

Fiats: Functional inference and training for surrogates

Damian Rouson¹, Dan Bonachea¹, Brad Richardson¹, Jordan A. Welsman¹, Jeremiah Bailey¹, Ethan D Gutmann², David Torres³, Katherine Rasmussen¹, Baboucarr Dibba¹, Yunhao Zhang¹, Kareem Weaver¹, Zhe Bai¹, and Tan Nguyen¹

¹ Lawrence Berkeley National Laboratory, United States ² NSF National Center for Atmospheric Research, United States ³ Northern New Mexico College, United States

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

Software

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

Editor: [Evan Spotte-Smith](#)

Reviewers:

- [@jwallwork23](#)
- [@niccolozanotti](#)

Submitted: 02 July 2025

Published: unpublished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

Fiats provides a platform for research on the training and deployment of neural-network surrogate models for computational science. Fiats also supports exploring, advancing, and combining functional, object-oriented, and parallel programming patterns in Fortran 2023 ([Fortran Standards Committee JTC1/SC22/WG5, Nov 2023](#)). As such, the Fiats name has dual expansions: “Functional Inference And Training for Surrogates” or “Fortran Inference And Training for Science.” Fiats inference and training procedures are pure and therefore satisfy a language constraint imposed on procedure invocations inside Fortran’s parallel loop construct: `do concurrent`. Furthermore, the Fiats training procedures are built around a `do concurrent` parallel reduction. Several compilers can automatically parallelize `do concurrent` on Central Processing Units (CPUs) or Graphics Processing Units (GPUs). Fiats thus aims to achieve performance portability through standard language mechanisms.

In addition to an example subdirectory with illustrative codes, the Fiats `demo/app` subdirectory contains three demonstration applications:

1. One trains a cloud-microphysics surrogate for the Berkeley Lab [fork](#) of the Intermediate Complexity Atmospheric Research (ICAR) model.
2. Another calculates input- and output-tensor statistics for ICAR’s physics-based microphysics models.
3. A third performs batch inference using an aerosols surrogate for the Energy Exascale Earth Systems Model ([E3SM](#)).

Ongoing research explores how Fiats can exploit multi-image execution, a set of Fortran features for Single-Program, Multiple-Data (SPMD) parallel programming with a Partitioned Global Address Space (PGAS) ([Numrich, 2018](#)), where the PGAS features center around “coarray” distributed data structures.

To explore how new language features and novel uses of longstanding features can power deep learning, Fiats contributors work to advance Fortran by

- Participating in the Fortran standardization process and
- Contributing to compiler development through
 - Writing unit tests ([Rasmussen et al., 2022](#)),
 - Studying performance ([Rouson, Bai, Bonachea, Ergawy, et al., 2025](#)),
 - Isolating and reporting compiler bugs and fixing front-end bugs,
 - Publishing and updating the Parallel Runtime Interface for Fortran (PRIF) ([Bonachea et al., 2024a, 2024b](#)), and
 - Developing the first PRIF implementation: [Caffeine](#) ([Bonachea et al., 2025](#); [Rouson & Bonachea, 2022](#)).

43 Fiats thus facilitates studying deep learning for science and studying programming paradigms
44 and patterns for deep learning in Fortran 2023.

45 Statement of need

46 Fortran 2008 introduced two forms of parallelism: `do concurrent` for loop-level parallelism and
47 multi-image execution for SPMD/PGAS parallelism in shared or distributed memory. Fortran
48 2018 and 2023 expanded and refined these features by adding, for example,

- 49 1. `Do concurrent` iteration locality specifiers, including reductions and
- 50 2. Collective subroutines, image teams, events (semaphores), atomic subroutines, and more,

51 which creates a need for libraries and frameworks that support users who adopt these features.
52 For example, one requirement impacting library design stems from the aforementioned language
53 constraint allowing only side-effect-free (pure) procedure invocations inside `do concurrent`.

54 All intrinsic functions defined in the Fortran 2023 standard are `simple`, an attribute that
55 implies pure plus additional constraints. Libraries that export pure procedures thus behave
56 like extensions of the language. To wit, the Fortran 2023 standard states: “It is expected that
57 most library procedures will conform to the constraints required of pure procedures, and so
58 can be declared pure and referenced in `do concurrent` constructs... and within user-defined
59 pure procedures.”

60 Conversely, multi-image execution in a library places a requirement on the client code. The
61 Fortran standard defines steps for synchronized image launch, synchronized normal termination,
62 and single-image initiation of global error termination. Multi-image execution in a library thus
63 requires support for multi-image execution in the main program.

64 A surrogate model’s utility hinges upon inference calculations executing faster than the physics-
65 based model the surrogate replaces. This commonly restricts the surrogate neural network to a
66 few thousand tunable parameters. For networks of modest size, useful insights can sometimes
67 be gleaned from visually inspecting the network parameters. Fiats therefore stores networks in
68 human-readable JavaScript Object Notation (JSON) format. The Fiats companion package
69 [Nexport](#) exports Fiats JSON files [PyTorch](#).

70 State of the field

71 At least six open-source software packages provide deep learning services to Fortran. Three
72 provide Fortran application programming interfaces (APIs) that wrap C++ libraries:

- 73 ■ [Fortran-TF-Lib](#) is a Fortran API for [TensorFlow](#),
- 74 ■ [FTorch](#) is a Fortran API for `libtorch`, the PyTorch back-end, and
- 75 ■ [TorchFort](#) is also a Fortran API for `libtorch`.

76 As of this writing, recursive searches in the root directories of the these three projects find
77 no pure procedures. Procedures are pure if declared as such or if declared `simple` or if
78 declared `elemental` without the `impure` attribute. Because any procedure invoked within a
79 pure procedure must also be pure, the absence of pure procedures precludes the use of these
80 APIs anywhere in the call stack inside a `do concurrent` construct. Also, as APIs backed by
81 C++ libraries, none use Fortran’s multi-image execution features.

82 Three packages supporting deep learning in Fortran are themselves written in Fortran:

- 83 ■ [Athena](#) ([Taylor, 2024](#))
- 84 ■ [Fiats](#) ([Rouson, Bai, Bonachea, Ergawy, et al., 2025](#))
- 85 ■ [neural-fortran](#) ([Curcic, 2019](#))

86 Searching the Athena, Fiats, and neural-fortran src subdirectories finds that over half of the
87 procedures in each are pure, including 75% of Fiats procedures. Included in these tallies are
88 procedures explicitly marked as pure along with simple procedures and elemental procedures
89 without the impure attribute. Athena, Fiats, and neural-fortran each employ do concurrent
90 extensively. Only Fiats, however, leverages the locality specifiers introduced in Fortran 2018
91 and expanded in Fortran 2023 to include parallel reductions.

92 Of the APIs and libraries discussed here, only neural-fortran and Fiats use multi-image features:
93 neural-fortran in its core library and Fiats in a demonstration application. Both use multi-image
94 features minimally, leaving considerable room for researching parallelization strategies.

95 Each of the Fortran deep learning APIs and libraries discussed in this paper is actively developed
96 except Fortran-TF-Lib. Fortran-TF-Lib's most recent commit was in 2023 and no releases
97 have been posted. The other mentioned projects have most-recent commits no older than two
98 months as of June 2025.

99 Recent research and scholarly publications

100 Fiats supports research in training surrogate models and parallelizing batch inference calculations
101 for atmospheric sciences. This research has generated two peer-reviewed paper submissions:
102 one accepted to appear in workshop proceedings (Rouson, Bai, Bonachea, Ergawy, et al., 2025)
103 and one in open review (Rouson, Bai, Bonachea, Dibba, et al., 2025).

104 Four programs in the Fiats repository played significant roles in these two papers:

- 105 1. [example/concurrent-inferences.f90](#),
- 106 2. [example/learn-saturated-mixing-ratio.f90](#),
- 107 3. [app/demo/infer-aerosols.f90](#), and
- 108 4. [app/demo/train-cloud-microphysics.f90](#).

109 Rouson, Bai, Bonachea, Ergawy, et al. (2025) used program 1 to study automatically
110 parallelizing batch inferences via do concurrent. Rouson, Bai, Bonachea, Dibba, et al. (2025)
111 used programs 2–4 to study neural-network training for cloud microphysics and inference
112 for atmospheric aerosols. The derived types in the Unified Modeling Language (UML) class
113 diagram in Figure 1 enabled these studies.

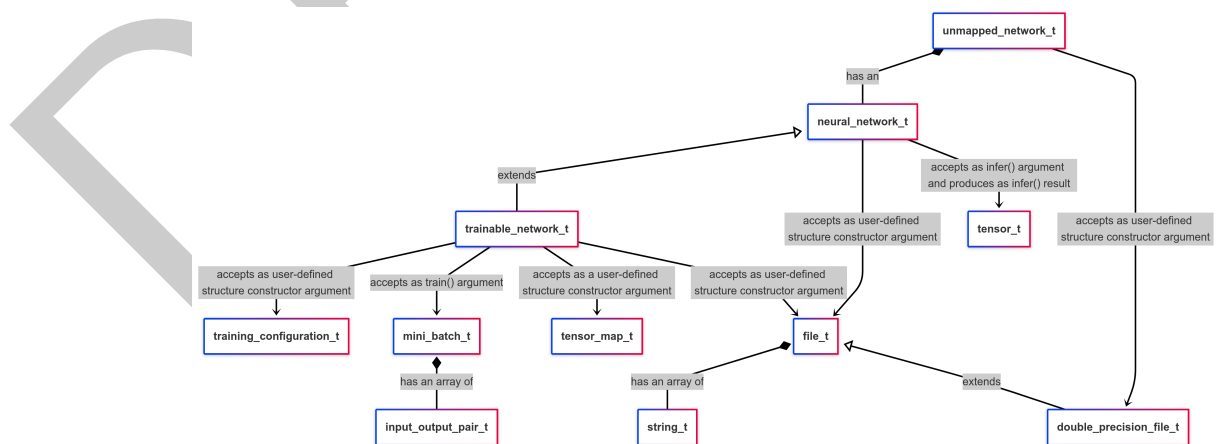


Figure 1: Class diagram: derived types (named in bordered white boxes), type relationships (connecting lines), type extension (open triangles), composition (solid diamonds), or directional relationship (arrows). Read relationships as sentences wherein the type named at the base of an arrow is the subject followed by an annotation (in an unbordered gray box) followed by the type named at the arrow's head as the object. Type extension reads with the type adjacent to the open triangle as the subject. Composition reads with the type adjacent to the closed diamond as the subject.

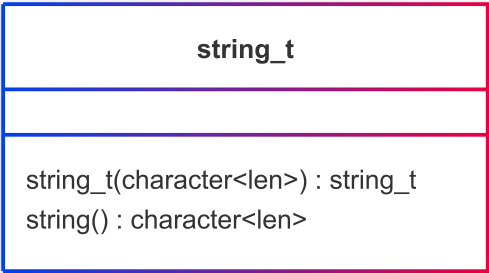


Figure 2: String class diagram



Figure 3: File class diagram

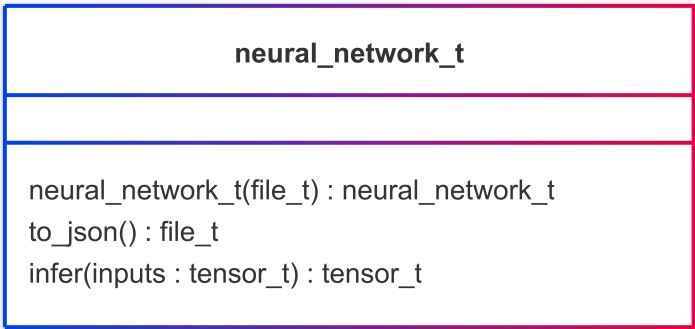


Figure 4: Neural network class diagram

114 Figure 1 includes two of the [Julienne](#) correctness-checking framework’s derived types, `string_t`
115 and `file_t`. These are included because other parts of the figure reference these types. The
116 rightmost four types in Figure 1 exist primarily to support inference. The leftmost six support
117 training. Because inference is considerably simpler, it makes sense to describe the right side of
118 the diagram before the left side.

119 The concurrent-inferences example program performs batch inference using the `string_t`,
120 `file_t`, and `neural_network_t` types. Figure 2 through Figure 4 show class diagrams with
121 more details on these types. Each detailed diagram displays a top panel listing the type name,
122 an empty middle panel with private components omitted, and a bottom panel listing public
123 procedure bindings.

124 The bottom panel also lists what the Fortran 2023 standard describes as user-defined structure
125 constructors, which are generic interfaces through which to invoke functions that define a
126 result of the named type ([Fortran Standards Committee JTC1/SC22/WG5, Nov 2023](#)). We
127 henceforth refer to these as “constructors.” From the bottom of the class hierarchy in Figure 1,
128 the concurrent-inferences program does the following:



Figure 5: Tensor class diagram

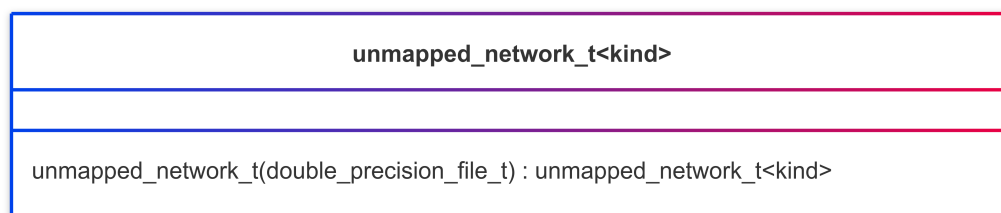


Figure 6: Unmapped network class diagram

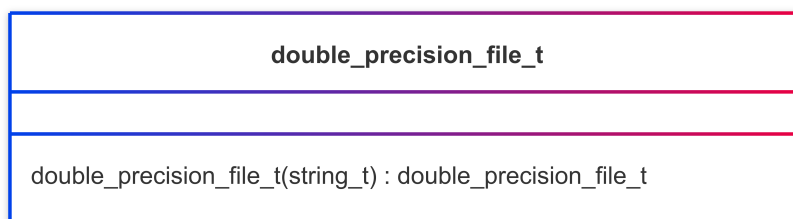


Figure 7: Double precision file class diagram

- 129 1. Gets a character file name from the command line,
- 130 2. Passes the name to a string_t constructor,
- 131 3. Passes the resulting string_t object to a file_t constructor, and
- 132 4. Passes the resulting file_t object to a neural_network_t constructor.

133 The program then repeatedly invokes the infer type-bound procedure on a three-dimensional
 134 (3D) array of tensor_t objects (see Figure 5) using OpenMP directives or do concurrent or
 135 an array statement. The array statement takes advantage of infer being elemental. Line 101
 136 of example/concurrent-inferences.f90 at git tag joss-line-references demonstrates
 137 neural-network construction from a file. Line 109 demonstrates using the network for inference.

138 The infer-aerosols program performs inferences by invoking double precision versions
 139 of the infer generic binding on an object of type unmapped_network_t (see Figure 6), a
 140 parameterized derived type (PDT) that has a kind type parameter. To match the expected
 141 behavior of the aerosol model, which was trained in PyTorch, the unmapped_network_t
 142 implementation ensures the use of raw network input and output tensors without the nor-
 143 malizations and remappings that are performed by default for a neural_network_t object.
 144 The double_precision_file_t (see Figure 7) type controls the interpretation of the JSON
 145 network file: JSON does not distinguish between categories of numerical values such as real,
 146 double precision, or even integer, so something external to the file must determine the
 147 interpretation of the numbers in a JSON file.

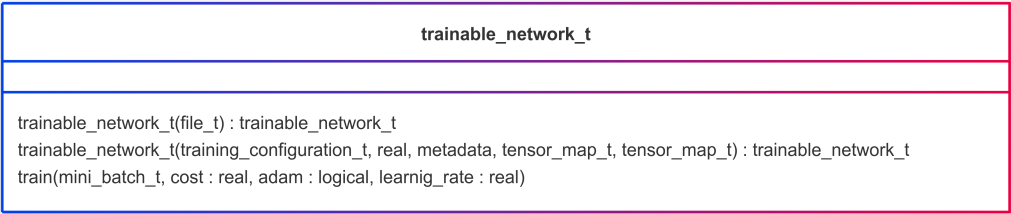


Figure 8: Trainable network class diagram

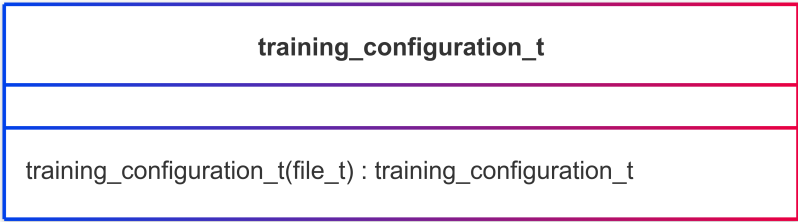


Figure 9: Training configuration class diagram

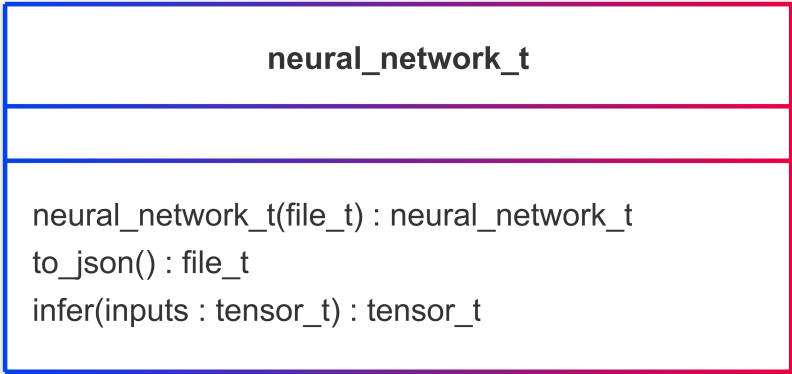


Figure 10: Mini-batch class diagram

148 The learn-saturated-mixing-ratio and train-cloud-microphysics programs focus on
149 using a trainable_network_t object (see Figure 8) for training. The former trains neural
150 network surrogates for a thermodynamic function from ICAR: the saturated mixing ratio, a
151 scalar function of temperature and pressure. The latter trains surrogates for the complete cloud
152 microphysics models in ICAR – models implemented in thousands of lines of code. Whereas
153 diagrammed relationships of neural_network_t reflect direct dependencies of only two types
154 (file_t and tensor_t), even describing the basic behaviors of trainable_network_t requires
155 showing dependencies on five types:

- 156 ■ A training_configuration_t object (see Figure 9), which holds hyperparameters such
157 as the learning rate and choice of optimization algorithms,
- 158 ■ A file_t object from which the training_configuration is read inside the
159 trainable_network_t constructor,
- 160 ■ A mini_batch_t object (see Figure 10) that stores an array of input_output_pair
161 objects (see Figure 11) from the training data set,
- 162 ■ Two tensor_map_t objects (see Figure 12) storing the linear functions that map inputs
163 to the training data range and map outputs from the training data range back to the

- 164 application range, and
165 ■ A parent `neural_network_t` object storing the network architecture, including weights,
166 biases, layer widths, etc.

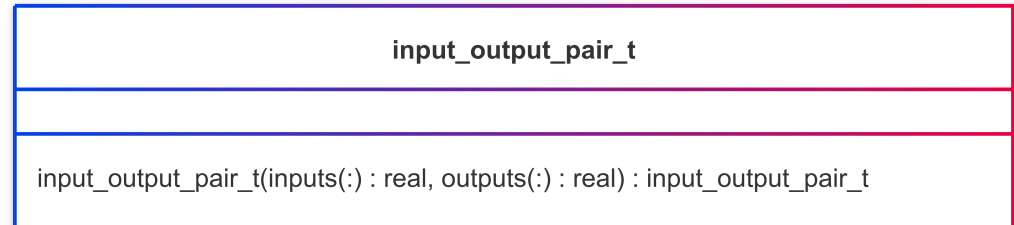


Figure 11: Input/Output tensor pair class diagram

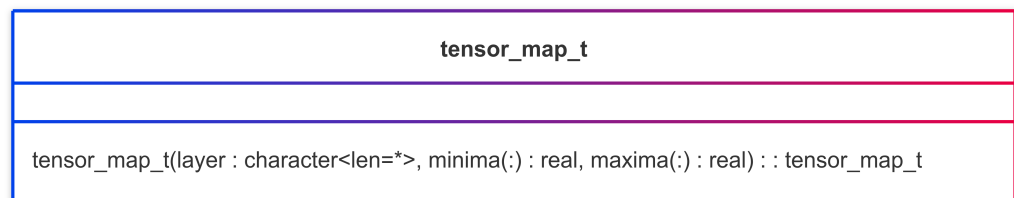


Figure 12: Tensor map class diagram

167 The `trainable_network_t` type stores a `workspace_t` (not shown) as a scratch-pad for
168 training purposes. The workspace is not needed for inference. During each training
169 step, a `trainable_network_t` object passes its `workspace_t` to a learn procedure binding
170 (not shown) on its parent `neural_network_t`. Lines 388–396 of `demo/app/train-cloud-`
171 `microphysics.f90` at git tag `joss-line-references` demonstrate:

- 172 1. A loop over epochs,
- 173 2. The shuffling of the `input_output_pair_t` objects at the beginning of each epoch,
- 174 3. The grouping of `input_output_pair_t` objects into `mini_batch_t` objects, and
- 175 4. The invocation of the train procedure for each mini-batch,

176 where steps 2 and 3 express deep learning’s stochastic gradient descent algorithm.

177 Acknowledgments

178 This material is based upon work supported by the U.S. Department of Energy, Office of
179 Science, Office of Advanced Scientific Computing Research and Office of Nuclear Physics, Sci-
180 entific Discovery through Advanced Computing (SciDAC) Next-Generation Scientific Software
181 Technologies (NGSST) programs under Contract No. DE-AC02-05CH11231. This material
182 is also based on work supported by Laboratory Directed Research and Development (LDRD)
183 funding from Lawrence Berkeley National Laboratory, provided by the Director, Office of
184 Science, of the U.S. DOE under Contract No. DE-AC02-05CH11231. This manuscript has
185 been authored by an author at Lawrence Berkeley National Laboratory under Contract No. DE-
186 AC02-05CH11231 with the U.S. Department of Energy. The U.S. Government retains, and the
187 publisher, by accepting the article for publication, acknowledges, that the U.S. Government
188 retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the
189 published form of this manuscript, or allow others to do so, for U.S. Government purposes.

References

- Bonachea, D., Rasmussen, K., Richardson, B., & Rouson, D. (2024a). Parallel Runtime Interface for Fortran (PRIF): A Multi-Image Solution for LLVM Flang. *Tenth Workshop on the LLVM Compiler Infrastructure in HPC (LLVM-HPC2024)*. <https://doi.org/10.25344/S4N017>
- Bonachea, D., Rasmussen, K., Richardson, B., & Rouson, D. (2024b). *Parallel runtime interface for Fortran (PRIF) specification (rev. 0.5)*. Lawrence Berkeley National Laboratory (LBNL), Berkeley, CA (United States). <https://doi.org/10.25344/S4CG6G>
- Bonachea, D., Rasmussen, K., Richardson, B., & Rouson, D. (2025). Caffeine: A parallel runtime library for supporting modern Fortran compilers. *Journal of Open Source Software*, 10(107), 7895. <https://doi.org/10.21105/joss.07895>
- Curcic, M. (2019). A parallel Fortran framework for neural networks and deep learning. *Acm Sigplan Fortran Forum*, 38, 4–21. <https://doi.org/10.1145/3323057.3323059>
- Fortran Standards Committee JTC1/SC22/WG5. (Nov 2023). *Information technology – programming languages – ISO/IEC 1539-1:2023*. International Organization for Standardization (ISO).
- Numrich, R. W. (2018). *Parallel Programming with Co-arrays*. CRC Press. <https://doi.org/doi:10.1201/9780429437182>
- Rasmussen, K., Rouson, D., George, N., Bonachea, D., Kadhem, H., & Friesen, B. (2022). *Agile acceleration of LLVM flang support for Fortran 2018 parallel programming*. Lawrence Berkeley National Laboratory (LBNL), Berkeley, CA (United States). <https://doi.org/10.25344/S4CP4S>
- Rouson, D., Bai, Z., Bonachea, D., Dibba, B., Gutmann, E., Rasmussen, K., Torres, D., Welsman, J., & Zhang, Y. (2025). Cloud microphysics training and aerosol inference with the Fiats deep learning library. *2025 Improving Scientific Software Conference*, (in review).
- Rouson, D., Bai, Z., Bonachea, D., Ergawy, K., Gutmann, E., Klemm, M., Rasmussen, K., Richardson, B., Shende, S., Torres, D., & Zhang, Y. (2025). Automatically parallelizing batch inference on deep neural networks using Fiats and Fortran 2023 “do concurrent.” *5th International Workshop on Computational Aspects of Deep Learning (CADL)*. <https://doi.org/10.25344/S4VG6T>
- Rouson, D., & Bonachea, D. (2022). Caffeine: CoArray Fortran Framework of Efficient Interfaces to Network Environments. *2022 IEEE/ACM Eighth Workshop on the LLVM Compiler Infrastructure in HPC (LLVM-HPC)*, 34–42. <https://doi.org/10.25344/S4459B>
- Taylor, N. T. (2024). ATHENA: A Fortran package for neural networks. *Journal of Open Source Software*, 9(99), 6492. <https://doi.org/10.21105/joss.06492>