

Empirical: A scientific software library for research, education, and public engagement

Anya Vostinar^{1,2,3}, Alexander Lalejini⁴, Charles Ofria^{1,2,3}, Emily Dolson^{1,2,3}, and Matthew Andres Moreno^{4,5}

¹ BEACON Center for the Study of Evolution in Action ² Computer Science and Engineering, Michigan State University ³ Ecology, Evolutionary Biology, and Behavior, Michigan State University ⁴ Computer Science, Carleton College ⁵ Ecology and Evolutionary Biology, University of Michigan ⁶ Center for the Study of Complex Systems, University of Michigan ⁷ Computer Science, Grand Valley State University

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

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

Editor: [↗](#)

Submitted: 14 December 2023

Published: unpublished

License

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

Summary

Empirical is a C++ library designed to promote open science and facilitate the development of scientific software that is efficient, reliable, and easily distributable to researchers and non-experts alike. Specifically, the library sets out to fulfill the following goals:

- Utility:** Empirical features a wide selection of helper tools to streamline common scientific computing tasks such as configuration, data management, random number generation, and mathematical manipulations. Documentation, examples, and project templates streamline developer onboarding.
- Efficiency:** Empirical tools emphasize efficiency to make heavy scientific computing workloads tractable. Where possible, computation is moved to compile-time, heap allocation is avoided, and memory-locality is respected.
- Reliability:** Scientific simulation and data analysis is meaningful insofar as it is correct. In conjunction with extensive unit testing, Empirical benefits from — and provides — sophisticated debug-mode instrumentation tools. These include audited memory management, drop-in replacements for standard C++ containers with safety checking, and self-documenting assertions.
- Distributability:** Scientific software should be easily used by any researcher, student, or citizen scientist who wants to experiment, explore, or replicate studies. As such, Empirical is highly portable, uses common data formats, facilitates flexible runtime configuration, and simplifies compilation to a performant web app.

In addition to many helpful utilities to improve the scientific programming experience, the core Empirical support library comprises four major features:

- debug-instrumented fundamental type and C++ standard library container wrappers,
- implementations of general-purpose data structures and algorithms,
- integrated, end-to-end frameworks for data and configuration management, and
- object-oriented bindings for Emscripten/WebAssembly GUI elements.

Statement of Need

Modern web-based interfaces give computational research the unique potential to embody open science objectives: they can make the scientific process more transparent with auditable and extensible code, clear and replicable methodologies, and production of accessible results (Woelfle et al., 2011). In practice, however, many scientific software applications are difficult to obtain, install, or use, and produce data in proprietary formats.

41 High quality open-science tools encourage researchers to follow effective software development
42 practices by simplifying development and helping them improve code quality, scientific rigor,
43 and ease of replication or extension. In the process, researchers get more out of their own code.
44 For example, adding a GUI not only helps other users of the software, but could also help the
45 researcher themselves gain “soft knowledge” of their system ([Dudley & Kristensson, 2018](#)).

46 Recent developments in web technology such as WebAssembly enable compilation of native
47 source code to browser-based interactive interfaces. However, many scientists lack web
48 development training and do not have the bandwidth to learn new languages and frameworks.
49 Empirical catalyzes progress toward open science ideals by streamlining the development of
50 in-browser software in C++.

51 To be reusable and extendable, software must remain robust at scales and in contexts that
52 were not originally envisioned. Empirical seeks to make writing correct, efficient scientific
53 software easier. Empirical’s debugging suite helps protect against common C++ programming
54 pitfalls such as iterator invalidation, memory leakage, and out-of-bounds indexing. Bundled
55 algorithms and data structures provide optimized, well-tested drop-in implementations for
56 common scientific computing tasks. Throughout, library design obviates trade-offs between
57 performance and safety; compile-time switches toggle safety checks for undefined or incorrect
58 behavior.

59 Empirical Features

60 Facilitating Better Code for Scientific Software

61 Software produced by academics, especially for one-off use, often foregoes rigorous programming
62 practices for the sake of expediency. By furnishing prepackaged components that address
63 common tasks for scientific software, the Empirical library helps scientific developers efficiently
64 write high quality C++ code. Utilities for common tasks empower users to craft more
65 readable and maintainable code. Bugs can be avoided by replacing one-off implementations
66 with Empirical components subjected to structured code review, unit testing (with coverage
67 tracking), and other modern software development practices detailed [in our documentation](#).
68 Sustained effort invested into optimization of the library’s utilities enables developers to
69 produce efficient software at far less effort. Finally, off-the-shelf solutions reduce barriers to
70 computational research, especially for developers new to scientific software design patterns.

71 To these ends, Empirical provides a comprehensive framework to manage runtime configuration
72 and flexible tools for data aggregation and recording. For example, Empirical’s configuration
73 framework includes utilities to

- 74 ▪ define and document default configuration values in a single line,
- 75 ▪ set configuration values via a combination of a configuration file or command line flags,
- 76 ▪ save configuration values to a file,
- 77 ▪ perform on-the-fly configuration adjustments, and
- 78 ▪ support multiple independent configuration subsystems.

79 Where appropriate, Empirical’s scientific software tools include features that integrate directly
80 into the browser environment. The configuration framework, for example, accepts input
81 via URL query parameters and ties in with a pre-built, in-browser GUI for setting-by-setting
82 adjustments.

83 In addition to the core C++ Library, we maintain a [template project](#) that streamlines laying
84 out crosscompilation boilerplate via a command-line wizard.

85 High-quality software cannot succeed without a thriving community of engineers. As detailed
86 [in our documentation](#), our development practices encourage the meaningful inclusion of all
87 interested contributors, and emphasize the importance of diverse backgrounds and perspectives
88 in the community.

Realizing the Promise of Emscripten-based Web UIs

Educational or outreach versions of scientific software can provide accessible windows into contemporary scientific work, promoting classroom learning, social media engagement, and citizen science. Scientific software projects that want to reach these additional audiences must typically maintain multiple codebases. As the codebases diverge, it becomes time consuming to synchronize changes across them, a problem that is only compounded when also maintaining crossplatform interfaces. These development costs preclude many scientific projects from providing easy access to the public. Even in better-resourced projects, this splintering effect absorbs limited developer hours and often leads to some versions of the code falling into neglect and drifting out of sync.

The Emscripten compiler promises to remedy this source splintering by enabling a single codebase to target web browsers alongside traditional native runtime environments (Zakai, 2011). Browser-based delivery can yield particularly effective public-facing apps due to widespread cross-platform compatibility, no-install access, and rich graphical interfaces. Native compilations are still required by most scientific applications, however, due to greater speeds and compatibility with high-performance computing environments.

Empirical amplifies the potential of Emscripten by fleshing out its rudimentary interface for interaction with browser elements. At the lowest level, Empirical provides tools for reciprocal data transfer between C++ code and the browser. DOM elements (such as `<button>`, `<div>`, and `<canvas>`) are given corresponding C++ objects (`emp::Button`, `emp::Div`, and `emp::Canvas`) and can be easily used from within C++ code. With these tools, users no longer need to manage JavaScript resources, and thus need much less preexisting web-programing knowledge. At a higher level of abstraction, Empirical packages pre-configured, pre-styled collections of DOM elements as prefabricated widgets (e.g., configuration managers, collapsible read-outs, modal messages, etc.). Empirical's tools aim to make generating a mobile-friendly, web-based GUI for existing software so trivial that the practice becomes ubiquitous.

Below, we give an example of Empirical's DOM interface in action. This example creates a button that increments an on-screen counter every time the button is clicked. You can view the resulting web page live at <https://devosoft.github.io/empirical-joss-demo/>.

C++ source:

```
#include "emp/web/web.hpp"

emp::web::Document doc("target");

int x = 5;
int main() {
    doc << "<h1>Hello World!</h1>";
    doc << "Original x = " << x << "<br>";
    doc << "Current x = " << emp::web::Live(x) << "<br>";

    // Create a button to modify x.
    emp::web::Button my_button( [](){ x+=5; doc.Redraw(); }, "Click me!" );
    doc << my_button;
}
```

HTML source:

```
<body>
  <div id="target"> </div>
</body>

<script src="https://code.jquery.com/jquery-1.11.2.min.js" integrity="sha256-Ls0pXS1b7AY"
<script type="text/javascript" src="main.js"></script>
```

120 A live demo of more sophisticated Empirical widgets, presented alongside their source C++
121 code, is available on our [prefab demos page](#).

122 Facilitating Runtime Efficiency

123 WebAssembly's runtime efficiency is a major driver of its increasing popularity for web app
124 development. WebAssembly is a virtual bytecode, but just-in-time compilation engines translate
125 critical sections into native machine code ([Haas et al., 2017](#)). This compilation model allows
126 WebAssembly to achieve at least 50% — and at times closer to 90% — of native performance,
127 providing an order of magnitude speed increase over JavaScript alone ([Jangda et al., 2019](#)).
128 These performance improvements open the door to entirely new possibilities for browser-based
129 scientific computation. For example, the Avida-ED web viewer at <https://avida-ed.msu.edu/>,
130 uses WebAssembly to simulate hundreds of thousands of generations of digital organisms within
131 the span of a class period. Such rich, intensive in-browser experiences necessitate efficient
132 source code.

133 More broadly, Empirical caters to the necessity for runtime efficiency across all scientific
134 computing, including in-browser applications, local runs, or high-performance computing
135 (HPC) cluster deployments. In some contexts, even modest performance gains can save
136 substantial hardware costs, meaningfully reduce energy use, and shave days or even weeks
137 off run times. Order-of-magnitude performance gains can meaningfully broaden the scope of
138 scientific questions that are tractable.

139 Empirical supports runtime efficiency in scientific computing by providing optimized tools
140 for performance-critical tasks. For example, `emp::BitSet` and `emp::BitVector` classes
141 are faster drop in replacements for their standard library equivalents (`std::bitset` and
142 `std::vector<bool>`) while providing extensive additional functionality for rapid bit manipula-
143 tions. Likewise, `emp::Random` wraps a cutting-edge high-performance pseudorandom number
144 generator algorithm ([Widynski, 2020](#)). Benchmark-informed development practices ensure that
145 optimizations translate into consistent performance enhancements. At a more fundamental
146 level, Empirical's header-only design prioritizes ease of use and runtime performance at the
147 cost of somewhat longer compilation times.

148 Facilitating Debugging

149 Identifying and correcting incorrect program behavior consumes a large fraction of developer
150 hours for any software project. Software bugs that slip through into production can inflict even
151 greater costs, especially in scientific contexts where the validity of generated data and analyses
152 is paramount.

153 In conjunction with unit tests and integration tests, runtime safety checks are commonly used
154 to flag potential bugs. Assert statements typify runtime safety checks. These statements abort
155 program execution at the point of failure with a helpful error message if an expected runtime
156 condition is not met. Runtime safety checks like `assert` don't necessarily oblige a performance
157 cost to compute the asserted runtime condition; these checks can be verified only in debug
158 mode and ignored in production mode to maximize performance.

159 Indeed, the C++ standard library's `assert` macro follows this paradigm. Empirical provides an
160 extended `emp_assert` macro that prints custom error messages with current values of specified
161 expressions, and dispatches a UI alert when triggered in a web environment.

162 In addition to user-defined asserts, most programming languages (Java, Python, Ruby, Rust,
163 etc.) provide built-in support to detect common runtime violations, such as out-of-bounds
164 indexing or bad type conversions. C++ does not in an effort to maximize performance. However,
165 standard library vendors — like GCC's `libstdc++`, Clang's `libc++`, and Microsoft's `stl` — do
166 provide some proprietary support for such safety checks. This support, however, is limited and

167 poorly documented¹. Empirical supplements vendors' runtime safety checking by providing
168 drop-in replacements for `std::array`, `std::optional`, and `std::vector` with stronger runtime
169 safety checks, but only while in debug mode. In addition, Empirical furnishes a safety-checked
170 pointer wrapper, `emp::Ptr`, that identifies memory leaks and invalid memory access in debug
171 mode while retaining the full speed of raw pointers in release mode.

172 Because of poor support for built-in runtime safety checks, C++ developers typically use an
173 external toolchain to detect and diagnose runtime violations. Popular tools include Valgrind,
174 GDB, and runtime sanitizers. Although this tooling is very mature and quite powerful, there
175 are fundamental limitations to the runtime violations it can detect. For example, Clang 12.0.0's
176 sanitizers cannot detect the iterator invalidation described above ([live example](#)). Additionally,
177 most of this tooling is not available when debugging WASM code compiled with Emscripten —
178 a core use case targeted by the Empirical library. Although Emscripten provides some [sanitizer](#)
179 [support](#) and [other debugging features](#), tooling limitations (such as the lack of a steppable
180 debugger) make runtime safety checking particularly critical.

181 Outlook and Future Plans

182 Empirical remains under active development. Current priorities include assembling higher-
183 level web widgets for common tasks, making existing classes more web-friendly (such as
184 file management and rich text handling), and adding more step-by-step tutorials to our
185 documentation.

186 We are committed to maintaining a stable interface for existing users. Last year, we took a major
187 step towards fulfilling this objective on an ongoing basis by completing a major reorganization
188 informed by best practices to expose sustainable, consistent API to our end-users. We maintain
189 an extensive suite of unit tests and integration tests to ensure that continuing development
190 retains backward compatibility. In addition, our software releases are archived on Zenodo in
191 order to guarantee uninterrupted, perpetual access to our software for those who depend on it
192 ([Ofria et al., 2020](#)).

193 Empirical has already been successfully incorporated into major projects within our research
194 group's primary domains: digital evolution, artificial life, and genetic programming. We
195 aim for potential utility across a much broader swath of the scientific software ecosystem,
196 particularly among projects that prioritize open science objectives. To this end, we look forward
197 to welcoming new collaborations and supporting a wider collection of end-users.

198 Related Software Packages

199 Software Addressing Related Needs

200 There are many existing software platforms that provide functionalities overlapping with
201 Empirical. However, most are not in C++, and there is value in this functionality being easily
202 available to C++ programmers. See the Non-C++ Comparable Software section for citations
203 to software platforms that provide some of Empirical's functionality in different languages.

204 RepastHPC

205 RepastHPC, accessible at <https://repast.github.io/>, is a C++ modeling framework targeted
206 at large computing clusters and supercomputers ([Collier & North, 2013](#); [North et al., 2013](#)).
207 A Java-based counterpart, Repast Symphony, provides interactive GUI support. As such,
208 simultaneous on-cluster and in-browser support requires maintenance of two separate code
209 bases.

¹For example, neither GCC 10.3 nor Clang 12.0.0 detect `std::vector` iterator invalidation when appending
to a `std::vector` happens to fall within existing allocated buffer space ([GCC live example](#); [Clang live example](#)).

210 Boost C++ Libraries

211 Boost C++ Libraries, available at <https://www.boost.org/>, provide an enormous range of
212 software components. Several of the Boost libraries have been already incorporated into the
213 C++ standard, and we build off of those with Empirical. However, Boost does not contain
214 libraries for web-based GUI tools, configuration management, or data management specifically
215 tailored to scientific software.

216 Emscripten

217 Emscripten is available at <https://emscripten.org/> (Zakai, 2011). It provides cross-compilation
218 from C++ to WebAssembly and we use it in Empirical. Empirical's tools build abstractions from
219 Emscripten intrinsics tailored to visualization and interactive control of scientific simulations.

220 Cheerp

221 Cheerp, another C++ to WebAssembly compiler, is available at [https://leaningtech.com/](https://leaningtech.com/cheerp/)
222 [cheerp/](https://leaningtech.com/cheerp/). Like Emscripten, Cheerp provides primarily low-level APIs for interaction with browser
223 GUI elements.

224 Non-C++ Comparable Software

- 225 ▪ [TinyGo](#)
- 226 ▪ [WebIO](#)
- 227 ▪ [GWT](#)
- 228 ▪ [yew](#)
- 229 ▪ Pyodide (Droettboom & Developers, 2021)
- 230 ▪ Shiny (Chang et al., 2020)

231 Projects Using the Software

- 232 ▪ [Aagos](#) (Gillespie et al., 2018)
233 – An interactive model for empirically studying how environmental change affects the
234 evolution of genetic architectures.
- 235 ▪ [conduit](#) (Moreno et al., 2020)
236 – A C++ library that wraps intra-thread, inter-thread, and inter-process communica-
237 tion in a uniform, modular, object-oriented interface, with a focus on asynchronous
238 high-performance computing applications.
- 239 ▪ [Dishtiny](#) (Moreno & Ofria, 2019)
240 – An interactive web demo of an artificial life platform built to study major transitions
241 in evolution.
- 242 ▪ [ecology in evolutionary computation explorer](#) (Dolson & Ofria, 2018)
243 – Interactive exploration of interaction networks in evolutionary computation under
244 different selection schemes.
- 245 ▪ [Symbulation](#) (Vostinar, 2017)
246 – An interactive artificial life model focused on the evolution of symbiosis between
247 parasitism to mutualism.
- 248 ▪ [SignalGP](#) (Lalejini & Ofria, 2018)
249 – SignalGP is a new GP technique designed to give evolution direct access to the
250 event-driven programming paradigm where computations are triggered response to
251 signals from the environment, from other agents, or that are internally generated.
- 252 ▪ [Model of cancer evolution on an oxygen gradient](#)
253 – A companion model to a series of wet lab experiments on cancer evolution in
254 spatially heterogenous environments

Acknowledgements

This research was supported in part by NSF grants DEB-1655715 and DBI-0939454, by the National Science Foundation Graduate Research Fellowship under Grant No. DGE-1424871, by Michigan State University through the computational resources provided by the Institute for Cyber-Enabled Research, and by the Eric and Wendy Schmidt AI in Science Postdoctoral Fellowship, a Schmidt Futures program. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

References

- Chang, W., Cheng, J., Allaire, J., Xie, Y., & McPherson, J. (2020). *Shiny: Web application framework for r*. <https://CRAN.R-project.org/package=shiny>
- Collier, N., & North, M. (2013). Parallel agent-based simulation with repast for high performance computing. *SIMULATION*, 89(10), 1215–1235. <https://doi.org/10.1177/0037549712462620>
- Dolson, E., & Ofria, C. (2018). Ecological theory provides insights about evolutionary computation. *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 105–106. <https://doi.org/10.1145/3205651.3205780>
- Droettboom, M., & Developers, P. (2021). *Shiny: Web application framework for r*. <https://pyodide.org/>
- Dudley, J. J., & Kristensson, P. O. (2018). A review of user interface design for interactive machine learning. *ACM Trans. Interact. Intell. Syst.*, 8(2). <https://doi.org/10.1145/3185517>
- Gillespie, L., Dolson, E., Lalejini, A., & Ofria, C. (2018). Changing environments drive the separation of genes and increased evolvability in NK-inspired landscapes. *Late Breaking Abstract at The 2018 Conference on Artificial Life*.
- Haas, A., Rossberg, A., Schuff, D. L., Titzer, B. L., Holman, M., Gohman, D., Wagner, L., Zakai, A., & Bastien, J. (2017). Bringing the web up to speed with WebAssembly. *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*, 185–200. <https://doi.org/10.1145/3062341.3062363>
- Jangda, A., Powers, B., Berger, E. D., & Guha, A. (2019). *Not so fast: Analyzing the performance of webassembly vs. Native code*. 2019 {USENIX} Annual Technical Conference ({USENIX}{ATC} 19), 107–120. ISBN: 9781939133038
- Lalejini, A., & Ofria, C. (2018). Evolving event-driven programs with SignalGP. *Proceedings of the Genetic and Evolutionary Computation Conference on - GECCO '18*, 1135–1142. <https://doi.org/10.1145/3205455.3205523>
- Moreno, M. A., & Ofria, C. (2019). Toward Open-Ended Fraternal Transitions in Individuality. *Artificial Life*, 25(2), 117–133. https://doi.org/10.1162/artl_a_00284
- Moreno, M. A., rodsan0, perryk12, & Badger, C. (2020). *mmore500/conduit: Initial release (Version v0.1.0)*. Zenodo. <https://doi.org/10.5281/zenodo.4118608>
- North, M. J., Collier, N. T., Ozik, J., Tatara, E. R., Macal, C. M., Bragen, M., & Sydelko, P. (2013). Complex adaptive systems modeling with repast simphony. *Complex Adaptive Systems Modeling*, 1(1), 1–26. <https://doi.org/10.1186/2194-3206-1-3>
- Ofria, C., Moreno, M. A., Dolson, E., Lalejini, A., rodsan0, Fenton, J., perryk12, Jorgensen, S., hoffmanriley, grenewode, & al., et. (2020). *Devosoft/empirical*. <https://doi.org/10.5281/zenodo.2575606>

- 300 Vostinar, A. E. (2017). *Suicide, signals, and symbionts: Evolving cooperation in agent-based*
301 *systems*. Michigan State University. ISBN: 978-0-355-07992-0
- 302 Widynski, B. (2020). Squares: A fast counter-based RNG. *arXiv Preprint arXiv:2004.06278*.
303 <https://arxiv.org/abs/2004.06278v3>
- 304 Woelfle, M., Oliaro, P., & Todd, M. H. (2011). Open science is a research accelerator. *Nature*
305 *Chemistry*, 3(10), 745–748. <https://doi.org/10.1038/nchem.1149>
- 306 Zakai, A. (2011). Emscripten: An LLVM-to-JavaScript compiler. *Proceedings of the ACM*
307 *International Conference Companion on Object Oriented Programming Systems Languages*
308 *and Applications Companion*, 301–312. <https://doi.org/10.1145/2048147.2048224>

DRAFT