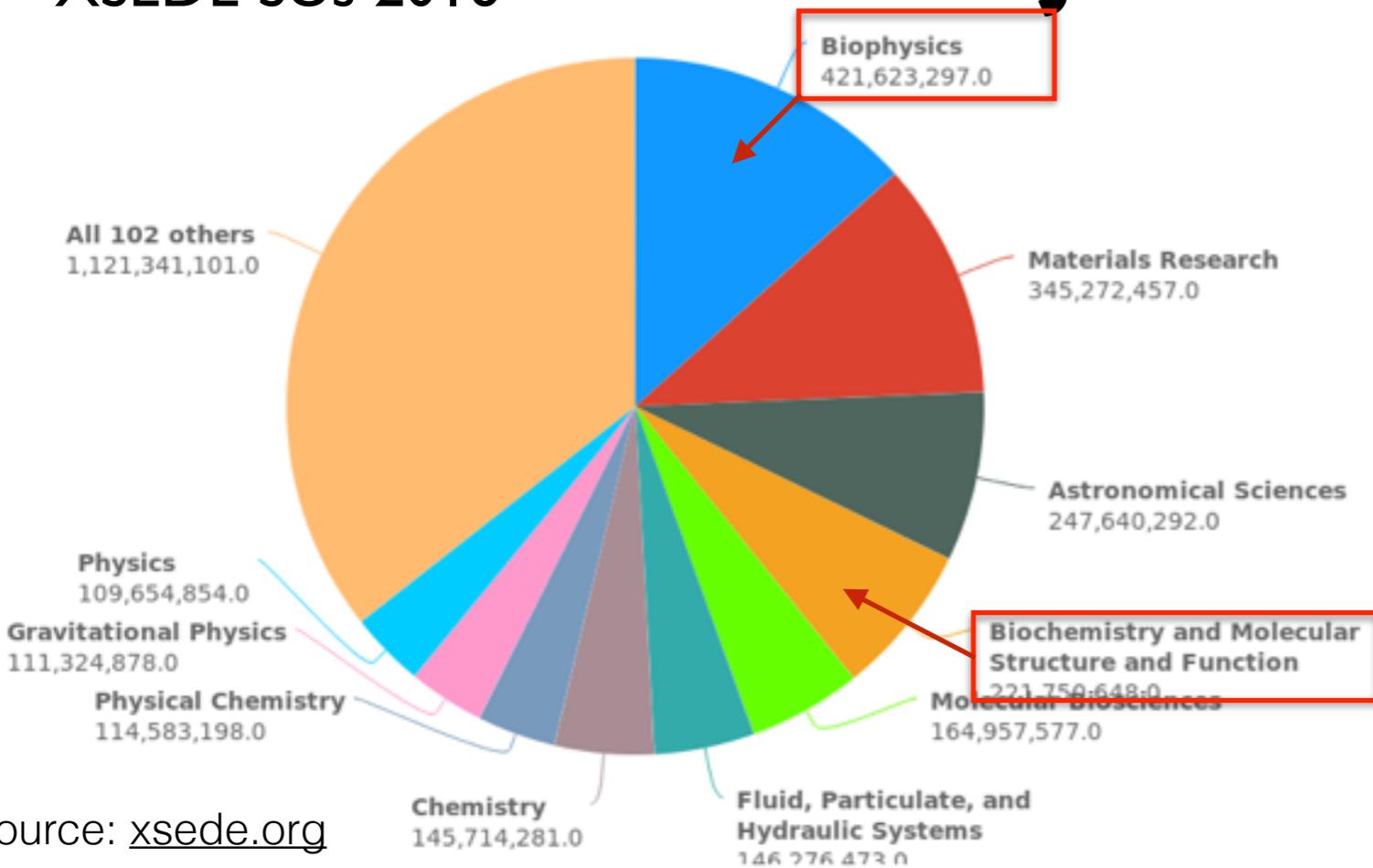
12 atoms

10296 atoms

**12421 atoms**

MD simulation of intestinal fatty acid binding protein. Rendered with VMD.

# MD trajectories can be big-ish

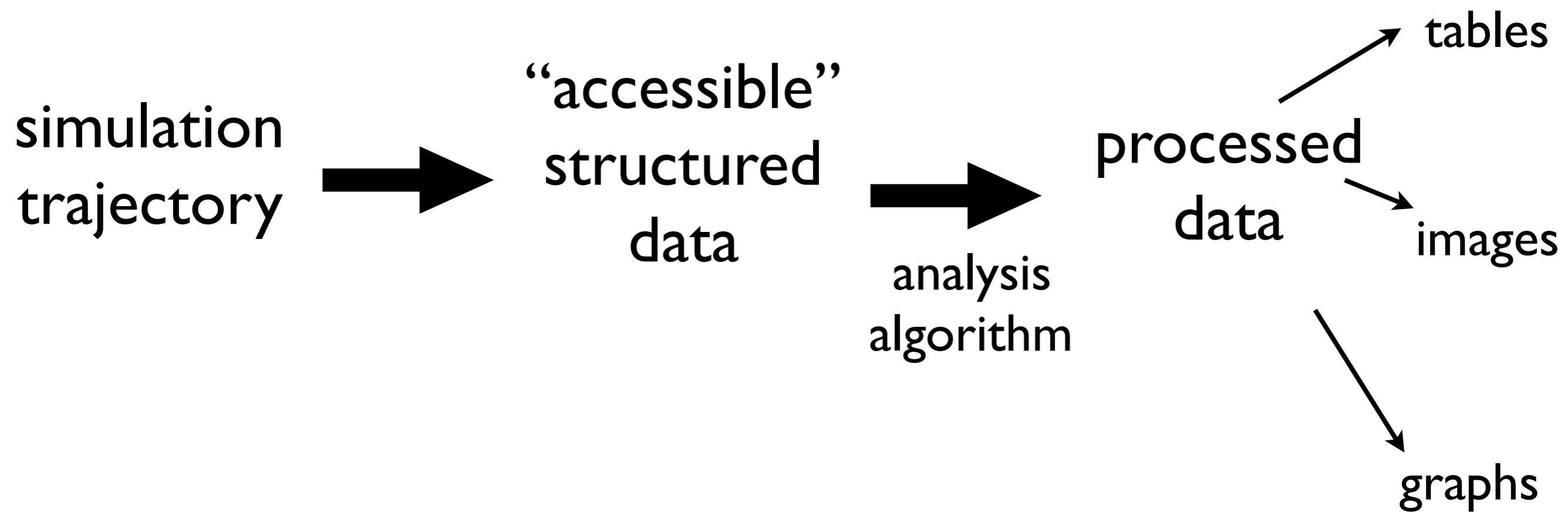


trajectory

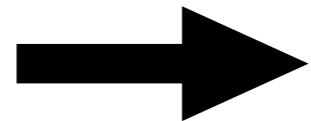
$$(\mathbf{r}_1(t), \dots, \mathbf{r}_N(t))$$

$$0 \leq t \leq \tau$$

	max (2016)	typical (2016)
atoms $N$	$\sim 10^7$	$\sim 10^5$
simulated time $\tau$	$\sim 10 \mu\text{s}$	$0.1\text{--}1 \mu\text{s}$
trajectory frames	$\sim 10^9$	$\sim 10^5$
trajectory size	< 10 TiB	150 GiB



simulation  
trajectory



“accessible”  
structured  
data



analysis  
algorithm

processed  
data

tables

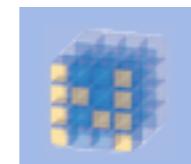
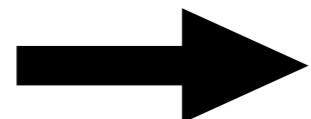
images

graphs

# Basic Idea: Use NumPy



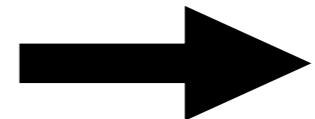
simulation  
trajectory



NumPy

python™

“accessible”  
structured  
data



analysis  
algorithm

processed  
data

tables

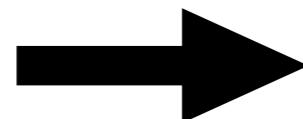
images

graphs

# Basic Idea: Use NumPy

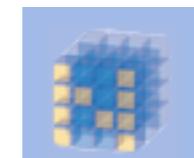


simulation  
trajectory



dcd, xtc, trr,  
ncdf, trj, pdb,  
pqr, gro, crd,  
dms, trz, mol2,  
xyz, config,  
history, gms, ...

psf, tpr,  
prmtop, dms,  
mol2, hoomd  
xml, ...

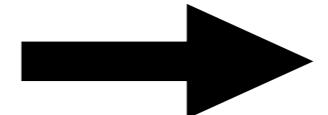


NumPy



python™

“accessible”  
structured  
data



analysis  
algorithm

processed  
data

tables

images

graphs

# Basic Idea: Use NumPy



simulation  
trajectory

dcd, xtc, trr,  
ncdf, trj, pdb,  
pqr, gro, crd,  
dms, trz, mol2,  
xyz, config,  
history, gms, ...

psf, tpr,  
prmtop, dms,  
mol2, hoomd  
xml, ...



“accessible”  
structured  
data

analysis  
algorithm

processed  
data

tables

images

graphs



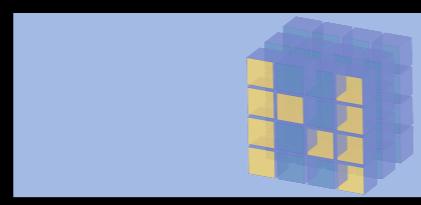
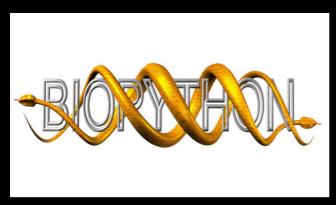
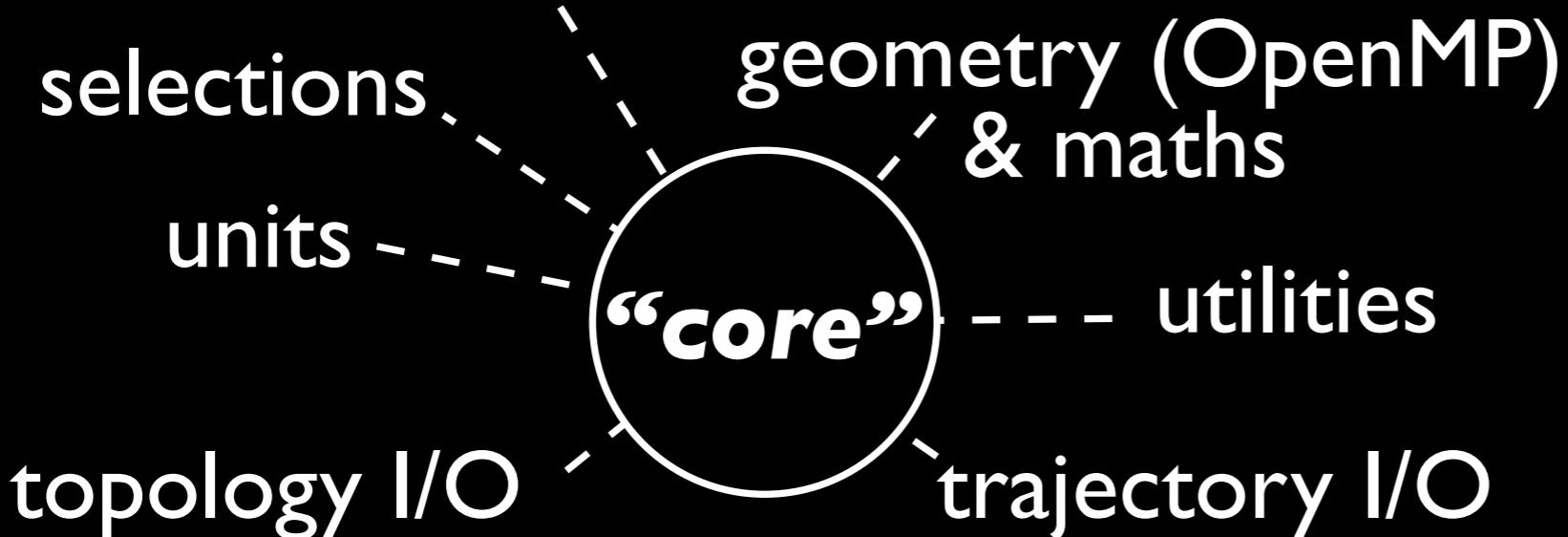
# MDAnalysis

<http://mdanalysis.org>

Universe  
AtomGroup  
*(main data structures in the user interface)*

MDAnalysis.analysis

MDAnalysis.visualization



NumPy

**Code base:**

- python 2.7
- py 3 (~80%)
- cython
- C
- ~42k LOC
- ~24k lines comments



# MDAnalysis

<http://mdanalysis.org>

Universe  
AtomGroup

*(main data  
structures in the  
user interface)*

MDAnalysis.  
analysis

MDAnalysis.  
visualization

selections

units

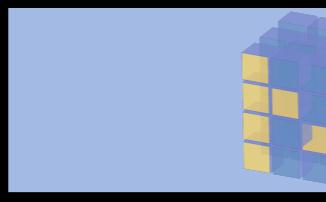
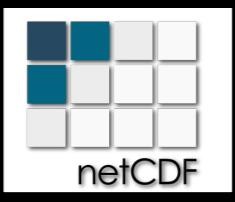
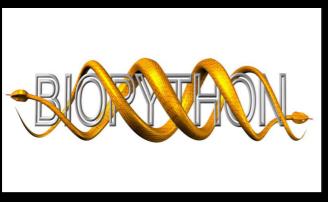
topology I/O

geometry (OpenMP)  
& maths

“core”

utilities

trajectory I/O



NumPy

Code base:

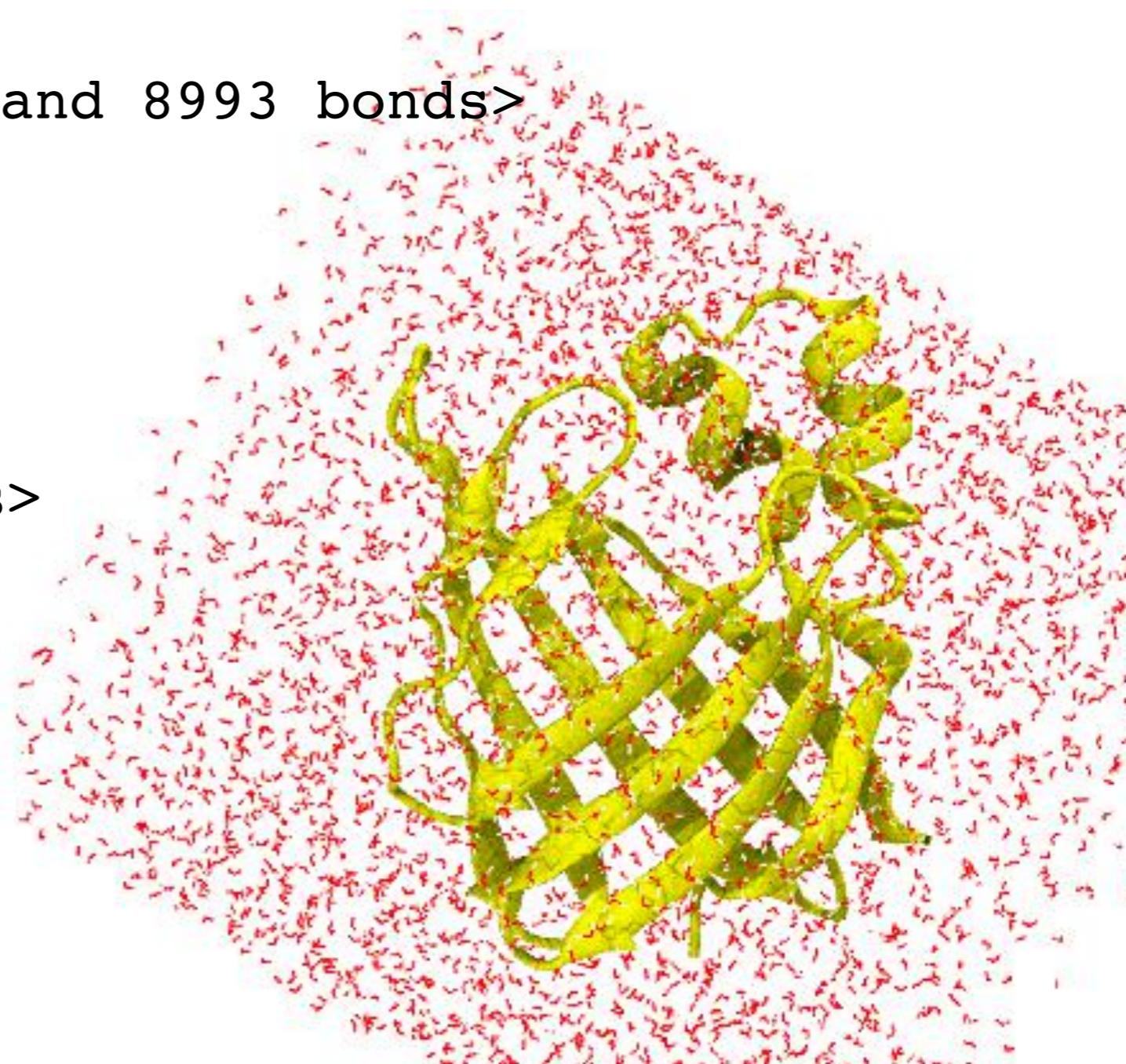
- python 2.7
- py 3 (~80%)
- cython
- C
- ~42k LOC
- ~24k lines comments

# Fundamental data structures: *Universe*

```
import MDAnalysis as mda
u = mda.Universe('topol.tpr', 'traj.trr')

print(u)
<Universe with 12421 atoms and 8993 bonds>

print(u.atoms)
<AtomGroup with 12421 atoms>
```



# Fundamental data structures: *AtomGroup*

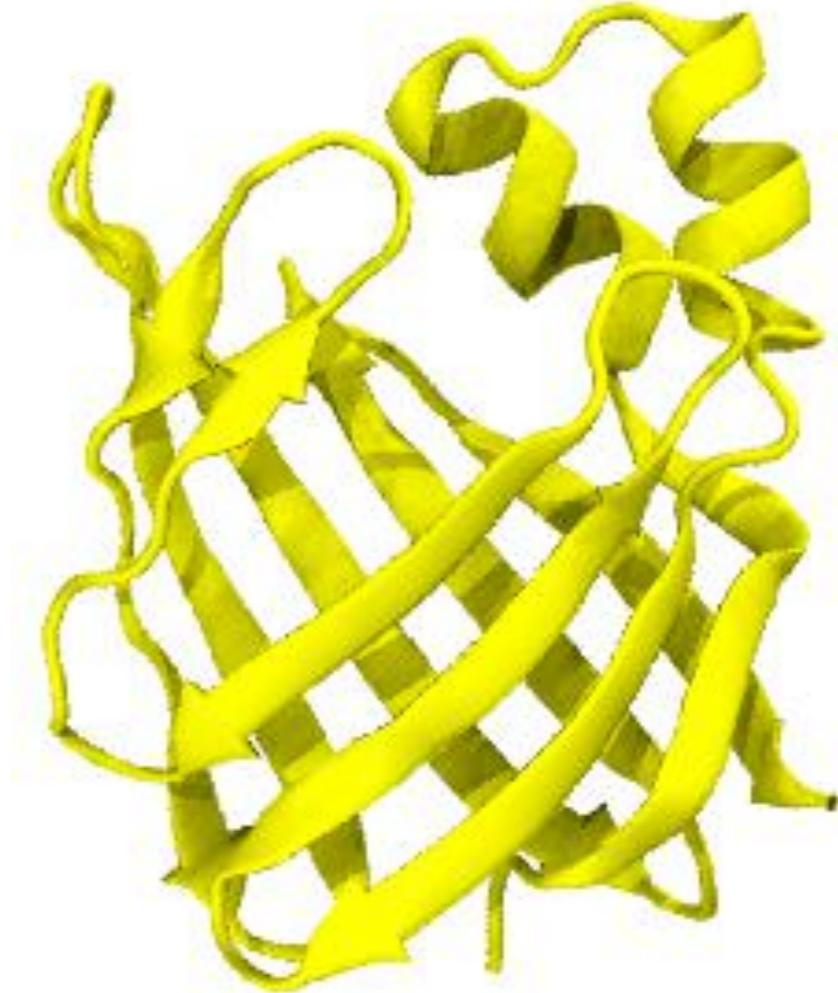
```
protein =  
u.atoms.select_atoms(  
    "protein")
```

```
print(protein)
```

```
<AtomGroup with 2113 atoms>
```

```
print(list(protein[:5]))
```

```
[<Atom 1: N of type NH3 of resname ALA, resid 1 and segid IFAB>,  
<Atom 2: HT1 of type HC of resname ALA, resid 1 and segid IFAB>,  
<Atom 3: HT2 of type HC of resname ALA, resid 1 and segid IFAB>,  
<Atom 4: HT3 of type HC of resname ALA, resid 1 and segid IFAB>,  
<Atom 5: CA of type CT1 of resname ALA, resid 1 and segid IFAB>]]
```

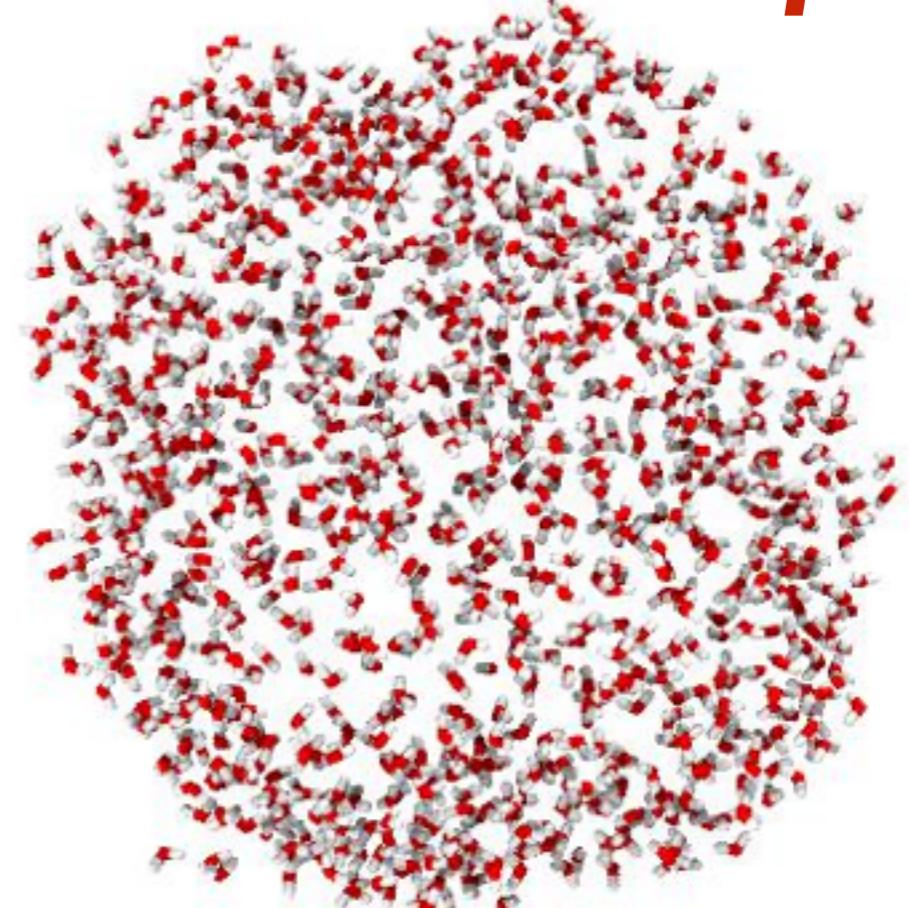


# Fundamental data structures: *AtomGroup*

```
solvshell =  
    u.atoms.select_atoms(  
        "resname SOL and  
        around 5.0 protein")
```

```
print(solvshell)
```

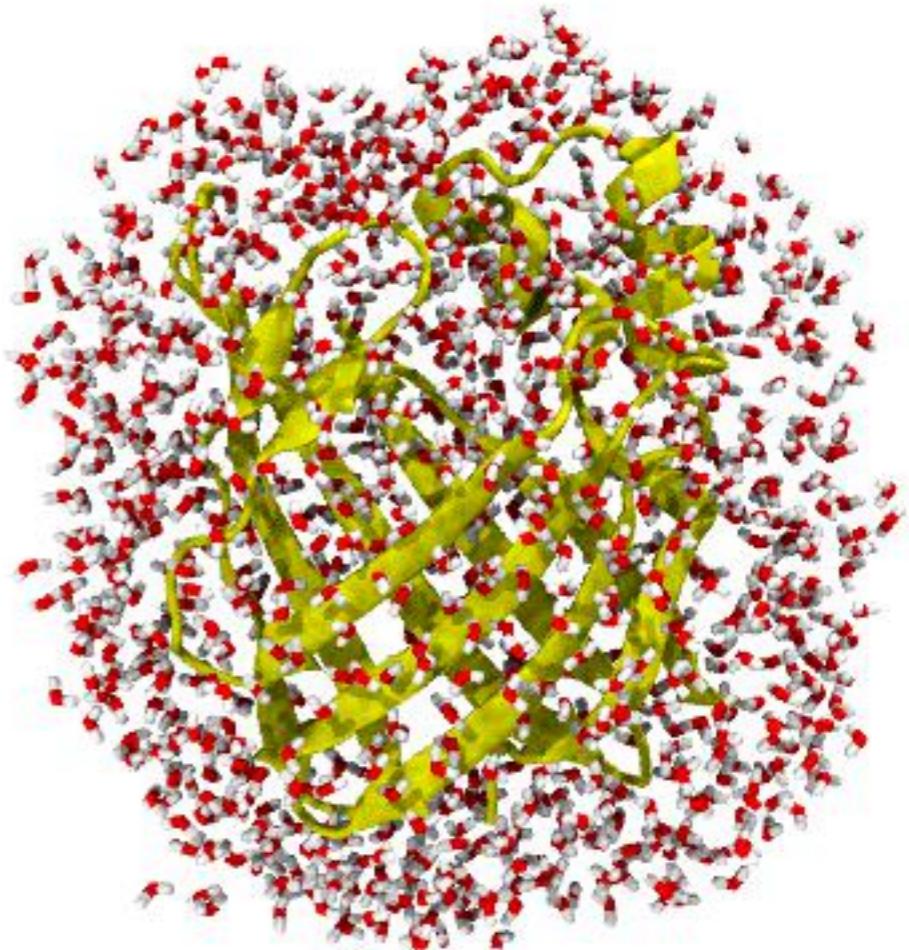
```
<AtomGroup with 2792 atoms>
```



```
ag = protein + solvshell
```

```
print(ag)
```

```
<AtomGroup with 4905 atoms>
```



# AtomGroups: Residues and Segments

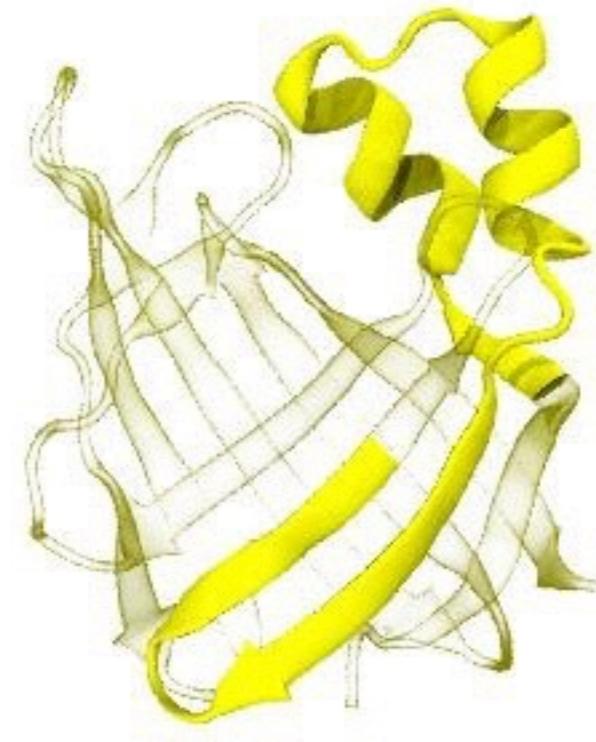
```
protein.residues[10:50]
```

```
list(protein.residues[10:50])
```

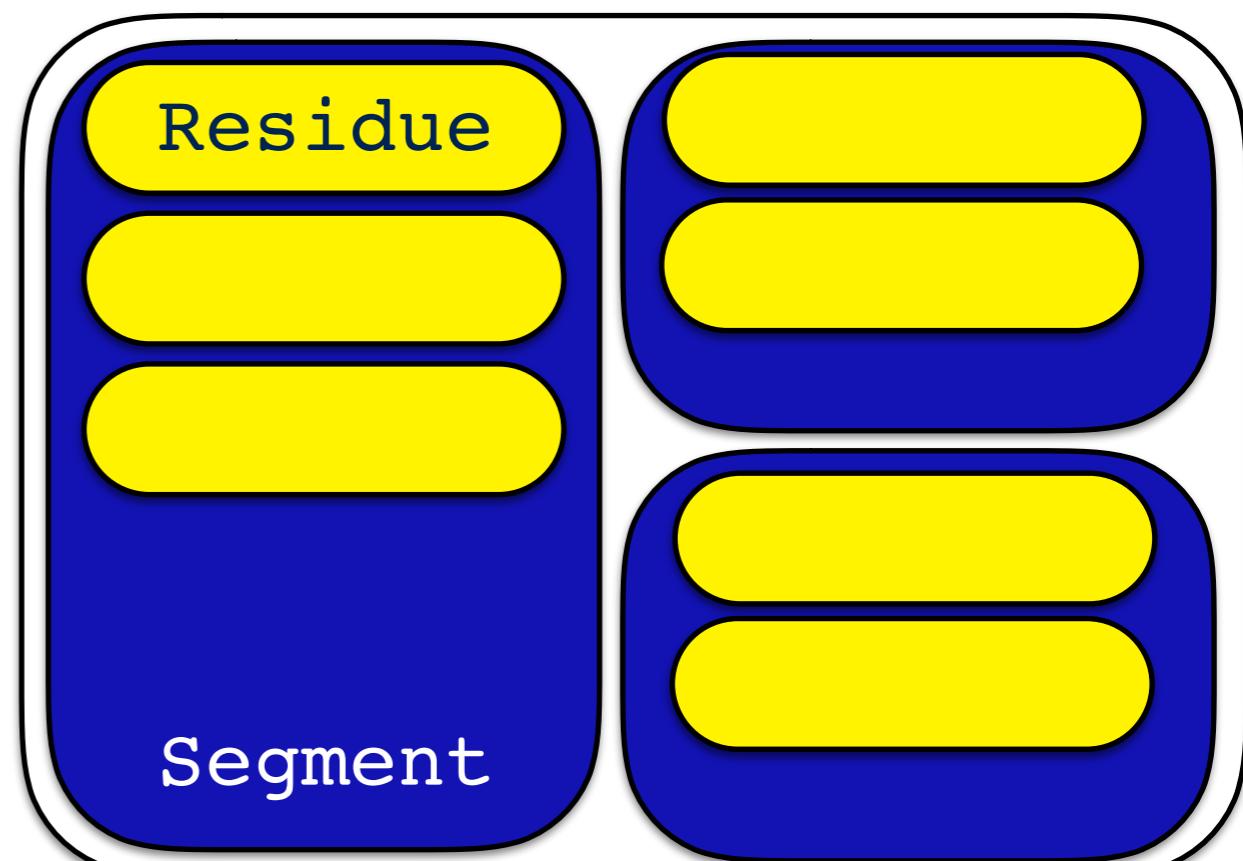
```
[<Residue ASN, 11>, <Residue GLU, 12>,
<Residue ASN, 13>, <Residue TYR, 14>,
<Residue GLU, 15>, ...]
```

```
list(protein.segments)
```

```
[<Segment IFAB>]
```



Universe



# Atom data as NumPy arrays

AtomGroups contain particles (“atoms”).

Properties of all particles are **NumPy arrays**:

`ag.names`                           `array(['N', 'HT1', 'HT2', ..., 'OH2', 'H1', 'H2'],  
  dtype='|S4')`

`ag.charges`                       `array([-0.3 , 0.33 , 0.33 ,  
  ...,  
   -0.834, 0.417, 0.417])`

`ag.positions`                       `array([[ -12.57699966, 10.42199993, -5.22900009],  
   [-13.59200001, 10.19900036, -5.19299984],  
   [-12.31599998, 10.22900009, -6.21700001],  
   ...,  
   [ -5.02600002, -12.31200027, 13.30200005],  
   [ -5.45100021, -11.82499981, 12.59500027],  
   [ -4.14099979, -12.47900009, 12.97900009]],  
  dtype=float32)`

`ag.velocities`

`ag.forces`

... and many more

# Basic analysis pattern: Iterate over frames

- trajectories contain *frames*: one snapshot of all particles at a specific time (*positions*[, *velocities*[, *forces*]])
- `Universe.trajectory` is iterable

```
for ts in u.trajectory[::10]:  
    analyze(ag.positions)
```

updates  
every step

- random access to a specific frame  
`u.trajectory[72]`
- number of frames  
`len(u.trajectory)`

# Example analysis: Per-residue RMSF

- root mean square fluctuation  $\rho_i$  (RMSF) measures local flexibility of amino acid  $i$
- standard quantity to compute for protein simulations

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$

- use the “C-alpha” atom in each residue to characterise the motion of the whole residue:  $\mathbf{x}_i$  (position of  $C_{\alpha,i}$ )

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")
means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)
for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)
rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)
for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)
rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

u = mda.Universe("topol.tpr", "trj.xtc")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

# $C_\alpha$ RMSF

$$\rho_i = \sqrt{\langle (\mathbf{x}_i(t) - \langle \mathbf{x}_i \rangle)^2 \rangle}$$



```
import numpy as np
import MDAnalysis as mda

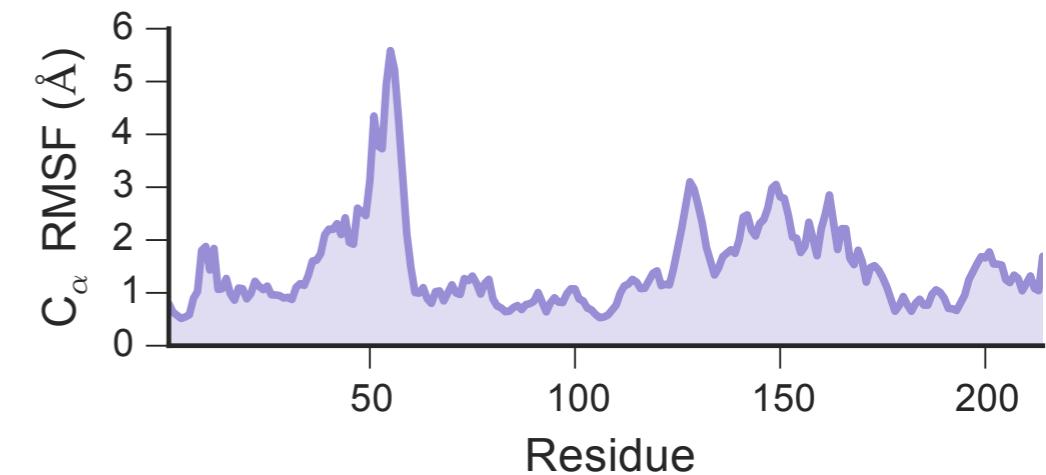
u = mda.Universe("topol.tpr", "trj")
ca = u.select_atoms("name CA")

means = np.zeros((len(ca), 3))
sumsq = np.zeros_like(means)

for k, ts in enumerate(u.trajectory):
    sumsq += (k/(k+1.0)) * (ca.positions - means)**2
    means[:] = (k*means + ca.positions)/(k+1.0)

rmsf = np.sqrt(sumsq.sum(axis=1)/(k+1.0))

matplotlib.pyplot.plot(ca.residues.resids, rmsf)
```

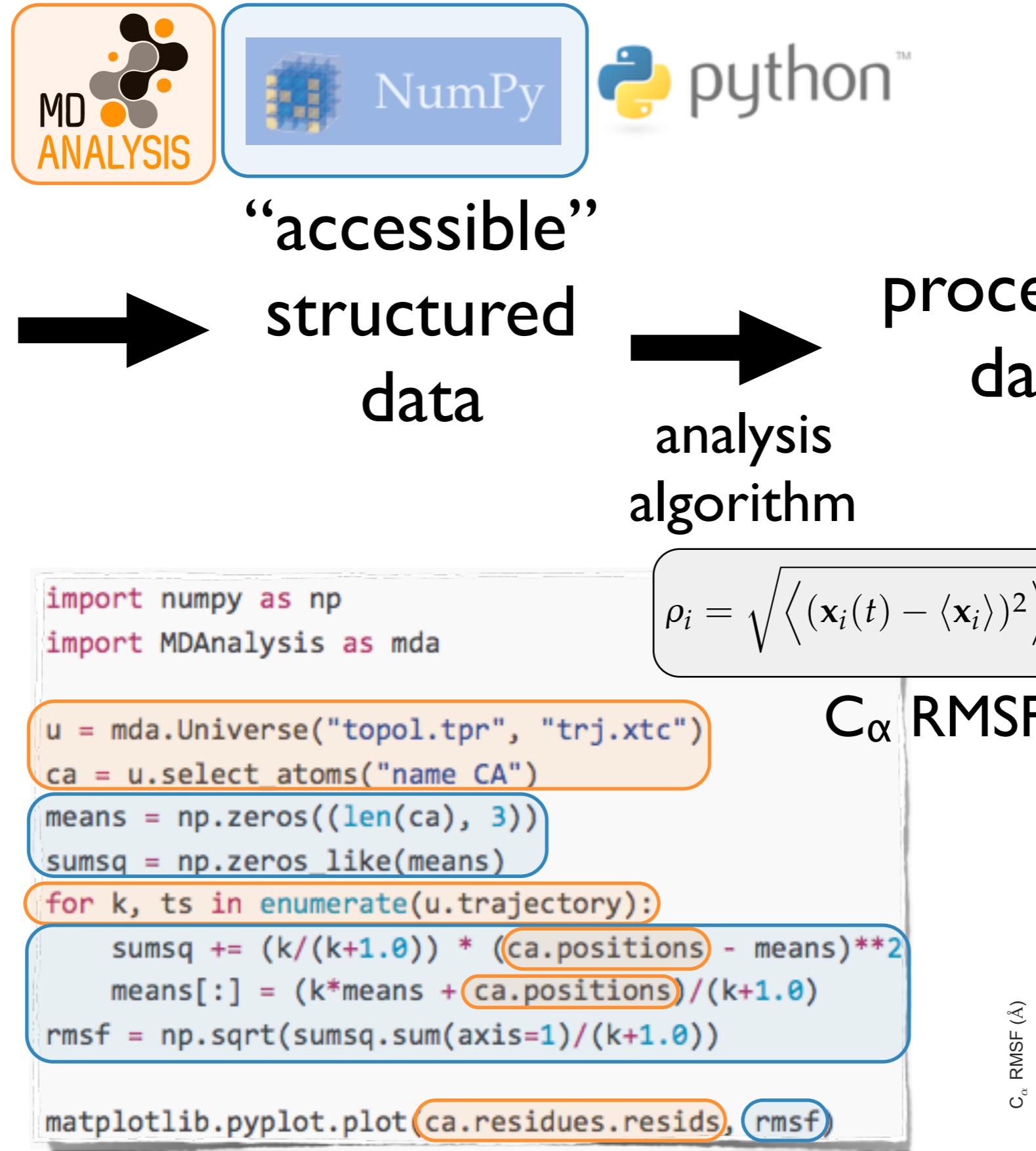


# simulation trajectory

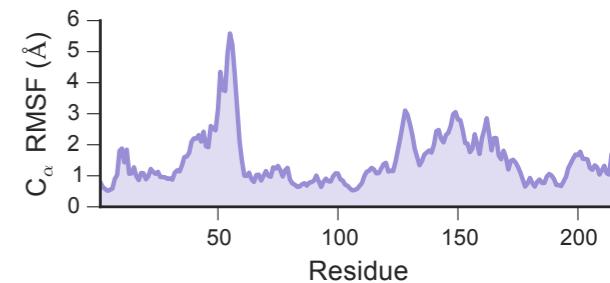
dcd, xtc, trr,  
ncdf, trj, pdb,  
pqr, gro, crd,  
dms, trz, mol2,  
xyz, config,  
history, gms, ...

psf, tpr,  
prmtop, dms,  
mol2, hoomd  
xml, ...

25 different  
formats —  
one analysis  
script



# RESULTS!



# Interactive use

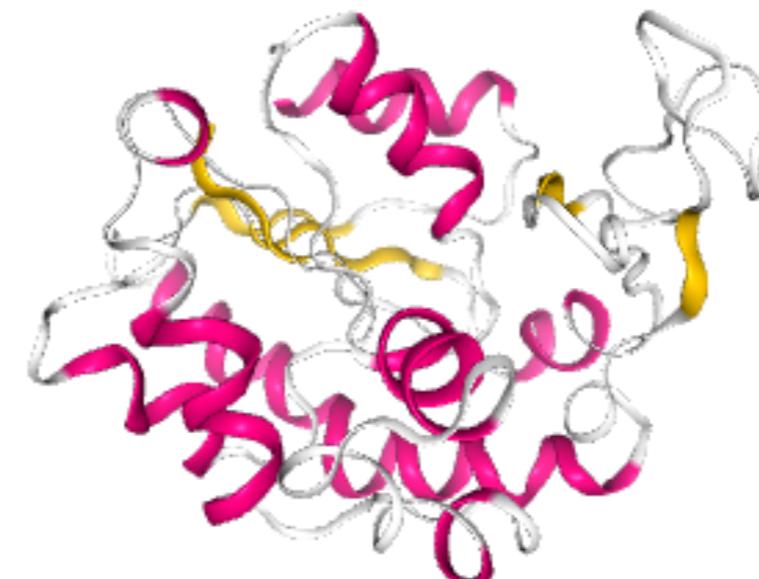
- Jupyter notebooks (+ pandas, ...)
- visualisation with nglviewer

```
In [1]: import nglview as nv
import MDAnalysis as mda
from MDAnalysisTests.datafiles import PSF, DCD
```

```
In [2]: u = mda.Universe(PSF, DCD)
protein = u.select_atoms('protein')
```

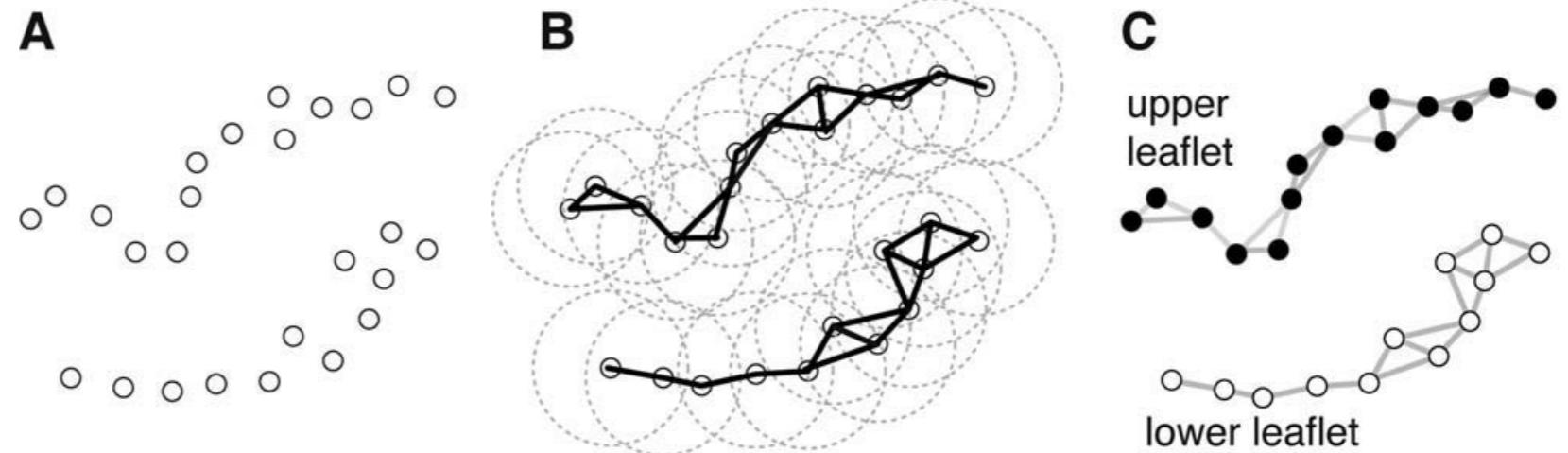
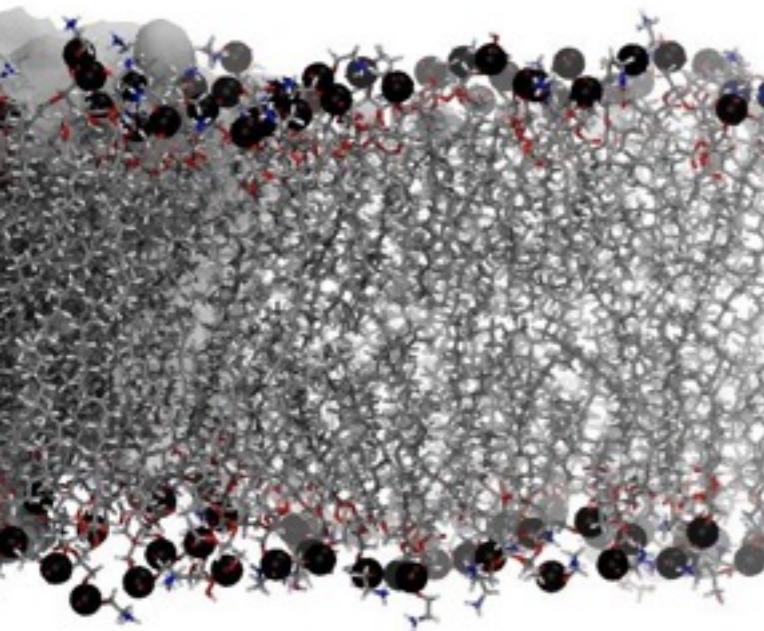
```
In [4]: w = nv.show_mdanalysis(protein)
w
```

x



play

# Interoperability with the Python Ecosystem



*LeafletFinder* algorithm

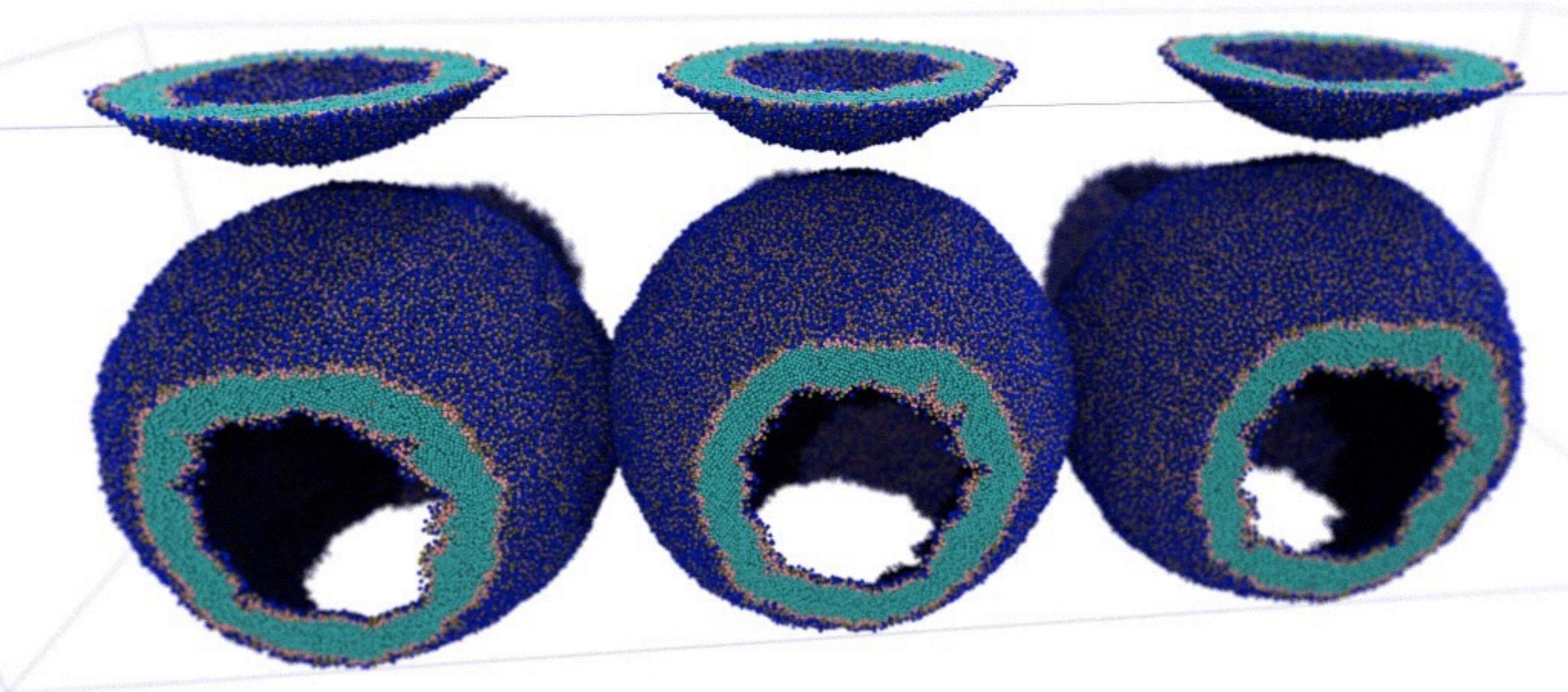
```
import MDAnalysis as mda
import networkx as nx
from MDAnalysis.lib.distances import distance_array
```

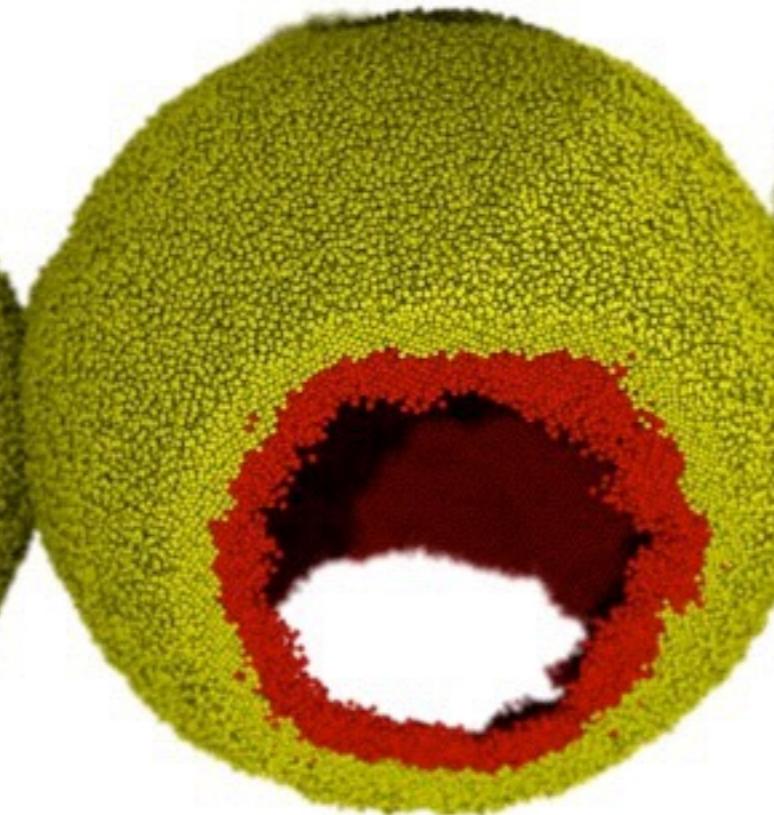
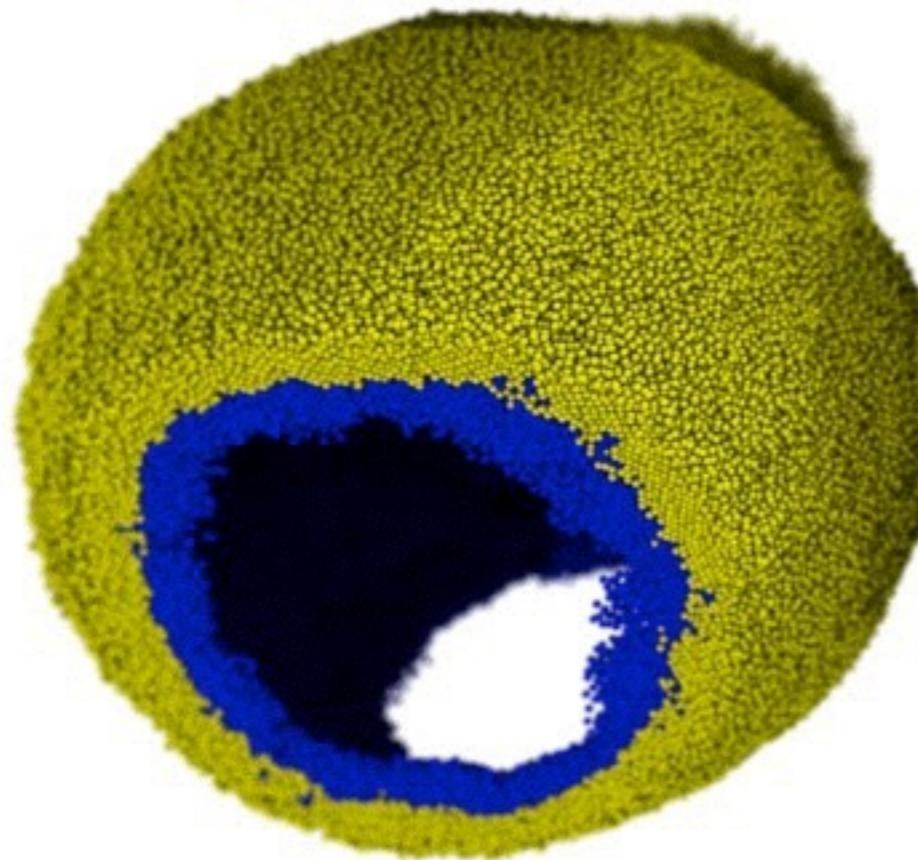
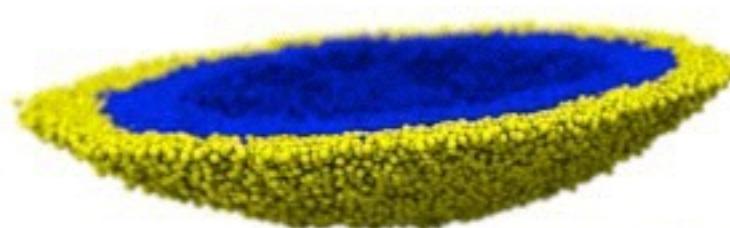
A  
B  
C

```
u = mda.Universe(pdb, xtc)
headgroup_atoms = u.select_atoms("name P*")
x = headgroup_atoms.positions
```

```
adj = (distance_array(x, x) < 12)
leaflets = sorted(nx.connected_components(nx.Graph(adj)), key=len, reverse=True)
```

```
A_lipids = headgroup_atoms[leaflets[0]].residues
B_lipids = headgroup_atoms[leaflets[1]].residues
```







# Analysis module: *MDAnalysis.analysis*

- standard analysis functionality (RMSD, RMSF, distances, density, hydrogen bonds, native contacts...)
- unique capabilities: Path Similarity Analysis, LeafletFinder, water dynamics, HOLE, diffusion map
- moving towards a standard API with an *AnalysisBase* class

```
from MDAnalysis.analysis.rms import RMSF
A = RMSF(protein).run()
A.plot()
print(A.rmsf)

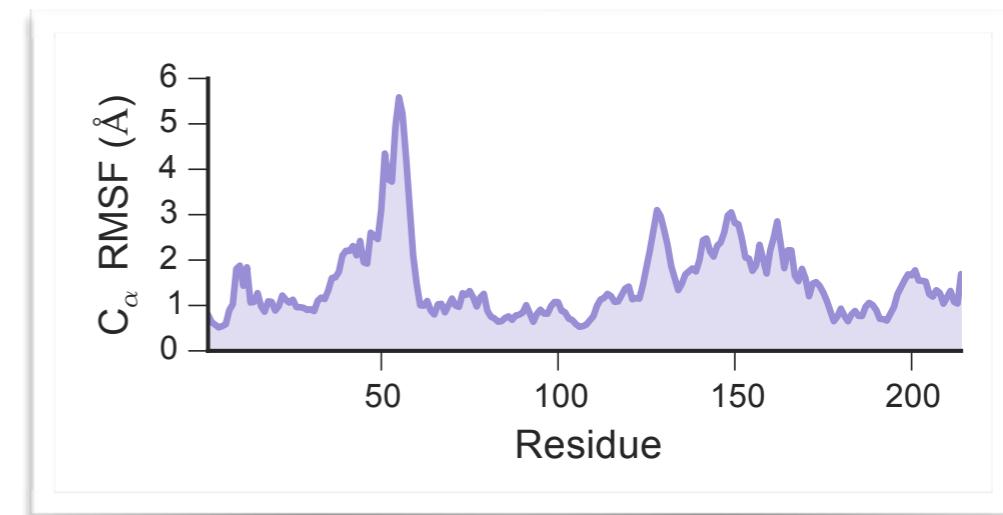
array([ 2.55843466,  3.07868589,
       2.37143333, ...,  3.60916556,
       4.05652126,  3.88028668])
```

# Analysis module: *MDAnalysis.analysis*

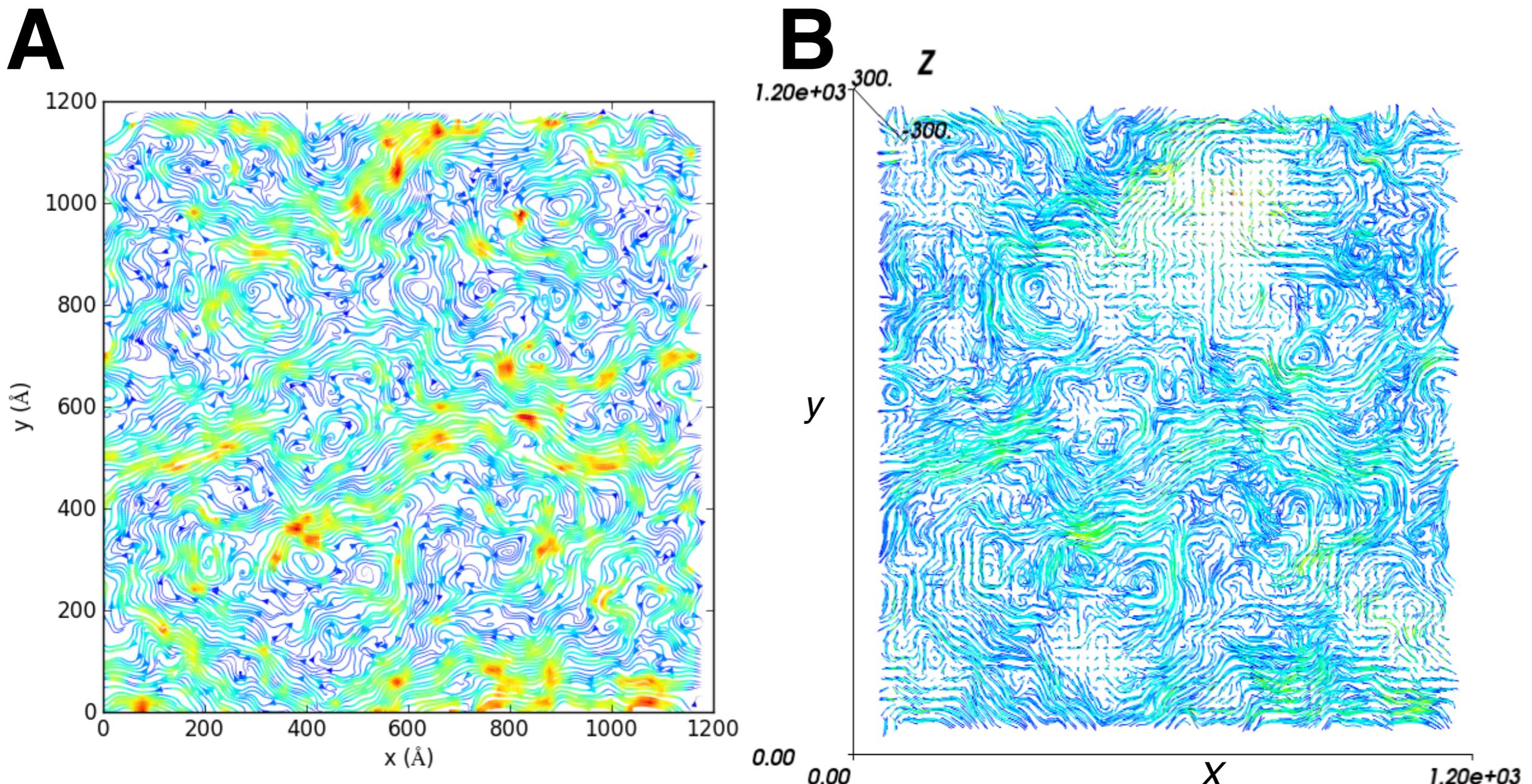
- standard analysis functionality (RMSD, RMSF, distances, density, hydrogen bonds, native contacts...)
- unique capabilities: Path Similarity Analysis, LeafletFinder, water dynamics, HOLE, diffusion map
- moving towards a standard API with an *AnalysisBase* class

```
from MDAnalysis.analysis.rms import RMSF
A = RMSF(protein).run()
A.plot()
print(A.rmsf)

array([ 2.55843466,  3.07868589,
       2.37143333, ...,  3.60916556,
       4.05652126,  3.88028668])
```



# Visualisation module: *MDAnalysis.visualization*



- **streamlines: lipid flow**

M. Chavent, T. Reddy, J. Goose, A. C. E. Dahl, J. E. Stone, B. Jobard, and M. S. P. Sansom. Methodologies for the analysis of instantaneous lipid diffusion in MD simulations of large membrane systems. *Faraday Discuss.*, 169:455–475, 2014. doi: 10.1039/C3FD00145H.

# Recent performance improvements

- originally: *AtomGroup* as a Python list of *Atom* instances
- now:
  - NumPy array of integer atom indices
  - additional arrays for Atom properties
  - Atom only manages access to arrays (but is seldom used)
  - *Topology* instance manages all arrays

Array of Structs → Struct of Arrays

# Recent performance improvements

subselecting AtomGroup from existing one

# atoms	v0.15.0	v0.16.0	speed up
1.75 M	19 ms	0.45 ms	42
3.50 M	18 ms	0.54 ms	33
10.1 M	17 ms	0.45 ms	38

```
ag = u.atoms[np.arange(3, 574393, 7)]
```

accessing attributes

# atoms	v0.15.0	v0.16.0	speed up
1.75 M	250 ms	35 ms	7.1
3.50 M	490 ms	72 ms	6.8
10.1 M	1500 ms	300 ms	5.0

```
resnames = u.atoms.resnames
```

# Recent performance improvements

loading topology file

# atoms	v0.15.0	v0.16.0	speed up
1.75 M	18 s	5 s	3.6
3.50 M	36 s	11 s	3.3
10.1 M	105 s	31 s	3.4

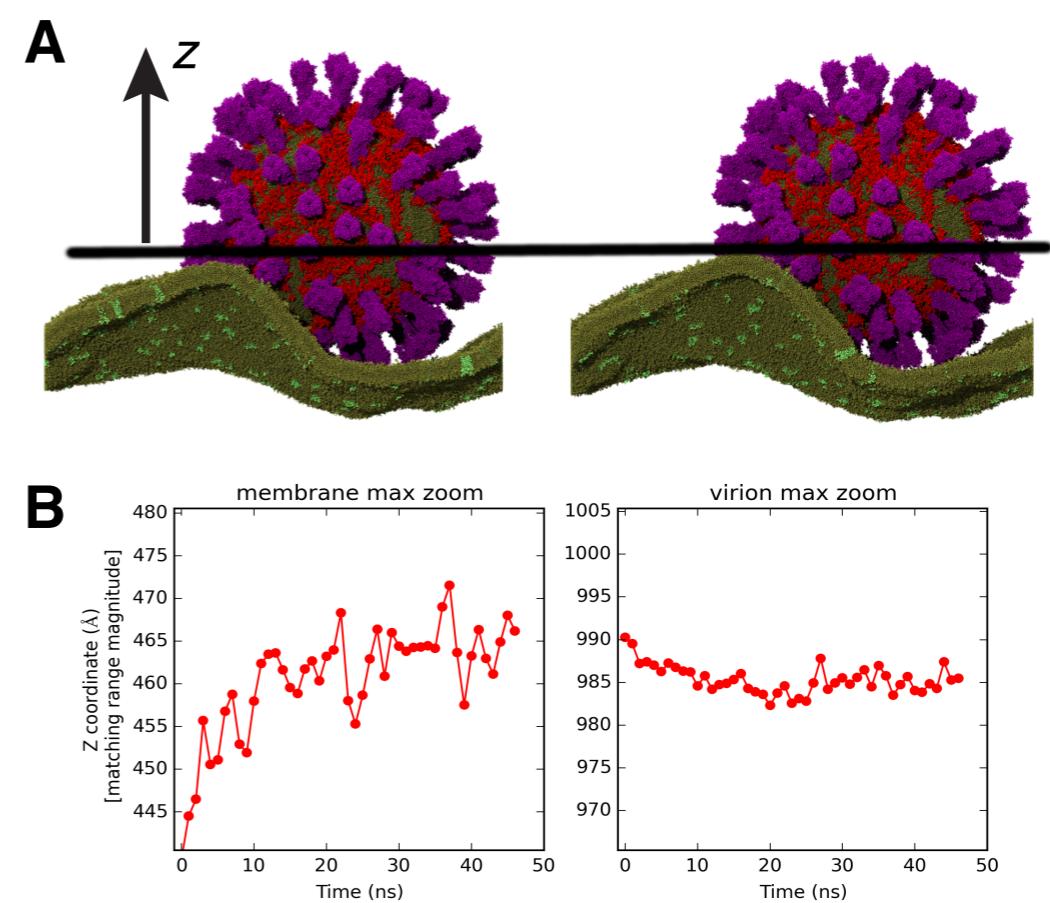
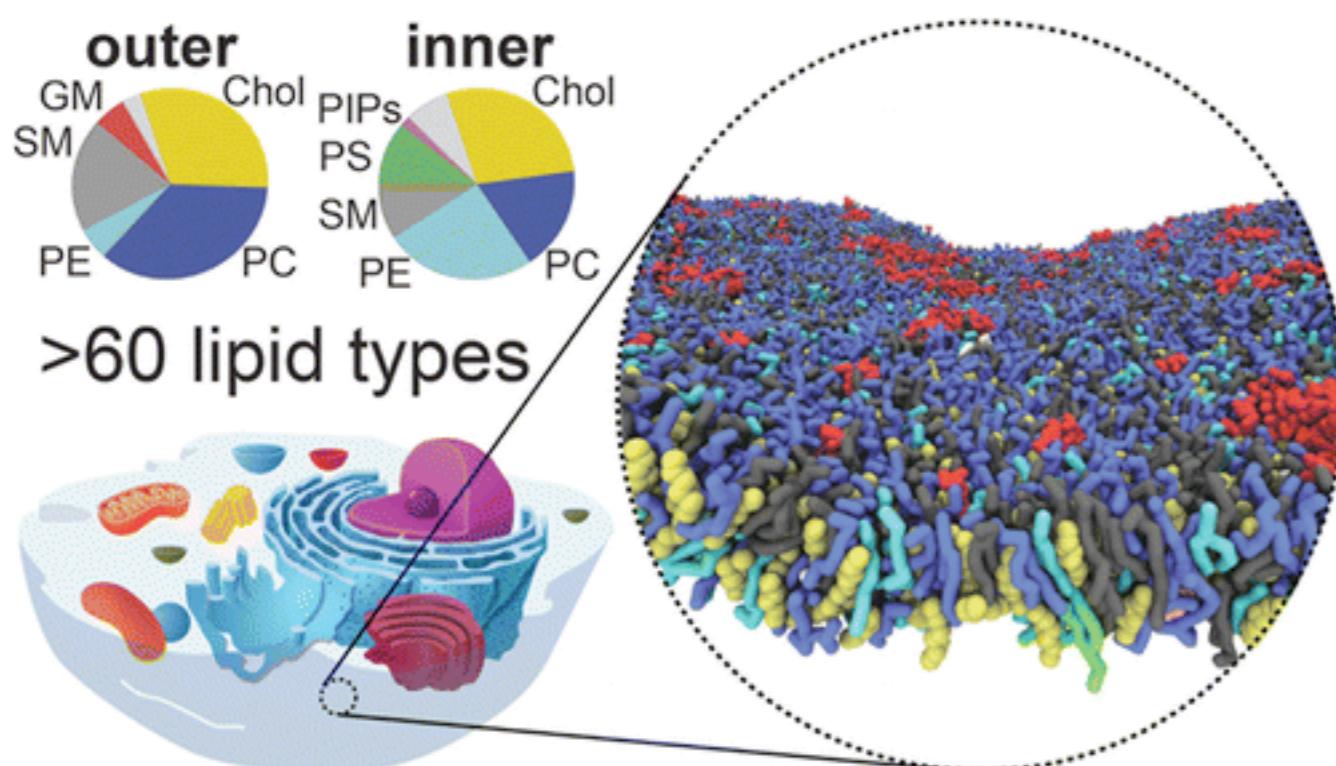
```
u = Universe("system.gro")
```

```
MDAnalysis.topology.GROParser.GROParser("system.gro").parse()
```

... + reduced memory footprint down to 1/3

# Very large systems

- MDAnalysis is capable to analyse very large simulation systems (millions of particles, many  $\mu$ s of data)
- E.g., lipid membranes, virions



H. I. Ingólfsson, M. N. Melo, F. J. V. Eerden, C. Arnarez, C. A. López, T. A. Wassenaar, X. Periole, A. H. D. Vries, D. P. Tieleman, and S. J. Marrink. Lipid Organization of the Plasma Membrane Lipid Organization of the Plasma Membrane. *J Am Chem Soc*, 136(41):14554–14559, 2014.

12.7 M influenza A + membrane  
(TJE Reddy, unpublished)

# Installation

- Python Package Index

```
pip install --upgrade MDAnalysis
```

- Conda

```
conda config --add channels MDAnalysis  
conda install mdanalysis
```

- from source (develop or master == release)

```
git clone https://github.com/MDAnalysis/mdanalysis.git  
cd mdanalysis/package  
python setup.py install
```



## Runs on

- Linux
- Mac OS X

## Open source

- GPL v2
- [github.com/MDAnalysis](https://github.com/MDAnalysis)

MDAnalysis  
MDAnalysis is an object-oriented python toolkit to analyze molecular dynamics  
<http://www.mdanalysis.org>

Repositories 21 People 21 Teams 6 Settings

Filters Find a repository... + New repository

**mdanalysis** Python ★ 38 ⚡ 18

MDAnalysis is a Python library to analyze molecular dynamics trajectories.  
Updated an hour ago

## Development process

- pull request / review / merge
- continuous integration with > 3,500 unit tests



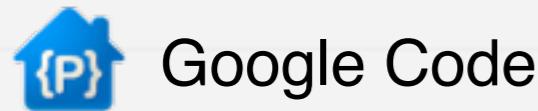
build passing coverage 88%

★ COVERALLS

quantifiedcode

# Impact of “Social Coding”?

MDAnalysis repository commit history



Author: Naveen Michaud-Agrawal  
Date: Thu Jan 31 11:42:33 2008 +0000  
initial import



Software News and Updates  
MDAnalysis: A Toolkit for the Analysis of Molecular  
Dynamics Simulations

GitHub



NAVEEN MICHAUD-AGRAWAL,<sup>1</sup> ELIZABETH J. DENNING,<sup>1,2</sup> THOMAS B. WOOLF,<sup>1,3</sup> OLIVER BECKSTEIN<sup>3,4</sup>

Received 23 October 2010; Revised 6 February 2011; Accepted 12 February 2011

DOI 10.1002/jcc.21787

Published online 15 April 2011 in Wiley Online Library (wileyonlinelibrary.com).

J Comput Chem 32: 2319–2327, 2011

(>200 citations on GoogleScholar (July 2016))

# MDA



# USERS

(... and developers, of course! )

# **MDAnalysis is not just code, it's also a community**

Open and inclusive:

- questions are answered (mailing list)
- extensive docs
- pull requests very welcome!

Reduce upgrade pain for users:

- semantic versioning
- unavoidable API breaks
  - deprecation cycles
  - fixer script based on standard lib2to3

**users → contributors → developers**

- 44 contributing authors (July 2016)
- 7 core developers

**MDA**   
**USERS**



- forms a bridge between most of the commonly encountered MD trajectory formats and NumPy arrays
- provides an object-oriented interface to data on particles
- comes with many analysis (`MDAnalysis.analysis`) and visualisation (`MDAnalysis.vizualiation`) classes/functions
- has a welcoming and active community

MDA   
USERS

Naveen Michaud-Agrawal, Elizabeth J. Denning, Joshua Adelman, **Jonathan Barnoud**, Balasubramanian, Christian Beckstein (logo), Alejandro Bernardin, Bart Bruininks, Sébastien Buchoux, David Caplan, Matthieu Chavent, **John Detlefs**, **David L. Dotson**, Xavier Deupi, Jan Domański, Lennard van der Feltz, Philip Fowler, Joseph Goose, **Richard J. Gowers**, Lukas Grossar, Benjamin Hall, Kyle J. Huston, Joe Jordan, **Max Linke**, Jinju Lu, Robert McGibbon, **Fiona Naughton**, Alex Nesterenko, **Manuel Nuno Melo**, Caio S. Souza, Hai Nguyen, Mattia F. Palermo, Danny Parton, Joshua L. Phillips, **Tyler Reddy**, Paul Rigor, Utkarsh Saxena, Sean L. Seyler, Andy Somogyi, Lukas Stelzl, Gorman Stock, Isaac Virshup, Zhuyi Xue, Carlos Yáñez S, and **Oliver Beckstein**

Join us at

[mdanalysis.org](http://mdanalysis.org)

[github.com/MDAnalysis](https://github.com/MDAnalysis)

