

# funsies: A minimalist, distributed and dynamic workflow engine

Cyrille Lavigne<sup>\*1</sup> and Alán Aspuru-Guzik<sup>1, 2, 3, 4</sup>

<sup>1</sup> Department of Computer Science, University of Toronto, 40 St. George St, Toronto, Ontario M5S 2E4, Canada <sup>2</sup> Chemical Physics Theory Group, Department of Chemistry, University of Toronto, 80 St. George St, Toronto, Ontario M5S 3H6, Canada <sup>3</sup> Vector Institute for Artificial Intelligence, 661 University Ave Suite 710, Toronto, Ontario M5G 1M1, Canada <sup>4</sup> Lebovic Fellow, Canadian Institute for Advanced Research (CIFAR), 661 University Ave, Toronto, Ontario M5G, Canada

DOI: [10.21105/joss.03274](https://doi.org/10.21105/joss.03274)

## Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Patrick Diehl](#) 

## Reviewers:

- [@gflfst](#)
- [@vijaysm](#)

Submitted: 28 April 2021

Published: 26 July 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

Large-scale, high-throughput computational investigations are increasingly common in chemistry and physics. Until recently, computational chemistry was primarily performed using all-in-one monolithic software packages ([Aprà et al., 2020](#); [Aquilante et al., 2020](#); [Barca et al., 2020](#); [Kühne et al., 2020](#); [Romero et al., 2020](#); [Smith et al., 2020](#)). However, the limits of individual programs become evident when tackling complex multifaceted problems. As such, it is increasingly common to use multiple disparate software packages in a single computational pipeline, often stitched together using shell scripts in languages such as Bash, or using Python and other interpreted languages.

These complex computational pipelines are difficult to scale and automate, as they often include manual steps and significant “human-in-the-loop” tuning. Shell scripting errors are often undetected, which can compromise scientific results. Conversely, exception-based error handling, the standard approach in Python, can readily bring a computational workflow to a halt when exceptions are not properly caught ([Weimer & Necula, 2008](#)).

**funsies** is a set of python programs and modules to describe, execute and analyze computational workflows, with first-class support for shell scripting. It includes a lightweight, decentralized workflow engine backed by a NoSQL store. Using **funsies**, external program and python-based computations are easily mixed together. Errors are detected and propagated throughout computations. Automatic, transparent incremental computing (based on a hash tree data structure) provides a convenient environment for iterative prototyping of computationally expensive workflows.

## Statement of need

Modern workflow management programs used in the private sector, such as Apache Airflow and Uber’s Cadence, are robust and extremely scalable, but are difficult to deploy. Scientific workflow management systems, many of which are compiled in ([Tommaso, 2021](#)) and systematically reviewed in ([Mölder et al., 2021](#)), are easier to set up on high-performance computing clusters, but are tuned to the needs of specific disciplines, such as bioinformatics or machine learning. This includes, for example, the use of configuration file formats (YAML, JSON, etc.), packaging tools (for example, conda or Docker), locked-in compute providers (Amazon Web Services, Google Cloud) and storage formats that may be common in specific scientific fields but not throughout the greater community.

<sup>\*</sup>Corresponding author.

40 For our own group's research program, we wanted to have available a lightweight workflow  
41 management system that could be readily deployed to new and varied computational facilities  
42 and local workstations with minimal effort. This system had to support our existing shell-  
43 based and Python-based scripts, and be flexible enough for rapid prototyping all the way to  
44 large-scale computational campaigns, and provide an embeddable solution that can be bundled  
45 within other software (Lavigne et al., 2020). Finally, we were looking for a tool that could  
46 integrate data generation and storage, to avoid the common practice of transforming the  
47 filesystem into what is effectively a schema-less database. We developed `funsies` to address  
48 those needs.

## 49 Features and Implementation

50 `funsies` is a Python library and a set of associated command-line tools. Using the `funs`  
51 `ies` library, general computational workflows are described in lazily evaluated Python code.  
52 Operations in `funsies` are taken to be pure, that is, all operation outputs are entirely and  
53 solely determined by their inputs. Workflows are orchestrated using python by manipulating  
54 pointers to yet-to-be-calculated data. Workflow instructions are transparently translated and  
55 saved as graph elements in a Redis database.

56 Computational evaluation is initiated by the user asking for specific output value. The task  
57 graph from these final outputs is walked back all the way to those operations with no de-  
58 pendencies. These initial operations are then queued for execution. Lightweight worker pro-  
59 cesses, instantiated from the command line on local or remote machines, connect to the Redis  
60 database and start executing the workflow. For each operation, the worker checks if outputs  
61 are already cached, and if not, executes the associated function and saves its outputs. It then  
62 enqueues any dependents for execution, by itself or by other workers. In this way, the entire  
63 computational graph is evaluated in a distributed, decentralized fashion without any scheduler  
64 or manager program. Errors in workflows are handled using a functional approach inspired  
65 by Rust (Klabnik & Nichols, 2019). Specifically, exceptions are propagated through workflow  
66 steps, canceling dependent tasks, without interrupting valid workflow branches. This provides  
67 both easy error tracing and a high degree of fault tolerance.

68 The main distinguishing feature of `funsies` is the hash tree structure that is used to encode  
69 all operations and their inputs. The causal hashing approach used in `funsies` can also be  
70 found in Snakemake (Mölder et al., 2021) as an optional component and the (now defunct)  
71 Koji workflow system (Maymounkov, 2018), as part of the Nix package manager (Dolstra et  
72 al., 2004) and in the Git version control system (Chacon & Straub, 2014). In `funsies`, we  
73 replace all filesystem operations with hash addressed operations; that is all I/O operations  
74 and dependencies are tracked.

75 Every operation has a hash address that is computed from the hash values of its dependencies  
76 and a hashed identifier for the associated operation on data. In this way, the consistency of  
77 data dependencies is strongly enforced. Changes to data and operations are automatically  
78 and transparently propagated, as changing a single dependency will cause a rehash of all  
79 its dependents, effectively producing a new workflow with no associated data that needs to  
80 be recomputed. Alternatively, if data already exists at a specific hash address, then it was  
81 generated from the same operations that produced that hash. In this way, the hash tree  
82 structure enables transparent and automatic incremental recomputing.

83 Using hash addresses also enables decentralization, as we can rely on the unlikeliness of hash  
84 collisions (Stevens et al., 2017) to eliminate centralized locks. An important advantage of  
85 this approach is that it allows worker processes to generate their own workflows of tasks  
86 dynamically. Results from these dynamic workflows can be collected and used further in the  
87 workflow description, provided they can be reduced to a number of outputs known at compile  
88 time, a technique similar to MapReduce (Dean & Ghemawat, 2004).

89 As of now, we have published one project ([Pollice et al., 2021](#)) that used an earlier iteration of f  
90 unsies, and are using it in multiple ongoing inquiries. We provide several sample workflows on  
91 Github, with a focus on computational chemistry, quantum computing, and high-performance  
92 computing infrastructure.

93 We intend to maintain funsies and of course welcome collaborations from contributors  
94 around the world.

## 95 Acknowledgements

96 We acknowledge testing by early users Cher-Tian Ser ([\(\(chertianser?\)\)](#)), Kjell Jorner  
97 ([\(\(kjelljorner?\)\)](#)) and Gabriel dos Passos Gomes ([\(\(gabegomes?\)\)](#)). CL also thanks Chris  
98 Crebolder ([\(\(ccrebolder?\)\)](#)) for help setting up documentation pages. We acknowledge the  
99 Defense Advanced Research Projects Agency (DARPA) under the Accelerated Molecular  
100 Discovery Program under Cooperative Agreement No. HR0011920027 dated August 1,  
101 2019. The content of the information presented in this work does not necessarily reflect the  
102 position or the policy of the Government. A. A.-G. thanks Dr. Anders G. Frøseth for his  
103 generous support. A. A.-G. also acknowledges the generous support of Natural Resources  
104 Canada and the Canada 150 Research Chairs program. We thank Compute Canada for  
105 providing computational resources.

## 106 References

- 107 Aprà, E., Bylaska, E. J., Jong, W. A. de, Govind, N., Kowalski, K., Straatsma, T. P., Valiev,  
108 M., Dam, H. J. J. van, Alexeev, Y., Anchell, J., Anisimov, V., Aquino, F. W., Atta-Fynn,  
109 R., Autschbach, J., Bauman, N. P., Becca, J. C., Bernholdt, D. E., Bhaskaran-Nair, K.,  
110 Bogatko, S., ... Harrison, R. J. (2020). NWChem: Past, present, and future. *Journal of*  
111 *Chemical Physics*, 152(18), 184102. <https://doi.org/10.1063/5.0004997>
- 112 Aquilante, F., Autschbach, J., Baiardi, A., Battaglia, S., Borin, V. A., Chibotaru, L. F., Conti,  
113 I., De Vico, L., Delcey, M., Fdez. Galván, I., Ferré, N., Freitag, L., Garavelli, M., Gong,  
114 X., Knecht, S., Larsson, E. D., Lindh, R., Lundberg, M., Malmqvist, P. Å., ... Veryazov,  
115 V. (2020). Modern quantum chemistry with [open]molcas. *Journal of Chemical Physics*,  
116 152(21), 214117. <https://doi.org/10.1063/5.0004835>
- 117 Barca, G. M. J., Bertoni, C., Carrington, L., Datta, D., De Silva, N., Deustua, J. E., Fedorov,  
118 D. G., Gour, J. R., Gunina, A. O., Guidez, E., Harville, T., Irle, S., Ivanic, J., Kowalski,  
119 K., Leang, S. S., Li, H., Li, W., Lutz, J. J., Magoulas, I., ... Gordon, M. S. (2020). Recent  
120 developments in the general atomic and molecular electronic structure system. *Journal of*  
121 *Chemical Physics*, 152(15), 154102. <https://doi.org/10.1063/5.0005188>
- 122 Chacon, S., & Straub, B. (2014). *Pro git*. Apress. ISBN: 978-1-4842-0076-6
- 123 Dean, J., & Ghemawat, S. (2004). MapReduce: Simplified data processing on large clusters.  
124 *OSDI'04: Sixth Symposium on Operating System Design and Implementation*, 137–150.
- 125 Dolstra, E., Jonge, M. de, & Visser, E. (2004). *Nix: A safe and policy-free system for software*  
126 *deployment*. 14. <https://edolstra.github.io/pubs/nspfssd-lisa2004-final.pdf>
- 127 Klabnik, S., & Nichols, C. (2019). *The rust programming language (covers rust 2018)*. No  
128 Starch Press. ISBN: 978-1-71850-044-0
- 129 Kühne, T. D., Iannuzzi, M., Del Ben, M., Rybkin, V. V., Seewald, P., Stein, F., Laino, T.,  
130 Khaliullin, R. Z., Schütt, O., Schiffmann, F., Golze, D., Wilhelm, J., Chulkov, S., Bani-  
131 Hashemian, M. H., Weber, V., Borštnik, U., Taillefumier, M., Jakobovits, A. S., Lazzaro,  
132 A., ... Hutter, J. (2020). CP2K: An electronic structure and molecular dynamics software

- 133 package - quickstep: Efficient and accurate electronic structure calculations. *Journal of*  
 134 *Chemical Physics*, 152(19), 194103. <https://doi.org/10.1063/5.0007045>
- 135 Lavigne, C., Passos Gomes, G. dos, Pollice, R., & Aspuru-Guzik, A. (2020). *Automatic discov-*  
 136 *ery of chemical reactions using imposed activation.* [https://doi.org/10.26434/chemrxiv.](https://doi.org/10.26434/chemrxiv.13008500.v2)  
 137 [13008500.v2](https://doi.org/10.26434/chemrxiv.13008500.v2)
- 138 Maymounkov, P. (2018). Koji: Automating pipelines with mixed-semantics data sources.  
 139 *arXiv:1901.01908 [cs]*. <http://arxiv.org/abs/1901.01908>
- 140 Mölder, F., Jablonski, K., Letcher, B., Hall, M., Tomkins-Tinch, C., Sochat, V., Forster, J.,  
 141 Lee, S., Twardziok, S., Kanitz, A., Wilm, A., Holtgrewe, M., Rahmann, S., Nahnsen, S.,  
 142 & Kster, J. (2021). Sustainable data analysis with snakemake [version 2; peer review:  
 143 1 approved, 1 approved with reservations]. *F1000Research*, 10(33). [https://doi.org/10.](https://doi.org/10.12688/f1000research.29032.2)  
 144 [12688/f1000research.29032.2](https://doi.org/10.12688/f1000research.29032.2)
- 145 Pollice, R., Friederich, P., Lavigne, C., Gomes, G. dos P., & Aspuru-Guzik, A. (2021). Organic  
 146 molecules with inverted gaps between first excited singlet and triplet states and appreciable  
 147 fluorescence rates. *Matter*. <https://doi.org/10.1016/j.matt.2021.02.017>
- 148 Romero, A. H., Allan, D. C., Amadon, B., Antonius, G., Applencourt, T., Baguet, L.,  
 149 Bieder, J., Bottin, F., Bouchet, J., Bousquet, E., Bruneval, F., Brunin, G., Caliste, D.,  
 150 Côté, M., Denier, J., Dreyer, C., Ghosez, P., Giantomassi, M., Gillet, Y., ... Gonze, X.  
 151 (2020). ABINIT: Overview and focus on selected capabilities. *Journal of Chemical Physics*,  
 152 152(12), 124102. <https://doi.org/10.1063/1.5144261>
- 153 Smith, D. G. A., Burns, L. A., Simmonett, A. C., Parrish, R. M., Schieber, M. C., Galvelis, R.,  
 154 Kraus, P., Kruse, H., Di Remigio, R., Alenaizan, A., James, A. M., Lehtola, S., Misiewicz,  
 155 J. P., Scheurer, M., Shaw, R. A., Schriber, J. B., Xie, Y., Glick, Z. L., Sirianni, D. A.,  
 156 ... Sherrill, C. D. (2020). PSI4 1.4: Open-source software for high-throughput quantum  
 157 chemistry. *Journal of Chemical Physics*, 152(18), 184108. [https://doi.org/10.1063/5.](https://doi.org/10.1063/5.0006002)  
 158 [0006002](https://doi.org/10.1063/5.0006002)
- 159 Stevens, M., Bursztein, E., Karpman, P., Albertini, A., & Markov, Y. (2017). The first  
 160 collision for full SHA-1. In J. Katz & H. Shacham (Eds.), *Advances in cryptology –*  
 161 *CRYPTO 2017* (pp. 570–596). Springer International Publishing. [https://doi.org/10.](https://doi.org/10.1007/978-3-319-63688-7_19)  
 162 [1007/978-3-319-63688-7\\_19](https://doi.org/10.1007/978-3-319-63688-7_19)
- 163 Tommaso, P. D. (2021). A curated list of awesome pipeline toolkits inspired by awe-  
 164 some sysadmin. In *GitHub repository*. GitHub. [https://github.com/pditommaso/](https://github.com/pditommaso/awesome-pipeline)  
 165 [awesome-pipeline](https://github.com/pditommaso/awesome-pipeline)
- 166 Weimer, W., & Nacula, G. C. (2008). Exceptional situations and program reliability. *ACM*  
 167 *Transactions on Programming Languages and Systems*, 30(2), 8:1–8:51. [https://doi.org/](https://doi.org/10.1145/1330017.1330019)  
 168 [10.1145/1330017.1330019](https://doi.org/10.1145/1330017.1330019)