# RustSASA: A Rust Crate for Accelerated Solvent Accessible Surface Area Calculations

**Maxwell J. Campbell** [1]

**1** University of California, San Francisco, United States

## Summary

Solvent accessible surface area (SASA) calculations are fundamental for understanding protein structure, function, and dynamics in computational biology. These calculations quantify the surface area of biomolecules accessible to solvent molecules, providing insights into protein folding, stability, and intermolecular interactions. The Shrake-Rupley algorithm has served as the standard for SASA calculations since 1973, but existing implementations often become computational bottlenecks when analyzing large protein datasets. As proteomics datasets continue to grow with initiatives like AlphaFold producing hundreds of millions of predicted protein structures, the need for efficient SASA calculation tools has increased dramatically. RustSASA addresses this challenge by providing a high-performance implementation of the Shrake-Rupley algorithm written in pure Rust, delivering a 5× speed improvement over FreeSASA while maintaining calculation accuracy and providing interfaces for multiple programming languages and frameworks (i.e. MDAnalysis). The source code is freely available at https://github.com/maxall41/RustSASA.

## Statement of need

Current SASA calculation tools represent a significant computational bottleneck in structural biology workflows, particularly for molecular dynamics simulations and high-throughput analyses. Popular implementations such as those in Biopython and FreeSASA, while accurate, become prohibitively slow when processing large protein datasets. RustSASA addresses this performance gap by leveraging Rust's efficient parallelization abstractions (via Rayon) and readily available SIMD instructions (via Pulp). These optimizations enable RustSASA's performance advantage over the simpler C implementation of the same algorithm in FreeSASA.

Benchmarking on representative protein structures demonstrates that RustSASA achieves a 5× improvement over FreeSASA, and a 63× performance improvement over Biopython. This performance advantage reduces computational costs for high-throughput structural analyses and makes large-scale comparative studies feasible. Furthermore, RustSASA's multi-language support (Rust and Python), command-line interface, and MDAnalysis package (Gowers et al., 2016; Michaud-Agrawal et al., 2011) ensure broad accessibility across the computational biology community.

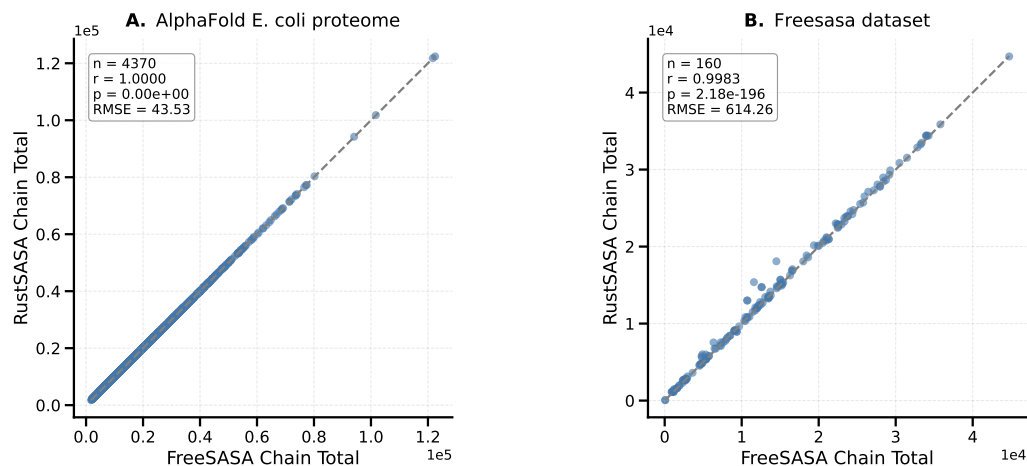<sup>34</sup> # Results

<sup>35</sup> # Calculation Quality



**Figure 1: A.** Comparing RustSASA against FreeSASA on *E. coli* proteome at the chain level. **B.** Comparing RustSASA against FreeSASA on FreeSASA comparison dataset at the chain level.

<sup>36</sup> To evaluate the accuracy of RustSASA calculations, we compared results to FreeSASA (Mit-
<sup>37</sup> ternacht, 2016) on both the predicted *E. coli* proteome from AlphaFold DB (Jumper et al.,
<sup>38</sup> 2021; Varadi et al., 2021) and the FreeSASA evaluation dataset. RustSASA produces SASA
<sup>39</sup> values that closely match those from FreeSASA, achieving Pearson correlation coefficients >
<sup>40</sup> 0.99 on both datasets (Figure 1).

## Performance



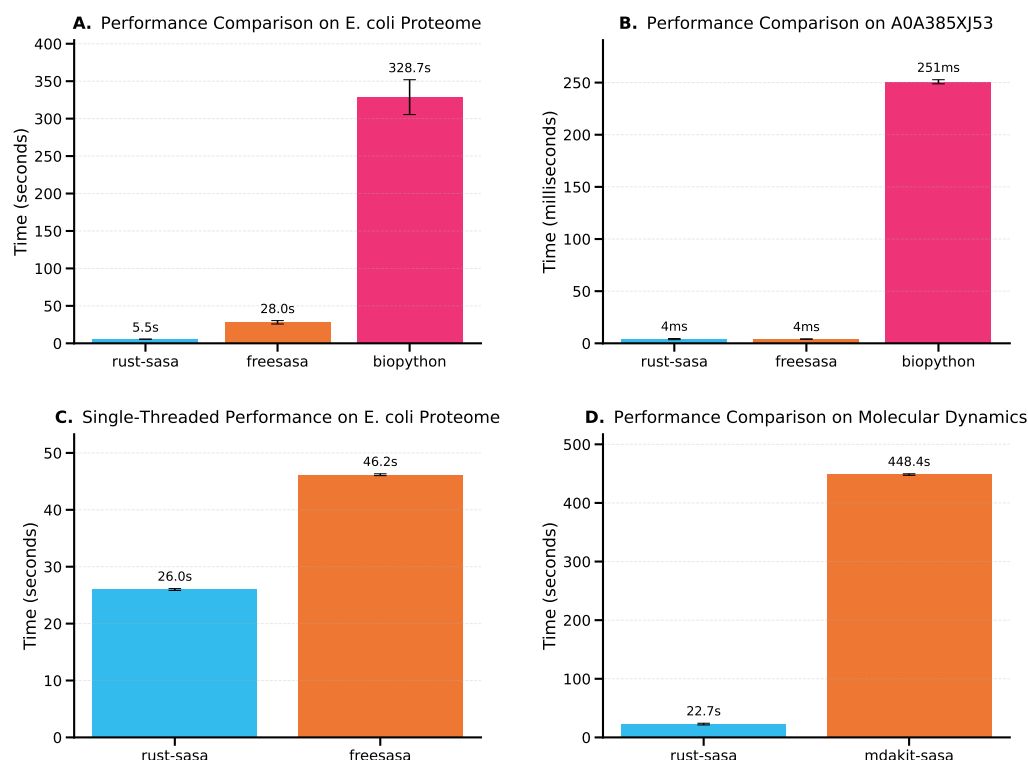**Figure 2: A.** Comparing the multi-threaded performance of RustSASA, FreeSASA, and Biopython on the full AlphaFold *E. coli* proteome. **B.** Comparing the multi-threaded performance of RustSASA, FreeSASA, and Biopython on A0A385XJ53, a protein randomly selected from the AlphaFold *E. coli* proteome. **C.** Comparing the single-threaded performance of RustSASA and FreeSASA on the full AlphaFold *E. coli* proteome. **D.** Comparing the performance of RustSASA and mdakit-sasa (based on freesasa) on a molecular dynamics trajectory

We evaluated the performance of FreeSASA, RustSASA, and Biopython ([Cock et al., 2009](#)) across four evaluations. First, we performed multi-threaded SASA calculations for all proteins in the *E. coli* proteome. Second, we evaluated the performance of these methods on a single randomly selected protein (A0A385XJ53) from the AlphaFold *E. coli* proteome. Third, we evaluated the single-threaded performance of RustSASA and FreeSASA on the *E. coli* proteome, Biopython was excluded from this benchmark due to its poor performance hindering timely evaluation. Fourth, we evaluated the performance of RustSASA and mdakit-sasa, which uses freesasa internally, on a molecular dynamics trajectory, specifically trajectory 10824 (4IAQ, 5HT receptor) from GPRCmd ([Rodríguez-Espigares et al., 2020](#)).

For the full proteome benchmarks (Figure 2A), we used Hyperfine ([Peter, 2023](#)) with 3 runs and 3 warmup iterations. All methods utilized parallel processing across eight cores. GNU parallel ([Tange, 2011](#)) was used to parallelize FreeSASA and Biopython, while RustSASA utilized its internal parallelization. RustSASA processed the entire proteome in ~5 seconds compared to ~28 seconds for FreeSASA and ~328 seconds for Biopython, representing 5× and 63× speed improvements, respectively.

For the single protein benchmark (Figure 2B), we used Hyperfine with 3 warmup iterations and 25 runs. RustSASA processed the protein in 4.0ms (±0.5), FreeSASA processed the protein in 4.0ms (±0.2), and Biopython processed the protein in 250.8ms (±2.0). On the single-threaded benchmark (Figure 2C), RustSASA processed the proteome in 26.0 seconds compared to 46.2 seconds for FreeSASA, representing a ~43% performance improvement, demonstrating that

62 RustSASA's performance advantage is not solely due to multi-threading.

63 For the molecular dynamics trajectory benchmark (Figure 2D), we used Hyperfine with 3
64 runs. RustSASA processed the trajectory in 22.7 seconds ($\pm$1.4), and mdakit-sasa processed
65 the trajectory in 448.4 seconds ($\pm$1.3), representing a ~20x performance improvement. The
66 magnitude of the improvement can be attributed to the inefficiency with which mdakit-sasa
67 utilizes FreeSASA.

## Methods

69 RustSASA computes solvent-accessible surface areas (SASA) using the Shrake–Rupley algorithm
70 (Shrake & Rupley, 1973). In this algorithm each atom is represented as a sphere with a radius
71 equal to its atomic van der Waals radius plus the radius of a spherical solvent probe; the
72 sphere surface is sampled with a dense quasi-uniform distribution of test points and a point
73 is considered solvent accessible if it is not occluded by any neighboring atom sphere. For all
74 calculations reported here a solvent probe radius of 1.4 Å (the approximate radius of a water
75 molecule) was used.

76 Atomic radii were assigned according to the ProtOr parameter set introduced by Tsai et
77 al. (Tsai et al., 1999). These radii were applied to all non-hydrogen heavy atoms present
78 in the structures; hydrogen atoms, when present in input files, were ignored for the SASA
79 computations to maintain consistency with common practice for protein SASA estimation.
80 Additionally, all HETATM records in the input—non-standard amino acids and ligands—were
81 ignored.

82 To ensure a fair comparison of RustSASA and FreeSASA in the single-threaded benchmark, a
83 C++ script was utilized to call the FreeSASA C API for all input proteins in a given folder.
84 This approach ensures that the command-line overhead is not responsible for RustSASA's
85 performance advantage. Furthermore, in all experiments, FreeSASA was configured to use the
86 Shrake-Rupley algorithm over its default algorithm, Lee & Richards, to ensure an accurate
87 comparison between the methods. Proteome-scale structure models for Escherichia coli
88 were obtained from the AlphaFold DB (entry UP000000625_83333_ECOLI_v6, available
89 at https://alphafold.ebi.ac.uk/download). All experiments were conducted on a 2024 Apple
90 MacBook Air with an M3 processor and 24GB of unified memory.

## Conclusion

92 RustSASA provides a significant advancement in SASA calculation performance while main-
93 taining accuracy, addressing a bottleneck in computational structural biology. The $5\times$ speed
94 improvement over current standards enables previously intractable analyses of large protein
95 datasets and molecular dynamics simulations. By providing interfaces for multiple programming
96 languages alongside a command-line tool and MDAnalysis package, RustSASA ensures broad
97 accessibility across the research community. As structural biology datasets continue to expand,
98 efficient computational tools like RustSASA become essential for advancing our understanding
99 of protein structure and function.

## Acknowledgements

# References

Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., & others. (2009). Biopython: Freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics*, *25*(11), 1422–1423. https://doi.org/10.1093/bioinformatics/btp163

Gowers, R. J., Linke, M., Barnoud, J., Reddy, T. J. E., Melo, M. N., Seyler, S. L., Domański, J., Dotson, D. L., Buchoux, S., Kenney, I. M., & Beckstein, O. (2016). MDAnalysis: A python package for the rapid analysis of molecular dynamics simulations. *SciPy 2016*. https://doi.org/10.25080/Majora-629e541a-00e

Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., Bridgland, A., Meyer, C., Kohl, S. A. A., Ballard, A. J., Cowie, A., Romera-Paredes, B., Nikolov, S., Jain, R., Adler, J., … Hassabis, D. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature*, *596*(7873), 583–589. https://doi.org/10.1038/s41586-021-03819-2

Michaud-Agrawal, N., Denning, E. J., Woolf, T. B., & Beckstein, O. (2011). MDAnalysis: A toolkit for the analysis of molecular dynamics simulations. *Journal of Computational Chemistry*, *32*(10), 2319–2327. https://doi.org/https://doi.org/10.1002/jcc.21787

Mitternacht, S. (2016). FreeSASA: An open source c library for solvent accessible surface area calculations. *F1000Research*, *5*, 189. https://doi.org/10.12688/f1000research.7931.1

Peter, D. (2023). *Hyperfine*. https://github.com/sharkdp/hyperfine

Rodríguez-Espigares, I., Torrens-Fontanals, M., Tiemann, J. K. S., Aranda-García, D., Ramírez-Anguita, J. M., Stepniewski, T. M., Worp, N., Varela-Rial, A., Morales-Pastor, A., Medel-Lacruz, B., Pándy-Szekeres, G., Mayol, E., Giorgino, T., Carlsson, J., Deupi, X., Filipek, S., Filizola, M., Gómez-Tamayo, J. C., Gonzalez, A., … Selent, J. (2020). GPCRmd uncovers the dynamics of the 3D-GPCRome. *Nature Methods*, *17*(8), 777–787. https://doi.org/10.1038/s41592-020-0884-y

Shrake, A., & Rupley, J. A. (1973). Environment and exposure to solvent of protein atoms. Lysozyme and insulin. *Journal of Molecular Biology*, *79*(2), 351–371. https://doi.org/https://doi.org/10.1016/0022-2836(73)90011-9

Tange, O. (2011). GNU parallel - the command-line power tool;*Login: The USENIX Magazine*, *36*(1), 42–47. http://www.gnu.org/s/parallel

Tsai, J., Taylor, R., Chothia, C., & Gerstein, M. (1999). The packing density in proteins: Standard radii and volumes11Edited by j. M. thornton. *Journal of Molecular Biology*, *290*(1), 253–266. https://doi.org/https://doi.org/10.1006/jmbi.1999.2829

Varadi, M., Anyango, S., Deshpande, M., Nair, S., Natassia, C., Yordanova, G., Yuan, D., Stroe, O., Wood, G., Laydon, A., Žídek, A., Green, T., Tunyasuvunakool, K., Petersen, S., Jumper, J., Clancy, E., Green, R., Vora, A., Lutfi, M., … Velankar, S. (2021). AlphaFold protein structure database: Massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Research*, *50*(D1), D439–D444. https://doi.org/10.1093/nar/gkab1061