

1 **THE F-1 ALGORITHM FOR EFFICIENT COMPUTATION OF THE**
2 **HESSIAN MATRIX OF AN OBJECTIVE FUNCTION DEFINED**
3 **IMPLICITLY BY THE SOLUTION OF A STEADY-STATE PROBLEM***

4 BENOÎT PASQUIER[†] AND FRANÇOIS PRIMEAU[†]

5 **Abstract.** Steady-state systems of nonlinear partial differential equations (PDEs) are common
6 in engineering and the biogeosciences. These systems are typically controlled by parameters that can
7 be estimated efficiently using second-order optimization algorithms. However, computing the gradi-
8 ent vector and Hessian matrix of a given objective function defined implicitly by the solution of large
9 PDE systems is seldom economical. Here we present a fast and easy-to-use algorithm for computing
10 the gradient and Hessian of an objective function implicitly constrained by a steady-state PDE sys-
11 tem. We call the new algorithm, which is based on the use of hyperdual numbers, the F-1 algorithm,
12 because it requires only one factorization of the constraint-equation Jacobian. Careful examination
13 of the relationships that arise from differentiating the PDE system reveal analytical shortcuts that
14 the F-1 algorithm leverages. We benchmark the F-1 algorithm against five numerical differentiation
15 schemes in the context of optimizing a global steady-state model of the marine phosphorus cycle that
16 depends explicitly on $m = 6$ parameters. In this context, the F-1 algorithm computes the Hessian 16
17 to 100 times faster than other algorithms, allowing for the entire optimization procedure to be per-
18 formed 4 to 26 times faster. This is because other algorithms require $\mathcal{O}(m)$ to $\mathcal{O}(m^2)$ factorizations,
19 which suggests even greater speedups for larger problems. To facilitate reproducibility and future
20 benchmarks, all the code developed for this study was implemented as open-source Julia packages.

21 **Key word.** Global, Ocean, Marine biogeochemistry, Biogeochemistry, Model, Modeling, In-
22 verse, Inverse model, Optimization, Uncertainty, Sensitivity, Differential equations, Partial differen-
23 tial equations, PDE, steady-state, Newton's method, Factorization, Automatic differentiation, Al-
24 gorithmic differentiation, Autodiff, AD, Finite differences, Complex-step, Dual numbers, Hyperdual
25 numbers, Nutrients, Nutrient cycles, Biogeochemical cycles, gradient, Hessian, Jacobian, Julia, Open
26 source

27 **AMS subject classifications.** 15Axx, 15A06, 15A09, 15A23, 15A24, 15A29, 15A66, 15A69,
28 15A99, 26B10, 26A24, 13P99, 90C53, 90C90

29 **1. Introduction.** The geosciences are rich with problems involving spatial data
30 that can be modeled using partial differential equations (PDEs). In cases where
31 steady-state or time-mean fields are of specific interest, such problems can be ex-
32 pressed generically as

33 (1.1)
$$\mathbf{F}(\mathbf{x}, \mathbf{p}) = 0,$$

34 where \mathbf{x} is the model state vector comprising one or more discretized field variables
35 and \mathbf{p} is a vector of adjustable parameters (see, e.g., [22, 50, 11, 24, 25, 44]).

36 A major modeling goal is then to find the value of \mathbf{x} and \mathbf{p} that are in the best
37 possible agreement with available observational data while satisfying (1.1). Mathe-
38 matically, this translates into the generic constrained optimization problem

39 (1.2)
$$\begin{cases} \underset{\mathbf{x}, \mathbf{p}}{\text{minimize}} & f(\mathbf{x}, \mathbf{p}) \\ \text{subject to} & \mathbf{F}(\mathbf{x}, \mathbf{p}) = 0, \end{cases}$$

40 where $f(\mathbf{x}, \mathbf{p})$ is some measure of how far the state and parameter vectors are from
41 the data and/or some assumed prior values. Here, we restrict ourselves to the case

*Submitted to the the SIAM Journal of Scientific Computing for review on 2019-07-09.

Funding: This work was funded by the US Department of Energy grant DE-SC0016539 and the National Science Foundation grant 1658380.

[†]Department of Earth System Science, University of California, Irvine, CA, United States (pasquieb@uci.edu, fprimeau@uci.edu).

42 where the solution to (1.1) defines \mathbf{x} as an implicit function of \mathbf{p} , which we denote by
 43 $\mathbf{s}(\mathbf{p})$, the steady-state solution. The problem defined by (1.2) is then equivalent to
 44 finding the minimum of the objective function defined by

$$45 \quad (1.3) \quad \hat{f}(\mathbf{p}) \equiv f(\mathbf{s}(\mathbf{p}), \mathbf{p}).$$

46 In a Bayesian formulation of the parameter estimation problem, \hat{f} would cor-
 47 respond to the negative logarithm of the posterior probability distribution. Solving
 48 $\hat{\mathbf{p}} = \arg \min_{\mathbf{p}} \hat{f}(\mathbf{p})$ is then equivalent to finding the most probable parameter val-
 49 ues. Efficient algorithms for minimizing $\hat{f}(\mathbf{p})$ in multidimensional parameter spaces
 50 make use of the gradient, $\nabla \hat{f}(\mathbf{p})$, and Hessian, $\nabla^2 \hat{f}(\mathbf{p})$, to select the most promising
 51 search directions. Furthermore, in parameter estimation problems, the Hessian ma-
 52 trix, $\nabla^2 \hat{f}(\mathbf{p})$, is of direct interest because its inverse evaluated at $\hat{\mathbf{p}}$ can be used to
 53 construct a useful approximation to the error covariance matrix for the parameters,
 54 which provides a useful summary of the parameter uncertainties (e.g., [51, 52, 55]).

55 The present study focuses on PDE problems with discretization schemes that
 56 lead to a sparse Jacobian matrix, $\nabla_{\mathbf{x}} \mathbf{F}(\mathbf{x}, \mathbf{p})$, that can be factored and stored in
 57 computer memory. For such problems Newton’s method can be used to efficiently
 58 solve (1.1). Here, we show how to take advantage of this fact together with the
 59 application of recently developed hyperdual numbers (e.g., [12, 14, 13]) to simplify
 60 and greatly reduce the computational cost of evaluating $\nabla \hat{f}(\mathbf{p})$ and $\nabla^2 \hat{f}(\mathbf{p})$.

61 The typical procedure for minimizing the objective function defined by (1.3) in-
 62 volves two nested iterative processes, as illustrated in Figure 1. The inner solver
 63 finds the steady-state solution, $\mathbf{s}(\mathbf{p})$, by iteratively updating the state, \mathbf{x} , until the
 64 norm of the state function, $\mathbf{F}(\mathbf{x}, \mathbf{p})$, is sufficiently small. This is indicated by the
 65 “ $\mathbf{F}(\mathbf{x}, \mathbf{p}) \approx 0$?” condition, which determines when the inner-solver loop terminates.
 66 On the outside, the optimizer iteratively searches for a minimum of \hat{f} . The opti-
 67 mizer loop updates the parameters, \mathbf{p} , and terminates when the norm of the gradi-
 68 ent of the objective function, $\nabla \hat{f}(\mathbf{p})$, is sufficiently small, which is indicated by the
 69 “ $\nabla \hat{f}(\mathbf{p}) \approx 0$?” condition. The outer optimizer problem, like the inner solver problem,
 70 can be solved using Newton’s method provided the search direction,

$$71 \quad (1.4) \quad \Delta \mathbf{p} \equiv - \left[\nabla^2 \hat{f}(\mathbf{p}) \right]^{-1} \nabla \hat{f}(\mathbf{p}),$$

72 can be computed.

73 However, computing the derivatives required to evaluate $\nabla \hat{f}(\mathbf{p})$ and $\nabla^2 \hat{f}(\mathbf{p})$ an-
 74 alytically is laborious, prone to errors, and potentially computationally expensive —
 75 see (2.5) for example, which involves five large third-order tensors. The evaluation of
 76 the gradient vector and Hessian matrix is therefore typically done using finite differ-
 77 ences applied directly to \hat{f} . However, finite-difference approximations for computing
 78 the Hessian matrix is also computationally expensive when m is moderately large and
 79 suffers from both round-off and truncation errors [29], which can have detrimental
 80 effects on the convergence rate of the optimizer.

81 A recently developed alternative to finite differences, which does not suffer from
 82 round-off or truncation errors is the application of dual numbers to efficiently compute
 83 numerical derivatives (see, e.g., [37]). Dual numbers, like complex numbers, extend
 84 the real numbers by introducing a new unit, denoted ε , but with $\varepsilon^2 \equiv 0$ rather than
 85 $i^2 = -1$ as is the case for the imaginary unit. A more detailed description of dual
 86 numbers is given in subsection SM4.2 and references therein. Implementations of dual

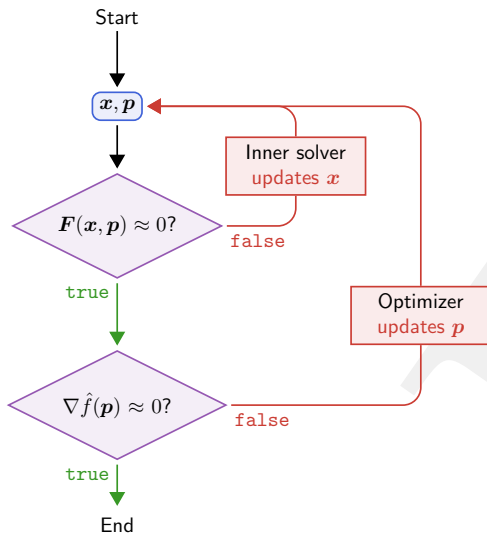


FIG. 1. Schematic diagram representing the optimization procedure. Starting at the top with an initial choice for the state, \mathbf{x} , and parameters, \mathbf{p} , the procedure goes through two nested iterative loops. The outer loop is the optimizer, which iterates on \mathbf{p} until it lies at the minimum of the objective function, \hat{f} (i.e., until $\nabla \hat{f}(\mathbf{p}) \approx 0$). Nested inside the optimizer is the inner solver, which is executed for each update of \mathbf{p} . The inner-solver loop updates \mathbf{x} until it approximately satisfies the steady-state condition, $\mathbf{F}(\mathbf{x}, \mathbf{p}) \approx 0$, with is equivalent to $\mathbf{x} = \mathbf{s}(\mathbf{p})$. We note that the conditional statements to terminate the loops are approximations because of finite precision.

87 numbers for efficiently and accurately computing derivatives are available for several
 88 scientific computing languages (see, e.g., [12, 14, 13, 56, 48, 37]).

89 To see how to compute the gradient using dual numbers let \mathbf{e}_j be the j th vector
 90 of the natural basis of \mathbb{R}^m , i.e., a vector of m zeros except for a 1 in the j th entry.
 91 ($j = 1, \dots, m$ indexes the m dimensions of the parameter space.) Then, the Taylor
 92 expansion of the objective function at \mathbf{p} in the $\varepsilon \mathbf{e}_j$ direction is given by

$$93 \quad (1.5) \quad \hat{f}(\mathbf{p} + \varepsilon \mathbf{e}_j) = \hat{f}(\mathbf{p}) + \varepsilon \nabla \hat{f}(\mathbf{p}) \mathbf{e}_j,$$

94 where we express the gradient, $\nabla \hat{f}(\mathbf{p})$, as a row vector, so that the product $\nabla \hat{f}(\mathbf{p}) \mathbf{e}_j$
 95 yields its scalar j th entry. Thus, each entry of the gradient can be computed
 96 evaluating the objective function with the dual-valued parameters $\mathbf{p} + \varepsilon \mathbf{e}_j$ and taking
 97 the dual part of the result. Rearranging each entry into a row vector gives a formula
 98 to compute the gradient in m dual-valued evaluations of the objective function,

$$99 \quad (1.6) \quad \nabla \hat{f}(\mathbf{p}) = \mathfrak{D} \left(\begin{bmatrix} \hat{f}(\mathbf{p} + \varepsilon \mathbf{e}_1) \\ \hat{f}(\mathbf{p} + \varepsilon \mathbf{e}_2) \\ \vdots \\ \hat{f}(\mathbf{p} + \varepsilon \mathbf{e}_m) \end{bmatrix}^\top \right).$$

100 Note that the dual-step algorithm cannot be naively applied recursively to com-
 101 pute second-order derivatives because $\varepsilon^2 = 0$ ensures that terms of order two (and
 102 above) vanish in the Taylor expansion. To compute second-order derivatives, Fike
 103 and Alonso [12] have developed hyperdual numbers. Two distinct hyperdual units,

104 ε_1 and ε_2 , are introduced, such that $\varepsilon_1^2 = \varepsilon_2^2 = 0$ but such that $\varepsilon_1\varepsilon_2 \neq 0$. Just
 105 like the dual unit, the hyperdual units play the role of infinitesimally small num-
 106 bers. However, because they are independent and do not cancel each other out, they
 107 can propagate infinitesimal perturbations in two directions simultaneously. For more
 108 details on hyperdual numbers, see [subsection SM4.3](#) and references therein.

109 By definition, hyperdual-valued Taylor expansions only extend to second order
 110 terms. In particular, for any pair $(\mathbf{e}_j, \mathbf{e}_k)$ of directions in parameter space (with j
 111 and k spanning the dimensions 1 to m of the parameter space), we have that

$$112 \quad (1.7) \quad \hat{f}(\mathbf{p} + \varepsilon_1\mathbf{e}_j + \varepsilon_2\mathbf{e}_k) = \hat{f}(\mathbf{p}) + \varepsilon_1\nabla\hat{f}(\mathbf{p})\mathbf{e}_j + \varepsilon_2\nabla\hat{f}(\mathbf{p})\mathbf{e}_k + \varepsilon_1\varepsilon_2\mathbf{e}_j^\top\nabla^2\hat{f}(\mathbf{p})\mathbf{e}_k,$$

113 where the product $\mathbf{e}_j^\top\nabla^2\hat{f}(\mathbf{p})\mathbf{e}_k$ yields the entry in the j th row and the k th column
 114 of the Hessian matrix and where the product $\nabla\hat{f}(\mathbf{p})\mathbf{e}_j$ yields the j th entry of the
 115 gradient. Thus, one can compute the Hessian matrix with $m(m+1)/2$ hyperdual-
 116 valued evaluations of the objective function. Specifically, denoting the hyperdual
 117 parameters by $\mathbf{p}_{jk} \equiv \mathbf{p} + \varepsilon_1\mathbf{e}_j + \varepsilon_2\mathbf{e}_k$, the Hessian is given by

$$118 \quad (1.8) \quad \nabla^2\hat{f}(\mathbf{p}) = \mathfrak{H} \left(\begin{bmatrix} \hat{f}(\mathbf{p}_{11}) & \hat{f}(\mathbf{p}_{12}) & \cdots & \hat{f}(\mathbf{p}_{1m}) \\ \hat{f}(\mathbf{p}_{12}) & \hat{f}(\mathbf{p}_{22}) & \cdots & \hat{f}(\mathbf{p}_{2m}) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{f}(\mathbf{p}_{1m}) & \hat{f}(\mathbf{p}_{2m}) & \cdots & \hat{f}(\mathbf{p}_{mm}) \end{bmatrix} \right),$$

119 where $\mathfrak{H}(x)$ is the $\varepsilon_1\varepsilon_2$ coefficient of x .

120 Although they provide an attractive alternative to the fully analytical approach,
 121 the numerical algorithms listed above come at a price. Indeed, in practice, these nu-
 122 merical methods suffer large computational costs on top of potential implementation
 123 pitfalls. Computing the gradient, $\nabla\hat{f}(\mathbf{p})$, via (1.6) seems straightforward and compu-
 124 tationally efficient at face value but we note that each evaluation of \hat{f} will generate
 125 a call to the inner solver. That is, each call will need to find the dual-valued steady-
 126 state solution, $\mathbf{s}(\mathbf{p} + \varepsilon\mathbf{e}_j)$, thus forcing the inner solver, which uses Newton's method,
 127 to perform at least one computationally-expensive factorization of each dual-valued
 128 matrix $\nabla_{\mathbf{x}}\mathbf{F}(\mathbf{x}, \mathbf{p} + \varepsilon\mathbf{e}_j)$ and potentially multiple such factorizations. Similarly, com-
 129 puting the Hessian, $\nabla^2\hat{f}(\mathbf{p})$, via (1.8) will generate at a minimum $m(m+1)/2$ calls
 130 to the inner solver to find the hyperdual-valued steady-state solution, $\mathbf{s}(\mathbf{p}_{jk})$, thus
 131 forcing the inner solver to perform an even larger number of expensive factorizations
 132 of hyperdual-valued $\nabla_{\mathbf{x}}\mathbf{F}(\mathbf{x}, \mathbf{p}_{jk})$ matrices. The additional calls to the inner solver
 133 also expose the user to potential implementation pitfalls if for some reason the solver
 134 does not handle non-real numbers properly. This would be the case for example if the
 135 inner solver invoked operations with non-real numbers internally in the first place,
 136 or did not check for convergence of non-real parts. (These pitfalls are discussed in
 137 [section 6](#).)

138 Optimization problems defined generically in the form of (1.2) are common in
 139 physical sciences and engineering and practical solutions have been investigated and
 140 documented. For example, for aerospace engineering, Rumpfkeil and Mavriplis [49]
 141 suggested an efficient solution to a similar optimization problem to improve airfoil
 142 aerodynamism. They showed that taking the adjoint of the derivatives of their steady-
 143 state problem combined with algorithmic differentiation could lead to an optimally-
 144 efficient algorithm for computing the Hessian matrix. Here, in a similar approach,
 145 we show that a careful refactoring of the algorithm for computing $\nabla\hat{f}$ and $\nabla^2\hat{f}$ using

146 an adjoint formulation and hyperdual numbers can avoid all the calls to the inner
147 solver. This leads to an algorithm, which we call F-1, for computing the gradient and
148 Hessian that is simultaneously easy-to-use, fast, and accurate. The name, F-1, of the
149 new algorithm relates to the fact that it is fast and to the fact that it requires only 1
150 factorization of the large Jacobian matrix for the PDE constraint equation

151 By using dual and hyperdual numbers, the accuracies of the gradient and Hessian
152 computed by the F-1 algorithm are close to machine precision. Furthermore, the F-1
153 algorithm requires only a single factorization of the real-valued matrix $\mathbf{A} \equiv \nabla_{\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p})$
154 followed by $m + 1$ forward and back substitutions to compute both the gradient and
155 Hessian — the minimum possible. (We use \mathbf{s} instead of $\mathbf{s}(\mathbf{p})$ for brevity throughout.)
156 Additionally, it does not require any call to the inner solver, avoiding the pitfalls of
157 autodifferentiating through an iterative solver. Finally, the F-1 algorithm requires no
158 analytical derivatives with respect to \mathbf{p} , making it simple to use.

159 We emphasize that an important requirement for the F-1 algorithm to be applica-
160 ble is that it must be possible to create, store, and factorize the Jacobian matrix, \mathbf{A} .
161 In other words, a generalization of the F-1 algorithm to problems for which the state is
162 too large for the Jacobian to be factored or problems with millions of parameters is be-
163 yond the scope of this study. To illustrate the efficiency of the F-1 algorithm, we apply
164 it to the optimization of a global marine phosphorus-cycling model. (In this model,
165 described in detail in [subsection 3.2](#), the Jacobian, \mathbf{A} , is a sparse $400\,320 \times 400\,320$
166 matrix with 3 800 846 non-zero entries, and \mathbf{p} is of length $m = 6$.)

167 Global nutrient-cycling models play a key role in our understanding of the Earth
168 system. Photosynthetic microbes living in the sunlit upper ocean continuously remove
169 dissolved carbon dioxide (CO_2) and nutrients from surface waters to produce organic
170 matter that gets exported to depth in the form of sinking particles. The downward
171 flux of these particles supplies the carbon and energy that sustain life in the dark
172 subsurface waters. The respiration at depth of the exported organic matter particles
173 also maintains a vertical gradient of CO_2 in the ocean against the tendency of mixing
174 and overturning circulation to homogenize the concentration of dissolved constituents.
175 As such, this “biological carbon pump” [[54](#), [46](#), [1](#)] sets the partitioning of CO_2 be-
176 tween the atmosphere and ocean with important consequences for the climate of the
177 Earth. The strength of the biological pump is strongly regulated by the availability
178 of nutrients such as phosphate (PO_4).

179 The biogeochemical mechanisms that control the cycling of nutrients, which in-
180 volve hundreds of thousands of species, are complex. Additionally, the ocean is hard
181 to sample with sufficient spatial and temporal coverage, hindering our capacity to
182 understand and quantify marine processes. Oceanographers therefore often rely on
183 mathematical models with biogeochemical parameterizations that are calibrated by
184 requiring the model to reproduce available observations. Estimating these parameters
185 through objective optimization is a task whose importance is increasingly recognized
186 in marine biogeochemistry (see, e.g., [[24](#), [25](#), [15](#), [8](#), [10](#), [44](#)]) but whose implementation
187 remains a formidable challenge for state-of-the-art Earth System Models.

188 We demonstrate the performance of the F-1 algorithm in the context of optimizing
189 $m = 6$ parameters of a global marine biogeochemistry model of the phosphorus cycle,
190 for which the state has size $n = 400\,320$. We benchmark the F-1 algorithm against
191 other numerical differentiation algorithms by recording computation times in the same
192 conditions (i.e., for the same state function, \mathbf{F} , and objective function, \hat{f} , with the
193 same parameters, \mathbf{p} , on the same computer, and so forth). We show that the F-1
194 algorithm affords significant speedups. In fact, we show that the F-1 algorithm can
195 compute gradient and Hessian at virtually no added cost relative to a fully analytical

196 approach, as has been suggested in [49].

197 In the phosphorus-cycling-model optimization context, the F-1 algorithm com-
 198 putes the Hessian matrix from 16 to 100 times faster than other algorithms, affording
 199 4- to 26-fold computational speedups overall. Based on simple time-complexity argu-
 200 ments, we expect the computational cost benefits of the F-1 algorithm to scale with
 201 the size of the problem, because current state-of-the-art numerical differentiation al-
 202 gorithms require $\mathcal{O}(m)$ to $\mathcal{O}(m^2)$ factorizations, compared to a single factorization
 203 for the F-1 algorithm. In fact, for fixed n , we expect that computing the Hessian
 204 using the F-1 algorithm would be 3 orders of magnitude faster than finite-differences
 205 for $m \sim 20$ parameters, and 5 orders of magnitude for $m \sim 200$ parameters.

206 We start by describing the F-1 algorithm and its derivation in section 2, where
 207 we also give a short description of six other numerical differentiation algorithms. We
 208 describe our implementation in section 3 and show the results of the optimization of
 209 the phosphorus-cycling model and of the benchmarks of the F-1 algorithm in section 4.
 210 We conclude in section 5 and further discuss in section 6.

211 2. Theory.

212 **2.1. Analytical formulas.** The gradient, $\nabla \hat{f}(\mathbf{p})$, which we express as a row
 213 vector, is obtained by differentiating (1.3) via the chain rule:

$$214 \quad (2.1) \quad \underbrace{\nabla \hat{f}(\mathbf{p})}_{1 \times m} = \underbrace{\nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})}_{1 \times n} \underbrace{\nabla \mathbf{s}(\mathbf{p})}_{n \times m} + \underbrace{\nabla_{\mathbf{p}} f(\mathbf{s}, \mathbf{p})}_{1 \times m},$$

215 where $\nabla_{\mathbf{x}} f$ and $\nabla_{\mathbf{p}} f$ are the partial derivatives of $f(\mathbf{x}, \mathbf{p})$ with respect to \mathbf{x} and \mathbf{p} ,
 216 respectively, and $\nabla \mathbf{s}(\mathbf{p})$ is the derivative of \mathbf{s} with respect to \mathbf{p} . The matrix size of
 217 each derivative, which defines the row/column orientation of vectors and matrices,
 218 is indicated below each term. Here, strictly speaking, the row vector $\nabla \hat{f}(\mathbf{p})$ is the
 219 transpose of the gradient of $\hat{f}(\mathbf{p})$, which we usually take to be an $m \times 1$ column vector.
 220 However, we refer to $\nabla \hat{f}(\mathbf{p})$ as the gradient (of the objective function) throughout
 221 for simplicity. The $\nabla \mathbf{s}(\mathbf{p})$ term in (2.1) is obtained by differentiating the steady-state
 222 equation, (1.1), and gives

$$223 \quad (2.2) \quad \underbrace{\mathbf{A}}_{n \times n} \underbrace{\nabla \mathbf{s}(\mathbf{p})}_{n \times m} + \underbrace{\nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})}_{n \times m} = 0,$$

224 where $\nabla_{\mathbf{p}} \mathbf{F}$ is the partial derivative of \mathbf{F} with respect to \mathbf{p} , and the Jacobian matrix
 225 $\mathbf{A} = \nabla_{\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p})$ is the partial derivative of $\mathbf{F}(\mathbf{x}, \mathbf{p})$ with respect to \mathbf{x} evaluated at
 226 $\mathbf{x} = \mathbf{s}(\mathbf{p})$ and \mathbf{p} . A closed formula for the gradient, $\nabla \hat{f}(\mathbf{p})$ is then obtained by
 227 inserting the solution of (2.2) into (2.1), giving

$$228 \quad (2.3) \quad \underbrace{\nabla \hat{f}(\mathbf{p})}_{1 \times m} = - \underbrace{\nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})}_{1 \times n} \underbrace{\mathbf{A}^{-1}}_{n \times n} \underbrace{\nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})}_{n \times m} + \underbrace{\nabla_{\mathbf{p}} f(\mathbf{s}, \mathbf{p})}_{1 \times n}.$$

229 The computation time for evaluating (2.3) strongly depends on the size of the
 230 state, n , through the need to solve for the $n \times 1$ vector \mathbf{s} and because of the need to
 231 evaluate and invert the $n \times n$ matrix \mathbf{A} . We note, however, that once the factors of
 232 \mathbf{A} are available, they can be used to evaluate the first term on the right hand side for
 233 relatively little additional cost by first evaluating $\mathbf{A}^{-\top} \nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})^{\top}$, which can be done
 234 with a single forward and backward substitution, and then multiplying its transpose
 235 by $\nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})$, rather than first evaluating $\mathbf{A}^{-1} \nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})$, which would consist of
 236 solving m linear systems instead of one.

237 We now turn to the Hessian, for which an analytical expression is obtained by
 238 differentiating (2.1). Using the compact tensor-product notation of Manton [27], we
 239 get

$$\begin{aligned}
 \nabla^2 \hat{f}(\mathbf{p}) &= \nabla_{\mathbf{x}\mathbf{x}} f(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \nabla \mathbf{s}) + \nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}}) \\
 &+ \nabla_{\mathbf{p}\mathbf{x}} f(\mathbf{s}, \mathbf{p}) (\mathbf{I}_{\mathbf{p}} \otimes \nabla \mathbf{s}) + \nabla_{\mathbf{p}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) (\mathbf{I}_{\mathbf{p}} \otimes \mathbf{I}_{\mathbf{p}}) \\
 &+ \nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p}) \nabla^2 \mathbf{s},
 \end{aligned}
 \tag{2.4}$$

241 where we have omitted the \mathbf{p} argument of $\nabla \mathbf{s}(\mathbf{p})$ and $\nabla^2 \mathbf{s}(\mathbf{p})$ for brevity. Here we
 242 give a brief explanation on how to interpret the tensor-product notation, and refer the
 243 interested reader to [27] for more details. In (2.4), the tensor products can be under-
 244 stood merely as separating the arguments for each combination of directions, which
 245 second-order derivatives are applied to. For example, evaluating $\nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}})$
 246 is done by applying the $n \times m$ tensor, $\nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p})$, to the combined direction $(\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}})$.
 247 That is, one must contract the first dimension of $\nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p})$ on the first dimension of
 248 $\nabla \mathbf{s}$ (which is a $n \times m$ matrix) and its second dimension on the first dimension of $\mathbf{I}_{\mathbf{p}}$
 249 (which is the $m \times m$ identity). Effectively, $\nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}})$ results in a $m \times m$
 250 matrix, which can be understood as the matrix product $\nabla \mathbf{s}^T \nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p})$.

251 Note that the evaluation of the expression on the right-hand side of (2.4) makes
 252 use of the second derivative of the steady-state solution, i.e., the $n \times m \times m$ tensor,
 253 $\nabla^2 \mathbf{s}(\mathbf{p})$. It is obtained by differentiating (2.2), which yields

$$\begin{aligned}
 0 &= \nabla_{\mathbf{x}\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \nabla \mathbf{s}) + \nabla_{\mathbf{x}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}}) \\
 &+ \nabla_{\mathbf{p}\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p}) (\mathbf{I}_{\mathbf{p}} \otimes \nabla \mathbf{s}) + \nabla_{\mathbf{p}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) (\mathbf{I}_{\mathbf{p}} \otimes \mathbf{I}_{\mathbf{p}}) \\
 &+ \mathbf{A} \nabla^2 \mathbf{s}.
 \end{aligned}
 \tag{2.5}$$

255 In (2.5), note that every term in the sum is a $n \times m \times m$ tensor (none of which
 256 can be represented in matrix form, justifying our use of the tensor notation of [27]).
 257 For example, $\nabla_{\mathbf{x}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) (\nabla \mathbf{s} \otimes \mathbf{I}_{\mathbf{p}})$ is the double contraction of the second dimension
 258 of the $n \times n \times m$ tensor $\nabla_{\mathbf{x}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})$ on the first dimension of $\nabla \mathbf{s}$ and of its third
 259 dimension on the first dimension of $\mathbf{I}_{\mathbf{p}}$. Also note that $\mathbf{A} \nabla^2 \mathbf{s}$ must be understood
 260 as \mathbf{A} multiplying each of the m^2 column vectors of dimension $n \times 1$ contained in the
 261 $n \times m \times m$ tensor $\nabla^2 \mathbf{s}$.

262 In principle, (2.5), can be substituted into the adjoint of (2.4) to compute the
 263 Hessian matrix by solving a single linear system involving the Jacobian matrix \mathbf{A} .
 264 Thus, at least one factorization of \mathbf{A} and one forward and back substitutions are
 265 required to compute the Hessian with a closed analytical formula. However, numeri-
 266 cally computing each term in (2.5) is tedious at best. Instead, a better solution is to
 267 compute the premultiplied terms arising from inserting (2.5) into (2.4) so as to avoid
 268 computing the third-order tensors, as suggested by [49].

269 **2.2. The F-1 algorithm.** For the gradient, the F-1 algorithm first uses (2.2) to
 270 compute $\nabla \mathbf{s}(\mathbf{p})$, which requires m forward and back substitutions and the factoriza-
 271 tion of \mathbf{A} . Once computed, $\nabla \mathbf{s}(\mathbf{p})$ is then inserted into (2.1) to evaluate the gradient.
 272 Importantly, the partial derivatives of \mathbf{F} and f with respect to the parameters, \mathbf{p} , are
 273 computed numerically using dual numbers. Specifically, $\nabla_{\mathbf{p}} f(\mathbf{s}, \mathbf{p})$ and $\nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p})$ are
 274 computed in m dual-valued evaluations, via

$$\nabla_{\mathbf{p}} f(\mathbf{s}, \mathbf{p}) = \mathcal{D}[f(\mathbf{s}, \mathbf{p} + \varepsilon \mathbf{e}_1), \dots, f(\mathbf{s}, \mathbf{p} + \varepsilon \mathbf{e}_m)]
 \tag{2.6}$$

276 and

$$\nabla_{\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) = \mathcal{D}[\mathbf{F}(\mathbf{s}, \mathbf{p} + \varepsilon \mathbf{e}_1), \dots, \mathbf{F}(\mathbf{s}, \mathbf{p} + \varepsilon \mathbf{e}_m)],
 \tag{2.7}$$

278 respectively.

279 For the Hessian, the F-1 algorithm uses hyperdual numbers, but exploits a com-
 280 bination of (2.4) and (2.5) that provides an optimally efficient analytical shortcut,
 281 which reduces the cost of computing the Hessian down to a single forward and back
 282 substitution. For $\mathbf{p}_{jk} = \mathbf{p} + \varepsilon_1 \mathbf{e}_j + \varepsilon_2 \mathbf{e}_k$, let $\mathbf{x}_{jk} \equiv \mathbf{s} + \varepsilon_1 \nabla \mathbf{s} \mathbf{e}_j + \varepsilon_2 \nabla \mathbf{s} \mathbf{e}_k$ denote a
 283 carefully chosen corresponding hyperdual-valued state. (We have dropped the \mathbf{p}
 284 dependency on $\mathbf{s}(\mathbf{p})$ and on $\nabla \mathbf{s}(\mathbf{p})$ for brevity again.) Then, the entries of the Hessian
 285 matrix of the objective function are given by

$$286 \quad (2.8) \quad [\nabla^2 \hat{f}(\mathbf{p})]_{jk} = \mathfrak{H}[f(\mathbf{x}_{jk}, \mathbf{p}_{jk})] - \mathfrak{H}[\mathbf{F}(\mathbf{x}_{jk}, \mathbf{p}_{jk})^\top] \mathbf{A}^{-\top} \nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})^\top,$$

287 where $\mathfrak{H}(x)$ is the $\varepsilon_1 \varepsilon_2$ coefficient of x . (A formal derivation of (2.8) is given later in
 288 this section.) With (2.8), a single forward and back substitution for $\mathbf{A}^{-\top} \nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})^\top$
 289 is required for all the entries of the Hessian because it is independent of j and k .
 290 Additionally, because the Hessian is symmetric, only $m(m+1)/2$ hyperdual-valued
 291 evaluations of f and \mathbf{F} are necessary.

292 Hence, the F-1 algorithm computes the gradient in a single factorization of \mathbf{A} ,
 293 m forward and back substitutions, and $\mathcal{O}(m)$ function evaluations, and computes the
 294 Hessian in a single forward and back substitution and $\mathcal{O}(m^2)$ function evaluations.
 295 Additionally, the F-1 algorithm requires only f , \mathbf{F} , $\nabla_{\mathbf{x}} f$, and $\nabla_{\mathbf{x}} \mathbf{F}$ from the user.
 296 We note that although the F-1 algorithm requires the user to supply the Jacobian,
 297 $\nabla_{\mathbf{x}} \mathbf{F}$, as well as $\nabla_{\mathbf{x}} f$, these may sometimes be easily derived analytically or computed
 298 numerically, as is the case for our phosphorus-cycling model (details in [section SM2](#)).
 299 Thus in some cases, with little extra work, the F-1 algorithm provides an automatic
 300 differentiation tool that requires no derivatives from the user — only \mathbf{F} and f .

301 We note that (2.8) is similar to Equation (13) in the work of Rumpfkeil and
 302 Mavriplis [49], who suggest to premultiply the third-order tensors in (2.5) and use
 303 an algorithmic differentiation tool to compute the directional derivatives for each (j, k)
 304 direction. Here, we accomplish the same thing by simply evaluating f and \mathbf{F} with
 305 appropriately chosen hyperdual-valued arguments. Thus, the F-1 algorithm provides
 306 an easy-to-implement and similarly optimally-efficient alternative.

307 We now derive (2.8). The hyperdual-valued Taylor expansion of f at (\mathbf{s}, \mathbf{p}) in the
 308 $(\varepsilon_1 \nabla \mathbf{s} \mathbf{e}_j + \varepsilon_2 \nabla \mathbf{s} \mathbf{e}_k, \varepsilon_1 \mathbf{e}_j + \varepsilon_2 \mathbf{e}_k)$ direction gives exactly the cross terms of (2.4) as its
 309 $\varepsilon_1 \varepsilon_2$ coefficient. That is,

$$310 \quad (2.9) \quad \mathfrak{H}[f(\mathbf{x}_{jk}, \mathbf{p}_{jk})] = \mathbf{e}_j^\top \nabla \mathbf{s}^\top \nabla_{\mathbf{x}\mathbf{x}} f(\mathbf{s}, \mathbf{p}) \nabla \mathbf{s} \mathbf{e}_k + \mathbf{e}_j^\top \nabla \mathbf{s}^\top \nabla_{\mathbf{x}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) \mathbf{e}_k \\ + \mathbf{e}_j^\top \nabla_{\mathbf{p}\mathbf{x}} f(\mathbf{s}, \mathbf{p}) \nabla \mathbf{s} \mathbf{e}_k + \mathbf{e}_j^\top \nabla_{\mathbf{p}\mathbf{p}} f(\mathbf{s}, \mathbf{p}) \mathbf{e}_k,$$

311 where $\mathfrak{H}(x)$ is the $\varepsilon_1 \varepsilon_2$ coefficient of x . Mathematically, $\mathfrak{H}[f(\mathbf{x}_{jk}, \mathbf{p}_{jk})]$ is simply
 312 the second-order directional derivative of f at (\mathbf{s}, \mathbf{p}) in the combined $(\nabla \mathbf{s} \mathbf{e}_j, \mathbf{e}_j)$ and
 313 $(\nabla \mathbf{s} \mathbf{e}_k, \mathbf{e}_k)$ directions.

314 The entry-wise version of (2.4) can be rearranged and expressed as

$$315 \quad (2.10) \quad [\nabla^2 \hat{f}(\mathbf{p})]_{jk} = \mathfrak{H}[f(\mathbf{x}_{jk}, \mathbf{p}_{jk})] + \nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p}) [\nabla^2 \mathbf{s}]_{jk},$$

316 where $[\nabla^2 \mathbf{s}]_{jk}$ is the $n \times 1$ column vector given by the second partial derivative of
 317 $\mathbf{s}(\mathbf{p})$ with respect to the j th and k th parameters (which contracts on the $1 \times n$ row
 318 vector multiplied to its left, $\nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})$, as per the tensor notation of [27] resulting in
 319 a scalar entry).

320 Similarly, the hyperdual-valued Taylor expansion of \mathbf{F} , also taken at (\mathbf{s}, \mathbf{p}) and
 321 in the same $(\varepsilon_1 \nabla \mathbf{s} \mathbf{e}_j + \varepsilon_2 \nabla \mathbf{s} \mathbf{e}_k, \varepsilon_1 \mathbf{e}_j + \varepsilon_2 \mathbf{e}_k)$ direction, gives exactly the cross terms

322 of (2.5) as its $\varepsilon_1\varepsilon_2$ coefficient. That is,

$$323 \quad (2.11) \quad \mathfrak{H}[\mathbf{F}(\mathbf{x}_{jk}, \mathbf{p}_{jk})] = \mathbf{e}_j^\top \nabla \mathbf{s}^\top \nabla_{\mathbf{x}\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p}) \nabla \mathbf{s} \mathbf{e}_k + \mathbf{e}_j^\top \nabla \mathbf{s}^\top \nabla_{\mathbf{x}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) \mathbf{e}_k \\ + \mathbf{e}_j^\top \nabla_{\mathbf{p}\mathbf{x}} \mathbf{F}(\mathbf{s}, \mathbf{p}) \nabla \mathbf{s} \mathbf{e}_k + \mathbf{e}_j^\top \nabla_{\mathbf{p}\mathbf{p}} \mathbf{F}(\mathbf{s}, \mathbf{p}) \mathbf{e}_k,$$

324 which can be interpreted as the second-order directional derivative of \mathbf{F} at (\mathbf{s}, \mathbf{p}) in
325 the combined $(\nabla \mathbf{s} \mathbf{e}_j, \mathbf{e}_j)$ and $(\nabla \mathbf{s} \mathbf{e}_k, \mathbf{e}_k)$ directions. Thus, the entry-wise version of
326 (2.5) can be recast as

$$327 \quad (2.12) \quad [\nabla^2 \mathbf{s}]_{jk} = -\mathbf{A}^{-1} \mathfrak{H}[\mathbf{F}(\mathbf{x}_{jk}, \mathbf{p}_{jk})].$$

328 Inserting (2.12) into (2.10) and taking the adjoint yields (2.8).

329 **3. Implementation.** We chose the Julia language [3] for the implementation of
330 this study because it affords, among many other features, (i) state-of-the-art factoriza-
331 tion of sparse matrices and solving of sparse linear systems using the standard-library
332 `LinearAlgebra`, `SparseArrays`, and `SuiteSparse` packages [5, 6], (ii) state-of-the-art op-
333 timizations using the `Optim` package [30, 31], and (iii) efficient implementations of
334 the dual- and hyperdual-number types using the `DualNumbers` and `HyperDualNum-
335 bers` packages. We emphasize that different implementations are possible, so that
336 one could use another scientific computing language or other packages. However, this
337 combined choice of Julia and relevant packages allows for a simple yet fast, modern,
338 and open-source implementation of the algorithms benchmarked in this study.

339 3.1. Algorithms for the gradient and Hessian.

340 **3.1.1. F-1 algorithm.** The Julia implementation of the F-1 algorithm is pub-
341 licly available online as the `F1Method` package [42], which was developed for this study
342 and used in the optimization benchmarks of section 4. Thanks to the expressivity and
343 syntax of the Julia language, the code for the entire `F1Method` package requires just
344 a few lines of code that closely match (2.1), (2.2), and (2.6)–(2.8). The `F1Method`
345 package essentially defines five functions, which are called to update the steady-state
346 solution, update a memory cache, and compute the objective, its gradient, and its
347 Hessian, respectively.

348 The memory cache, denoted `mem`, used by the `F1Method` package, is an instance of
349 a custom Julia type, `Mem`, which is used to store the results of reusable computations,
350 i.e., the steady-state solution, $\mathbf{s}(\mathbf{p})$, the factors of the Jacobian, \mathbf{A} , the derivatives
351 $\nabla \mathbf{s}(\mathbf{p})$ and $\nabla_{\mathbf{x}} f(\mathbf{s}, \mathbf{p})$, and the corresponding parameter values, \mathbf{p} . When calling either
352 of the objective, the gradient, or the Hessian functions, the contents of `mem` are only
353 updated if the parameters, \mathbf{p} , are modified.

354 **3.1.2. Other algorithms.** We chose to benchmark the F-1 algorithm against
355 five algorithms that apply their respective numerical scheme to either the analytical
356 gradient function, $\nabla \hat{f}$, as defined by (2.3), or directly to the objective function, \hat{f} , as
357 defined by (1.3), taken as black boxes. The algorithms that use the analytical gradi-
358 ent formula are (i) the FD1 algorithm, which applies a second-order finite-difference
359 scheme, (ii) the CSD algorithm, which applies the complex-step scheme, and (iii) the
360 DUAL algorithm, which applies the dual-step scheme. The algorithms that directly
361 use the objective function to compute both gradient and Hessian are (iv) the FD2
362 algorithm, which applies second-order finite-difference schemes for both the gradi-
363 ent and Hessian, and (v) the HYPER algorithm, which applies the dual-step scheme
364 for the gradient and the hyperdual-step scheme for the Hessian. These algorithms,

365 together with the F-1 algorithm, are collected in [Table 1](#), with a link to their Julia
 366 implementation and a short description with an estimate of their computational costs.
 367 Overall, we thus benchmark six algorithms that are run in the same conditions and
 368 on the same machine, allowing for a fair set of comparisons.

TABLE 1

Collection of the algorithms benchmarked in this study. The FD1, CSD, and DUAL algorithms compute the Hessian, $\nabla^2 \hat{f}(\mathbf{p})$, by numerically differentiating the gradient, $\nabla \hat{f}$, using the analytical formula, (2.3). They require the user to supply f , \mathbf{F} , $\nabla_{\mathbf{x}} f$, $\nabla_{\mathbf{x}} \mathbf{F}$, $\nabla_{\mathbf{p}} f$, and $\nabla_{\mathbf{p}} \mathbf{F}$. The F-1 algorithm requires f , \mathbf{F} , $\nabla_{\mathbf{x}} f$, and $\nabla_{\mathbf{x}} \mathbf{F}$. The FD2 and HYPER algorithms compute both gradient and Hessian directly from the objective function and thus require only \hat{f} . The “cost” column shows the number of factorizations as a function of the number of parameters, m . Names in the first column are clickable URLs that point to their implementation in GitHub.

Algorithm	Cost	Definition
F-1	$\mathcal{O}(1)$	The algorithm presented in this study. Computes $\nabla \hat{f}(\mathbf{p})$ and $\nabla^2 \hat{f}(\mathbf{p})$ using the analytical shortcuts described in subsection 2.2 . Overall requires 1 factorization and $m + 1$ forward and back substitutions.
FD1	$\mathcal{O}(m)$	Central finite-differences algorithm. Computes $\nabla^2 \hat{f}(\mathbf{p})$ as the Jacobian of $\nabla \hat{f}$. Executes $2m$ calls to $\nabla \hat{f}$ (i.e., $2m$ factorizations) and to the inner solver (about 1 iteration each time). Overall requires about $4m$ additional factorizations.
CSD	$\mathcal{O}(m)$	Complex-step differentiation algorithm. Computes $\nabla^2 \hat{f}(\mathbf{p})$ by evaluating $\nabla \hat{f}$ with complex parameters. Executes m calls to the complex $\nabla \hat{f}$ (i.e., m complex factorizations) and to the inner solver (2 iterations each time). Overall requires about $2m$ additional complex factorizations. (Note that complex-valued operations can be more expensive than real-valued ones.)
DUAL	$\mathcal{O}(m)$	Dual-step differentiation algorithm. Computes $\nabla^2 \hat{f}(\mathbf{p})$ by evaluating $\nabla \hat{f}$ with dual parameters. Executes m calls to the dual $\nabla \hat{f}$ and to the inner solver (2 iterations each time). Overall requires about $2m$ additional dual factorizations. (Note the DUAL method applies (SM5.1) for solving dual-valued linear systems).
FD2	$\mathcal{O}(m^2)$	Second-order finite-differences algorithm. Computes $\nabla \hat{f}(\mathbf{p})$ and $\nabla^2 \hat{f}(\mathbf{p})$ directly from \hat{f} . Executes $2m$ calls to the inner solver for the gradient and $2m^2$ for the Hessian. Overall requires about $2m^2 + 2m$ additional factorizations. (Note that the low accuracy of the FD2 algorithm slows the convergence of the optimizer.)
HYPER	$\mathcal{O}(m^2)$	Dual- and Hyperdual-step differentiation algorithm. Computes $\nabla \hat{f}(\mathbf{p})$ and $\nabla^2 \hat{f}(\mathbf{p})$ by evaluating \hat{f} with dual-valued and hyperdual parameters, respectively. Executes m calls to the inner solver with dual parameters (2 iterations each time), and $m(m + 1)/2$ calls with hyperdual-valued parameters (3 iterations each time). Overall requires about m additional dual factorizations and $m(m + 1)/2$ additional hyperdual factorizations. (Note that the HYPER algorithm applies (SM5.1) and (SM5.2) for dual and hyperdual linear systems.)

369 Importantly, we note that the CSD, DUAL, and HYPER algorithms invoke the
 370 inner solver with non-real number types. Thus, the inner solver must be able to check
 371 the convergence of the imaginary, dual, and hyperdual parts. In practice, the solver
 372 that we use applies the Shamanskii method, a quasi-Newton method that computes
 373 Newton steps [\[19, 20\]](#). We thus add a conditional statement on the relative size of
 374 the non-real parts of the Newton step for the inner-solver loop to terminate.

375 Additionally, for the CSD, DUAL, and HYPER algorithms, the inner-solver New-
 376 ton steps require solving complex-valued, dual-valued, or hyperdual-valued linear sys-
 377 tems. While UMFPACK [\[5\]](#), which is the C package called by Julia’s [SuiteSparse](#)
 378 package [\[6\]](#) for the LU factorization of unsymmetric sparse matrices (like \mathbf{A} in our
 379 phosphorus-cycling model), can handle complex numbers, it cannot deal with dual or

380 hyperdual valued number types. Thus an important advantage of the F-1 algorithm in
 381 this regard is that despite using dual and hyperdual numbers, the F-1 algorithm does
 382 not require the solution to any dual-valued or hyperdual-valued linear systems. For
 383 the DUAL and HYPER algorithms — specifically, for the inner solver to handle dual-
 384 valued and hyperdual-valued factorizations and forward and back substitutions — we
 385 developed two Julia packages, `DualMatrixTools` [39] and `HyperDualMatrixTools` [40].
 386 These packages afford efficient factorization of dual and hyperdual sparse matrices, as
 387 well as solving dual and hyperdual linear systems with a minimum number of forward
 388 and back substitutions. Both packages rely on analytical identities for the inverse of
 389 dual-valued and hyperdual-valued matrices that are derived in [section SM5](#).

390 **3.2. Phosphorus cycling model.** To generate the global marine phosphorus
 391 cycling model, we use the `AIBECS` package (for Algebraic Implicit Biogeochemical
 392 Elemental Cycling System, [41]), which was developed in parallel to this study. In
 393 the AIBECS, steady-state problems are built as objects of the `SteadyStateProblem`
 394 type as defined by the `DiffEqBase` package [45] in Julia. The steady-state solution is
 395 computed via a state-of-the-art quasi-Newton algorithm implementing the Shamanskii
 396 method, derived from the work of [19, 20] and coded inside the AIBECS. We reiterate
 397 that the solver invoked by the AIBECS has been carefully designed to handle real-,
 398 complex-, dual-, and hyperdual-valued state and parameters, and integrates the `Dual-`
 399 `MatrixTools` [39] and `HyperDualMatrixTools` [40] packages for dual- and hyperdual-
 400 valued factorizations and forward and back substitutions of sparse linear systems.

401 Our phosphorus-cycling model consists of two marine tracers, dissolved inorganic
 402 phosphorus (DIP), i.e., phosphate, and particulate organic phosphorus (POP). DIP is
 403 transported by water currents and turbulent eddies, and is taken up by phytoplank-
 404 ton in the euphotic layer (where light is available for photosynthesis to occur) and
 405 converted to sinking POP. As it sinks, POP is remineralized into DIP. The sequence
 406 of uptake, sinking, and remineralization, provides the downward transport mechanism
 407 (the biological pump [54, 46, 1]). The system is in steady state when the circulation,
 408 which brings nutrients back to the surface, balances the biological pump. The model
 409 only explicitly tracks DIP and POP, thus the state of the system is entirely deter-
 410 mined by the concentration fields of DIP and POP, denoted by $x_{\text{DIP}}(\mathbf{r})$ and $x_{\text{POP}}(\mathbf{r})$
 411 at location \mathbf{r} .

412 The evolution of the system is determined by the coupled mass-conservation equa-
 413 tions for DIP and POP. They are given by

$$414 \quad (3.1) \quad \begin{cases} (\partial_t + \nabla_{\mathbf{r}} \cdot [\mathbf{u} - \mathbf{K}\nabla_{\mathbf{r}}]) x_{\text{DIP}} = -U(x_{\text{DIP}}) + R(x_{\text{POP}}) \\ (\partial_t + \nabla_{\mathbf{r}} \cdot \mathbf{w}) x_{\text{POP}} = +U(x_{\text{DIP}}) - R(x_{\text{POP}}) \end{cases},$$

415 where $\nabla_{\mathbf{r}}$ is the classical three-dimensional gradient operator. (We have omitted the
 416 \mathbf{r} dependency of the DIP and POP fields for brevity.) [Equation \(3.1\)](#) defines a system
 417 of two coupled PDEs where we assume no-flux (Neumann) boundary conditions at
 418 the land–ocean and atmosphere–ocean interfaces. On the left hand sides of [\(3.1\)](#), \mathbf{u}
 419 is the 3D water velocity vector field, \mathbf{K} is the 3×3 eddy-diffusivity tensor, and \mathbf{w} is the
 420 particulate sinking velocity. The flux divergence of DIP due to the ocean’s currents
 421 and eddies is thus represented by the action of the advective–eddy-diffusive transport
 422 operator, $\nabla_{\mathbf{r}} \cdot [\mathbf{u} - \mathbf{K}\nabla_{\mathbf{r}}]$, on x_{DIP} .

423 The flux divergence of sinking POP is represented by the action of $\nabla_{\mathbf{r}} \cdot \mathbf{w}$ on x_{POP} .
 424 The remineralization profile of POP is assumed to follow a power-law with depth after
 425 the observations of Martin [28]. Following [23], this is equivalent to assuming that

TABLE 2

Parameters for the global marine phosphorus-cycling model with their prior and posterior means.

Symbol	Description	Prior mean	Posterior mean	Unit
x^{geo}	Mean DIP concentration	2.17	2.12	mmol m^{-3}
k	Half-saturation constant (Michaelis-Menten)	10.00	6.62	$\mu\text{mol m}^{-3}$
w_0	Sinking velocity at surface	1.00	0.64	m d^{-1}
w'	Vertical gradient of sinking velocity	0.22	0.13	d^{-1}
κ	Dissolution rate constant (POP to DIP)	0.19	0.19	d^{-1}
τ	Maximum uptake rate timescale	30.00	236.52	d

the magnitude of the sinking velocity, w , increases linearly with depth, an approach we adopt here, such that $w = w'z + w_0$, where w' and w_0 are optimizable parameters.

On the right hand sides of (3.1), U and R are the local uptake and remineralization rates, respectively. The specific phosphate uptake by phytoplankton in the euphotic layer is modeled according to a simple Monod term [32] with maximum set by the timescale τ and half-saturation rate constant k , and the POP remineralization is modeled after a first order reaction depending only on the POP concentration. Hence, U and R are defined by

$$(3.2) \quad \begin{cases} U(x_{\text{DIP}}) \equiv \frac{x_{\text{DIP}}}{\tau} \frac{x_{\text{DIP}}}{x_{\text{DIP}} + k} & \text{where } z \leq z_0, \\ R(x_{\text{POP}}) \equiv \kappa x_{\text{POP}} \end{cases},$$

where τ , k , and κ are optimizable parameters. (For our discrete model grid, the depth of the bottom of the euphotic layer, z_0 , lies at the bottom of the second layer, i.e., at about 73 m, below which $U \equiv 0$.)

Because (3.1) does not contain external sources and sinks to the system, the global means are not constrained and could be chosen arbitrarily (see, e.g., [25]). We prescribe the global mean phosphate concentration by slowly restoring the DIP concentration everywhere to a mean value, x^{geo} , with a timescale of $\tau_{\text{geo}} = 1 \text{ Myr}$ that is larger than the typical timescale for a tracer to be homogeneously mixed in the ocean. Thus, in practice, $(x^{\text{geo}} - x_{\text{DIP}})/\tau_{\text{geo}}$ is added to the right hand side of the DIP equation in (3.1), where x^{geo} is an optimizable parameter. We note that estimating the value of x^{geo} is of interest because the total inventory of DIP in the ocean, which is uncertain, is given by x^{geo} multiplied by the total volume of the ocean. The $m = 6$ optimizable parameters are collected in Table 2.

The continuous equations in (3.1) are discretized onto a 3D grid of the ocean. Specifically, the steady-state version of (3.1) is recast into (1.1) by rearranging the 3D concentration fields of DIP and POP into a state vector $\mathbf{x} = \begin{bmatrix} \mathbf{x}_{\text{DIP}} \\ \mathbf{x}_{\text{POP}} \end{bmatrix}$ and by replacing the linear operators for the flux divergences by large sparse matrices. For the advective-diffusive transport operator, we use the Ocean Circulation Inverse Model (OCIM1, [9, 7]), which defines the 3D grid of the ocean. With two tracers and the 200 160 boxes of the OCIM1 grid, the state vector, \mathbf{x} , has length $n = 400\,320$. (More details on creating the discrete model and on OCIM1 are given in section SM1.)

For the optimization, we define the objective function as the sum of the squared mismatch of the modeled state against observations and the posterior parameters against their prior mean. Specifically, the objective function is defined by

$$(3.3) \quad f(\mathbf{x}, \mathbf{p}) \equiv \frac{\omega_{\mathbf{x}}}{2} \delta \mathbf{x}^{\text{T}} \boldsymbol{\Omega}_{\mathbf{x}} \delta \mathbf{x} + \frac{\omega_{\mathbf{p}}}{2} \delta \boldsymbol{\lambda}^{\text{T}} \boldsymbol{\Omega}_{\boldsymbol{\lambda}} \delta \boldsymbol{\lambda},$$

460 where $\omega_{\mathbf{x}}$ and $\omega_{\mathbf{p}}$ are hyper parameters that control the relative weights of the state
 461 and the parameters, respectively. In (3.3), $\delta\mathbf{x}$ is the difference between the modeled
 462 and observed DIP concentrations, where the observations are from the World Ocean
 463 Atlas (WOA18, [16, 43]). The diagonal matrix $\mathbf{\Omega}_{\mathbf{x}}$ is taken as the inverse of the
 464 covariance matrix from regridding the WOA18 data onto the OCIM1 grid. For $\delta\mathbf{\lambda}$,
 465 assuming prior log-normal distributions for the parameters, we use $\delta\mathbf{\lambda} \equiv \log(\mathbf{p}) - \boldsymbol{\mu}$,
 466 where $\boldsymbol{\mu}$ is the prior mean in logarithmic space. In practice, we feed $\boldsymbol{\lambda} = \log(\mathbf{p})$ to the
 467 optimizer instead of \mathbf{p} so that the parameters remain positive throughout. (We note
 468 that while our parameters are necessarily positive, they need not be positive in general
 469 and in such a case one would forgo the logarithmic transformation.) The diagonal
 470 matrix $\mathbf{\Omega}_{\boldsymbol{\lambda}}$ is taken as the inverse of the prior covariance matrix in logarithmic space.
 471 (The non-logarithmic prior variances of the parameters are prescribed as the square
 472 of their non-logarithmic prior means.)

473 **3.3. Optimizer.** For the optimization, we use the Trust-Region Newton algo-
 474 rithm of Julia’s `Optim` package [30, 31], which we use to optimize the parameters in
 475 logarithmic space. For the initial choice of parameters, we chose their prior means as
 476 collected in Table 2. The optimizer is deemed to have converged when the norm of
 477 the gradient of the objective function is less than 10^{-8} . The initial state is chosen to
 478 be equal to x^{geo} everywhere. Accurate measurements of computation times for the
 479 benchmarks are performed with the `BenchmarkTools` and `TimerOutputs` packages.

480 4. Results.

481 **4.1. Optimized model.** The prior values of the parameters, which we use as
 482 the initial guess in the optimization, are given in Table 2, along with the posterior
 483 values (i.e., the optimized values). In our phosphorus-cycling model, at the start of the
 484 optimization, for the first steady-state solution, \mathbf{s}_0 , such that $\mathbf{F}(\mathbf{s}_0, \mathbf{p}_0) = 0$, the DIP
 485 field has a large mismatch with observations, with a root-mean-square error (RMSE)
 486 relative to the mean observed DIP of about 41 %. Our careful choice of the weights $\omega_{\mathbf{x}}$
 487 and $\omega_{\mathbf{p}}$ ensures that the optimization effectively reduces the mismatch of the modeled
 488 state, i.e., the first term of (3.3), such that, at the end of the optimization, the RMSE
 489 of the modeled DIP field of the optimal steady-state solution, $\mathbf{s}(\hat{\mathbf{p}})$, for the optimal
 490 parameters $\hat{\mathbf{p}}$, is of only about 6 %. Figure 2 shows a number of diagnostics of the
 491 DIP field of this optimized steady-state solution.

492 The DIP field of the optimized solution is shown at a depth of about 919 m in
 493 Figure 2a. For the same depth, the relative mismatch with observations is shown
 494 in Figure 2b, revealing the mismatch range of approximately $\pm 20\%$. Most of the
 495 mismatch lies within $\pm 5\%$ although there are some significant positive biases in the
 496 Arctic and negative biases at low-latitudes.

497 In order to investigate the mismatch at different depths, Figure 2c shows the
 498 horizontally-averaged DIP concentrations of the optimized solutions and of the ob-
 499 servations for each of the Atlantic (ATL), Pacific (PAC), and Indian (IND) basins.
 500 The excellent fit shows that the optimized model captures the global vertical gradi-
 501 ents, which quantify the strength of the biological pump, remarkably well. This is
 502 confirmed by Figure 2d, which allows us to evaluate the DIP mismatch at every loca-
 503 tion by showing the cost-weighted cumulative joint probability density function of the
 504 modeled and observed DIP fields. The joint distribution, which is concentrated on
 505 the 1 : 1 line, shows that the optimized model strongly agrees with the observations
 506 over the entire ocean.

507 We emphasize that the phosphorus-cycling model we use is simplistic in the sense

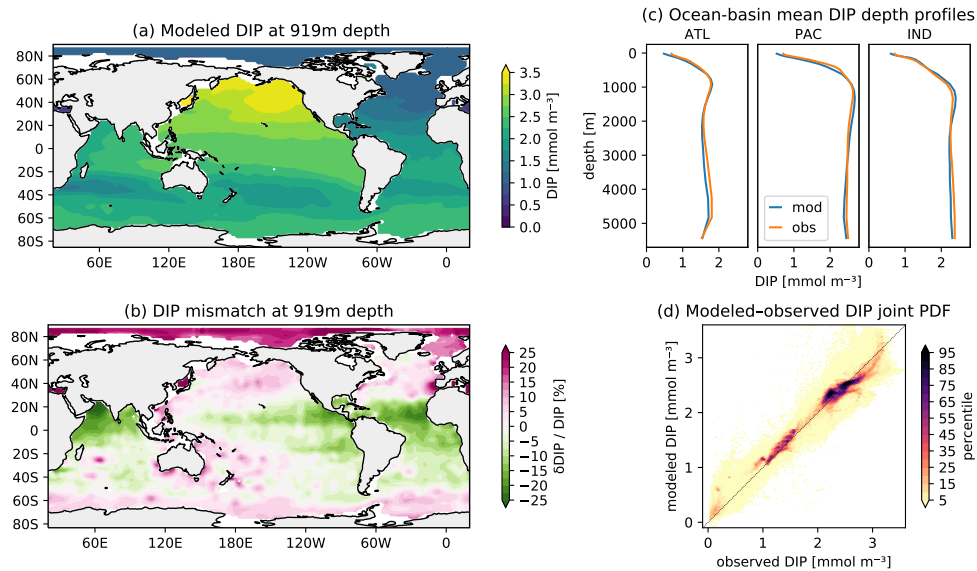


FIG. 2. (a) Modeled DIP at about 919 m depth. (b) Relative mismatch between modeled DIP and observations at the same depth as (a). (c) Horizontally averaged (volume-weighted) concentrations of the modeled and observed DIP for the Atlantic (ATL), Pacific (PAC), and Indian (IND) oceans. (d) Volume- and inverse-variance-weighted cumulative joint probability density function of the modeled and observed DIP. (*i*% of the data lies outside of the contour of the *i*th percentile.)

508 that it does not contain any information on light availability or other nutrient limi-
 509 tations, which are important controls on the distribution of DIP. Hence, the quality
 510 of the fit of the optimized steady-state solution to observations is remarkably good.

511 **4.2. Benchmarks.** We now compare the computation times afforded by using
 512 the F-1 algorithm against the other algorithms collected in Table 1 in the context of
 513 optimizing our global marine phosphorus-cycling model. We run the entire optimiza-
 514 tion procedure as illustrated in Figure 1. (Details of the implementation, such as the
 515 initial state and parameters, are described in section 3.)

516 Figure 3 shows the convergence rate of the optimization as the norm of the gra-
 517 dient, $\|\nabla \hat{f}(\mathbf{p})\|$, versus computation time. Recall that the gradient, $\nabla \hat{f}(\mathbf{p})$, must be
 518 equal to zero where the parameters are optimal. For all the algorithms, convergence is
 519 achieved in 9 iterations of the optimizer, except in the case of the FD2 algorithm, for
 520 which 10 iterations are needed. Using the FD2 algorithm requires more optimizer-loop
 521 iterations because both its gradient and Hessian are inaccurate. However, an accurate
 522 Hessian is not as important as an accurate gradient [18], so that the optimization run
 523 using the FD1 algorithm also converges in 9 iterations despite an inaccurate Hessian.

524 The optimization using the F-1 algorithm converges in about 7 min and is, by far,
 525 the fastest. Using the DUAL, FD1, and CSD algorithms, convergence is achieved in
 526 about 28 min, 40 min and 47 min, respectively. The HYPER algorithm converges in
 527 almost 2 h while the FD2 algorithm takes almost 3 h. The F-1 algorithm is about 24
 528 times faster than the FD2 algorithm, which is the most common numerical differen-
 529 tiation algorithm.

530 We note that each optimization run includes unavoidable computations, regardless

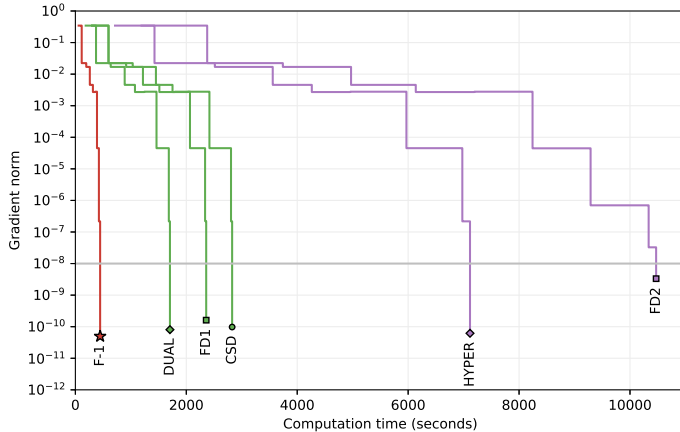


FIG. 3. Convergence for all the algorithms from Table 1 quantified by the gradient norm, $\|\nabla f(\mathbf{p})\|$, which is zero at the minimum, versus computation time. The algorithm is indicated at the end of each curve. The fastest algorithm is the F-1 algorithm (red star). Other Hessian-only computing algorithms (DUAL, FD1, and CSD) are shown in green, and gradient-and-Hessian-computing algorithms (HYPER and FD2) are shown in purple. The horizontal gray line indicates the tolerance of the optimizer, which terminates when $\|\nabla f(\mathbf{p})\| \leq 10^{-8}$.

531 of the algorithm used. For instance, every time the (real-valued) parameters, \mathbf{p} , are
 532 updated by the optimizer, a non-negligible fraction of the computation time is spent
 533 finding the corresponding steady-state solution, $\mathbf{s}(\mathbf{p})$, by invoking the inner solver.
 534 Below, we partition the computation time for the entire optimization into the time
 535 spent for the the gradient and Hessian. This allows to further untangle the differences
 536 in performance for each algorithm.

537 Figure 4 shows the partition of the time spent computing the gradient (in gray)
 538 and the Hessian (purple) during the same optimization runs as for Figure 3. In
 539 optimization problems, most of the time is usually spent computing the highest-order
 540 derivative, i.e., the Hessian in our context. This is true of all the algorithms except
 541 the F-1 algorithm. In fact, the optimization run using the F-1 algorithm is the only
 542 case where the time spent computing the Hessian is smaller than for computing the
 543 gradient.

544 Importantly, the cost of computing the Hessian only by the F-1 algorithm is spec-
 545 tacularly low. Specifically, the F-1 algorithm is about 16 times, 24 times, and 32
 546 times faster than the DUAL, FD1, and CSD algorithms, for computing the Hessian.
 547 Furthermore, based on the number of factorizations of these algorithms, which scales
 548 like $\mathcal{O}(m)$, one should expect these performance ratios to roughly scale with the num-
 549 ber of parameters, m . (Recall that $m = 6$ in our benchmarks.) For example, one
 550 would reasonably expect speedups of about two orders of magnitude with $m = 25$
 551 parameters, and of about three orders of magnitude with $m = 250$. This is remark-
 552 able, considering the F-1 algorithm only requires the partials $\nabla_{\mathbf{x}}\mathbf{F}$ and $\nabla_{\mathbf{x}}f$ from the
 553 user, compared to the DUAL, FD1, and CSD algorithms, which are all state-of-the-
 554 art (although naive) applications of numerical differentiation, and which additionally
 555 require the analytical formula for $\nabla_{\mathbf{p}}\mathbf{F}$ and $\nabla_{\mathbf{p}}f$. We note however that for large m ,
 556 the cost of function evaluations for the F-1 algorithm, which scales like $\mathcal{O}(m^2)$, may
 557 become predominant.

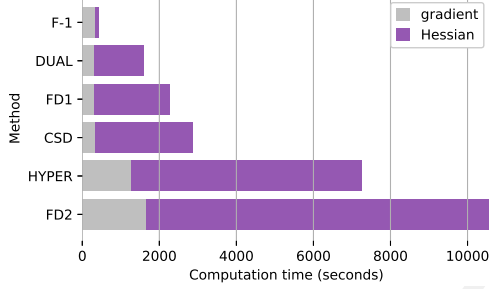


FIG. 4. Cumulated computation times for the entire optimization run for all the algorithms. The computation time for the gradient, $\nabla \hat{f}$, and the Hessian, $\nabla^2 \hat{f}$, are indicated in gray and purple, respectively. (Computation times for the objective function are negligible in the optimization context — see details in [section SM3](#).)

558 A fairer comparison is against the HYPER and FD2 algorithms. Normalized by
 559 the number of calls (which is greater for the FD2 algorithm), the F-1 algorithm is
 560 effectively 76 times faster than the HYPER algorithm, and 100 times faster than
 561 the FD2 algorithm, for computing the Hessian. This is also spectacular, even more
 562 so when considering that the number of factorizations required by the HYPER and
 563 FD2 algorithms scales as $\mathcal{O}(m^2)$. In fact, one should expect the F-1 algorithm to be
 564 about 3 orders of magnitude faster than the HYPER and FD2 algorithms for $m \sim 20$
 565 parameters, and 5 orders of magnitude faster for $m \sim 200$.

566 **5. Conclusions.** We presented a computationally efficient method, the F-1 al-
 567 gorithm, to numerically evaluate the gradient and Hessian of an objective function,
 568 $\hat{f}(\mathbf{p})$, which quantifies a model’s skill (its ability to match observations) as a function
 569 of its parameters, \mathbf{p} . The algorithm is applicable to steady-state problems represented
 570 by a system of discretized nonlinear PDEs, $\mathbf{F}(\mathbf{x}, \mathbf{p}) = 0$, for which the steady-state
 571 solution, $\mathbf{x} = \mathbf{s}(\mathbf{p})$, can be efficiently computed using a Newton-type solver. Addition-
 572 ally, the F-1 algorithm requires that the Jacobian matrix of the problem, $\nabla_{\mathbf{x}} \mathbf{F}(\mathbf{x}, \mathbf{p})$,
 573 can be created, stored, and factored. Requiring minimal input from the user, the
 574 F-1 algorithm performs significantly better than other state-of-the-art differentiation
 575 algorithms.

576 The F-1 algorithm relies on existing numerical differentiation schemes that are
 577 often-overlooked, even in advanced scientific applications. These techniques are based
 578 on the concepts of dual numbers and hyperdual numbers, which allow numerical
 579 differentiation of first and second derivatives, respectively, with machine-precision
 580 accuracy (see, e.g., [12, 14, 13, 37]). In addition to providing increased accuracy,
 581 using dual and hyperdual numbers is essential for the F-1 algorithm to be both fast
 582 and easy to implement.

583 While it builds on existing autodifferentiation tools [12, 14, 13, 37] and concepts
 584 [49], the F-1 algorithm elegantly combines them to leverage analytical shortcuts that
 585 we derive in this study. These shortcuts eliminate expensive calculations that are
 586 unavoidable when differentiating a black-box steady-state solver. In particular, they
 587 avoid redundant factorizations of the Jacobian. Specifically, the F-1 algorithm com-
 588 putes both gradient and Hessian in a single factorization, $\mathcal{O}(m)$ forward and back
 589 substitutions, and $\mathcal{O}(m^2)$ inexpensive function evaluations. Because factorizations
 590 are typically computationally expensive, the single-factorization feature of the F-1
 591 algorithm affords large computational savings.

592 Naturally, the computational costs of computing gradients and Hessians depend
593 on multiple variables, including the structure of the model itself. While we do not
594 make any definitive statement on how the time complexity depends on this structure,
595 it is reasonable to assume that it predominantly depends on the sparsity pattern of
596 the Jacobian given by $\nabla_{\mathbf{x}}\mathbf{F}$. Furthermore, it is also reasonable to assume that com-
597 putational costs scale as a function of the number of state variables and the number
598 of parameters, i.e., as a function of n and m . However, the algorithms benchmarked
599 in this study mainly differ by the number of factorizations they require, which is
600 essentially a function of m only. Thus, throughout, we qualitatively estimated the
601 expected computational costs as a function of m to explain the effective differences
602 in performance and to extrapolate our performance estimates to problems of different
603 sizes.

604 From our experience, the computational costs depend primarily on the number
605 of factorizations, secondarily on the number of forward and back substitutions, and
606 tertiarily on the number of function evaluations. Because the FD1, CSD, and DUAL
607 algorithms apply a numerical differentiation scheme to the gradient, they require
608 $\mathcal{O}(m)$ factorizations, forward and back substitutions, and function evaluations. Being
609 applied to the objective function, the HYPER and FD2 algorithms are more expen-
610 sive, requiring $\mathcal{O}(m^2)$ factorizations, forward and back substitutions, and function
611 evaluations. Furthermore, the computational costs of each of these algorithms likely
612 scale with the size of the state, n . Based on these considerations alone, the F-1 al-
613 gorithm, which already outperforms the other algorithms by a large margin (in the
614 context of our global marine phosphorus-cycling model), should outperform the other
615 algorithms even more for larger problems, i.e., for larger m and n .

616 We demonstrated the computational performance of the F-1 algorithm by bench-
617 marking it against five other state-of-the-art numerical differentiation algorithms (the
618 DUAL, FD1, CSD, HYPER, and FD2 algorithms) in the context of optimizing a global
619 marine phosphorus-cycling model, embedded in a global steady-state data-assimilated
620 ocean circulation. The performance of the F-1 algorithm benefits from being used in
621 the context of an optimization, during which previous computations, such as that
622 of the factors of the Jacobian matrix, \mathbf{A} , can be reused. Overall, optimizing these
623 $m = 6$ parameters takes about 7 minutes with the F-1 algorithm, against 28 to 47
624 minutes using the DUAL, CSD, or FD1 algorithms, about 2 hours using the HYPER
625 algorithm, and almost 3 hours using the FD2 algorithm.

626 We further investigated performance by recording the time spent computing the
627 gradient and Hessian. In particular, we found that for computing the Hessian only,
628 the F-1 algorithm is about 16 to 32 times faster than the DUAL, CSD, and FD1 algo-
629 rithms, which have access to the analytical gradient. Furthermore, the F-1 algorithm
630 is 76 and 100 times faster than the HYPER and FD2 algorithms, respectively, for
631 only $m = 6$ parameters.

632 We extrapolated our results to different problem sizes (different m) based on
633 qualitative algorithmic complexity arguments. Assuming computation times scale
634 primarily with the number of factorizations, we find that the performance of the F-1
635 algorithm would be amplified for larger problems. In particular, we expect the F-1
636 algorithm to be about 2 orders of magnitude faster than the DUAL, CSD, and FD1
637 algorithms with $m = 25$ parameters, and 3 orders of magnitude faster with $m = 250$.
638 Furthermore, we expect the F-1 algorithm to be about 3 and 5 orders of magnitude
639 faster than the HYPER and FD2 algorithms with $m = 20$ and $m = 200$, respectively.

640 In summary, we have presented an optimally efficient algorithm for computing the
641 gradient and Hessian of the objective function for problems defined implicitly by the

642 steady-state solution of a system of discretized nonlinear PDEs. Our algorithm out-
643 performs other state-of-the art algorithms for numerical differentiation, spectacularly
644 so in the context of optimization. This is because classical numerical differentia-
645 tion methods invoke nested iterative algorithms as black boxes, effectively repeating
646 redundant computations, which incur significant computational costs. Instead, our
647 algorithm leverages analytical shortcuts that are not available with naive black-box
648 approaches. The performance gains likely scale with the size, n , of the system of
649 PDEs, and with the number, m , of parameters, such that larger models, e.g., with
650 finer resolution or more detailed mechanisms, would benefit even more from our al-
651 gorithm than the benchmarks presented in this study.

652 The F-1 algorithm is ideally suited to a number of geoscientific model optimiza-
653 tions, provided the models can be represented by a steady-state PDE system of which
654 the Jacobian can be stored and factored. However, the F-1 algorithm could potentially
655 be extended to a larger scope of problems: (i) For very large m , where the optimizer
656 does not create the full Hessian matrix but instead uses a matrix-free approach (i.e.,
657 only evaluates matrix–vector products). This is the case when one wishes to optimize
658 a 3D field with a large number of entries, rather than a few scalar parameters. (ii) For
659 very large n , where the inner solver similarly does not create the Jacobian matrix \mathbf{A} ,
660 but only evaluates matrix–vector products, e.g., using a Newton–Krylov type of solver
661 (as in, e.g., [26, 49]). Other avenues of research include exploring potentially faster
662 strategies for constrained optimization problems, for which the solver is not nested
663 inside the optimizer, allowing for updates of the state, \mathbf{x} , outside of the manifold
664 of steady-state solutions (i.e., not satisfying the steady-state condition at every update
665 of the parameters, \mathbf{p}) [38]. Exploring the potential generalization of the F-1 algorithm
666 to non-steady problems, as in [49] is also a promising research direction. Finally, the
667 question remains whether the F-1 algorithm is applicable to problems that can lever-
668 age a distributed structure. In the case where the state vector can be separated into
669 chunks that the solver can update in parallel, the computation times would be set
670 by the sizes of the submatrices of the Jacobian now separated into smaller blocks. A
671 parallel solver and F-1 algorithm could be combined, e.g., in the offline optimization
672 of much larger models than the phosphorus-cycling model presented here [2]. Such
673 models are common in global marine biogeochemistry, e.g., in order to simulate a
674 large number of marine tracers, like the Biogeochemical Elemental Cycling (BEC)
675 model [33, 35, 34].

676 **6. Discussion.** Although our approach applies to a vast range of steady-state
677 models defined through the implicit solution of a discretized system of nonlinear PDEs,
678 it does not apply to all optimization problems of that form. Specifically, we focused
679 on the cases where one is interested in computing the Hessian, with a particular
680 focus on optimization algorithms that use quasi-Newton’s methods (see, e.g., [53]).
681 However, we should point out that there are different algorithms that can be used
682 to minimize \hat{f} that do not require the Hessian matrix. For example, some require
683 only evaluations of the objective function, like the Simulated Annealing (e.g., [21])
684 and Nelder-Mead [36] algorithms. Others, like the Broyden-Fletcher-Goldfarb-Shanno
685 (e.g., [38]), the gradient descent [4], and the conjugate gradient (e.g., [17]) algorithms,
686 require evaluations of the gradient. It might be the case that the problem at hand is
687 not suitable for a Newton-like method for the optimization algorithm, in which case
688 the F-1 algorithm would not be needed.

689 Because it does not invoke the inner solver to compute derivatives, the F-1 algo-
690 rithm avoids a number of implementation pitfalls that all other available algorithms

691 fall into. First, the F-1 algorithm allows for inner solvers that use non-real operations.
 692 For example, if the inner solver used the complex-step algorithm to compute
 693 the Jacobian, $\mathbf{A} = \nabla_{\mathbf{x}}\mathbf{F}(\mathbf{x}, \mathbf{p})$, then the CSD algorithm could not be naively applied
 694 to compute the Hessian matrix because of perturbation confusion. Similarly, if the
 695 inner solver used the dual-step algorithm to compute the Jacobian, then the DUAL
 696 and HYPER algorithms would not work either. We note that it is possible *in theory*
 697 to carefully chose the size of the complex, dual, or hyperdual steps so that the CSD,
 698 DUAL, or HYPER algorithms work with an inner solver that uses non-real operations
 699 internally. However, we do not recommend such an approach because it comes at the
 700 risk of failing silently. We also note that perturbation confusion can be avoided by
 701 carerul autodifferentiation implementations [48, 47].

702 Second, the F-1 algorithm avoids having to carefully chose the step sizes. This is
 703 particularly important in the case of finite-difference methods (e.g., FD1 and FD2),
 704 for which if the step size h is too small, the inner-solver loop may not execute a single
 705 iteration, potentially causing large errors. In our implementation, an optimal choice
 706 for h was about $[\mathbf{p}]_j/10^4$ for the FD1 algorithm, and $[\mathbf{p}]_j/10^2$ for the non-diagonal
 707 terms of the Hessian for the FD2 algorithm, as can be seen in the code, accessible
 708 from the URLs in Table 1. ($[\mathbf{p}]_j$ denotes the j th optimizable parameter value.) Such
 709 a large relative step is likely the reason for the slower convergence (in terms of number
 710 of optimizer iterations) of the FD2 algorithm.

711 Third, the F-1 algorithm avoids having to carefully chose the tolerances of the
 712 iterative algorithms. In the case of the DUAL, CSD, and HYPER algorithms, an
 713 additional tolerance for each non-real part must be added to the inner solver, as
 714 detailed in section SM6, and the choice of said tolerance will matter. Although one
 715 may get away with forgetting to set the non-real tolerances, in this case the choice of
 716 the tolerance on the real part will determine when the inner loop terminates and may
 717 cause large errors.

718 Fourth, the F-1 algorithm does not need an inner solver that can handle complex,
 719 dual, or hyper-valued parameter or state inputs. In contrast, this is the case for
 720 the CSD, DUAL, and HYPER algorithms. In order for these algorithms to work,
 721 the inner-solver Newton steps go through solving complex-, dual-, and hyperdual-
 722 valued linear systems. As mentioned in subsection 3.1.2, the calls to underlying
 723 UMFPACK allow for complex-valued linear systems, but fails on dual- and hyperdual-
 724 valued systems. These failures compelled us to develop the `DualMatrixTools` and
 725 `HyperDualMatrixTools` packages specifically designed to overcome this likely common
 726 shortcoming: there is no guarantee that dual or hyperdual-valued numbers will be
 727 handled correctly by underlying package that solves linear systems.

728 Naturally, one should always avoid repeating expensive computations. This is
 729 what is accomplished by the memory cache in the implementation of the F-1 algo-
 730 rithm, which stores, e.g., the factors of \mathbf{A} , for multiple subsequent forward and back
 731 substitutions. While the other algorithms do not store all the information that the
 732 F-1 algorithm stores, they still keep the real-valued steady-state solution, $\mathbf{s}(\mathbf{p})$, in
 733 memory. A common strategy in computer sciences that would benefit the other algo-
 734 rithms is memoization of the factorization function, such that the factors of \mathbf{A} would
 735 be computed only once, just like for the F-1 algorithm. However, while this seems
 736 like a good strategy at face value, it turns out that there is numerical noise as the
 737 state, \mathbf{x}_l , gets close to the theoretical steady-state solution, $\mathbf{s}(\mathbf{p})$. That is, the fac-
 738 tors of \mathbf{A} would be updated with a high probability at every update of either the
 739 state or the parameters, regardless of how small the change is. Additionally, storing
 740 a large number of factors of a large sparse matrix would likely cause memory issues.

741 In comparison, by leveraging exact analytical shortcuts, the F-1 algorithm provides
 742 a finely-tuned storage-and-reuse approach that avoids redundant computations in an
 743 optimal way.

744 **Acknowledgments.** FP and BP acknowledge funding from DOE grant DE-
 745 SC0016539 and NSF grant 1658380. The computations were performed using the
 746 Linux computational cluster Katana supported by the Faculty of Science, University
 747 of New South Wales, Australia.

748

REFERENCES

- 749 [1] D. E. ARCHER, G. ESHEL, A. WINGUTH, W. BROECKER, R. PIERREHUMBERT, M. TOBIS, AND
 750 R. JACOB, *Atmospheric pCO₂ sensitivity to the biological pump in the ocean*, *Global Bio-*
 751 *geochem. Cycles*, 14 (2000), pp. 1219–1230, <https://doi.org/10.1029/1999GB001216>.
- 752 [2] A. BARDIN, F. PRIMEAU, AND K. LINDSAY, *An offline implicit solver for simulating prebomb*
 753 *radiocarbon*, *Ocean Modelling*, 73 (2014), pp. 45–58, [https://doi.org/10.1016/j.ocemod.](https://doi.org/10.1016/j.ocemod.2013.09.008)
 754 [2013.09.008](https://doi.org/10.1016/j.ocemod.2013.09.008).
- 755 [3] J. BEZANSON, A. EDELMAN, S. KARPINSKI, AND V. SHAH, *Julia: A fresh approach to numerical*
 756 *computing*, *SIAM Review*, 59 (2017), pp. 65–98, <https://doi.org/10.1137/141000671>.
- 757 [4] A. CAUCHY, *Méthode générale pour la résolution des systemes d'équations simultanées*, *Comp.*
 758 *Rend. Sci. Paris*, 25 (1847), pp. 536–538.
- 759 [5] T. A. DAVIS, *Algorithm 832: UMFPACK v4.3 — an unsymmetric-pattern multifrontal method*,
 760 *ACM Trans. Math. Softw.*, 30 (2004), pp. 196–199, <https://doi.org/10.1145/992200.992206>.
- 761 [6] T. A. DAVIS, *Direct methods for sparse linear systems*, vol. 2, SIAM, 2006.
- 762 [7] T. DEVRIES, *The oceanic anthropogenic CO₂ sink: Storage, air-sea fluxes, and transports over*
 763 *the industrial era*, *Global Biogeochem. Cycles*, 28 (2014), pp. 631–647, [https://doi.org/10.](https://doi.org/10.1002/2013GB004739)
 764 [1002/2013GB004739](https://doi.org/10.1002/2013GB004739).
- 765 [8] T. DEVRIES, C. DEUTSCH, P. A. RAFTER, AND F. PRIMEAU, *Marine denitrification rates*
 766 *determined from a global 3-D inverse model*, *Biogeosciences*, 10 (2013), pp. 2481–2496,
 767 <https://doi.org/10.5194/bg-10-2481-2013>.
- 768 [9] T. DEVRIES AND F. PRIMEAU, *Dynamically and observationally constrained estimates of water-*
 769 *mass distributions and ages in the global ocean*, *J. Phys. Oceanogr.*, 41 (2011), pp. 2381–
 770 2401, <https://doi.org/10.1175/JPO-D-10-05011.1>.
- 771 [10] T. DEVRIES AND T. WEBER, *The export and fate of organic matter in the ocean: New con-*
 772 *straints from combining satellite and oceanographic tracer observations*, *Global Biogeo-*
 773 *chemical Cycles*, 31 (2017), pp. 535–555, <https://doi.org/10.1002/2016GB005551>.
- 774 [11] H. A. DIJKSTRA AND W. WEIJER, *Stability of the global ocean circulation: Basic bifurcation*
 775 *diagrams*, *Journal of Physical Oceanography*, 35 (2005), pp. 933–948, [https://doi.org/10.](https://doi.org/10.1175/JPO2726.1)
 776 [1175/JPO2726.1](https://doi.org/10.1175/JPO2726.1).
- 777 [12] J. FIKE AND J. ALONSO, *The development of hyper-dual numbers for exact second-derivative*
 778 *calculations*, in 49th AIAA Aerospace Sciences Meeting including the New Horizons Forum
 779 and Aerospace Exposition, 2011, p. 886.
- 780 [13] J. A. FIKE AND J. J. ALONSO, *Automatic differentiation through the use of hyper-dual num-*
 781 *bers for second derivatives*, in *Recent Advances in Algorithmic Differentiation*, S. Forth,
 782 P. Hovland, E. Phipps, J. Utke, A. Walther, T. J. Barth, M. Griebel, D. E. Keyes,
 783 R. M. Nieminen, D. Roose, and T. Schlick, eds., vol. 87 of *Lecture Notes in Computa-*
 784 *tional Science and Engineering*, Springer Berlin Heidelberg, 2012, pp. 163–173, https://doi.org/10.1007/978-3-642-30023-3_15.
- 785 [14] J. A. FIKE, S. JONGSMA, J. J. ALONSO, AND E. VAN DER WEIDE, *Optimization with gradient*
 786 *and Hessian information calculated using hyper-dual numbers*, in AIAA paper 2011-3807,
 787 29th AIAA Applied Aerodynamics Conference, 2011.
- 788 [15] M. FRANTS, M. HOLZER, T. DEVRIES, AND R. MATEAR, *Constraints on the global marine iron*
 789 *cycle from a simple inverse model*, *Journal of Geophysical Research: Biogeosciences*, 121
 790 (2016), pp. 28–51, <https://doi.org/10.1002/2015jg003111>.
- 791 [16] H. E. GARCIA, K. WEATHERS, C. R. PAVER, I. SMOLYAR, T. P. BOYER, R. A. LOCARNINI,
 792 M. M. ZWENG, A. V. MISHONOV, O. K. BARANOVA, D. SEIDOV, AND J. R. REAGAN,
 793 *World Ocean Atlas 2018*, Volume 4: Dissolved Inorganic Nutrients (phosphate, nitrate
 794 and nitrate+nitrite, silicate) (in preparation). A. Mishonov Technical Ed.

- 796 [17] W. W. HAGER AND H. ZHANG, *Algorithm 851: CG_DESCENT, a Conjugate Gradient Method*
797 *with Guaranteed Descent*, ACM Trans. Math. Softw., 32 (2006), pp. 113–137, [https://doi.](https://doi.org/10.1145/1132973.1132979)
798 [org/10.1145/1132973.1132979](https://doi.org/10.1145/1132973.1132979).
- 799 [18] C. T. KELLEY, *Iterative Methods for Optimization*, Frontiers in Applied Mathematics, Society
800 for Industrial and Applied Mathematics (SIAM, 3600 Market Street, Floor 6, Philadelphia,
801 PA 19104), 1999.
- 802 [19] C. T. KELLEY, *Solving Nonlinear Equations with Newton's Method*, SIAM, 2003, ch. 1. Intro-
803 duction, pp. 1–25, <https://doi.org/10.1137/1.9780898718898.ch1>.
- 804 [20] C. T. KELLEY, *Solving Nonlinear Equations with Newton's Method*, SIAM, 2003, ch. 2. Find-
805 ing the Newton Step with Gaussian Elimination, pp. 27–55, [https://doi.org/10.1137/1.](https://doi.org/10.1137/1.9780898718898.ch2)
806 [9780898718898.ch2](https://doi.org/10.1137/1.9780898718898.ch2).
- 807 [21] S. KIRKPATRICK, C. D. GELATT, AND M. P. VECCHI, *Optimization by simulated annealing*,
808 Science, 220 (1983), pp. 671–680, <https://doi.org/10.1126/science.220.4598.671>.
- 809 [22] O. KOLDITZ AND H.-J. DIERSCH, *Quasi-steady-state strategy for numerical simulation of*
810 *geothermal circulation in hot dry rock fractures*, International Journal of Non-Linear Me-
811 chanics, 28 (1993), pp. 467–481, [https://doi.org/10.1016/0020-7462\(93\)90020-L](https://doi.org/10.1016/0020-7462(93)90020-L).
- 812 [23] I. KRIEST AND A. OSCHLIES, *On the treatment of particulate organic matter sinking in large-*
813 *scale models of marine biogeochemical cycles*, Biogeosciences, 5 (2008), pp. 55–72, [https:](https://doi.org/10.5194/bg-5-55-2008)
814 [//doi.org/10.5194/bg-5-55-2008](https://doi.org/10.5194/bg-5-55-2008).
- 815 [24] E. Y. KWON AND F. PRIMEAU, *Sensitivity and optimization study of a biogeochemistry ocean*
816 *model using an implicit solver and in-situ phosphate data*, Global Biogeochem. Cycles, 20
817 (2006), GB4009, <https://doi.org/10.1029/2005GB002631>.
- 818 [25] E. Y. KWON AND F. PRIMEAU, *Optimization and sensitivity of a global biogeochemistry ocean*
819 *model using combined in situ DIC, alkalinity, and phosphate data*, Journal of Geophysical
820 Research: Oceans, 113 (2008), <https://doi.org/10.1029/2007JC004520>. C08011.
- 821 [26] X. LI AND F. W. PRIMEAU, *A fast Newton–Krylov solver for seasonally varying global ocean*
822 *biogeochemistry models*, Ocean Modelling, 23 (2008), pp. 13–20, [https://doi.org/10.1016/](https://doi.org/10.1016/j.ocemod.2008.03.001)
823 [j.ocemod.2008.03.001](https://doi.org/10.1016/j.ocemod.2008.03.001).
- 824 [27] J. H. MANTON, *Differential Calculus, Tensor Products and the Importance of Notation*, ArXiv
825 e-prints, (2012), <https://arxiv.org/abs/1208.0197v2>, <https://arxiv.org/abs/1208.0197>.
- 826 [28] J. W. MARTIN, G. A. KNAUER, D. M. KARL, AND W. W. BROENKOW, *VERTEX: Carbon*
827 *cycling in the NE Pacific*, Deep-Sea Research, 34 (1987), pp. 267–285.
- 828 [29] J. R. A. MARTINS, P. STURDZA, AND J. J. ALONSO, *The complex-step derivative approxima-*
829 *tion*, ACM Trans. Math. Softw., 29 (2003), pp. 245–262, [https://doi.org/10.1145/838250.](https://doi.org/10.1145/838250.838251)
830 [838251](https://doi.org/10.1145/838250.838251).
- 831 [30] P. K. MOGENSEN AND A. N. RISETH, *Optim: A mathematical optimization package for Julia*,
832 Journal of Open Source Software, 3 (2018), p. 615, <https://doi.org/10.21105/joss.00615>.
- 833 [31] P. K. MOGENSEN, J. M. WHITE, A. N. RISETH, T. HOLY, M. LUBIN, C. STOCKER, A. LEVITT,
834 C. ORTNER, B. JOHNSON, A. NOACK, Y. YU, K. CARLSSON, D. LIN, T. R. COVERT,
835 R. ROCK, J. REGIER, B. KUHN, A. WILLIAMS, RYAN, K. SATO, D. SMITH, R. ANANTHARA-
836 MAN, M. GOMEZ, J. REVELS, I. DUNNING, D. MACMILLEN, C. RACKAUCKAS, B. LEGAT,
837 AND A. STUKALOV, *JuliaNLSolvers/Optim.jl: Bugfix for Fminbox, better docs, and exper-*
838 *imental maximize function*, Sept. 2018, <https://doi.org/10.5281/zenodo.1412092>.
- 839 [32] J. MONOD, *Microbiologie: Recherches sur la croissance des cultures bactériennes. I*, Actualités
840 scientifiques et industrielles, Hermann & cie, 1942.
- 841 [33] J. MOORE, S. C. DONEY, J. A. KLEYPAS, D. M. GLOVER, AND I. Y. FUNG, *An intermedi-*
842 *ate complexity marine ecosystem model for the global domain*, Deep Sea Research Part
843 II: Topical Studies in Oceanography, 49 (2001), pp. 403 – 462, [https://doi.org/10.1016/](https://doi.org/10.1016/S0967-0645(01)00108-4)
844 [S0967-0645\(01\)00108-4](https://doi.org/10.1016/S0967-0645(01)00108-4). The US JGOFs Synthesis and Modeling Project: Phase 1.
- 845 [34] J. K. MOORE AND O. BRAUCHER, *Sedimentary and mineral dust sources of dissolved iron*
846 *to the world ocean*, Biogeosciences, 5 (2008), pp. 631–656, [https://doi.org/10.5194/](https://doi.org/10.5194/bg-5-631-2008)
847 [bg-5-631-2008](https://doi.org/10.5194/bg-5-631-2008).
- 848 [35] J. K. MOORE, S. C. DONEY, AND K. LINDSAY, *Upper ocean ecosystem dynamics and iron cycling*
849 *in a global three-dimensional model*, Global Biogeochem. Cycles, 18 (2004), GB4028, [https:](https://doi.org/10.1029/2004GB002220)
850 [//doi.org/10.1029/2004GB002220](https://doi.org/10.1029/2004GB002220).
- 851 [36] J. A. NELDER AND R. MEAD, *A simplex method for function minimization*, The Computer
852 Journal, 7 (1965), pp. 308–313, <https://doi.org/10.1093/comjnl/7.4.308>.
- 853 [37] M. NEUENHOFEN, *Review of theory and implementation of hyper-dual numbers for first and*
854 *second order automatic differentiation*, CoRR, abs/1801.03614 (2018), [http://arxiv.org/](http://arxiv.org/abs/1801.03614)
855 [abs/1801.03614](http://arxiv.org/abs/1801.03614), <https://arxiv.org/abs/1801.03614>.
- 856 [38] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media,
857 2006.

- 858 [39] B. PASQUIER, *DualMatrixTools: A Julia package for solving dual-valued linear systems*, Nov.
859 2018, <https://doi.org/10.5281/zenodo.1493571>.
- 860 [40] B. PASQUIER, *HyperDualMatrixTools: A Julia package for solving hyperdual-valued linear sys-*
861 *tems*, Nov. 2018, <https://doi.org/10.5281/zenodo.1670841>.
- 862 [41] B. PASQUIER, *AIBECs.jl: The ideal tool for exploring global marine biogeochemical cycles*,
863 May 2019, <https://doi.org/10.5281/zenodo.2864051>.
- 864 [42] B. PASQUIER, *F1Method.jl: A julia package for computing the gradient and Hessian of an*
865 *objective function defined implicitly by the solution to a steady-state problem*, May 2019,
866 <https://doi.org/10.5281/zenodo.2667835>.
- 867 [43] B. PASQUIER, *Worldoceanatlastools.jl: Download and process world ocean data in julia*, May
868 2019, <https://doi.org/10.5281/zenodo.2677666>.
- 869 [44] B. PASQUIER AND M. HOLZER, *Inverse-model estimates of the ocean's coupled phosphorus,*
870 *silicon, and iron cycles*, Biogeosciences, 14 (2017), pp. 4125–4159, [https://doi.org/10.5194/](https://doi.org/10.5194/bg-14-4125-2017)
871 [bg-14-4125-2017](https://doi.org/10.5194/bg-14-4125-2017).
- 872 [45] C. RACKAUCKAS AND P. G. NIE, *Differentialequations.jl – a performant and feature-rich ecosystem*
873 *for solving differential equations in julia*, Journal of Open Research Software, 5 (2017),
874 <https://doi.org/10.5334/jors.151>.
- 875 [46] J. A. RAVEN AND P. G. FALKOWSKI, *Oceanic sinks for atmospheric CO₂*, Plant, Cell & Envi-
876 ronment, 22 (1999), pp. 741–755, <https://doi.org/10.1046/j.1365-3040.1999.00419.x>.
- 877 [47] J. REVELS, T. BESARD, M. J. INNES, R. DEITS, M. PHIBELEHT, M. SCHAUER, L. WHITE, AND
878 K. FISCHER, *Cassette.jl*, Jan. 2019, <https://doi.org/10.5281/zenodo.2549715>.
- 879 [48] J. REVELS, M. LUBIN, AND T. PAPAMARKOU, *Forward-mode automatic differentiation in julia*,
880 arXiv:1607.07892 [cs.MS], (2016), <https://arxiv.org/abs/1607.07892>.
- 881 [49] M. P. RUMPFKEIL AND D. J. MAVRIPLIS, *Efficient hessian calculations using automatic differen-*
882 *tiation and the adjoint method with applications*, AIAA Journal, 48 (2010), pp. 2406–2417,
883 <https://doi.org/10.2514/1.J050451>.
- 884 [50] G. A. SCHMIDT AND L. A. MYSAK, *Can increased poleward oceanic heat flux explain the warm*
885 *cretaceous climate?*, Paleoceanography, 11 (1996), pp. 579–593, [https://doi.org/10.1029/](https://doi.org/10.1029/96PA01851)
886 [96PA01851](https://doi.org/10.1029/96PA01851).
- 887 [51] D. SIVIA AND J. SKILLING, *Data analysis: a Bayesian tutorial*, OUP Oxford, 2006.
- 888 [52] Y. TENG, F. W. PRIMEAU, J. K. MOORE, M. W. LOMAS, AND A. MARTINY, *Global-scale*
889 *variations of the ratios of carbon to phosphorus in exported marine organic matter*, Nature
890 Geosci., 7 (2014), pp. 895–898, <https://doi.org/10.1038/NGEO2303>.
- 891 [53] G. VENTER, *Review of optimization techniques*, Encyclopedia of aerospace engineering, (2010).
- 892 [54] T. VOLK AND M. HOFFERT, *Ocean carbon pumps: analysis of relative strengths and efficienc-*
893 *ies in ocean-driven atmospheric CO₂ changes.*, American Geophysical Union; Geophysical
894 Monograph 32, 1985, pp. 99–110.
- 895 [55] W.-L. WANG, C. LEE, AND F. W. PRIMEAU, *A bayesian statistical approach to inferring particle*
896 *dynamics from in-situ pump poc and chloropigment data from the mediterranean sea*,
897 Marine Chemistry, (2019), <https://doi.org/10.1016/j.marchem.2019.04.006>.
- 898 [56] W. YU AND M. BLAIR, *DNAD, a simple tool for automatic differentiation of Fortran codes*
899 *using dual numbers*, Computer Physics Communications, 184 (2013), pp. 1446–1452, <https://doi.org/10.1016/j.cpc.2012.12.025>.
- 900