

---

# AbstractDifferentiation.jl: Backend-Agnostic Differentiable Programming in Julia

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 No single Automatic Differentiation (AD) system is the optimal choice  
2 for all problems. This means informed selection of an AD system and  
3 combinations can be a problem-specific variable that can greatly impact  
4 performance. In the Julia programming language, the major AD systems  
5 target the same input and thus in theory can compose. Hitherto, switching  
6 between AD packages in the Julia Language required end-users to familiarize  
7 themselves with the user-facing API of the respective packages. Furthermore,  
8 implementing a new, usable AD package required AD package developers  
9 to write boilerplate code to define convenience API functions for end-users.  
10 As a response to these issues, we present AbstractDifferentiation.jl for  
11 the automatized generation of an extensive, unified, user-facing API for  
12 any AD package. By splitting the complexity between AD users and  
13 AD developers, AD package developers only need to implement one or  
14 two primitive definitions to support various utilities for AD users like  
15 Jacobians, Hessians and lazy product operators from native primitives such  
16 as pullbacks or pushforwards, thus removing tedious – but so far inevitable  
17 – boilerplate code, and enabling the easy switching and composing between  
18 AD implementations for end-users.

## 19 1 Introduction

20 Differentiable programming ( $\partial P$ ), i.e., the ability to differentiate general computer program  
21 structures, has enabled the efficient combination of existing packages for scientific computation  
22 and machine learning [Raissi et al., 2019, Rackauckas et al., 2020a, de Avila Belbute-Peres  
23 et al., 2018]. Black-box machine learning approaches are flexible but require a large amount  
24 of data. Incorporating scientific knowledge about the structure of a problem via  $\partial P$  reduces  
25 the amount of data needed. It allows the learning task to be simplified, for example,  
26 by focusing on learning only the parts of the model that are missing [Rackauckas et al.,  
27 2020b, Dandekar et al., 2020]. There are already many examples where such differentiable  
28 frameworks have provided performance and accuracy advantages over black-box approaches  
29 to machine learning, including but not limited to protein-folding [AlQuraishi, 2018, Ingraham  
30 et al., 2018], fluid dynamics [Schenck and Fox, 2018], robotics [Schenck and Fox, 2018], and  
31 quantum control [Schäfer et al., 2020, 2021].

32  $\partial P$  is (commonly) realized by automatic differentiation (AD), a family of techniques to  
33 efficiently and accurately differentiate numeric functions expressed as computer programs.  
34 Generally, besides forward- and reverse-mode AD, the two main AD branches, many software  
35 implementations with different pros and cons exist. Some AD software implementations  
36 work at a lower level code representation, possibly mixing in LLVM-level compiler passes, to

37 fully optimize scalar operations [Revels et al., 2016, Moses and Churavy, 2020] while others  
38 perform transformations at a higher level to keep linear algebra operations intact for optimal  
39 usage of BLAS primitives [Innes et al., 2019, Paszke et al., 2017]. The goal is to make the  
40 best choice of AD system in every part of the program without requiring users to extensively  
41 contort their code to the differing APIs.

42 The AD landscape of the Julia programming language is developed in a manner in which  
43 composability between the AD systems is possible. While many automatic differentiation  
44 systems require specific formulations of the code, for example PyTorch using an alterna-  
45 tive implementation of the NumPy API known as torch.numpy [Paszke et al., 2017] with  
46 torch.tensor and similarly for Jax with jax.numpy [Bradbury et al., 2018] each differing from  
47 the original NumPy [Oliphant, 2006] API in subtle ways with different numerical properties,  
48 the Julia AD systems generally act directly on the standard Julia syntax, with its standard  
49 library, array implementation, its standard GPU acceleration tools [Besard et al., 2018], and  
50 more. This has previously been shown to allow packages in Julia which were developed  
51 without knowledge of AD systems to be fully differentiable without modification by multiple  
52 different tools [Rackauckas et al., 2020a]. Furthermore, Julia has a common ground on  
53 which differentiation rules are defined, ChainRules.jl [White et al., 2021], which is shared  
54 amongst the AD packages. This empowers the idea of a “glue AD” system [Rackauckas]  
55 where software library authors define ChainRules overloads to add domain insight into the  
56 automatic differentiation process without tying to one particular AD system.

57 However, switching from one backend to another on the user side can still be tedious because  
58 the user has to learn and adapt the code towards the user-facing API of the new AD package.  
59 Similarly, for the author of the AD package defining an extensive API supporting every  
60 possible differentiation use case requires a lot of boilerplate code, e.g. to define the Jacobian  
61 function, Jacobian-vector product, Hessian, Hessian-vector product, etc. Defining all of these  
62 functions for each AD implementation is tedious and unnecessary since the relationship  
63 between these functions is abstract and not implementation-specific. Therefore, while in  
64 theory switching between AD systems can be trivially done, in practice the competing APIs  
65 of the various AD mechanisms has limited its use throughout the language’s ecosystem.

66 The Julia Language [Bezanson et al., 2012] has over a dozen automatic differentiation  
67 packages [White]. Different packages have different user interfaces and offer different  
68 tradeoffs. Popular systems include:

- 69 1. ForwardDiff.jl [Revels et al., 2016], an operator-overloading-based, forward-mode  
70 AD implementation, with many years of extensive use and thus very high reliability
- 71 2. ReverseDiff.jl [Revels, 2018], an operator-overloading-based, reverse-mode AD im-  
72 plementation, featuring several tape-based optimizations
- 73 3. Zygote.jl [Innes et al., 2019], a reverse-mode AD implementation that does source  
74 code transformation to generate the derivative’s code from the function’s code,  
75 operating at the level of Julia’s intermediate representation. Zygote is therefore able  
76 to handle arbitrary Julia code but is unable to handle mutation.
- 77 4. Enzyme.jl [Moses and Churavy, 2020], a reverse-mode AD implementation that runs  
78 by source code transformation at the LLVM level, with excellent performance on  
79 scalar operations, but at present lesser performance on large matrix operations.
- 80 5. Diffraction.jl [Fischer], a new source-to-source AD package promising high perfor-  
81 mance on both scalar and vector/tensor code

82 A more detailed summary of the strengths and limitations of different AD packages is given  
83 in Appendix A.

84 Each of these AD systems (and each of the many others) has its own unique set of advantages  
85 and disadvantages. Additionally, all of them only define API functions for a subset of all the  
86 possible differentiation use cases, often requiring users to do package-specific implementations  
87 of quantities like Jacobian-vector product or Hessian-vector product when needed. Beside  
88 the existing stable AD implementations, any new implementation may or may not be mature  
89 enough to handle perturbation confusion properly [Siskind and Pearlmutter, 2005, Manzyuk  
90 et al., 2019] which prevents one from doing general, higher-order AD correctly. A simple

91 workaround is to compose various AD packages for each level of differentiation, further giving  
92 rise to applications where changing between AD mechanisms is increasingly common.

93 As AD systems have different pros and cons, a software author will want to change AD systems  
94 depending on the problem and available hardware resources, see Appendix B. However, this  
95 is more challenging than it might seem. Changing AD systems results in forking the code,  
96 even though the nominal value of the software using the AD remains the same. To give  
97 some examples: Flux.jl<sup>1</sup> changed from using the Tracker.jl<sup>2</sup> to Zygote.jl [Innes et al., 2019].  
98 This resulted in a fork being created, viz. TrackerFlux.jl<sup>3</sup> for those who want to use the old  
99 AD system – even though conceptually Flux is a Neural Network library that should be  
100 abstracted away from the AD. PyMC4 was created as an attempt to move from Theano [Al-  
101 Rfou et al., 2016], as used in PyMC3 [Salvatier et al., 2016], to using TensorFlow [Abadi  
102 et al., 2015]. This attempt was eventually abandoned, in favor of keeping Theano but adding  
103 a Jax [Bradbury et al., 2018] backend [The PyMC Development Team]. Not only did the  
104 code need to be forked, but the overall attempt was not successful. Admittedly, this was a  
105 particularly complex case beyond just AD, with TensorFlow and Theano being more general  
106 computational frameworks with AD as just one feature. The work we present here aims to  
107 ensure that changing the AD system is accessible by providing consistent abstractions that  
108 the author of the  $\partial P$  algorithm implementation can use.

109 A similar but more complex problem was solved by the MathOptInterface.jl [Legat et al.,  
110 2020]. MathOptInterface.jl provides common abstractions across constrained mathematical  
111 optimizers such as IPOPT [Wachter, 2002], Cbc [Forrest et al., 2018], and Gurobi [Gurobi  
112 Optimization, LLC, 2021]. It in turn is used by mathematical optimization frameworks  
113 including JuMP [Dunning et al., 2017] and Convex.jl [Udell et al., 2014]. Each of the  
114 different mathematical optimizers has their own very unique internal set of abstractions, but  
115 MathOptInterface.jl exposes them all in the same way. An additional complication is that  
116 each supports different kinds of problems and so this too must be exposed. Further still, for  
117 some classes of problems they can be re-expressed as a different kind through a mathematical  
118 transformation, MathOptInterface exposes this through an extensible system of so-called  
119 ”bridges”, that will automatically perform these reformulations. This system is considerably  
120 more complicated than our setting as every AD system can perform all the operations, to  
121 varying degrees of efficiency. The MathOptInterface system has proven very successful, which  
122 supports the idea that this kind of abstraction is valuable and can be practically realized.

123 In light of the above, the authors believe it is necessary to have a backend-agnostic interface  
124 to provide objects like the function value, its gradient, Hessian, etc. as well as combining  
125 AD implementations together for higher-order AD. Such an interface can help us avoid a  
126 combinatorial explosion of code when supporting every differentiation package in Julia in  
127 every piece of software requiring gradients and/or Hessians. This is especially important  
128 for higher-order derivatives because one can combine any two differentiation backends to  
129 create a new higher-order backend. More generally for a  $k^{th}$  order derivative, the amount of  
130 code required to support  $n$  differentiation packages in  $m$   $\partial P$  algorithm implementations is  
131  $O(m \times n^k)$ .

132 In this paper, we present AbstractDifferentiation.jl [Anonymous], a package that:

- 133 • Defines an abstract, extensive API for differentiation in Julia enabling the de-  
134 velopment of algorithms requiring first and higher-order derivatives in an AD-  
135 implementation-agnostic way using a single, unified interface reducing the code  
136 complexity from  $O(m \times n^k)$  to  $O(m + n)$ .
- 137 • Automatically defines most of the extensive user-facing API for any new AD package  
138 from just one or two primitive API function definitions, thus making it easier for  
139 the AD package developer to support every possible use case without a great deal of  
140 boilerplate code.

---

<sup>1</sup><https://github.com/FluxML/Flux.jl>

<sup>2</sup><https://github.com/FluxML/Tracker.jl>

<sup>3</sup><https://github.com/ASTupidBear/TrackerFlux.jl>

141 **2 Levels of abstraction in Julia’s AD ecosystem**

142 In Figure 1, an overview of the levels of abstraction in Julia’s AD ecosystem with Abstract-  
143 Differentiation.jl is presented. At the bottom level, we have libraries of differentiation rules  
144 (DiffRules.jl and ChainRules.jl) for specific functions. These rules are either defined by AD  
145 developers for basic Julia constructs, or by AD users for specific user-defined functions with  
146 known analytic derivatives.

147 Sitting on top of the library of rules are all the AD package implementations. At this level,  
148 numerous design decisions and optimizations can be made giving a variety of different AD  
149 package implementations with different tradeoffs. Each AD package developer will then  
150 define a minimal set of primitives and a backend type extending AbstractDifferentiation.jl.  
151 These minimal definitions then enable AbstractDifferentiation.jl to automatically define an  
152 extensive set of user-facing API functions for AD users to use, e.g. derivative, Jacobian,  
153 Hessian, Jacobian-vector product, Hessian-vector product, etc.

154 At the top level, AD users can then use the relevant part of the AbstractDifferentiation.jl API  
155 to implement algorithms requiring  $\partial P$ . With this abstraction design, the amount of code  
156 needed to support all of  $n$  AD packages in  $m$  algorithms requiring  $k^{th}$  order derivatives is  
157 only  $O(m+n)$ , a significant reduction from the  $O(m \times n^k)$  without AbstractDifferentiation.jl.  
158 Additionally, the AD users and developers do not need to add unnecessary boilerplate code  
159 to extend an AD package’s API anymore, since AbstractDifferentiation.jl automatically does  
160 this for them.

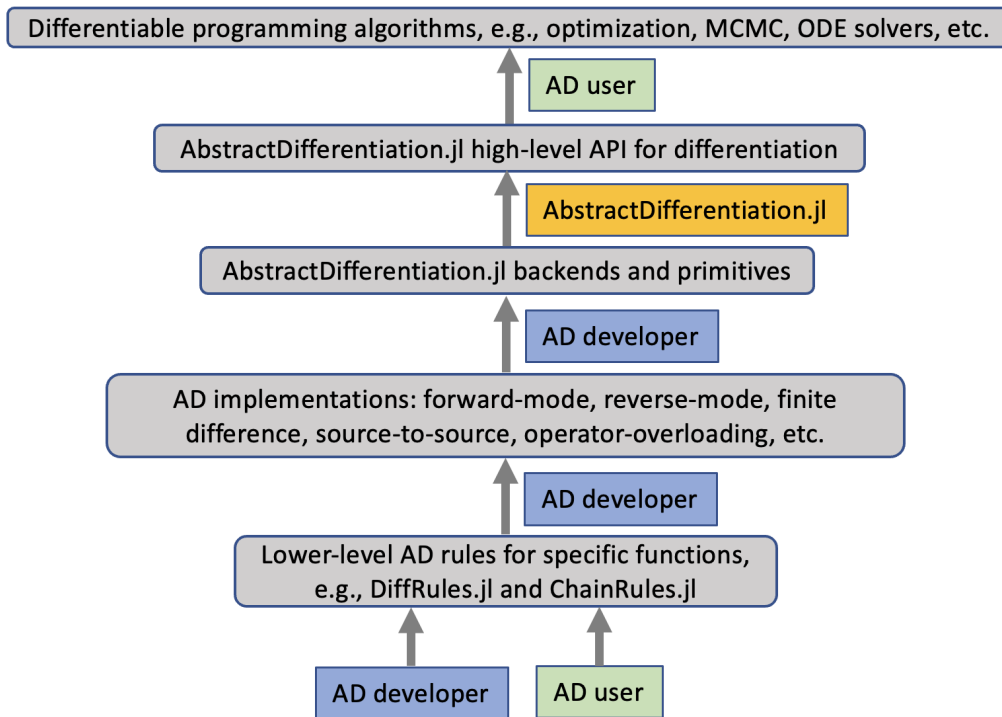


Figure 1: The levels of abstraction in Julia’s AD ecosystem.

161 **3 API description**

162 **Installation and loading** AbstractDifferentiation.jl is a registered Julia package and can  
163 be installed by the Julia package manager. The package can be loaded by

```

164
165 # alternatively: import AbstractDifferentiation as AD
166 using AbstractDifferentiation
167
168

```

170 Note that AbstractDifferentiation.jl exports “AD” as an alias for the AbstractDifferentiation  
171 module. This alias allows us to conveniently access names within AbstractDifferentiation.jl via  
172 AD instead of typing the full package name.

### 173 3.1 Backends and primitives

174 **Forward-mode, reverse-mode, and finite-difference backends** All functionalities  
175 in AbstractDifferentiation.jl are implemented based on an `ab::AbstractBackend` type.  
176 An AD package developer (or the AD user if necessary) first constructs a backend in-  
177 stance that subtypes `ab::AbstractForwardMode`, `ab::AbstractReverseMode`, or `ab::`  
178 `AbstractFiniteDifference`, which are themselves subtypes of `ab::AbstractBackend`.  
179 For example, backends that support ForwardDiff.jl or Zygote.jl are defined as follows:

```

180
181
182 ## ForwardDiff
183 struct ForwardDiffBackend <: AD.AbstractForwardMode end
184 const forwarddiff_backend = ForwardDiffBackend()
185
186 ## Zygote
187 struct ZygoteBackend <: AD.AbstractReverseMode end
188 const zygote_backend = ZygoteBackend()
189
190

```

191 By adding fields to the backend struct, we can control configurations of the differentiation  
192 package such as chunk sizes, compilation flags, or method choices. To use a finite differencing  
193 method at a central grid of 5 points as implemented in the FiniteDifferences.jl package, we  
194 write:

```

195
196 ## FiniteDifferences
197 struct FDMBackend{A} <: AD.AbstractFiniteDifference
198     alg::A
199 end
200 # 1 denotes the order of the derivative to estimate.
201 FDMBackend() = FDMBackend(central_fdm(5, 1))
202
203

```

205 **Higher-order backends** To compute higher-order derivatives, it may be desirable  
206 to combine different backends. We provide `AD.HigherOrderBackend` to implement  
207 higher-order backends. Let `ab_f` be a forward-mode automatic differentiation back-  
208 end and let `ab_r` be a reverse-mode automatic differentiation backend. To construct  
209 a higher-order backend for doing forward-over-reverse-mode automatic differentiation,  
210 one defines `AD.HigherOrderBackend((ab_f, ab_r))`. Analogously, higher-order back-  
211 end for doing reverse-over-forward-mode automatic differentiation is constructed via `AD.`  
212 `HigherOrderBackend((ab_r, ab_f))`.

213 **Jacobian, pushforward, and pullback as primitive operation** In addition to the  
214 definition of a backend, the AD package developer needs to define one of the following  
215 primitive operations:

```

216
217 AD.@primitive function jacobian(ab::backend, f, xs...)
218     return ..
219 end
220 AD.@primitive function pushforward_function(ab::backend, f, xs...)
221     return ..
222 end
223 AD.@primitive function pullback_function(ab::backend, f, xs...)
224     return ..
225 end
226
227

```

229 AbstractDifferentiation.jl then generates the other two primitive functions. For in-  
 230 stance, a source-to-source reverse-mode AD package developer can specify only `AD.`  
 231 `pullback_function` as the native primitive operation.

```

232
233
234 ## Zygote is source-to-source reverse-mode
235 AD.@primitive function pullback_function(ab::ZygoteBackend, f, xs...)
236     return function (vs)
237         # Supports only single output
238         _, back = Zygote.pullback(f, xs...)
239         if vs isa AbstractVector
240             return back(vs)
241         else
242             # vs isa Tuple
243             @assert length(vs) == 1
244             return back(vs[1])
245         end
246     end
247 end
248
  
```

250 In the case of operator overloading AD implementations, we require additionally the definition  
 251 of `AD.primal_value` returning the primal value of the forward pass.

### 252 3.2 Automatically provided functions

253 After these preparatory steps, AbstractDifferentiation.jl automatically defines various func-  
 254 tions for AD users making use of the primitives defined. Some of the most important  
 255 API functions provided are presented in the following. We refer the reader to the package  
 256 documentation for further details [Anonymous].

#### 257 Derivative, gradient, jacobian, hessian

```

258
259 ds = AD.derivative(ab::AD.AbstractBackend, f, xs::Number...)
260 gs = AD.gradient(ab::AD.AbstractBackend, f, xs...)
261 js = AD.jacobian(ab::AD.AbstractBackend, f, xs...)
262 h = AD.hessian(ab::AD.AbstractBackend, f, x)
263
264
  
```

#### 266 Value and derivative, gradient, jacobian, hessian

```

267
268
269 v, ds = AD.value_and_derivative(ab::AD.AbstractBackend, f, xs::Number
270     ...)
271 v, gs = AD.value_and_gradient(ab::AD.AbstractBackend, f, xs...)
272 v, js = AD.value_and_jacobian(ab::AD.AbstractBackend, f, xs...)
273 v, h = AD.value_and_hessian(ab::AD.AbstractBackend, f, x)
274 v, g, h = AD.value_gradient_and_hessian(ab::AD.AbstractBackend, f, x)
275
276
  
```

#### 277 Lazy operators

278 Finally, we implemented lazy versions of the derivative, gradient, Jacobian, and Hessian,

```

279
280
281 ld = lazy_derivative(ab::AbstractBackend, f, xs::Number...)
282 lg = lazy_gradient(ab::AbstractBackend, f, xs...)
283 lj = lazy_jacobian(ab::AbstractBackend, f, xs...)
284 lh = lazy_hessian(ab::AbstractBackend, f, x)
285
  
```

287 which are of interest in iterative solvers. For example, we compute a vector-Jacobian product  
 288 by multiplying a single (transposed) vector, or a tuple of an appropriate shape, from the left  
 289 to the lazy Jacobian operator.

290 **4  $\partial P$  use cases and an example**

291 Many numerical algorithms require the computation of constructs such as the ones described  
 292 in Section 3.2. Table 1 presents a rough summary linking some of the most widely adopted  
 293 routines across different domains to the quantities used in the respective iterative update  
 294 steps. As an example, we present a (simple, non-optimized) backend-agnostic implementation  
 295 of the Gauss-Newton algorithm to solve non-linear least squares problems in Appendix C.

296 We also expect AbstractDifferentiation.jl to be specifically handy for (future) AD package like  
 297 Diffactor.jl or Enzyme.jl where computing constructs like Jacobians or Hessians is technically  
 298 possible but not yet part of the public API due to abstractions or naming conventions made  
 299 in the package.

algorithms	required quantities
<b>root finding</b>	
Newton–Raphson	Jacobian
Jacobian-Free Newton Krylov	Jacobian-vector products
<b>optimization</b>	
ADAM	gradient
Newton	gradient, Hessian
Levenberg–Marquardt	Jacobian
Gauss-Newton	Jacobian
<b>differential equations</b>	
stiff ODE solvers	Jacobian
stiff ODE Jacobian-free solvers	Jacobian-vector products
forward sensitivity methods	Jacobian-vector products
adjoint sensitivity methods	vector-Jacobian product

Table 1: Domain-specific algorithms requiring derivatives, gradients, Jacobians, Hessians, vector-Jacobian products, Jacobian-vector products commonly computed by AD packages.

300 **5 Summary & Future work**

301 The ability to straightforwardly combine different packages in one workflow is one of the  
 302 most versatile and key features of Julia. Switching between different AD packages and  
 303 combining them for higher-order derivatives is a useful feature to have when selecting the  
 304 best AD implementation for a specific application. We have presented the AbstractDifferentiation.jl  
 305 package which makes this switching and combining of AD implementations as  
 306 painless as possible for end-users while also reducing the amount of necessary boilerplate  
 307 code per AD package to support all differentiation use cases.

308 In the future, we aim to support AbstractDifferentiation.jl in all of the AD packages in  
 309 Julia and remove lots of boilerplate code from popular Julia packages (e.g. in the SciML  
 310 and TuringLang organizations) that heavily employ AD.

311 **References**

312 M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis,  
 313 J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia,  
 314 R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore,  
 315 D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker,  
 316 V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke,  
 317 Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems,  
 318 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.

319 R. Al-Rfou, G. Alain, A. Almahairi, C. Angermueller, D. Bahdanau, N. Ballas, F. Bastien,  
 320 J. Bayer, A. Belikov, A. Belopolsky, Y. Bengio, A. Bergeron, J. Bergstra, V. Bisson,

- 321 J. Bleecher Snyder, N. Bouchard, N. Boulanger-Lewandowski, X. Bouthillier, A. de Brébis-  
322 son, O. Breuleux, P.-L. Carrier, K. Cho, J. Chorowski, P. Christiano, T. Cooijmans, M.-A.  
323 Côté, M. Côté, A. Courville, Y. N. Dauphin, O. Delalleau, J. Demouth, G. Desjardins,  
324 S. Dieleman, L. Dinh, M. Ducoffe, V. Dumoulin, S. Ebrahimi Kahou, D. Erhan, Z. Fan,  
325 O. Firat, M. Germain, X. Glorot, I. Goodfellow, M. Graham, C. Gulcehre, P. Hamel, I. Har-  
326 louchet, J.-P. Heng, B. Hidasi, S. Honari, A. Jain, S. Jean, K. Jia, M. Korobov, V. Kulkarni,  
327 A. Lamb, P. Lamblin, E. Larsen, C. Laurent, S. Lee, S. Lefrancois, S. Lemieux, N. Léonard,  
328 Z. Lin, J. A. Livezey, C. Lorenz, J. Lowin, Q. Ma, P.-A. Manzagol, O. Mastropietro,  
329 R. T. McGibbon, R. Memisevic, B. van Merriënboer, V. Michalski, M. Mirza, A. Orlandi,  
330 C. Pal, R. Pascanu, M. Pezeshki, C. Raffel, D. Renshaw, M. Rocklin, A. Romero, M. Roth,  
331 P. Sadowski, J. Salvatier, F. Savard, J. Schlüter, J. Schulman, G. Schwartz, I. V. Serban,  
332 D. Serdyuk, S. Shabanian, E. Simon, S. Spieckermann, S. R. Subramanyam, J. Sygnowski,  
333 J. Tanguay, G. van Tulder, J. Turian, S. Urban, P. Vincent, F. Visin, H. de Vries, D. Warde-  
334 Farley, D. J. Webb, M. Willson, K. Xu, L. Xue, L. Yao, S. Zhang, and Y. Zhang. Theano:  
335 A Python framework for fast computation of mathematical expressions. *arXiv e-prints*,  
336 abs/1605.02688, May 2016. URL <http://arxiv.org/abs/1605.02688>.
- 337 M. AlQuraishi. End-to-end differentiable learning of protein structure. *bioRxiv*, 2018. doi:  
338 10.1101/265231. URL [https://www.biorxiv.org/content/early/2018/02/14/](https://www.biorxiv.org/content/early/2018/02/14/265231)  
339 [265231](https://www.biorxiv.org/content/early/2018/02/14/265231).
- 340 Anonymous. AbstractDifferentiation.jl. URL [https://anonymous.4open.science/](https://anonymous.4open.science/r/AbstractDifferentiation/README.md)  
341 [r/AbstractDifferentiation/README.md](https://anonymous.4open.science/r/AbstractDifferentiation/README.md).
- 342 T. Besard, C. Foket, and B. De Sutter. Effective extensible programming: Unleashing Julia  
343 on GPUs. *IEEE Transactions on Parallel and Distributed Systems*, 2018. ISSN 1045-9219.  
344 doi: 10.1109/TPDS.2018.2872064.
- 345 J. Bezanson, S. Karpinski, V. B. Shah, and A. Edelman. Julia: A fast dynamic language for  
346 technical computing. *arXiv preprint arXiv:1209.5145*, 2012.
- 347 J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula,  
348 A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transfor-  
349 mations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- 350 R. T. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural ordinary differential  
351 equations. *arXiv preprint arXiv:1806.07366*, 2018.
- 352 R. Dandekar, C. Rackauckas, and G. Barbastathis. A machine learning-aided global diagnostic  
353 and comparative tool to assess effect of quarantine control in covid-19 spread. *Patterns*, 1  
354 (9):100145, 2020.
- 355 F. de Avila Belbute-Peres, K. Smith, K. Allen, J. Tenenbaum, and J. Z. Kolter. End-to-end  
356 differentiable physics for learning and control. *Advances in neural information processing*  
357 *systems*, 31:7178–7189, 2018.
- 358 I. Dunning, J. Huchette, and M. Lubin. Jump: A modeling language for mathematical  
359 optimization. *SIAM Review*, 59(2):295–320, 2017.
- 360 K. Fischer. Non-local compiler transformations in the presence of dynamic dispatch. URL  
361 <https://www.youtube.com/watch?v=mQnSRfseu0c>.
- 362 J. Forrest, T. Ralphs, S. Vigerske, LouHafer, B. Kristjansson, jpfasano, EdwinStraver,  
363 M. Lubin, H. G. Santos, rlougee, and M. Saltzman. coin-or/cbc: Version 2.9.9, July 2018.  
364 URL <https://doi.org/10.5281/zenodo.1317566>.
- 365 Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2021. URL [https:](https://www.gurobi.com)  
366 [//www.gurobi.com](https://www.gurobi.com).
- 367 J. Ingraham, A. Riesselman, C. Sander, and D. Marks. Learning protein structure with a  
368 differentiable simulator. In *International Conference on Learning Representations*, 2018.

- 369 M. Innes, A. Edelman, K. Fischer, C. Rackauckus, E. Saba, V. B. Shah, and W. Tebbutt.  
370 Zygote: A differentiable programming system to bridge machine learning and scientific  
371 computing. arXiv preprint arXiv:1907.07587, 2019. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1907.07587)  
372 **1907.07587**.
- 373 B. Legat, O. Dowson, J. D. Garcia, and M. Lubin. Mathoptinterface: a data structure  
374 for mathematical optimization problems, 2020. URL [https://arxiv.org/abs/2002.](https://arxiv.org/abs/2002.03447)  
375 **03447**.
- 376 O. Manzyuk, B. A. Pearlmutter, A. A. Radul, D. R. Rush, and J. M. Siskind. Perturbation  
377 confusion in forward automatic differentiation of higher-order functions. *Journal of*  
378 *Functional Programming*, 29, 2019.
- 379 W. Moses and V. Churavy. Instead of rewriting foreign code for machine learning, automati-  
380 cally synthesize fast gradients. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan,  
381 and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages  
382 12472–12485. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2020/file/9332c513ef44b682e9347822c2e457ac-Paper.pdf)  
383 **cc/paper/2020/file/9332c513ef44b682e9347822c2e457ac-Paper.pdf**.
- 384 T. E. Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- 385 A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison,  
386 L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.
- 387 C. Rackauckas. Glue AD for full language differentiable pro-  
388 gramming. URL [http://www.stochasticlifestyle.com/](http://www.stochasticlifestyle.com/glue-ad-for-full-language-differentiable-programming/)  
389 **glue-ad-for-full-language-differentiable-programming/**.
- 390 C. Rackauckas, Y. Ma, V. Dixit, X. Guo, M. Innes, J. Revels, J. Nyberg, and V. Ivaturi. A  
391 comparison of automatic differentiation and continuous sensitivity analysis for derivatives  
392 of differential equation solutions. *arXiv preprint arXiv:1812.01892*, 2018.
- 393 C. Rackauckas, A. Edelman, K. Fischer, M. Innes, E. Saba, V. B. Shah, and W. Tebbutt.  
394 Generalized physics-informed learning through language-wide differentiable programming.  
395 In *AAAI Spring Symposium: MLPS*, 2020a.
- 396 C. Rackauckas, Y. Ma, J. Martensen, C. Warner, K. Zubov, R. Supekar, D. Skinner, and  
397 A. Ramadhan. Universal differential equations for scientific machine learning. *arXiv*  
398 *preprint arXiv:2001.04385*, 2020b.
- 399 M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep  
400 learning framework for solving forward and inverse problems involving nonlinear partial  
401 differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- 402 J. Revels. ReverseDiff.jl, 2018. URL [http://github.com/JuliaDiff/ReverseDiff.](http://github.com/JuliaDiff/ReverseDiff.jl)  
403 **jl**.
- 404 J. Revels, M. Lubin, and T. Papamarkou. Forward-mode automatic differentiation in julia.  
405 *arXiv preprint arXiv:1607.07892*, 2016.
- 406 J. Salvatier, T. V. Wiecki, and C. Fonnesbeck. Probabilistic programming in python using  
407 PyMC3. *PeerJ Computer Science*, 2:e55, apr 2016. doi: 10.7717/peerj-cs.55. URL  
408 <https://doi.org/10.7717/peerj-cs.55>.
- 409 F. Schäfer. Abstractdifferentiation.jl for AD-backend agnostic code. URL [https://](https://frankschae.github.io/post/abstract_differentiation/)  
410 **frankschae.github.io/post/abstract\_differentiation/**.
- 411 F. Schäfer, M. Kloc, C. Bruder, and N. Lörch. A differentiable programming method for  
412 quantum control. *Machine Learning: Science and Technology*, 1(3):035009, 2020.
- 413 F. Schäfer, P. Sekatski, M. Koppenhöfer, C. Bruder, and M. Kloc. Control of stochastic  
414 quantum dynamics by differentiable programming. *Machine Learning: Science and*  
415 *Technology*, 2(3):035004, 2021.

- 416 C. Schenck and D. Fox. Spnets: Differentiable fluid dynamics for deep neural networks. In  
417 *Conference on Robot Learning*, pages 317–335. PMLR, 2018.
- 418 J. M. Siskind and B. A. Pearlmutter. Perturbation confusion and referential transparency:  
419 Correct functional implementation of forward-mode AD. 2005.
- 420 F. Srajer, Z. Kukulova, and A. Fitzgibbon. A benchmark of selected algorithmic differentiation  
421 tools on some problems in computer vision and machine learning. *Optimization Methods*  
422 *and Software*, 33(4-6):889–906, 2018.
- 423 The PyMC Development Team. The future of pymc3, or: Theano is  
424 dead, long live theano. URL [https://pymc-devs.medium.com/  
425 the-future-of-pymc3-or-theano-is-dead-long-live-theano-d8005f8a0e9b](https://pymc-devs.medium.com/the-future-of-pymc3-or-theano-is-dead-long-live-theano-d8005f8a0e9b).
- 426 M. Udell, K. Mohan, D. Zeng, J. Hong, S. Diamond, and S. Boyd. Convex optimization in  
427 Julia. *SC14 Workshop on High Performance Technical Computing in Dynamic Languages*,  
428 2014.
- 429 A. Wachter. *An interior point algorithm for large-scale nonlinear optimization with applica-*  
430 *tions in process engineering*. PhD thesis, Carnegie Mellon University, 2002.
- 431 L. White. Juliadiff website. URL <https://juliadiff.org/>.
- 432 L. White, M. Zgubic, M. Abbott, J. Revels, A. Arslan, S. Axen, S. Schaub, N. Robinson,  
433 Y. Ma, G. Dhingra, W. Tebbutt, N. Heim, A. D. W. Rosemberg, N. Schmitz, C. Rack-  
434 auckas, D. Widmann, R. Heintzmann, F. Schäfer, K. Fischer, A. Robson, M. Brzezinski,  
435 A. Zhabinski, M. Besançon, P. Vertechi, S. Gowda, A. Fitzgibbon, C. Lucibello,  
436 C. Vogt, D. Gandhi, and F. Chorney. Juliadiff/chainrules.jl: v1.11.5, Sept. 2021. URL  
437 <https://doi.org/10.5281/zenodo.5467874>.

438 **A Summary of the current state of AD packages in Julia as of**  
 439 **September 2021**

Table 2: This table summarizes the current state of popular Julia AD packages in September 2021. “Scalar” refers to scalar operations support including defining custom rules for scalar-valued functions of scalars. “Vector/tensor” refers to native vector/tensor support as a construct including the ability to define custom differentiation rules for vector/tensor-valued functions and/or functions of vectors/tensors. Similarly, “First class struct support” refers to the native support of Julia structs as a construct including the ability to define custom differentiation rules for struct-valued functions or functions of structs. “GPU” refers to the ability to differentiate functions of or returning GPU arrays. “GC” refers to supporting functions that call the Julia garbage collector. “Mutation” refers to the ability to support mutating arrays and structs. “Runtime branches” refers to the ability to support “piece-wise” functions with control flow such that the path that the function takes and ultimately the structure of the mathematical function differentiated depends on the values of the inputs to the function. “Maturity” refers to a subjective measure of how mature each package is in the eyes of the community as well as the feature maturity of the package.

Package	Scalar	Vector / tensor	First class struct support	GPU	GC	Mutation	Runtime branches	Maturity
ForwardDiff	✓	✗	✗	✓	✓	✓	✓	very high
ReverseDiff	slow	✓	✗	✗	✓	limited	✓	high
ReverseDiff with compiled tape	✓	✓	✗	✗	✓	limited	✗	high
Tracker	slow	✓	✗	✓	✓	limited	✓	high
Zygote	slow	✓	✓	✓	✓	✗	✓	high
Enzyme	✓	limited	wip	wip	wip	✓	✓	low
Difffractor	✓	✓	✓	✓	✓	✗	✓	low

440 Table 2 summarizes the current state of the most popular AD packages in the Julia ecosystem  
 441 as of the time of the writing of this paper.

442 **B AD performance can be problem-specific**

443 It is well known that for a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with  $n$  independent input variables and  $m$   
 444 dependent output variables, forward-mode AD is preferred to build the Jacobian when  $m \gg n$   
 445 while reverse-mode AD is preferred when  $n \gg m$ , i.e. as one increases the number of inputs  
 446 within the same problem, reverse-mode AD mode will eventually overtake forward-mode  
 447 AD in performance. This has been investigated in depth for differential equations when  
 448 applied to models relevant to biopharmacology, alongside various adjoint options [Rackauckas  
 449 et al., 2018]. This work shows that on sufficiently small ODEs (<100 ODEs + parameters),  
 450 forward-mode AD via ForwardDiff.jl is the most efficient method comparing against analytical  
 451 techniques and adjoint techniques using Tracker.jl, Enzyme.jl, and ReverseDiff.jl. When the  
 452 size of the ODEs+parameters is increased in a stiff partial differential equation, it was shown  
 453 that Enzyme.jl vector-Jacobian products mixed with a specific adjoint method was the most  
 454 efficient, outperforming the ForwardDiff.jl techniques long with ReverseDiff.jl and Tracker.jl.

455 In what follows, we demonstrate on two additional examples that the choice of the specific  
 456 reverse-mode AD package may also significantly impact the performance [Srajer et al., 2018].  
 457 These examples show ReverseDiff.jl in compiled mode outperforming Enzyme.jl under certain  
 458 circumstances. However, as ReverseDiff.jl is not compatible with GPUs and was shown to not  
 459 be performance competitive on other potential equations in scientific computing applications,  
 460 which allows Zygote.jl and Tracker.jl to be the most efficient. Together this shows in one  
 461 application that 5 AD systems could potentially be the optimal choice depending on user  
 462 inputs into the package code. This establishes that the optimal choice of AD mechanism  
 463 can be rather complex for users and package developers, and thus decreasing the cost of  
 464 performing such benchmarks is of value to many scientists.

465 **Example 1: Lotka–Volterra model**

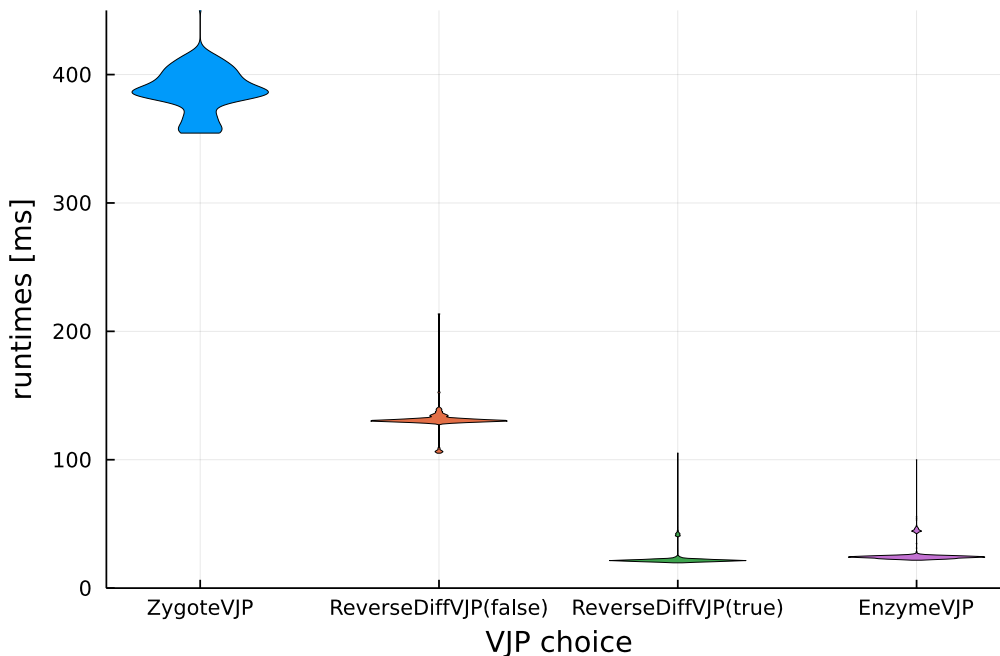


Figure 2: Benchmark 1: Lotka–Volterra model. In all cases, we use a checkpointed interpolating adjoint method [Rackauckas et al., 2020b] to compute the local sensitivities. ‘false’ and ‘true’ indicate if the tape in ReverseDiff.jl is precompiled.

466 Suppose that we have an instantaneous objective

$$l(x(t), y(t)) = x(t) + y(t) \tag{1}$$

467 for the Lotka–Volterra model

$$\dot{x} = \alpha x - \beta xy, \tag{2}$$

$$\dot{y} = \gamma xy - \delta y, \tag{3}$$

468 with initial conditions  $x(t = 0) = 1$  and  $y(t = 0) = 1$ . Let  $\xi$  denote any of the parameters  
469  $\alpha, \beta, \gamma, \delta$ . We are interested in the sensitivities  $\frac{\partial}{\partial \xi} \sum_i l(x(t_i), y(t_i))$  with respect to an equally  
470 spaced time grid between 0 and 10 with a grid spacing of 0.1.

471 Figure 2 shows a violin plot for the runtimes for four choices of the internally used AD  
472 system. This demonstrates that the vector-Jacobian products which use static compilation of  
473 the ODE function, ReverseDiff.jl with compilation enabled and Enzyme.jl, vastly outperform  
474 the other choices for small ODEs with a lot of scalar indexing, which is a common feature  
475 in many nonlinear physical and biochemical models. Note that all adjoint techniques were  
476 shown to be outperformed by ForwardDiff.jl on this example elsewhere [Rackauckas et al.,  
477 2018], but this example still confirms that in many scalar indexing cases the Zygote.jl system  
478 can perform rather poorly.

### 479 **Example 2: Neural ODE**

480 This example is the Spiral Neural ODE chosen from the Neural Ordinary Differential  
481 Equations manuscript [Chen et al., 2018]. It is an ODE defined as a neural network applied  
482 to the cubed states of the system:

$$\dot{u} = \text{NN}(u^3) \tag{4}$$

483 where  $\text{NN}(u)$  is a multilayer perceptron with one hidden layer of size 50 and a tanh activation  
484 function, and  $u \in \mathbb{R}^2$ . Figure 3 shows a violin plot for the runtimes for four choices of  
485 the internally used AD system. The results show that for direct differentiation on CPUs,  
486 ReverseDiffVJP with a compiled tape is the most efficient method. However, this has many  
487 caveats. One caveat is that ReverseDiff.jl’s tape-compiled form is only applicable if the code  
488 has no branching, and thus would be incompatible with activation functions like relu.

489 Additionally, by testing over various sizes of hidden layers, we established that a RTX  
490 2080 Super GPU outperformed a Ryzen 9 5950x CPU when the hidden layer size reached  
491 approximately 7,500 (note the crossover point could potentially be a lot smaller in many  
492 scenarios if the neural network is deeper since the first and last layer sizes are 2, matching  
493 the dimensionality of  $u$ ). At around this size of neural networks, the Zygote.jl and Tracker.jl  
494 strategies on GPUs become more efficient than the one of ReverseDiff.jl which is restricted  
495 to CPUs.

496 These two examples, in addition to the prior research, clearly demonstrate that the internal  
497 AD system must be carefully chosen based on the problem (and hardware resources) at hand.

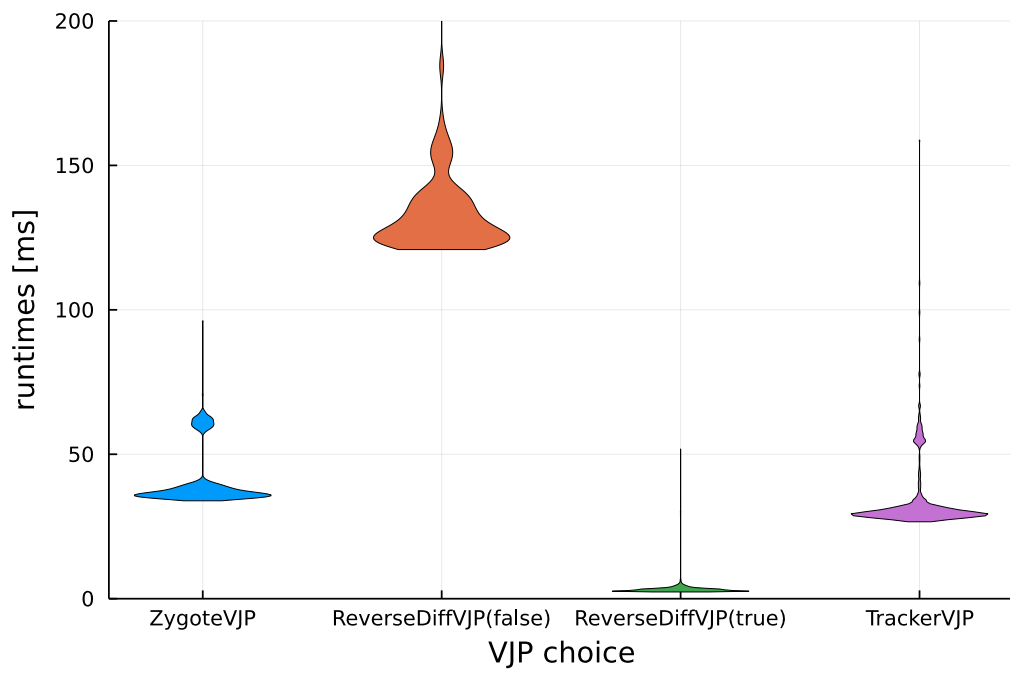


Figure 3: Benchmark 2: Spiral Neural ODE model. In all cases, we use a checkpointed interpolating adjoint method [Rackauckas et al., 2020b] to compute the local sensitivities. ‘false’ and ‘true’ indicate if the tape in ReverseDiff.jl is precompiled.

## 498 C Implementation of the Gauss-Newton algorithm

499 In this appendix, we use AbstractDifferentiation.jl for the implementation of the Gauss-Newton  
 500 algorithm for solving nonlinear least squares problems [Schäfer]. The Gauss-Newton  
 501 algorithm iteratively finds the value of the  $N$  variables  $\mathbf{x} = (x_1, \dots, x_N)$  minimizing the sum  
 502 of squares of  $M$  residuals  $(f_1, \dots, f_M)$

$$S(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^M f_i(\mathbf{x}; \mathbf{p})^2. \quad (5)$$

503 Starting from an initial guess  $\mathbf{x}_0$  for the minimum, the method runs through the iterations

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k (J^T J)^{-1} J^T f(\mathbf{x}^k; p), \quad (6)$$

504 where the residuals  $f(\mathbf{x}^k; p)$  depend on the current step  $\mathbf{x}^k$  and parameters  $p$ .  $J$  is the  
 505 Jacobian matrix at  $\mathbf{x}^k$ , and  $\alpha_k$  is the step length determined via a line search subroutine.

```

506
507
508 ## Gauss-Newton scheme
509 function GaussNewton!(xs, x, p backend; maxiter=100)
510     for i=1:maxiter
511         x = step(x, p, backend)
512         push!(xs, x)
513     end
514     return xs, x
515 end
516 function step(x, p, backend, α=1//1)
517     x2 = deepcopy(x)
518     while !done(x, x2, p) # line-search condition
519         # first return value of AD.jacobian is ∂f∂x
520         # second return value of AD.jacobian is ∂f∂p
521         J = AD.jacobian(backend, f, x, p)[1]
522         d = -inv(J'*J)*J'*f(x, p)
523         copyto!(x2, x + α*d)
524         α = α//2
525     end
526     return x2
527 end
528

```

530 Switching between different AD systems is then easily accomplished by passing different  
 531 backends (`zygote_backend`, `forwarddiff_backend`, etc.) as input to the `GaussNewton!`  
 532 function.