

Figure 5: A Folding@Home result showing the value of CBPR via Ebola

Background

- MDAnalysis is an open-source Python library capable of analyzing molecular dynamics(MD) simulations [1].
 - Creating a common accessible interface to access raw simulation data of any format
- PDB(Protein Data Bank) - An international repository of protein and molecular data
- MD simulations are **pivotal** in designing new approaches to **medicine, drug delivery, and biofuels.** [2]

Literature Review:

- PDB founded in 1971[3] by Berman et al., stores **208,000+ of structure files**[4] and data.
- 2019 - Formal recognition of PDBx/mmCIF as the master file format of the PDB. [5]
- RCSB PDB(a branch of the PDB) accounts for **1 million annual users** and **\$5.5 billion** in use value, alone [6]
- PDBx/mmCIF, despite the machine-readability, lack homogeneity between similar objects, making **parsing and interpretation difficult** [7].

Research Questions:

- How can users and the PDB design for more efficient and reliable data representations of diverse organic structures?
- Why might researchers prefer more varied methods of representing organic matter over a standardized format?
- Can we use precedents in linguistics, stemmatics, and biosemiotics to understand how different file formats impact stakeholders' perceptions of molecular dynamics simulations and the Protein Data Bank?

Hypothesis:

- Various PDB formats are motivated by scientific techniques and discoveries
- Design cannot be limited by current capabilities
- Need to design flexibility for an endearing design

Figure 1: An application of MDAnalysis on a trajectory file

```
import MDAnalysis
from MDAnalysis.tests.datafiles import PSF, DCD
import numpy.linalg

u = MDAnalysis.Universe(PSF,DCD)
nterm = u.select_atoms('segid 4AKE and name N')[0]
cterm = u.select_atoms('segid 4AKE and name C')[0]

bb = u.select_atoms('protein and backbone')

for ts in u.trajectory:
    r = cterm.position - nterm.position
    d = numpy.linalg.norm(r)
    rgyr = bb.radius_of_gyration()
    print("frame = {0}: d = {1} A, Rgyr = {2} A".format(
        ts.frame, d, rgyr))
```

Materials and Methods:

- Compiled literature examining the interplay between the PDB and its users using keywords such as molecular dynamics, PDB, computational chemistry, languages, stemmatics, and linguistics. The research papers were collected from e-libraries such as Elsevier, Science Direct, Lens, and the MIT libraries.
- Interviewed users, such as Dr. Beckstein [8], and PDB board members, including Dr. Berman [9], Dr. Vallat [9], and Dr. Peisach [10, 11].

Results & Discussion

- Experimental methods, knowledge of proteins, data processing methods** [12] and **user bias/error** [13] are major factors that influence the “dialects” of PDB files
 - Dr. Berman: In 2011, there were 27 different “dialects” of PDB files

Legacy-PDB											
ATOM	16	N	LEU	A	2	28.555	13.855	-10.636	1.00	27.76	N
ATOM	17	CA	LEU	A	2	28.797	15.269	-10.398	1.00	25.21	C
ATOM	18	C	LEU	A	2	29.492	15.983	-11.585	1.00	24.21	C
ATOM	19	O	LEU	A	2	30.250	15.240	-12.306	1.00	23.80	O
ATOM	20	CB	LEU	A	2	29.688	15.470	-9.152	1.00	24.30	C
ATOM	21	CG	LEU	A	2	29.884	15.416	-7.751	1.00	22.96	C
ATOM	22	CD1	LEU	A	2	28.730	13.988	-7.390	1.00	22.03	C
ATOM	23	CD2	LEU	A	2	30.085	16.008	-6.776	1.00	21.94	C

PDBx/mmCIF												
loop_	_atom_site.group_PDB											
	_atom_site.id											
	_atom_site.type_symbol											
	_atom_site.label_atom_id	ATOM	16	N	LEU	A	2	28.555	13.855	-10.636	1.00 27.76 ? 2 LEU A N 1	
	_atom_site.label_alt_id	ATOM	17	C	CA	LEU	A	2	28.797	15.269	-10.398	1.00 25.21 ? 2 LEU A CA 1
	_atom_site.label_comp_id	ATOM	18	C	C	LEU	A	2	29.492	15.983	-11.585	1.00 24.21 ? 2 LEU A C 1
	_atom_site.label_entity_id	ATOM	19	O	O	LEU	A	2	30.250	15.240	-12.306	1.00 23.80 ? 2 LEU A O 1
	_atom_site.label_seq_id	ATOM	20	C	CB	LEU	A	2	29.688	15.470	-9.152	1.00 24.30 ? 2 LEU A CB 1
	_atom_site.label_seq_id	ATOM	21	C	CG	LEU	A	2	29.884	15.416	-7.751	1.00 22.96 ? 2 LEU A CG 1
	_atom_site.pdbx_ins_code	ATOM	22	C	CD1	LEU	A	2	28.730	13.988	-7.390	1.00 22.03 ? 2 LEU A CD1 1
	_atom_site.cartn_x	ATOM	23	C	CD2	LEU	A	2	30.085	16.008	-6.776	1.00 21.94 ? 2 LEU A CD2 1
	_atom_site.cartn_y											
	_atom_site.cartn_z											
	_atom_site.pdbx_occupancy											
	_atom_site.is_disordered											
	_atom_site.pdbx_formal_charge											
	_atom_site.auth_seq_id											
	_atom_site.auth_comp_id											
	_atom_site.auth_atom_id											
	_atom_site.pdbx_PDB_model_num											

Figure 2: The difference between the legacy PDB and PDBx/mmCIF file format. Smith, R.D et. al, 1997)

- The PDB adheres to a lateral approach, **equally** valuing each stakeholder and lets **scientific advancement** determine its actions.

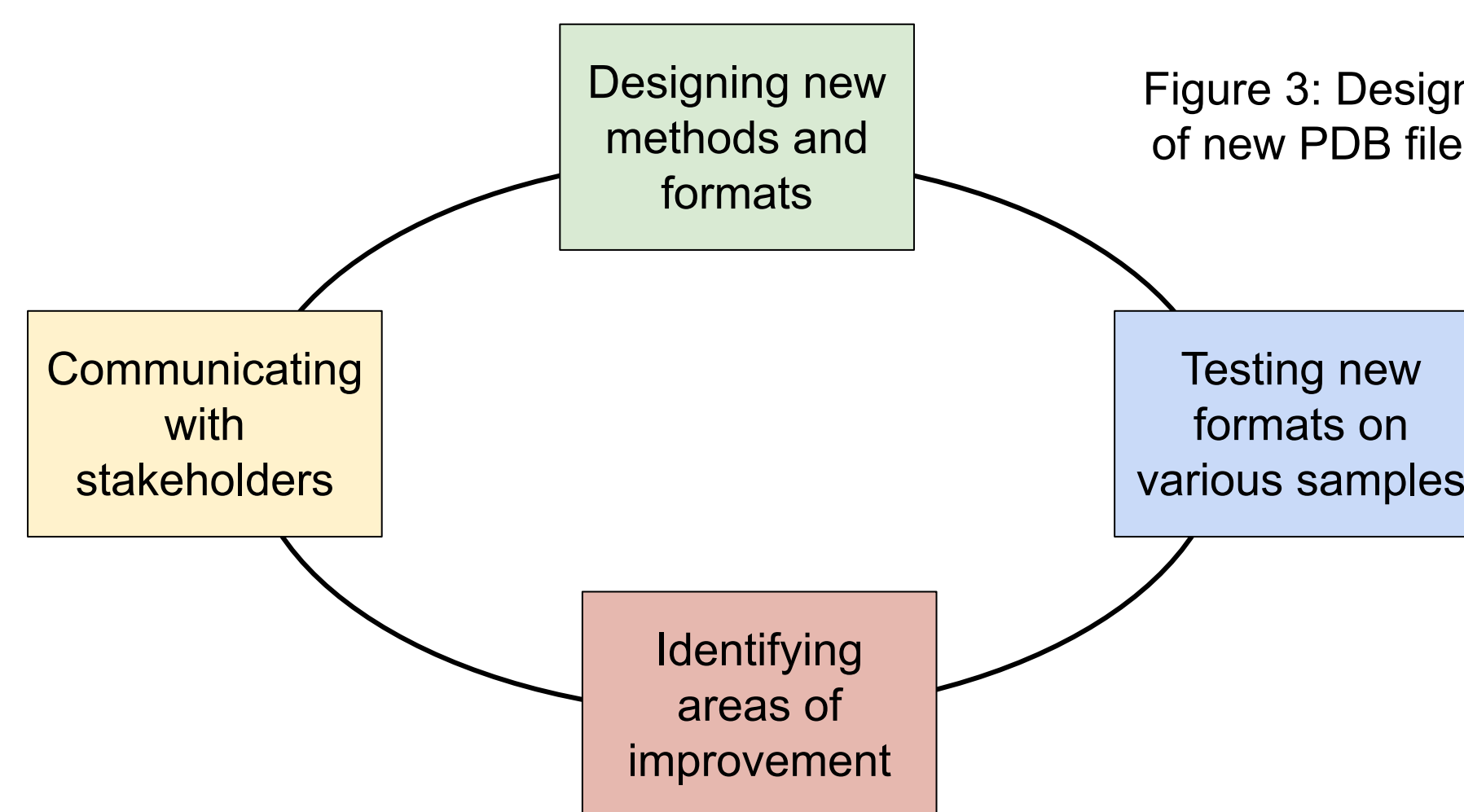


Figure 3: Design process of new PDB file formats.

- Accurate results need accurate data; missteps may cost years of research.
 - MDAnalysis stakeholders work in fields that need accurate data.

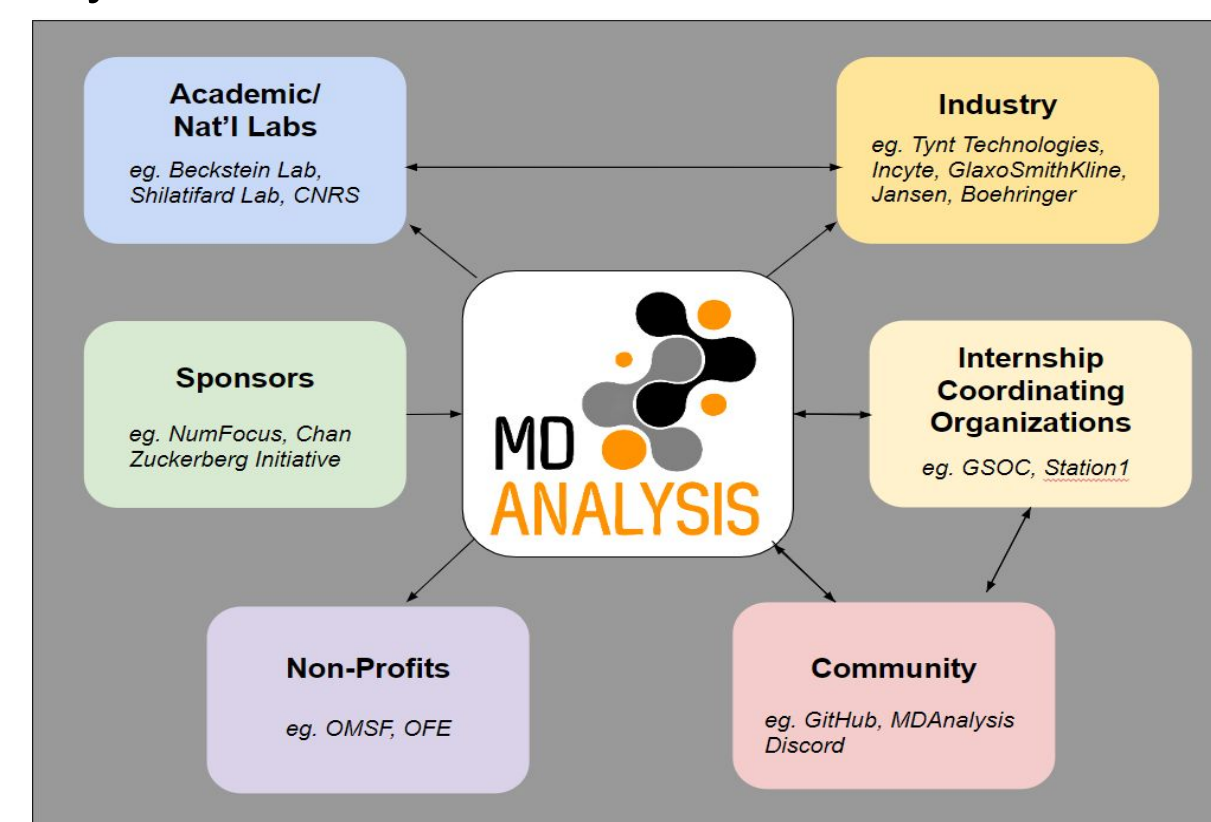


Figure 4: MDAnalysis's stakeholders map.

Social Impact Discussion

- Open-source information and software allows for community-based participatory research (ie. Folding@Home).
- Design for the Margins [14]: increase diversity of users and professionals (developers, etc.).
- Non-profit with **400,000+ downloads** - creates a **strong and sustainable community.**[15]

Future Work

MDAnalysis can:

- Implement **resilient** approaches to utilizing PDB data aided by our findings
- Continue support for various file formats that **balance human and machine readability**
- Continue to invest in methods of better automation in submission techniques for efficient resource distribution
- Support raw files to expand methodology and integrity
- Continued **support for scientific flexibility**

PDB stakeholders can:

- Adhere to the definitions of the PDBx/mmCIF dictionary
- Continued commitment to academic **integrity** and iterative **experimentation**

Integrative Conclusions

These insights provide valuable information on the social implications of both repositories like the Protein Data Bank and Python libraries like MDAnalysis. Understanding the creation of different file formats and their impact on both professionals and other users helps broaden the accessibility of such tools and, ultimately, provides a deeper understanding of their cognitive implications.

Acknowledgements

We'd like to thank all Station1 instructors, Dr. Gowers, and Dr. MacDermott-Opeskin for their efforts this summer. We'd also like to thank Dr. Berman, Dr. Beckstein, Dr. Vallat, and Dr. Peisach for setting time aside to speak to us. Finally, we'd like to thank the MDAnalysis community and the GSOC interns for their continued support.

References

[1] Gowers, Richard, Max Link, Jonathan Barnoud, Tyler Reddy, Manuel Melo, Sean Seyer, Jan Domanski, et al. "MDAnalysis: A Python Package for the Rapid Analysis of Molecular Dynamics Simulations." 98-105. Austin, Texas, 2016. <https://doi.org/10.25080/Majors-629e641a-00a>

[2] Perilla, Juan R, Boon Chong Goh, C Keith Cassidy, Bo Liu, Rafael C Bernardi, Til Rudack, Hang Yu, Zhe Wu, and Klaus Schulten. "Molecular Dynamics Simulations of Large Macromolecular Complexes." Current Opinion in Structural Biology 31 (April 2016): 64-74. <https://doi.org/10.1016/j.sbspro.2015.05.007>

[3] Berman, Helen M., John Westbrook, Zukang Feng, Gary Gilliland, T. N. Bhat, Helge Weissig, Ilya N. Shindyalov, and Philip E. Bourne. "The Protein Data Bank." Nucleic Acids Research 28, no. 1 (January 1, 2000): 235-42. <https://doi.org/10.1093/nar/28.1.235>

[4] PDB statistics: Overall growth of released structures per year. RCSB PDB. (n.d.). <https://www.rcsb.org/structure/growth-released-structures>

[5] Adams, P. D., P. V. Kovalek, K. Bakker, N. M. Berman, J. Bernstein, G. Brocas, D. G. Brown, et al. "Announcing Mandatory Submission of PDBx/mmCIF Format Files for Crystallographic Depositions to the Protein Data Bank (PDB)." Acta Crystallographica Section D: Structural Biology 75, no. 4 (April 1, 2019): 451-54. <https://doi.org/10.1107/S2029798319004522>

[6] Sullivan, Kevin P., Peggy Brennan-Tonetta, and Lucas J. Marxen. "Economic Impacts of the Research Collaboratory for Structural Bioinformatics (RCSB) Protein Data Bank." RCSB Protein Data Bank, May 1, 2017. <https://doi.org/10.22106/pdb-2017-05-01>

[7] D'Amico, Nancy, Deborah Giordano, Bernardina Scafuri, Angelo Facchiano, and Anna Marabotti. "Standardizing Macromolecular Structure Files: Further Efforts Are Needed." Trends in Biochemical Sciences 48, no. 7 (July 2023): 590-98. <https://doi.org/10.1016/j.tics.2023.03.002>

[8] Oliver R. Beckstein. Interview by Karen Bekhazi, John Ong, July 28, 2023. Interview 1, recording. <https://drive.google.com/file/d/1nK9E9QaYsJectHUCoWBF50EFC4U4v9w2uq/view>

[9] Helen M. Berman. Brenda Vallat. Interview by Karen Bekhazi, John Ong, July 28, 2023. Interview 1, recording. https://drive.google.com/file/d/1H1L1G78P8B28713p7SYu4EzYz4M8Vw/view?usp=drive_link

[10] Ezra Peisach. Interview by Karen Bekhazi, John Ong, July 28, 2023. Interview 1, recording. <https://drive.google.com/file/d/1nK9E9QaYsJectHUCoWBF50EFC4U4v9w2uq/view>

[11] Bhat, T. N., Philip Bourne, Zukang Feng, Gary Gilliland, Shri Jau, Veerasamy Ravichandran, Bohdan Schneider, et al. "The PDB Data Uniformity Project." Nucleic Acids Research 29, no. 1 (January 1, 2001): 214-18. <https://doi.org/10.1093/nar/29.1.214>

[12] Dauter, Zbigniew, Alexander Wlodawer, Wlodek Minor, Mariusz Jaskolski, and Bernhard Rupp. "Avoidable Errors in Deposited Macromolecular Structures: An Impediment to Efficient Data Mining." IUCr 1, no. Pt 3 (April 14, 2014): 179-93. <https://doi.org/10.1107/S2029225614005442>

[14] McQuweil, M.T. "Comp. Design Principles, 2018 Pdf - LA Data." Accessed August 2, 2023. <https://www.improbook.com/compdesign/McQuweil%20DesignPrinciples2018.pdf>

[15] MDAnalysis. Anaconda.org. (n.d.). <https://anaconda.org/conda-forge/mdanalysis>