

Differentiable Programming for Differential Equations: A Review

Facundo Sapienza^{*1}, Jordi Bolibar², Frank Schäfer³, Patrick Heimbach^{6,7}, Giles Hooker⁴,
Fernando Pérez¹, Per-Olof Persson⁵, Christopher Rackauckas^{8,9}, and Victor Boussange¹⁰

¹*Department of Statistics, University of California, Berkeley (USA)*

²*TU Delft, Department of Geosciences and Civil Engineering, Delft (Netherlands)*

³*CSAIL, Massachusetts Institute of Technology, Cambridge (USA)*

⁴*Department of Statistics and Data Science, University of Pennsylvania (USA)*

⁵*Department of Mathematics, University of California, Berkeley (USA)*

⁶*Oden Institute for Computational Engineering and Sciences, University of Texas at Austin (USA)*

⁷*Jackson School of Geosciences, University of Texas at Austin (USA)*

⁸*Massachusetts Institute of Technology, Cambridge (USA)*

⁹*JuliaHub, Cambridge (USA)*

¹⁰*Swiss Federal Research Institute WSL, Birmensdorf (Switzerland)*

March 18, 2024

^{*}Corresponding author: fsapienza@berkeley.edu

Abstract

The differentiable programming paradigm has become a central component of modern machine learning techniques. A long tradition of this paradigm exists in the context of scientific computing, in particular in differential equation-constrained, gradient-based optimization. The recognition of the strong conceptual synergies between inverse methods and machine learning offers the opportunity to lay out a coherent framework applicable to both fields. For models described by differential equations, the calculation of sensitivities and gradients requires careful algebraic and numeric manipulations of the underlying dynamical system. Here, we provide a comprehensive review of existing techniques to compute gradients of numerical solutions of differential equation systems. We first discuss the importance of gradients of solutions of ODEs in a variety of scientific domains, covering computational fluid dynamics, electromagnetism, geosciences, meteorology, oceanography, climate science, flux inversion, glaciology, solid earth geophysics, biology and ecology, and quantum physics. Second, we lay out the mathematical foundations of the various approaches and compare them with each other. Finally, we delve into the computational considerations and explore the solutions available in modern scientific software.

To the community, by the community. *This manuscript was conceived with the goal of shortening the gap between developers and practitioners of differentiable programming applied to modern scientific machine learning. With the advent of new tools and new software, it is important to create pedagogical content that allows the broader community to understand and integrate these methods into their workflows. We hope this encourages new people to be an active part of the ecosystem, by using and developing open-source tools. This work was done under the premise **open-science from scratch**, meaning all the contents of this work, both code and text, have been in the open from the beginning and that any interested person can contribute to the project. You can contribute directly to the GitHub repository github.com/ODINN-SciML/DiffEqSensitivity-Review.*

Contents

1	Introduction	5
2	Scientific motivation: A domain science perspective	7
2.1	Domain-specific applications	7
2.1.1	Computational physics and optimal design	8
2.1.1.1	Quantum physics	8
2.1.2	Geosciences	8
2.1.2.1	Meteorology	9
2.1.2.2	Oceanography	9
2.1.2.3	Climate science	10
2.1.2.4	Glaciology	10
2.1.3	Biology and ecology	10
3	Methods: A mathematical perspective	11
3.1	Preliminaries	12
3.1.1	Numerical solvers for differential equations	13
3.1.2	What to differentiate and why?	14
3.1.3	Optimization	15
3.1.4	Sensitivity matrix	15
3.2	Finite differences	16
3.3	Automatic differentiation	16
3.3.1	Forward mode	17
3.3.1.1	Dual numbers	17
3.3.1.2	Computational graph	17
3.3.2	Reverse mode	18
3.3.3	AD connection with JVPs and VJPs	18
3.4	Complex step differentiation	20
3.5	Symbolic differentiation	21
3.6	Sensitivity equations	21
3.7	Discrete adjoint method	22
3.7.1	Adjoint state equations	22
3.7.2	Simple linear system	23
3.8	Continuous adjoint method	25
3.9	Mathematical comparison of the methods	26
3.9.1	Forward AD and complex step differentiation	26
3.9.2	AD and symbolic differentiation with sparsity	26
3.9.3	Discrete adjoints and reverse AD	27
3.9.4	Consistency: forward AD and sensitivity equations	28
3.9.5	Consistency: discrete and continuous adjoints	29
4	Implementation: A computer science perspective	30
4.1	Direct methods	30
4.1.1	Finite differences	30
4.1.2	Automatic differentiation	31
4.1.2.1	Forward AD based on dual numbers	31
4.1.2.2	Reverse AD based on computational graph	32
4.1.2.3	Checkpointing	33
4.1.3	Complex step differentiation	33
4.2	Solver-based methods	34
4.2.1	Sensitivity equation	34
4.2.1.1	Computing VJPs inside the solver	34

4.2.2	Adjoint methods	35
4.2.2.1	Discrete adjoint method	35
4.2.2.2	Continuous adjoint method	35
4.2.2.3	Solving the quadrature	36
5	Recommendations	36
5.1	Chaotic systems	37
6	Conclusions	37
	Appendices	38
A	Lagrangian formulation of the adjoint method	38
A.1	Discrete adjoint	38
A.2	Continuous adjoint	38
A.3	The adjoint from a functional analysis perspective	39
B	When AD is algorithmically correct but numerically wrong	40
C	Supplementary code	42
	References	43

Plain language summary

Differential equations are mathematical tools that explicitly describe the processes and dynamics within various systems, based on prior knowledge. They are fundamental in many scientific disciplines for modeling phenomena such as physical processes, population dynamics, social interactions, and chemical reactions. By contrast, data-driven models do not necessarily require a detailed understanding of the underlying physical processes, and learn patterns and relationships directly from data. Data-driven models are particularly useful in scenarios where the underlying processes are poorly understood or too complex to be captured by traditional mathematical models. The combination of mechanistic models with data-driven models is becoming increasingly common in many scientific domains. To achieve this, these models need to leverage both domain knowledge and data, to have an accurate representation of the underlying dynamics. Being able to determine which model parameters are most influential and further compute derivatives of such a model is key to correctly assimilating and learning from data, but a myriad of sensitivity methods exist to do so. In this review, we present an overview of the different sensitivity methods that exist, providing (i) guidelines on the best use cases for different scientific domain problems, (ii) detailed mathematical analyses of their characteristics, and (iii) computational implementations on how to solve them efficiently.

1 Introduction

Evaluating how the value of a function changes with respect to its arguments and parameters plays a central role in optimization, sensitivity analysis, Bayesian inference, inverse methods, and uncertainty quantification, among many.¹⁹⁶ Modern machine learning applications require the use of gradients to efficiently exploit the high-dimensional space of parameters to be inferred or learned (e.g., the weights of a neural network). When optimizing an objective function, gradient-based methods (for example, gradient descent and its many variants²⁰⁰) are more efficient at finding a minimum and converge faster to them than gradient-free methods. When numerically computing the posterior of a probabilistic model, gradient-based sampling strategies are better at estimating the posterior distribution than gradient-free methods. Hessians further help to improve the convergence rates of these algorithms and enable uncertainty quantification around parameter values.³⁸ Furthermore, the *curse of dimensionality* renders gradient-free optimization and sampling methods computationally intractable for most large-scale problems.¹⁷⁶

A gradient serves as a compass in modern data science: it tells us in which direction in the vast, open ocean of parameters we should move towards in order to increase our chances of success.

Models based on differential equations arising in simulation-based science, which play a central role in describing the behaviour of systems in natural and social sciences, are not an exception to the rule.⁸¹ The solution of differential equations can be seen as functions that map parameter and initial conditions to state variables, similar to machine learning models. Some authors have recently suggested differentiable programming as the bridge between modern machine learning and traditional scientific models.^{77,190,193,214} Being able to compute gradients or sensitivities of dynamical systems opens the door to more complex data assimilation models that leverage in strong physical priors at the same time they offer flexibility to adapt to observations. This is very appealing in fields like computational physics, geophysics, and biology, to mention a few, where there is a broad literature on physical models and a long tradition in numerical methods. The first goal of this work is to introduce some of the applications of this emerging technology and to motivate its incorporation for the modelling of complex systems in the natural and social sciences.

Question 1. *What are the scientific applications of differentiable programming for dynamical systems?*

Sensitivity analysis corresponds to any method aiming to calculate how much the output of a function or program changes when we vary one of the function (or model) parameters. This task is performed in different ways by different communities when working with dynamical systems. In statistics, the sensitivity equations enable the computation of gradients of the likelihood of the model with respect to the parameters of the dynamical system, which can be later used for inference.¹⁹² In numerical analysis, sensitivities quantify how the solution of a differential equation fluctuates with respect to certain parameters. This is particularly useful

in optimal control theory,⁸⁶ where the goal is to find the optimal value of some control (e.g. the shape of a wing) that minimizes a given loss function. In recent years, there has been an increasing interest in designing machine learning workflows that include constraints in the form of differential equations. Examples of this include methods that numerically solve differential equations, such as physics-informed neural networks,¹⁹¹ as well as methods that augment and learn parts of the differential equation, such as universal differential equations,^{51,189} which also includes the case of neural ordinary differential equations⁴⁵ and neural stochastic differential equations.¹⁴¹

However, when working with differential equations, the computation of gradients is not an easy task, both regarding the mathematical framework and software implementation involved. Except for a small set of particular cases, most differential equations require numerical methods to approximate their solution. This means that solutions cannot be directly differentiated and require special treatment to compute first or second-order derivatives. Furthermore, numerical solutions introduce approximation errors. These errors can be propagated and amplified during the computation of the gradient. Alternatively, there is a broad literature on numerical methods for solving differential equations.^{97,248} Although each method provides different guarantees and advantages depending on the use case, this means that the tools developed to compute gradients when using a solver need to be universal enough in order to be applied to all or at least to a large set of them. As coined by Uwe Naumann, *the automatic generation of optimal (in terms of robustness and efficiency) adjoint versions of large-scale simulation code is one of the great open challenges in the field of High-Performance Scientific Computing*.¹⁶⁹ The second goal of this article is to review different methods that exist to achieve this goal.

Question 2. *How to efficiently compute the gradient of a function that depends on the numerical solution of a differential equation?*

The broader set of tools known as automatic or algorithmic differentiation (AD) aims at computing derivatives by sequentially applying the chain to the sequence of unit operations that constitute a computer program.^{95,169} The premise is simple: every computer program is ultimately an algorithm described by a nested concatenation of elementary algebraic operations, such as addition and multiplication, that are individually easy to differentiate and their composition is easy to differentiate by using the chain rule.⁸² More broadly than AD, differentiable programming encapsulates the set of software tools that allows to compute efficient and robust gradients through complex algorithms, including numerical solvers.¹¹¹ Although many modern differentiation tools use AD to some extent, there is also a family of methods that compute the gradient by relying on an auxiliary set of differential equations and/or compute an intermediate adjoint. Furthermore, it is important to be aware that when using AD or any other technique we are differentiating the algorithm used to lead to the numerical solution, not the numerical solution itself, which can lead to wrong results.⁵⁷

The differences between methods to compute sensitivities arise both from their mathematical formulation and their computational implementation. The first provides different guarantees on the method returning the actual gradient or a good approximation thereof. The second involves how theory is translated to software, and what are the data structures and algorithms used to implement it. Different methods have different computational complexities depending on the total number of parameters and size of the differential equation system, and these complexities are also balanced between total execution time and required memory. The third goal of this work, then, is to illustrate the different strengths and weaknesses of these methods, and how to use them in modern scientific software.

Question 3. *What are the advantages and disadvantages of different differentiation methods and how can I incorporate them in my research?*

Differentiable programming is opening new ways of doing research across sciences. Arguably, its potential has so far been under-explored but is being rediscovered in the age of data-driven science. In order to realize its full potential, we need close collaboration between domain scientists, methodological scientists, computational scientists, and computer scientists in order to develop successful, scalable, practical, and efficient frameworks for real world applications. As we make progress in the use of these tools, new methodological questions start to emerge. How do these methods compare? How can they be improved? In this review we present a comprehensive list of the methods that exists in the intersection of differentiable programming and differential equation modelling.

This review paper is structured in three main sections, looking at differentiable programming for differential equations from three different perspectives: a domain science perspective (Section 2), a mathematical perspective (Section 3) and a computer science perspective (Section 4).

2 Scientific motivation: A domain science perspective

Mechanistic (or process-based) models have played a central role in a wide range of scientific disciplines. They consist of precise mathematical descriptions of physical mechanisms, that include the modelling of causal interactions, feedback loops and dependencies between components of the system under consideration.¹⁸⁹ These mathematical representation typically take the form of differential equations. Together with the numerical methods to approximate their solutions, differential equations led to fundamental advances in the understanding and prediction of physical and biological systems. Differential equation models depend on parameters that determine the processes represented. The parameter values have traditionally been estimated independently of the model, which poses several problems.⁹⁸ First, the independent estimation of parameters and processes rapidly becomes impossible as the number of state variables increases, especially when considering highly non-linear processes. Second, the measurement of certain parameters are intrinsically difficult.²¹¹ Third, parameter values estimated from laboratory experiments may be specific to the experimental setting and resulting simulations are likely to diverge from observations in the field.

Due to the difficulty of estimating parameter values in mechanistic models, and accompanying the massive growth of data upon which they depend, statistical (or machine learning) models have lead the modelling field in the past decades.⁴⁹ Advances in the field of machine learning, and particularly in deep learning,¹³⁶ allowed statistical models to learn at multiple levels of abstraction and capture extremely complex nonlinear patterns and information hidden in large datasets. However, the use of machine learning models for predictions has been heavily criticized, as they critically assume that patterns contained in observed data will repeat in the future, which may not be the case.^{16,54} In contrast to purely statistical models, the process knowledge embedded in the structure of mechanistic models renders them more robust for predicting dynamics under different conditions.

The fields of mechanistic modelling and statistical modelling have mostly evolved independently,²⁵⁸ due to several reasons. On the one hand, domain scientists have often been reluctant in learning about machine learning methods, judging them as opaque black boxes, unreliable, and not respecting domain-established knowledge.⁴⁸ On the other hand, the field of machine learning has mainly been developed around data-driven applications, without including any *a priori* physical knowledge. However, there has been an increasing interest in making mechanistic models more flexible, as well as introducing domain-specific or physical constraints and interpretability in machine learning models.^{2,29,35,42,50,71,80,163,195,198,201,212,231,254} Inverse modelling is an attempt to bridge the statistical and mechanistic modelling fields.^{199,252} Differentiable programming is key in this process.

2.1 Domain-specific applications

Arguably, the notion of differentiable programming has a long tradition in computational physics which is founded on solving and/or inverting models based on differential equation. The overarching goal of inverse modelling is to find a set of optimal model parameters that minimizes an objective or cost function quantifying the misfit between observations and the simulated state. Depending on the nature of the inversion, we may distinguish between the following cases.

- **Initial conditions.** Inverting for uncertain initial conditions, which, when integrated using the model, leads to an optimal match between the observations and the simulated state (or diagnostics thereof); variants thereof are used for optimal forecasting.
- **Boundary conditions.** Inverting for uncertain surface (e.g., interface fluxes), bottom (e.g., bed properties), or lateral (e.g., open boundaries of a limited domain) boundaries, which, when used in the model, produces an optimal match of the observations; variants thereof are used in tracer or boundary (air-sea) flux inversion problems, e.g., related to the global carbon cycle.

- **Model parameters.** Inverting for uncertain model parameters amounts to an optimal model calibration problem. As a *learning of optimal parameters from data* problem, it is the closest to machine learning applications. Parametrization is a special case of parameter inversion, where a parametric function (e.g., a neural network) is used to approximate processes.^{beucler2024, boussange2024, 195}
- **Model selection.** Estimating the model evidence for each candidate model to discriminate between competing hypotheses.⁴⁰ Involves the calculation of the marginal likelihood for each model, or the calibration of each model and evaluation of their score in out-of-samples predictions.

Besides the use of sensitivity methods for optimization, inversion, estimation, or learning, gradients have also proven powerful tools for computing comprehensive sensitivities of quantities of interest; computing optimal perturbations (in initial or boundary conditions) that lead to maximum, non-normal amplification of specific norms of interest; and characterizing and quantifying uncertainties by way of second derivative (Hessian) information.

In recent years the use of machine learning methods has become more popular in many scientific domains. Differential equations can be used to describe a large variety of dynamical systems, while data-driven regression models (e.g., neural networks, Gaussian processes, reduced-order models, basis expansions) have been demonstrated to act as universal approximators, learning any possible function if enough data is available.⁹² This combined flexibility can be exploited by many different domain-specific problems to tailor modelling needs to both dynamics and data characteristics.

There is a long tradition of differentiable programming methods applied across the scientific spectrum. Without aiming to have a complete and fully comprehensive overview of applications to all scientific and engineering applications, as well as a perfect classification of them, the following are a few selected examples where these techniques had been used by different communities.

2.1.1 Computational physics and optimal design

There is a long tradition of computational physics models based on adjoint methods and automatic differentiation pipelines. These include examples in

particle physics⁵³ or quantum chemistry¹³
optimal design in nanophotonics¹⁶² optimal design in electromagnetism⁷⁹

There is a long tradition of computational models based on sensitivity methods for optimal design and optimal control.^{11,143,184} This includes applications to stellarator coil design;¹⁵⁷ fluid dynamics;^{86,161} supersonic aircraft design,^{66,108} biology²²⁵

2.1.1.1 Quantum physics

Quantum optimal control has diverse applications spanning a broad spectrum of quantum systems. Optimal control methods have been used to optimize pulse sequences, enabling the design of high-fidelity quantum gates and the preparation of complex entangled quantum states. Typically, the objective is to maximize the fidelity to a target state or unitary operation, accompanied by additional constraints or costs specific to experimental demands. The predominant control algorithms are gradient-based optimization methods, such as gradient ascent pulse engineering (GRAPE), and rely on the computation of derivatives for solutions of the differential equations modeling the time evolution of the quantum system. In cases where the analytical calculation of a gradient is impractical, numerical evaluation using AD becomes a viable alternative.^{3,4,89,116,117,139,207} Specifically, AD streamlines the adjustment to diverse objectives or constraints, and its efficiency can be enhanced by employing custom derivative rules for the time propagation of quantum states as governed by solutions to the Schrödinger equation.⁸⁹ Moreover, sensitivity methods for differential equations facilitate the design of feedback control schemes necessitating the differentiation of solutions to stochastic differential equations.²⁰⁸

2.1.2 Geosciences

Many geoscientific phenomena are governed by global and local conservation laws along with a set of empirical constitutive laws and subgrid-scale parametrization schemes. Together, they enable efficient description

of the system’s spatio-temporal evolution in terms of a set of partial differential equations (PDEs). Example are geophysical fluid dynamics,²³⁷ describing geophysical properties of many Earth systems, such as the atmosphere, oceans, land surface, and glaciers. In such models, calibrating model parameters is extremely challenging, due to datasets being sparse in both space and time, heterogeneous, and noisy; and computational models involving high-dimensional parameter spaces, often on the order of $O(10^3) - O(10^8)$. Moreover, many existing mechanistic models can only partially describe observations, with many detailed physical processes being ignored or poorly parameterized.

In the following, we sketch how differentiable programming has been used in different disciplines of geosciences, and how new concepts are emerging of combining inverse modeling and machine learning approaches where differentiable programming provides a key computational enabling framework, something recently coined as scientific machine learning.

2.1.2.1 Meteorology

Numerical weather prediction (NWP) is among the most prominent fields where adjoint methods have played an important role.⁶⁰ Adjoint methods was introduced to infer initial conditions that minimize the misfit between simulations and weather observations,^{47,226} with the value of second-derivative information also being recognized.⁵² This led to the development of the so-called *four-dimensional variational* (4D-Var) data assimilation (DA) technique^{186,187} at the European Centre for Medium-Range Weather Forecasts (ECMWF) as one of the most advanced DA approaches, and which contributed substantially to the *quiet revolution* in NWP.¹⁹ Related, within the framework of transient non-normal amplification or optimal excitation,^{62,63} the adjoint method has been used extensively to infer patterns in initial conditions that over time contribute to maximum uncertainty growth in forecasts^{39,180} and to infer the so-called *Forecast Sensitivity-based Observation Impact* (FSOI).¹³⁰ Except in very few instances and for experimental purposes,⁸⁴ automatic differentiation has not been used in the development of adjoint models in NWP. Instead, the adjoint code was derived and implemented manually.

2.1.2.2 Oceanography

The recognition of the benefit of adjoint methods for use in data assimilation in the ocean coincided roughly with that in meteorology.^{227,229} The first application appeared soon thereafter in the context of a basin-scale general circulation model.^{233–235} An important detail is that their work already differed from the “4D-Var” problem of NWP in that sensitivities were computed not only with respect to initial conditions but also with respect to surface boundary conditions, i.e., air-sea fluxes of buoyancy and momentum. Again, the role of the second-derivative for uncertainty quantification was readily realized.²²⁸ Similar to the work on calculating singular vectors in the atmosphere based on tangent linear and adjoint versions of a GCM to solve a generalized eigenvalue problem, the question of El Niño predictability invited model-based singular vector computations in models of the Tropical Pacific Ocean.^{164,165} Such model-based singular vectors were also later computed for optimal excitations of the North Atlantic thermohaline circulation.^{255–257} Notably in the context of this review, the consortium for Estimating the Circulation and Climate of the Ocean (ECCO)²²¹ set out in around 1999 to develop a parameter and state estimation framework, whereby a state-of-the-art ocean general circulation model is fit to diverse observations by way of PDE-constrained, gradient-based optimization, with the adjoint model of the GCM computing the gradient. Importantly, the adjoint model of the MIT general circulation model (MITgcm) is generated using source-to-source automatic differentiation,^{104,153} initially using the *Tangent linear and Adjoint Model Compiler* (TAMC)⁸² and then its commercial successor *Transformation of Algorithms in Fortran* (TAF).⁸⁴ Rigorous exploitation of AD enabled the simulation framework to be significantly extended over time in terms of vastly improved model numerics⁶⁷ and coupling other Earth system components, including biogeochemistry,⁵⁶ sea-ice,¹⁰² and sub-ice shelf cavities.¹⁰⁵ Unlike NWP-type 4D-Var, the use of AD also enabled extension of the framework to the problem of parameter calibration from observations.^{65,144,222} Arguably, this work heralded much of today’s efforts in online learning of parameterization schemes, where the functional representation between the parameters and the learning data are provided by the numerical implementation of a PDF rather than by a neural network. The desire to make AD for Earth system models written in Fortran (to date the vast majority) has also spurred the development of alternative AD tools with powerful reverse modes, notably

OpenAD²³⁶ and most recently Tapenade.^{72,73,99} There is enormous potential to seamlessly integrate the inverse-modeling and machine-learning based approaches through the concept of differentiable programming.

2.1.2.3 Climate science

The same goals that have driven the use of sensitivity information in numerical weather prediction (optimal initial conditions for forecasts) or ocean science (state and parameter estimation) apply in the world of climate modeling. The recognition that good initial conditions (e.g., such that are closest to the real or observed system) will lead to improved forecasts on subseasonal, seasonal, interannual, or even decadal time scales has driven major community efforts.¹⁵⁸ However, there has been a lack so far in exploiting the use of gradient information to achieve optimal initialization for coupled Earth system models.⁷⁰ One conceptual challenge is the presence of multiple timescales in the coupled system and the utility of gradient information beyond many synoptic time scales in the atmosphere and ocean.^{133,135} Nevertheless, efforts are underway to enable adjoint-based parameter estimation of coupled atmosphere-ocean climate models, with AD again playing a crucial role in generating the corresponding adjoint model.^{25,148,220} Complementary, recognizing the power of differentiable programming, efforts are also targeting the development of *neural atmospheric general circulation models* in JAX, which combine a differentiable dynamical core with neural operators as surrogate models of unresolved physics.¹²⁹

2.1.2.4 Glaciology

Due to the difficulty of having direct observations of internal and basal rheological processes of glaciers, adjoint methods have been widely used to study them, with a first paper three decades ago.¹⁵⁰ Since then, the adjoint method has been applied to many different studies investigating parameter and state estimation,⁹¹ ice volume sensitivity to basal, surface and initial conditions,¹⁰³ inversion of initial conditions¹⁶⁷ or inversion of basal friction.¹⁶⁶ All these studies derived the adjoint with a manual implementation. Additionally, the use of AD has become increasingly widespread in glaciology, paving the way for more complex modelling frameworks.^{100,145} Recently, differentiable programming has also facilitated the development of hybrid frameworks, combining numerical methods with data-driven models by means of universal differential equations.³⁰ Alternatively, some other approaches have dropped the use of numerical solvers in favour of different flavours of physics-informed neural networks, exploring the inversion of rheological properties of glaciers²⁴⁷ and to accelerate ice thickness inversions and simulations by leveraging GPUs.^{120,121}

2.1.3 Biology and ecology

Differential equation models have been broadly used in biology and ecology to model the dynamics of genes and alleles,¹⁷⁹ the ecological and evolutionary dynamics of biological units from bacteria to ecological communities,^{8,31,32,43,71,142,238,240} and biomass and energy fluxes and transformation at ecosystem levels.^{69,75,211,250} Estimating parameters of biological models from laboratory experiments is costly, may be extremely difficult per se²¹¹ and may result in simulations failing in capturing real biological dynamics.²⁴⁹ As a consequence, statistical models have been the main modelling paradigm in biology,²⁶² but inverse modelling methods are increasingly advocated to inform mechanistic models^{12,98,182} and e.g. provide robust forecasts of ecosystem responses to global change.¹² Inverse modelling can take the form of parameter estimation²¹¹ or model selection,¹¹⁸ both involving the use of inference methods to estimate the probability of the empirical data given the model parameters (i.e., the likelihood). Provided that they are inferred together with uncertainties, parameters can be interpreted to better understand the strengths and effects of the processes considered.¹⁸⁵ For instance,^{50,88,106} infer the parameters of dynamic community models to understand the processes involved in ecosystem functions. In model selection, candidate models embedding competing hypotheses about causal processes are derived, and the relative support of each model given the data is computed to discriminate between the hypotheses.^{12,118} For instance, using inverse modelling and eco-evolutionary models embedding competing evolutionary speed hypotheses,²¹⁸ shows that temperature-dependent evolutionary speed most likely explains variations in biodiversity patterns. The computation of the most probable model parameter values, or the computation of the different model supports, critically involves sensitivity methods which must adequately handle the typically larger number of parameters and the nonlinearities of biological models.⁷¹ While inverse modelling in biology has been mostly agnostic to AD techniques, their potential have recently

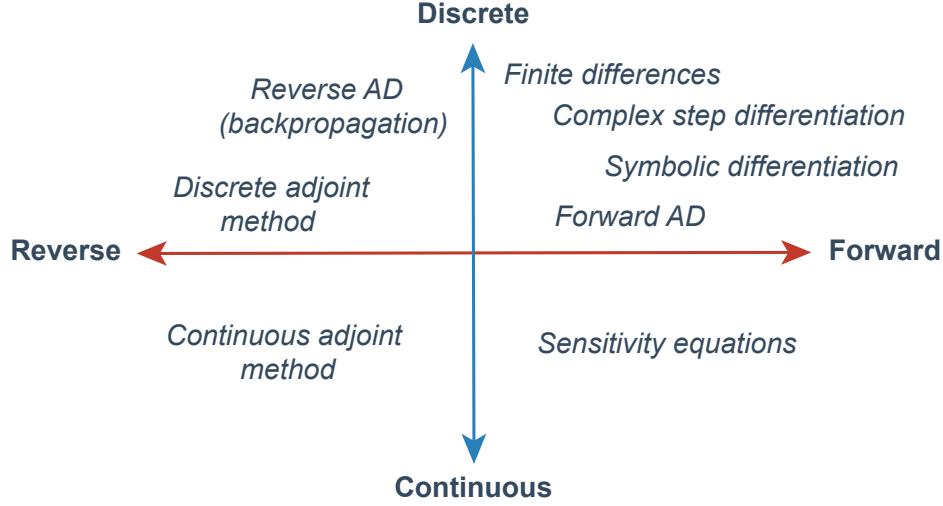


Figure 1: Schematic representation of the different methods available for differentiation involving differential equation solutions. These can be classified depending if they find the gradient by solving a new system of differential equations (*continuous*) or if instead they manipulate unit algebraic operations (*discrete*). Additionally, these methods can be categorized based on their alignment with the direction of the numerical solver. If they operate in the same direction as the solver, they are referred to as *forward* methods. Conversely, if they function in the opposite direction, they are known as *reverse* methods.

been highlighted^{12,68} and new approaches to accommodate for the specificity of biological and ecological models are increasingly proposed.^{33,183,254} Given that key ecological processes are not accurately represented in biological models,^{43,98,211} AD-based inverse modelling methods are particularly relevant to support hybrid approaches where neural networks are used as data-driven parametrization^{33,195} in differential equation models. Increasingly available biological dataset following the development of monitoring technologies such as paleo-time series,¹² environmental DNA,²⁰⁴ remote sensing,¹¹⁵ bioacoustics,⁷ and citizen observations,⁷⁴ are additional opportunities.

3 Methods: A mathematical perspective

There is a large family of methods for computing gradients and sensitivities of systems of differential equations. Depending on the number of parameters and the complexity of the differential equation we are trying to solve, they have different mathematical, numerical, and computational advantages. These methods can be roughly classified as follows.¹⁴⁹

- *Continuous* vs *discrete* methods
- *Forward* vs *reverse* methods

Figure 1 displays a classification of some methods under this two-fold classification.

The *continuous* vs *discrete* distinction is one of mathematical and numerical nature. When solving for the gradient of a differential equation, one needs to derive both a mathematical expression for the gradient (the *differentiation step*) and solve the equations using a numerical solver (the *discretization step*).^{34,178,216,259} Depending on the order of these two operations, we are going to talk about discrete methods (discretize-then-differentiate) or continuous methods (differentiate-then-discretize). In the case of *discrete* methods, gradients are computed based on simple function evaluations of the solutions of the numerical solver (finite differences, complex step differentiation) or by manipulation of atomic operations inside a numerical solver (AD, symbolic differentiation, discrete adjoint method). In the case of *continuous* methods, a new set of differential equations is derived for the sensitivity (sensitivity equations) or the adjoint (continuous adjoint

method) of the system, both quantities that allow the calculation of the desired gradient. When comparing between discrete and continuous methods, more than talking about computational efficiency we are focusing on the mathematical consistency of the method, that is, *is the method estimating the right gradient?*. When using discrete methods, we may have an algorithmically correct method, meaning that is computing the gradient of the solver discretization, but is not approximating the gradient of the true solution of the differential equation. For example, methods like automatic differentiation may compute the exact derivative of the numerical approximation of a loss function and yet not give a good approximation of the exact derivative of the loss.²⁴¹ However, one has to keep in mind that AD computes the exact derivative of an approximation of the objective and may not yield an approximation to the exact derivatives of the objective.

The *forward* vs *reverse* distinction regards when the gradient is computed, if this happens during the forward pass of the numerical solver or in a later recalculation.⁹⁵ In all *forward* methods the solution of the differential equation is solved sequentially and simultaneously with the gradient during the forward pass of the numerical solver. On the contrary, *reverse* methods compute the gradient tracking backwards the forward model by resolving a new problem that moves in the opposite direction as the original numerical solver. For systems of ordinary differential equations (ODEs) and initial value problems (IVPs), most numerical methods solve the differential equation progressively moving forward in time, reason why forward methods solve the gradient moving *forward* in time – or, instead, they solve a new system that goes *backwards* in time.

As we will discuss in the following sections, forward methods are very efficient for problems with a small number of parameters we want to differentiate with respect to, while backwards methods are more efficient for a large number of parameters but they come with a larger memory cost which needs to be overcome using different performance tricks. With the exception of finite differences and complex step differentiation, the rest of the forward methods (i.e. forward AD, sensitivity equations, symbolic differentiation) compute the full sensitivity of the differential equation, that is, how the full solution of the ODEs changes when we change the parameters of the model. This can be computationally expensive for large systems. Conversely, reverse methods are based on the computation of intermediate variables, known as the adjoint or dual variables, that cleverly avoid the unnecessary calculation of the full sensitivity at expenses of larger memory cost.⁸⁷ For this reason, reverse methods can be also labeled as adjoint methods.¹⁴⁹

One extra distinction between methods is with regards to how computationally entangled the numerical solver and the differentiation machinery are. With the exception of the discrete adjoint methods, this coincides with discrete-continuous classification. However, the construction of the discrete adjoint (which surprisingly is one of the most popular in the literature) is based on the numerical solver, something that does not happen with the other discrete methods. While this might not have big conceptual implications, it is an important consideration when using software that integrates numerical solvers and differentiation, a distinction that will help in the discussion in Section 4.

The rest of this section is organized as follows. We will first introduce some basic mathematical notions that are going to facilitate the discussion of the sensitivity methods (Section 3.1). Then we will embark in the mission of mathematically introducing each one of methods listed in Figure 1. We will finalize the discussion in Section 3.9 with an comparison of some of mathematical foundations of these methods.

3.1 Preliminaries

Consider a system of first order ordinary differential equations (ODEs) given by

$$\frac{du}{dt} = f(u, \theta, t), \quad (1)$$

where $u \in \mathbb{R}^n$ is the unknown solution; $f : \mathbb{R}^n \times \mathbb{R}^p \times \mathbb{R} \mapsto \mathbb{R}^n$ is a function that depends on the state u , some vector parameter $\theta \in \mathbb{R}^p$, and potentially the independent variable t which we will refer as time; and with initial condition $u(t_0) = u_0$. Here n denotes the total number of ODEs and p the dimension of a parameter embedded in the functional form of the differential equation. Although we here consider the case of ODEs, that is, when the derivatives are just with respect to the time variable t , the ideas presented here can be extended to the case of partial differential equations (PDEs; for example, via the method of lines¹⁴) and differential algebraic equations (DAE). In fact, PDEs play an essential role when formulating equations of motion via local conservation (and constitutive) laws in physics-based simulations. Except for a minority of functions $f(u, \theta, t)$, solutions to Equation (1) need to be computed using a numerical solver.

3.1.1 Numerical solvers for differential equations

Numerical solvers for the solution of ODEs or initial value problems can be classified as one-step methods, among which Runge-Kutta methods are the most widely used, and multi-step methods.⁹⁷ With a long historical record in numerical analysis, Runge-Kutta methods generalize quadrature rules to iterative solve for the time discretization $u^n \approx u(t_n)$. Given an integer s , a s -stage Runge-Kutta methods is defined by generalizing numerical integration quadrature rules as follows

$$\begin{aligned} u^{n+1} &= u^n + \Delta t_n \sum_{i=1}^s b_i k_i \\ k_i &= f \left(u^n + \sum_{j=1}^s a_{ij} k_j, \theta, t_n + c_i \Delta t_n \right) \quad i = 1, 2, \dots, s. \end{aligned} \quad (2)$$

where $u^n \approx u(t_n)$ approximates the solution at time t_n ; timesteps $\Delta t_n = t_{n+1} - t_n$; and coefficients a_{ij} , b_i , and c_j , with $i, j = 1, 2, \dots, s$, usually represented in the form of a tableau. Different choices of number of stages and coefficients give different orders of convergence of the numerical scheme.^{Butcher_Wanner_1996, Butcher_2001}

On the contrary, multisteps linear solvers are of the form

$$\sum_{i=0}^{K_1} \alpha_{ni} u^{n-i} + \Delta t_n \sum_{i=0}^{K_2} \beta_{ni} f(u^{n-i}, \theta, t_{n-i}) = 0. \quad (3)$$

with coefficients α_{ni} and β_{nj} . Notice that multisteps linear methods are linear in the function f , which is not the case in Runge-Kutta methods with intermediate evaluations.¹⁴ For multistep methods, solving the differential equation implies to be able to solve the system of constraints

$$g_i(u_i; \theta) = u_i - h \beta_{n0} f(u_i, \theta, t_i) - \alpha_i = 0 \quad (4)$$

where α_i has includes the information of all the past iterations. This system can be solved sequentially, by solving for u_i in increasing order of index using Newton method,

$$u_i^{(j+1)} = u_i^{(j)} - \left(\frac{\partial g_i}{\partial u_i}(u_i^{(j)}; \theta) \right)^{-1} g(u_i^{(j)}; \theta). \quad (5)$$

There are many considerations at the moment of picking a numerical solver. One of the most important ones is the stiffness of the differential equation we are trying to solve. Although stiffness is a known phenomena in the study of differential equation solver, different definitions (and types) of instability exist in the literature. This is due to historical reasons^{Dahlquist_1985} as well as the fact that different stiff equations suffer from different types of instabilities. Among them we select

- Stiff differential equations are characterized by dynamics with different time scale,¹²⁷ also characterized by the phenomena of increasing oscillations...
- Stiff equations are equations for which explicit methods do not work and implicit methods work better²⁴⁸

Stability properties can be achieved by the use of implicit methods over explicit methods. Explicit methods are characterized by $\beta_{n,0} = 0$ for the multistep and $a_{ij} = 0$ if $i \leq j$ for Runge-Kutta methods, otherwise, the method is implicit.

Another important consideration is how to pick the time-steps Δt_i in a numerical solver.⁹⁷ Modern solvers include stepsize controllers that pick Δt_i as large as possible to minimize the total number of steps at the same time that they control for large errors in the numerical solution controlled by adjustable relative and absolute tolerances (see Appendix B).

It is important to remark that direct methods to solve higher order systems of ODEs exist (Nystrom), but we can usually reduce them to a extended systems of ODEs.⁹⁷

3.1.2 What to differentiate and why?

We are interested in computing the gradient of a given function $L(u(\cdot, \theta))$ with respect to the parameter θ . This formulation is very general and allows to include many different applications, including the following.

- **Loss function and empirical risk function.** This is usually a real-valued function that quantifies the level of agreement between the model prediction and observations. Examples of loss functions include the squared error

$$L(\theta) = \frac{1}{2} \|u(t_1; \theta) - u^{\text{target}}(t_1)\|_2^2, \quad (6)$$

where $u^{\text{target}}(t_1)$ is the desired target observation at some later time t_1 . More generally, we can evaluate the loss function at points of the time series for which we have observations,

$$L(\theta) = \frac{1}{2} \sum_{i=1}^N \omega_i \|u(t_i; \theta) - u^{\text{target}}(t_i)\|_2^2. \quad (7)$$

with ω_i some arbitrary non-negative weights. More generally, misfit functions used in optimal estimation and control problems map from the model's state space, in this case the solution $u(t)$, to the observation space define by a new variable $y(t) = H(u(t, \theta))$, where $H : \mathbb{R}^n \mapsto \mathbb{R}^o$ is a given function mapping the latent state to observational space.³⁷ In these cases, the lost function is instead

$$L(\theta) = \frac{1}{2} \sum_{i=1}^N \omega_i \|H(u(t_i; \theta)) - y^{\text{target}}(t_i)\|_2^2. \quad (8)$$

We can also consider the continuous evaluated loss function of the form

$$L(u(\cdot, \theta)) = \int_{t_0}^{t_1} h(u(t; \theta), \theta) dt, \quad (9)$$

with h being a function that quantifies the contribution of the error term at every time $t \in [t_0, t_1]$. Defining a loss function where just the empirical error is penalized is known as trajectory matching.¹⁹² Other methods like gradient matching and generalized smoothing the loss depends on smooth approximations of the trajectory and their derivatives.

- **Likelihood function.** From a statistical perspective, it is common to assume that observations correspond to noisy observations of the underlying dynamical system, $y_i = H(u(t_i; \theta)) + \varepsilon_i$, with ε_i errors or residual that are independent of each other and of the trajectory $u(\cdot; \theta)$.¹⁹² When H is the identity, each y_i corresponds to the noise observation of the state $u(t_i; \theta)$. If $p(Y|t, \theta)$ is the probability distribution of $Y = (y_1, y_2, \dots, y_N)$, maximum likelihood estimation consists in finding the maximum a posteriori (MAP) estimate of the parameter θ as

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \ell(Y|\theta) = \prod_{i=1}^n p(y_i|\theta, t_i). \quad (10)$$

When $\varepsilon_i \sim N(0, \sigma_i^2 \mathbb{I})$ is the isotropic multivariate normal distribution, the maximum likelihood principle is the same as minimizing $-\log \ell(Y|\theta)$ which coincides with the mean squared error of Equation (8),¹⁰¹

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \{-\log \ell(Y|\theta)\} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \frac{1}{2\sigma_i^2} \|y_i - H(u(t_i; \theta))\|_2^2. \quad (11)$$

Provided with a prior distribution $p(\theta)$ for the parameter θ , we can further compute a posterior distribution for θ given the observations Y following Bayes theorem

$$p(\theta|Y) = \frac{p(Y|\theta)p(\theta)}{p(Y)}. \quad (12)$$

In practice, the posterior is difficult to evaluate and needs to be approximated using Markov chain Monte Carlo (MCMC) sampling methods.⁷⁸ Being able to further compute gradients of the likelihood allows to design more efficient sampling methods, such as Hamiltonian MCMC.²²

- **Quantity of interest.** Another important example is when L returns the value of the solution at one or many points, which is useful when we want to know how the solution itself changes as we move the parameter values.
- **Diagnosis of the solution.** In many cases we are interested in optimizing the value of some variable that is a function of the solution of a differential equation. This is the case in design control theory, a popular approach in aerodynamics modelling where goals include maximizing the speed of an airplane or the lift of a wing given the solution of the flow equation for a given geometry profile.^{85,113,160}

In the rest of the manuscript we will use letter L to emphasize that in many cases this will be a loss function, but without loss of generality this includes the richer class of functions included in the previous examples.

3.1.3 Optimization

In the context of optimization, the gradient of the loss allows performing gradient-based updates on the parameter θ by

$$\theta^{k+1} = \theta^k - \alpha_k \frac{dL}{d\theta^k}. \quad (13)$$

Gradient-based methods tend to outperform gradient-free optimization schemes, as they are not prone to the curse of dimensionality.²¹¹ While a direct implementation of gradient descent is prone to converge to a local minimum and slow down in a neighborhood of saddle points, variants employing more advanced updating strategies have been proposed to avoid convergence to local minima.²⁰⁰ These methods include Newton-type methods,^{second-order-optimization} quasi-Newton methods, acceleration techniques,^{JMLR:v22:20-207} and natural gradient descent methods.^{doi:10.1137/22M1477805} For instance, ADAM is an adaptive, momentum-based algorithm that remembers the solution update at each iteration, and determines the next update as a linear combination of the gradient and the previous update.¹²⁸ It has been widely adopted to train highly parametrized neural networks (up to the order of 10^8 parameters²³⁹). Other widely employed algorithms are the Broyden–Fletcher–Goldfarb–Shanno (BFGS) and its limited-memory version algorithm (L-BFGS), which determine the descent direction by preconditioning the gradient with curvature information. ADAM is less prone to converge to a local minimum, while (L-)BFGS has a faster converge rate. Using ADAM for the first iterations followed by (L-)BFGS proves to be a successful strategy to minimize a loss function with best accuracy.

3.1.4 Sensitivity matrix

In the general case, we are going to work with loss functions of the form $L(\theta) = L(u(\cdot, \theta), \theta)$. Using the chain rule we can derive

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial \theta} + \frac{\partial L}{\partial \theta}. \quad (14)$$

The two partial derivatives of the loss function on the right-hand side are usually easy to evaluate. For example, for the loss function in Equation (6) this are simply given by

$$\frac{\partial L}{\partial u} = u - u^{\text{target}}(t_1) \quad \frac{\partial L}{\partial \theta} = 0. \quad (15)$$

Just as in this last example, in most applications the loss function $L(\theta)$ will depend on θ just through u , meaning $\frac{\partial L}{\partial \theta} = 0$. The complicated term to compute is the matrix of derivatives $\frac{\partial u}{\partial \theta}$, usually referred to as the *sensitivity* s , and represents how much the full solution u varies as a function of the parameter θ ,

$$s = \frac{\partial u}{\partial \theta} = \begin{bmatrix} \frac{\partial u_1}{\partial \theta_1} & \cdots & \frac{\partial u_1}{\partial \theta_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial u_n}{\partial \theta_1} & \cdots & \frac{\partial u_n}{\partial \theta_p} \end{bmatrix} \in \mathbb{R}^{n \times p}. \quad (16)$$

The sensitivity s defined in Equation (16) is what is called a *Jacobian*, that is, a matrix of first derivatives for general vector-valued functions. Some of the methods we will discuss here will directly compute the sensitivity, while others will only deal with Jacobian-vector products (JVPs) of the form $\frac{\partial u}{\partial \theta} v$, for some

vector $v \in \mathbb{R}^p$. The product $\frac{\partial u}{\partial \theta} v$ is the directional derivative of the function $u(\theta)$, also known as the Gateaux derivative of $u(\theta)$ in the direction v , given by

$$\frac{\partial u}{\partial \theta} v = \lim_{h \rightarrow 0} \frac{u(\theta + hv) - u(\theta)}{h}, \quad (17)$$

representing how much the function u changes when we perturb θ in the direction of v .

3.2 Finite differences

The simplest way of evaluating a derivative is by computing the difference between the evaluation of the function at a given point and a small perturbation of the function. In the case of the function $L : \mathbb{R}^p \mapsto \mathbb{R}$, we can approximate

$$\frac{dL}{d\theta_i}(\theta) = \frac{L(\theta + \varepsilon e_i) - L(\theta)}{\varepsilon} + \mathcal{O}(\varepsilon), \quad (18)$$

with e_i the i -th canonical vector and ε the stepsize. Even better, the centered difference scheme leads to

$$\frac{dL}{d\theta_i}(\theta) = \frac{L(\theta + \varepsilon e_i) - L(\theta - \varepsilon e_i)}{2\varepsilon} + \mathcal{O}(\varepsilon^2). \quad (19)$$

While Equation (18) gives the derivative to an error of magnitude $\mathcal{O}(\varepsilon)$, the centered differences schemes improves the accuracy to $\mathcal{O}(\varepsilon^2)$.¹⁵

However, there are a series of problems associated with this approach. The first one is due to how this scales with the dimension p of parameter vector θ . Each directional derivative requires the evaluation of the loss function L twice. For the centered differences approach in Equation (19), this requires a total of $2p$ function evaluations which demands solving the differential equation each time for a new set of parameters.

A second problem is due to rounding errors. Every computer ultimately stores and manipulates numbers using floating point arithmetic.⁹⁰ Equations (18) and (19) involve the subtraction of two numbers that are very close to each other, which leads to large cancellation errors for small values of ε that are amplified by the division by ε . On the other hand, large values of the stepsize give inaccurate estimations of the gradient. Finding the optimal value of ε that balances these two effects is sometimes called the *stepsize dilemma*, for which some heuristics and algorithms have been introduced.^{17,107,156} Some of these methods require some a priori knowledge about the function to be differentiated, and others are based on arbitrary historical rules. If many analytical functions, like polynomials and trigonometric functions, can be computed with machine precision, numerical solutions of differential equations have errors larger than machine precision, which leads to inaccurate estimations of the gradient when ε is too small. We will further emphasize this point in Section 4.

Despite all these caveats, finite differences can be useful when computing Jacobian-vector products (JVPs). Given a Jacobian matrix $J = \frac{\partial f}{\partial u}$ (or the sensitivity $s = \frac{\partial u}{\partial \theta}$) and a vector v , the product Jv corresponding to the directional derivative and can be approximated as

$$Jv \approx \frac{f(u + \varepsilon v, \theta, t) - f(u, \theta, t)}{\varepsilon} \quad (20)$$

This approach is used in numerical solvers based on Krylov methods, where linear systems are solved by iteratively solving matrix-vectors products.¹¹²

3.3 Automatic differentiation

Automatic differentiation (AD) is a technology that generates new code representing derivatives of a given parent code. Examples are code representing the tangent linear or adjoint operator of the parent code.⁹⁵ The names *algorithmic* and *computational* differentiation had also been used in the literature, emphasizing the algorithmic rather than automatic nature of AD.^{95,152} The basis of all AD systems is the notion that complicated functions included in any computer program can be reduced to a sequence of simple algebraic operations that have straightforward derivative expressions, based upon elementary rules of differentiation.¹²² The derivatives of the outputs of the computer program (dependent variables) with respect to their inputs

(independent variables) are then combined using the chain rule. One advantage of AD systems is to automatically differentiate programs that include control flow, such as branching, loops or recursions. This is because any program can be reduced to a trace of input, intermediate and output variables.²⁰

Depending if the concatenation of these gradients is done as we execute the program (from input to output) or in a later instance where we trace-back the calculation from the end (from output to input), we refer to *forward* or *reverse* AD, respectively. Neither forward nor reverse mode is more efficient in all cases,⁹³ as we will discuss in Section 3.3.3.

3.3.1 Forward mode

Forward mode AD can be implemented in different ways depending on the data structures we use at the moment of representing a computer program. Examples of these data structures include dual numbers and computational graphs.²⁰

3.3.1.1 Dual numbers

Dual numbers extend the definition of a numerical variable that takes a certain value to also carry information about its derivative with respect to a certain parameter.⁴⁶ We define a dual number based on two variables: a *value* coordinate x_1 that carries the value of the variable and a *derivative* (also known as partial or tangent) coordinate x_2 with the value of the derivative $\frac{\partial x_1}{\partial \theta}$. Just as complex number, we can represent dual numbers as an ordered pair (x_1, x_2) , sometimes known as Argand pair, or in the rectangular form

$$x_\epsilon = x_1 + \epsilon x_2, \quad (21)$$

where ϵ is an abstract number called a perturbation or tangent, with the properties $\epsilon^2 = 0$ and $\epsilon \neq 0$. This last representation is quite convenient since it naturally allow us to extend algebraic operations, like addition and multiplication, to dual numbers.¹²⁴ For example, given two dual numbers $x_\epsilon = x_1 + \epsilon x_2$ and $y_\epsilon = y_1 + \epsilon y_2$, it is easy to derive using the fact $\epsilon^2 = 0$ that

$$x_\epsilon + y_\epsilon = (x_1 + y_1) + \epsilon(x_2 + y_2) \quad x_\epsilon y_\epsilon = x_1 y_1 + \epsilon(x_1 y_2 + x_2 y_1). \quad (22)$$

From these last examples, we can see that the derivative component of the dual number carries the information of the derivatives when combining operations. For example, suppose than in the last example the dual variables x_2 and y_2 carry the value of the derivative of x_1 and x_2 with respect to a parameter θ , respectively.

Intuitively, we can think of ϵ as being a differential in the Taylor series expansion, fact that we can observe in how the output of any scalar functions is extended to a dual number output:

$$\begin{aligned} f(x_1 + \epsilon x_2) &= f(x_1) + \epsilon x_2 f'(x_1) + \epsilon^2 \cdot (\dots) \\ &= f(x_1) + \epsilon x_2 f'(x_1). \end{aligned} \quad (23)$$

When computing first order derivatives, we can ignore everything of order ϵ^2 or larger, which is represented in the condition $\epsilon^2 = 0$. This implies that we can use dual numbers to implement forward AD through a numerical algorithm. In Section 4 we will explore how this is implemented.

Multidimensional dual number generalize dual number to include a different dual variable ϵ_i for each variable we want to differentiate with respect to.^{170,197} A multidimensional dual number is then defined as $x_\epsilon = x + \sum_{i=1}^p x_i \epsilon_i$, with the property that $\epsilon_i \epsilon_j = 0$ for all pairs i and j . Incorrect implementations of this aspect can lead to *perturbation confusion*,^{151,217} an existing problem in some AD software where dual variables corresponding to different variables result indistinguishable, especially in the case of nested functions.¹⁵¹ This problem can be further been overcome by computing the full gradient as the combination of independent directional derivatives (see Section 3.3.3) Another extension of dual numbers that should not be confused with multidimensional dual numbers are hyper-dual numbers, which allow to compute higher-order derivatives of a function.⁶⁶

3.3.1.2 Computational graph

A useful way of representing a computer program is via a computational graph with intermediate variables that relate the input and output variables. Most scalar functions of interest can be represented in this factorial

form as a acyclic directed graph with nodes associated to variables and edges to atomic operations,^{93,95} known as Kantorovich graph¹²³ or its linearized representation via Wengert trace/tape.^{18,95,251} We can define $v_1, v_2, \dots, v_p = \theta_1, \theta_2, \dots, \theta_p$ the input set of variables; v_{p+1}, \dots, v_{m-1} the set of all the intermediate variables, and finally $v_m = L(\theta)$ the final output of a computer program. This can be done in such a way that the order is strict, meaning that each variable v_i is computed just as a function of the previous variables v_j with $j < i$. Once the graph is constructed, we can compute the derivative of every node with respect to other (a quantity known as the tangent) using the Bauer formula^{18,177}

$$\frac{\partial v_j}{\partial v_i} = \sum_{\substack{\text{paths } w_0 \rightarrow w_1 \rightarrow \dots \rightarrow w_K \\ \text{with } w_0 = v_i, w_K = v_j}} \prod_{k=0}^{K-1} \frac{\partial w_{k+1}}{\partial w_k}, \quad (24)$$

where the sum is calculated with respect to all the directed paths in the graph connecting the input and target node. Instead of evaluating the last expression for all possible path, a simplification is to increasingly evaluate $j = p + 1, \dots, m$ using the recursion

$$\frac{\partial v_j}{\partial v_i} = \sum_{w \text{ such that } w \rightarrow v_j} \frac{\partial v_j}{\partial w} \frac{\partial w}{\partial v_i} \quad (25)$$

Since every variable node w such that $w \rightarrow v_j$ is an edge of the computational graph have index less than j , we can iterate this procedure as we run the computer program and solve for both the function and its gradient. This is possible because in forward mode the term $\frac{\partial w}{\partial v_i}$ has been computed in a previous iteration, while $\frac{\partial v_j}{\partial w}$ can be evaluated at the same time the node v_j is computed based on only the value of the parent variable nodes. The only requirement for differentiation is being able to compute the derivative/tangent of each edge/primitive and combine these using the recursion defined in Equation (25).

3.3.2 Reverse mode

Reverse mode AD is also known as the adjoint of cotangent linear mode, or backpropagation in the field of machine learning. The reverse mode of automatic differentiation has been introduced in different contexts⁹⁴ and materializes the observation made by Phil Wolfe that if the chain rule is implemented in reverse mode, then the ratio between the computation of the gradient of a function and the function itself can be bounded by a constant that does not depend of the number of parameters to differentiate,^{93,253} a point known as the *cheap gradient principle*.⁹⁴ Given a directional graph of operations defined by a Wengert list,²⁵¹ we can compute gradients of any given function in the same fashion as Equation (25) but in reverse mode as

$$\bar{v}_i = \frac{\partial \ell}{\partial v_i} = \sum_{w: v \rightarrow w \in G} \frac{\partial w}{\partial v_i} \bar{w}. \quad (26)$$

In this context, the notation $\bar{w} = \frac{\partial \ell}{\partial w}$ is introduced to signify the partial derivative of the output variable, here associated to the loss function, with respect to input and intermediate variables. This derivative is often referred to as the adjoint, dual, or cotangent, and its connection with the discrete adjoint method will be made more explicitly in Section 3.9.3.

Since in reverse-mode AD the values of \bar{w} are being updated in reverse order, in general we need to know the state value of all the argument variables v of w in order to evaluate the terms $\frac{\partial w}{\partial v}$. These state values (required variables) need to be either stored in memory during the evaluation of the function or recomputed on the fly in order to be able to evaluate the derivative. Checkpointing schemes exist to limit and balance the amount of storing versus recomputation (see section 4.1.2.3).

3.3.3 AD connection with JVPs and VJPs

When working with unit operations that involve matrix operations dealing with vectors of different dimensions, the order in which we apply the chain rule matters.⁸³ When computing a gradient using AD, we can encounter vector-Jacobian products (VJPs) or Jacobian-vector products (JVP). As their name indicates, the difference between them is that the quantity we are interested in is described by the product of a Jacobian

times a vector on the left side (VJP) or the right (JVP). Furthermore, both forward and reverse AD can be thought as a way of computing directional derivatives associated with JVPs (see Equation (17)) and VJPs, respectively. In other words, given a function $g : \mathbb{R}^{d_1} \mapsto \mathbb{R}^{d_2}$ that is evaluated during the forward mode of given program, AD will carry terms of the form $Dh(x) \cdot \dot{x}$ (JVP) in forward mode and $\bar{y}^T \cdot Dh(x)$ (VJP) in reverse mode.⁹⁵

Let us consider for example the case of a nested loss function $L : \mathbb{R}^p \mapsto \mathbb{R}$ taking a total of p arguments as inputs that can be decomposed as $L(\theta) = \ell \circ g_k \circ \dots \circ g_2 \circ g_1(\theta)$, with $\ell : \mathbb{R}^{d_k} \mapsto \mathbb{R}$ the final evaluation of the loss function after we apply in order a sequence of intermediate functions $g_i : \mathbb{R}^{d_{i-1}} \mapsto \mathbb{R}^{d_i}$, where we define $d_0 = p$ for simplicity. The final gradient is computed as the chain product of vectors and Jacobians as

$$\nabla_{\theta} L = \nabla \ell \cdot Dg_k \cdot Dg_{k-1} \cdot \dots \cdot Dg_2 \cdot Dg_1, \quad (27)$$

with Dg_i the Jacobian of each nested function evaluated at the intermediate values $g_{i-1} \circ g_{i-2} \circ \dots \circ g_i(\theta)$. Notice that in the last equation, $\nabla \ell \in \mathbb{R}^{d_k}$ is a vector. In order to compute $\nabla_{\theta} L$, we can solve the multiplication starting from the right side, which will correspond to multiplying the Jacobians forward from Dg_1 to Dg_k , or from the left side, moving backwards. The important aspect of the backwards case is that we will always be computing VJPs, since $\nabla \ell$ is a vector. Since VJPs are easier to evaluate than full Jacobians, the reverse mode will in general be faster when $1 \ll p$. This example is illustrated in Figure 2. For general rectangular matrices $A \in \mathbb{R}^{d_1 \times d_2}$ and $B \in \mathbb{R}^{d_2 \times d_3}$, the cost of the matrix multiplication AB is $\mathcal{O}(d_1 d_2 d_3)$. It is worth noticing that if well more efficient methods for matrix-matrix multiplication based on Strassen's recursive algorithm and its variants exist, these are not extensively used in most scientific applications.^{109,215} This implies that forward AD requires a total of

$$d_2 d_1 p + d_3 d_2 p + \dots + d_k d_{k-1} p + d_k p = \mathcal{O}(kp) \quad (28)$$

operations, while backwards mode AD requires

$$d_k d_{k-1} + d_{k-1} d_{k-2} + \dots + d_2 d_1 + d_1 p = \mathcal{O}(k + p) \quad (29)$$

operations.

In the general case of a function $L : \mathbb{R}^p \mapsto \mathbb{R}^q$ with multiple outputs and a total of k intermediate functions, the cost of forward AD is $\mathcal{O}(pk + q)$ and the cost of reverse is $\mathcal{O}(p + kq)$. When the function to differentiate has a larger input space than output ($q \ll p$), AD in reverse mode is more efficient as it propagates the chain rule by computing VJPs, the reason why reverse-mode AD is more used in modern machine learning. However, notice that backwards mode AD requires us to save intermediate variables through the forward run in order to run backwards afterwards,²¹ leading to performance overhead that makes forward AD more efficient when $p \lesssim q$.^{20,93,152} In other words, backwards AD is really more efficient when $q \ll p$. We discuss in section 4.1.2.3 how this problem can be overcome with a good checkpointing scheme.

In a practical sense, many AD systems are reduced to the computation of only directional derivatives (VJPs) and JVPs.⁹⁵ Full Jacobians $J \in \mathbb{R}^{n \times p}$ (e.g., the sensitivity $s = \frac{\partial u}{\partial \theta} \in \mathbb{R}^{n \times p}$) can be fully reconstructed by the independent computation of the p columns of J via the JVPs $J e_i$, with $e_i \in \mathbb{R}^p$ the canonical vectors; or by the calculation of the m rows of J via the VJPs $e_j^T J$, with $e_j \in \mathbb{R}^n$. An important observation here is then how to efficiently compute sparse Jacobians, which are commonplace in large-scale nonlinear systems, discretized PDEs, etc., and are often a major computational bottleneck for solving those problems. Consider the example of a Jacobian J_{sparse} with known sparsity pattern given by

$$J_{\text{sparse}} = \begin{bmatrix} \bullet & & & & \\ & \bullet & \bullet & & \\ & & & \bullet & \\ \bullet & \bullet & & & \bullet \\ & & & & \bullet \end{bmatrix}, \quad (30)$$

where \bullet denotes the non-zero elements of the Jacobian. For cases with known sparsity pattern, *colored AD* can be used to chunk multiple JVPs or VJPs using the colored Jacobian.⁷⁶ More concretely, we can color

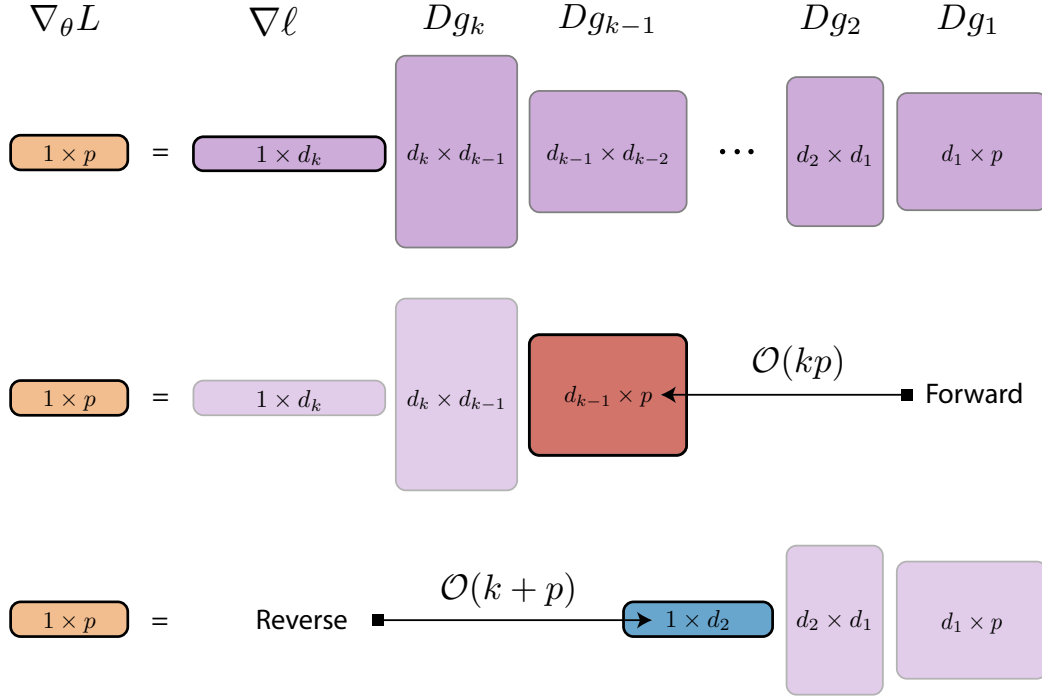


Figure 2: Comparison between forward and backward AD. Changing the order of Jacobian multiplications changes the total number of floating-point operations, which leads to different computational complexities between forward and reverse mode. When the multiplication is carried from the right side of the mathematical expression for $\nabla_{\theta} L$, each matrix simplification involves a matrix with size p , giving a total complexity of $\mathcal{O}(kp)$. This is the opposite of what happens when we carried the VJP from the left side of the expression, where the matrix of size $d_1 \times p$ has no effect in the intermediate calculations, making all the intermediate calculations $\mathcal{O}(1)$ with respect to p and a total complexity of $\mathcal{O}(k+p)$.

the above matrix as follows:

$$J_{\text{sparse}}^{(\text{col})} = \begin{bmatrix} \blacktriangle & & & & \\ & \blacksquare & \blacktriangle & & \\ & & & \blacktriangle & \\ \blacktriangle & \blacksquare & & & \\ & & & & \blacklozenge \\ & & & & \blacklozenge \end{bmatrix} \quad J_{\text{sparse}}^{(\text{row})} = \begin{bmatrix} \blacksquare & & & & \\ & \blacksquare & \blacksquare & & \\ & & & \blacksquare & \\ \blacklozenge & \blacklozenge & & & \blacklozenge \\ & & & & \blacksquare \end{bmatrix}. \quad (31)$$

To compute $J_{\text{sparse}}^{(\text{col})}$, we just need to perform three JVPs,

$$J_{\text{sparse}}^{(\text{col})} \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} \blacktriangle \\ \blacktriangle \\ \blacktriangle \\ \blacktriangle \\ \blacklozenge \end{bmatrix}, \quad J_{\text{sparse}}^{(\text{col})} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix}, \quad J_{\text{sparse}}^{(\text{col})} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \blacklozenge \\ \blacklozenge \\ \blacklozenge \\ \blacklozenge \\ \blacklozenge \end{bmatrix}, \quad (32)$$

compared to five JVPs for a 5×5 dense Jacobian. Similarly, since reverse mode materializes the Jacobian one row at a time, we need two VJP (once each for \blacksquare , and \blacklozenge) compared to five VJP for the dense counterpart.

3.4 Complex step differentiation

An alternative to finite differences that avoids subtractive cancellation errors is based on complex variable analysis. The first proposals originated in 1967 using the Cauchy integral theorem involving the numerical

evaluation of a complex-valued integral.^{146,147} A new approach recently emerged that uses the complex generalization of a real function to evaluate its derivatives.^{154,219} Assuming that the function $L(\theta)$ admits an holomorphic extension (that is, it can be extended to a complex-valued function that is analytical and differentiable²²³), the Cauchy-Riemann conditions can be used to evaluate the derivative with respect to one single scalar parameter $\theta \in \mathbb{R}$ as

$$\frac{dL}{d\theta} = \lim_{\varepsilon \rightarrow 0} \frac{\text{Im}(L(\theta + i\varepsilon))}{\varepsilon}, \quad (33)$$

where i is the imaginary unit satisfying $i^2 = -1$. The order of this approximation can be found using the Taylor expansion of a function,

$$L(\theta + i\varepsilon) = L(\theta) + i\varepsilon \frac{dL}{d\theta} - \frac{1}{2}\varepsilon^2 \frac{d^2L}{d\theta^2} + \mathcal{O}(\varepsilon^3). \quad (34)$$

Computing the imaginary part $\text{Im}(L(\theta + i\varepsilon))$ leads to

$$\frac{dL}{d\theta} = \frac{\text{Im}(L(\theta + i\varepsilon))}{\varepsilon} + \mathcal{O}(\varepsilon^2). \quad (35)$$

The method of *complex step differentiation* consists then in estimating the gradient as $\text{Im}(L(\theta + i\varepsilon))/\varepsilon$ for a small value of ε . Besides the advantage of being a method with precision $\mathcal{O}(\varepsilon^2)$, the complex step method avoids subtracting cancellation error and then the value of ε can be reduced to almost machine precision error without affecting the calculation of the derivative. However, a major limitation of this method is that it works just for complex analytical functions¹⁵⁴ and do not outperforms AD. Extension to higher order derivatives can be done by introducing multicomplex variables.¹³¹

3.5 Symbolic differentiation

In symbolic differentiation, functions are represented algebraically instead of algorithmically, reason why many symbolic differentiation tools are included inside computer algebra systems (CAS).^{Symbolics_jl_2022} Instead of numerically evaluating the final value of a derivative, symbolic systems define *algebraic* objects, including variable names, expressions, operations, and literals. For example, the relation $y = x^2$ is interpreted as expression with two variables, x and y , and the symbolic system need to generate the derivative $y' = 2 \times x$ with 2 a numeric literal, \times a binary operation, and x the same variable assignment than in the original expression. When the function to differentiate is large, symbolic differentiation can lead to *expression swell*, that is, exponentially large or complex symbolic expressions.²⁰ Here, an important piece of CAS is simplification routines that reduce the size and complexity of algebraic expressions by finding common sub-expressions. This can make symbolic differentiation very efficient when computing derivatives multiple times and for different input values.⁵⁵

It is important to remark on the close relationship between AD and symbolic differentiation. There is no agreement as to whether symbolic differentiation should be classified as AD^{58,122,132} or as a different method.²⁰ Both are equivalent in the sense that they perform the same operations but the underlying data structure is different.¹³² Here, expression swell is a consequence of the underlying representation when this does not allow for common sub-expressions. This can also be understood as if AD is symbolic differentiation performed by a compiler,⁵⁸ meaning that different AD can be classified based on the level of integration with the underlying source language.¹²²

3.6 Sensitivity equations

An easy way to derive an expression for the sensitivity s is by deriving the sensitivity equations,¹⁹² a method also referred to as continuous local sensitivity analysis (CSA). If we consider the original system of ODEs given by Equation (1) and we differentiate with respect to θ , we then obtain

$$\frac{d}{d\theta} \left(\frac{du}{dt} - f(u(\theta), \theta, t) \right) = 0. \quad (36)$$

Assuming that an unique solution exists and both $\frac{\partial f}{\partial u}$ and $\frac{\partial f}{\partial \theta}$ are continuous in the neighbour of the solution; or under the guarantee of interchangeability of the derivatives,^{gronwall1919note} for example by assuming that both $\frac{du}{dt}$ and $\frac{du}{d\theta}$ are differentiable,¹⁸¹ we can derive

$$\frac{d}{d\theta} \frac{du}{dt} = \frac{d}{d\theta} f(u(\theta), \theta, t) = \frac{\partial f}{\partial \theta} + \frac{\partial f}{\partial u} \frac{\partial u}{\partial \theta}. \quad (37)$$

Identifying the sensitivity matrix $s(t)$ defined in Equation (16), we obtain the *sensitivity differential equation*

$$\frac{ds}{dt} = \frac{\partial f}{\partial u} s + \frac{\partial f}{\partial \theta}. \quad (38)$$

Both the original system of n ODEs and the sensitivity equation of np ODEs are solved simultaneously, which is necessary since the sensitivity differential equation directly depends on the value of $u(t)$. This implies that as we solve the ODEs, we can ensure the same level of numerical precision for the two of them inside the numerical solver.

In contrast to the methods previously introduced, the sensitivity equations find the gradient by solving a new set of continuous differential equations. Notice also that the obtained sensitivity $s(t)$ can be evaluated at any given time t . This method can be labeled as forward, since we solve both $u(t)$ and $s(t)$ as we solve the differential equation forward in time, without the need of backtracking any operation though the solver. By solving the sensitivity equation and the original differential equation for $u(t)$ simultaneously, we ensure that by the end of the forward step we have calculated both $u(t)$ and $s(t)$.

3.7 Discrete adjoint method

Also known as the adjoint state method, it is another example of a discrete method that aims to find the gradient by solving an alternative system of linear equations, known as the *adjoint equations* simultaneously with the original system of linear equations defined by the numerical solver. These methods are extremely popular in optimal control theory in fluid dynamics, for example for the design of geometries for vehicles and airplanes that optimize performance.^{59,86}

The idea of the adjoint method is to treat the differential equation as a constraint in an optimization problem and then differentiate an objective function subject to that constraint. Mathematically speaking, this can be treated both from a duality or Lagrangian perspective.⁸⁶ In agreement with other authors, we prefer to derive the equation using the former as it may give better insights to how the method works and it allows generalization to other user cases.⁸⁷ The derivation of adjoint methods using the Lagrangian formulation can be found in Appendix A.

3.7.1 Adjoint state equations

The derivation of the discrete adjoint equations is carried out once the numerical scheme for solving Equation (1) has been specified (see Section 3.1.1). Given a discrete sequence of timesteps t_0, t_1, \dots, t_N , we aim to find approximate numerical solutions $u_i \approx u(t_i; \theta)$. Any numerical solver, including the ones discussed in Section 3.1.1, can be understood as solving the (in general nonlinear) system of equations defined by $G(U; \theta) = 0$, where U is the super-vector $U = (u_1, u_2, \dots, u_N) \in \mathbb{R}^{nN}$, and we had combine the systems of equations defined by the iterative solver as $G(U; \theta) = (g_1(u_1; \theta), \dots, g_N(u_N; \theta)) = 0$ (see Equation (4)).

We are interested in differentiating an objective or loss function $L(U, \theta)$ with respect to the parameter θ . Since here U is the discrete set of evaluations of the solver, examples of loss functions now include

$$L(U, \theta) = \frac{1}{2} \sum_{i=1}^N \|u_i - u_i^{\text{obs}}\|^2, \quad (39)$$

with u_i^{obs} the observed time-series. Now, same as Equation (14) we have

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial \theta} + \frac{\partial L}{\partial U} \frac{\partial U}{\partial \theta}. \quad (40)$$

We further need to impose the constraint that the solution satisfies the algebraic equation $G(U; \theta) = 0$, which gives

$$\frac{dG}{d\theta} = \frac{\partial G}{\partial \theta} + \frac{\partial G}{\partial U} \frac{\partial U}{\partial \theta} = 0 \quad (41)$$

and which is equivalent to

$$\frac{\partial U}{\partial \theta} = - \left(\frac{\partial G}{\partial U} \right)^{-1} \frac{\partial G}{\partial \theta}. \quad (42)$$

If we replace this last expression into equation (40), we obtain

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial \theta} - \frac{\partial L}{\partial U} \left(\frac{\partial G}{\partial U} \right)^{-1} \frac{\partial G}{\partial \theta}. \quad (43)$$

The important trick in the adjoint state methods is to observe that in this last equation, the right-hand side can be resolved as a vector-Jacobian product (VJP), with $\frac{\partial L}{\partial U}$ being the vector. Instead of computing the product of the matrices $\left(\frac{\partial G}{\partial U} \right)^{-1}$ and $\frac{\partial G}{\partial \theta}$, it is computationally more efficient first to compute the resulting vector from the VJP operation $\frac{\partial L}{\partial U} \left(\frac{\partial G}{\partial U} \right)^{-1}$ and then multiply this by $\frac{\partial G}{\partial \theta}$. This leads to the definition of the adjoint $\lambda \in \mathbb{R}^{n_N}$ as the solution of the linear system of equations

$$\left(\frac{\partial G}{\partial U} \right)^T \lambda = \left(\frac{\partial L}{\partial U} \right)^T, \quad (44)$$

or equivalently,

$$\lambda^T = \frac{\partial L}{\partial U} \left(\frac{\partial G}{\partial U} \right)^{-1}. \quad (45)$$

Finally, if we replace Equation (45) into (43), we obtain

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial \theta} - \lambda^T \frac{\partial G}{\partial \theta}. \quad (46)$$

The important trick used in the adjoint method is the rearrangement of the multiplicative terms involved in equation (43). Computing the full Jacobian/sensitivity $\partial U / \partial \theta$ will be computationally expensive and involves the product of two matrices. However, we are not interested in the calculation of the Jacobian, but instead in the VJP given by $\frac{\partial L}{\partial U} \frac{\partial U}{\partial \theta}$. By rearranging these terms and relying in the intermediate variable $G(U; \theta)$, we can make the same computation more efficient. These ideas are summarized in the diagram in Figure 3, where we can also see an interesting interpretation of the adjoint as being equivalent to $\lambda^T = -\frac{\partial L}{\partial G}$.

Notice that the algebraic equation of the adjoint λ in Equation (44) is a linear system of equations even when the original system $G(U; \theta) = 0$ was not necessarily linear in U . This means that while the forward mode may require multiple iterations in order to solve the non-linear system $G(U) = 0$ (e.g., by using Krylov methods), the backwards step to compute the adjoint is one single linear system of equations.

3.7.2 Simple linear system

To gain further intuition about the discrete adjoint method, let us consider the simple case of the explicit linear one-step methods, where at every step we need to solve the equation $u_{i+1} = g_i(u_i; \theta) = A_i(\theta) u_i + b_i(\theta)$, where $A_i(\theta) \in \mathbb{R}^{n \times n}$ and $b_i(\theta) \in \mathbb{R}^n$ are defined by the numerical solver.¹¹⁹ This condition can be written in a more compact manner as $G(U) = A(\theta)U - b(\theta) = 0$, that is

$$A(\theta)U = \begin{bmatrix} \mathbb{I}_{n \times n} & 0 & & & \\ -A_1 & \mathbb{I}_{n \times n} & 0 & & \\ & -A_2 & \mathbb{I}_{n \times n} & & \\ & & & \ddots & \\ & & & & -A_{N-1} & \mathbb{I}_{n \times n} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_N \end{bmatrix} = \begin{bmatrix} A_0 u_0 + b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_{N-1} \end{bmatrix} = b(\theta), \quad (47)$$

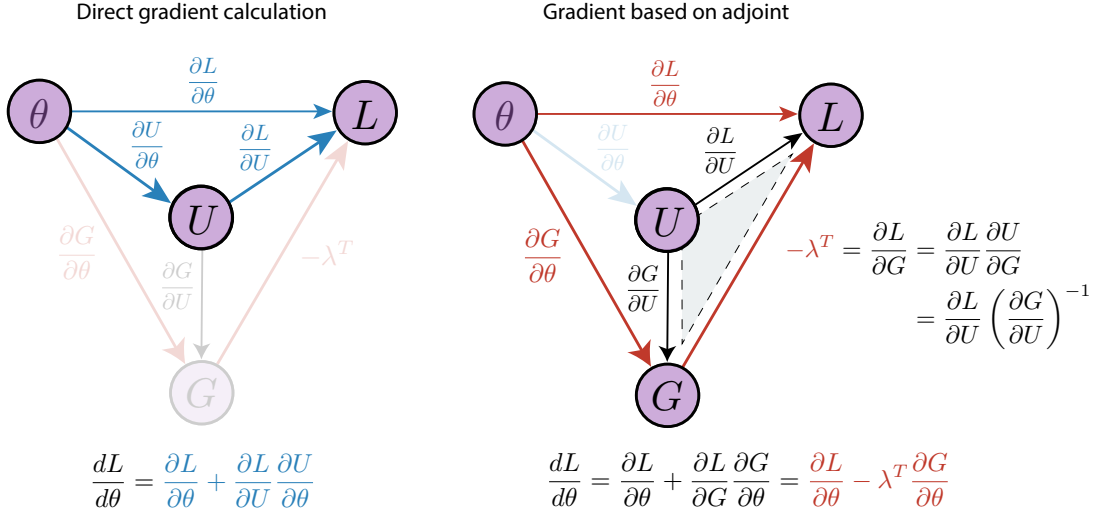


Figure 3: Diagram showing how gradients are computed using discrete adjoints. On the left, we see how gradients will be computed if we use the chain rule applied to the directed triangle defined by the variables θ , U , and L (blue arrows). However, we can define the intermediate vector variable $G = G(U; \theta)$, which satisfies $G = 0$ as long as the discrete system of differential equations are satisfied, and apply the chain rule instead to the triangle defined by θ , G , and L (red arrows). In the red diagram, the calculation of $\frac{\partial L}{\partial G}$ is done by pivoting in U as shown in the right diagram (shaded area). Notice that the use of adjoints avoids the calculation of the sensitivity $\frac{\partial U}{\partial \theta}$. The adjoint is defined as the partial derivative $\lambda^T = -\frac{\partial L}{\partial G}$ representing changes in the loss function due to variations in the discrete equation $G(U; \theta) = 0$.

with $\mathbb{I}_{n \times n}$ the identity matrix of size $n \times n$. Notice that in most cases, the matrix $A(\theta)$ is quite large and mostly sparse. While this representation of the discrete differential equation is convenient for mathematical manipulations, when solving the system we rely on iterative solvers that save memory and computation.

For the linear system of discrete equations $G(U; \theta) = A(\theta)U - b(\theta) = 0$, we have

$$\frac{\partial G}{\partial \theta} = \frac{\partial A}{\partial \theta} U - \frac{\partial b}{\partial \theta}, \quad (48)$$

so the desired gradient in Equation (46) can be computed as

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial \theta} - \lambda^T \left(\frac{\partial A}{\partial \theta} U - \frac{\partial b}{\partial \theta} \right) \quad (49)$$

with λ the discrete adjoint obtained by solving the linear system in Equation (44),

$$A(\theta)^T \lambda = \begin{bmatrix} \mathbb{I}_{n \times n} & -A_1^T & & & \\ 0 & \mathbb{I}_{n \times n} & -A_2^T & & \\ & 0 & \mathbb{I}_{n \times n} & -A_3^T & \\ & & & \ddots & -A_{N-1}^T \\ & & & 0 & \mathbb{I}_{n \times n} \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \vdots \\ \lambda_N \end{bmatrix} = \begin{bmatrix} u_1 - u_1^{\text{obs}} \\ u_2 - u_2^{\text{obs}} \\ u_3 - u_3^{\text{obs}} \\ \vdots \\ u_N - u_N^{\text{obs}} \end{bmatrix} = \frac{\partial L}{\partial U}^T. \quad (50)$$

This is a linear system of equations with the same size of the original $A(\theta)U = b(\theta)$, but involving the adjoint matrix A^T . Computationally this also means that if we can solve the original system of discretized equations then we can also solve the adjoint at the same computational cost (e.g., by using the LU factorization of $A(\theta)$). Another more natural way of finding the adjoints λ is by noticing that the system of equations (50) is equivalent to the final value problem

$$\lambda_i = A_i^T \lambda_{i+1} + (u_i - u_i^{\text{obs}}) \quad (51)$$

with final condition λ_N . This means that we can solve the adjoint equation backwards, i.e., in reverse mode, starting from the final state λ_N and computing the values of λ_i in decreasing index order. Unless the loss function L is linear in U , this procedure requires to know the value of u_i (or some equivalent form of it) at any given timestep.

3.8 Continuous adjoint method

The continuous adjoint method, also known as continuous adjoint sensitivity analysis (CASA), operates by defining a convenient set of new differential equations for the adjoint variable and using this to compute the gradient in a more efficient manner. We encourage the interested reader to make the effort of following how the continuous adjoint method follows the same logic than the discrete adjoint method, but where the discretization of the differential equation does not happen until the very last step, when the solutions of the differential equations involved need to be numerically evaluated. The Lagrangian derivation of the continuous adjoint method can also be found in Appendix A.

Consider an integrated loss function defined in Equation (9) of the form

$$L(u; \theta) = \int_{t_0}^{t_1} h(u(t; \theta), \theta) dt \quad (52)$$

and its derivative with respect to the parameter θ given by the following integral involving the sensitivity matrix $s(t)$:

$$\frac{dL}{d\theta} = \int_{t_0}^{t_1} \left(\frac{\partial h}{\partial \theta} + \frac{\partial h}{\partial u} s(t) \right) dt. \quad (53)$$

Just as in the case of the discrete adjoint method, the complicated term to evaluate in the last expression is the sensitivity (Equation (16)). Again, the trick is to evaluate the VJP $\frac{\partial h}{\partial u} \frac{\partial u}{\partial \theta}$ by first defining an intermediate adjoint variable. The continuous adjoint equation now is obtained by finding the dual/adjoint equation of the sensitivity equation using the weak formulation of Equation (38). The adjoint equation is obtained by writing the sensitivity equation in the form

$$\int_{t_0}^{t_1} \lambda(t)^T \left(\frac{ds}{dt} - f(u, \theta, t) s - \frac{\partial f}{\partial \theta} \right) dt = 0, \quad (54)$$

where this equation must be satisfied for every function $\lambda(t)$ in order for Equation (83) to be true. The next step is to get rid of the time derivative applied to the sensitivity $s(t)$ using integration by parts:

$$\int_{t_0}^{t_1} \lambda(t)^T \frac{ds}{dt} dt = \lambda(t_1)^T s(t_1) - \lambda(t_0)^T s(t_0) - \int_{t_0}^{t_1} \frac{d\lambda^T}{dt} s(t) dt. \quad (55)$$

Replacing this last expression into Equation (54) we obtain

$$\int_{t_0}^{t_1} \left(-\frac{d\lambda^T}{dt} - \lambda(t)^T f(u, \theta, t) \right) s(t) dt = \int_{t_0}^{t_1} \lambda(t)^T \frac{\partial f}{\partial \theta} dt - \lambda(t_1)^T s(t_1) + \lambda(t_0)^T s(t_0). \quad (56)$$

At first glance, there is nothing particularly interesting about this last equation. However, both Equations (53) and (56) involve a VJP with $s(t)$. Since Equation (56) must hold for every function $\lambda(t)$, we can pick $\lambda(t)$ to make the terms involving $s(t)$ in Equations (53) and (56) to perfectly coincide. This is done by defining the adjoint $\lambda(t)$ to be the solution of the new system of differential equations

$$\frac{d\lambda}{dt} = -f(u, \theta, t)^T \lambda - \frac{\partial h}{\partial u} \quad \lambda(t_1) = 0. \quad (57)$$

Notice that the adjoint equation is defined with the final condition at t_1 , meaning that it needs to be solved backwards in time. The definition of the adjoint $\lambda(t)$ as the solution of this last ODE simplifies Equation (56) to

$$\int_{t_0}^{t_1} \frac{\partial h}{\partial u} s(t) dt = \lambda(t_0)^T s(t_0) + \int_{t_0}^{t_1} \lambda(t)^T \frac{\partial f}{\partial \theta} dt. \quad (58)$$

Finally, replacing this inside the expression for the gradient of the loss function we have

$$\frac{dL}{d\theta} = \lambda(t_0)^T s(t_0) + \int_{t_0}^{t_1} \left(\frac{\partial h}{\partial \theta} + \lambda^T \frac{\partial f}{\partial \theta} \right) dt \quad (59)$$

The full algorithm to compute the full gradient $\frac{dL}{d\theta}$ can be described as follows:

1. Solve the original differential equation $\frac{du}{dt} = f(u, t, \theta)$;
2. Solve the backwards adjoint differential equation given by Equation (57);
3. Compute the gradient using Equation (59).

More general recipes for deriving continuous adjoint methods exists, including generalizations for PDEs. The adjoint methods can be formulated as ...

3.9 Mathematical comparison of the methods

In Sections 3.2-3.8 we focused in merely introducing each one of the sensitivity methods classified in Figure 1 as separate methods, postponing the discussion about their points in common. In this section, we are going to compare one-to-one these methods and show parallelism across them. We hope this enlightens the discussion on sensitivity methods and helps demystified misconceptions around the sometimes apparent differences across methods. This comparison will be strengthened again later in Section 4, where we will see how even small differences between methods can be translated to different software implementations with different advantages.

3.9.1 Forward AD and complex step differentiation

Notice that both AD based on dual number and complex-step differentiation introduce an abstract unit (ϵ and i , respectively) associated with the imaginary part of the extender value that carries forward the numerical value of the gradient. Although these methods seem similar, we emphasize that AD gives the exact gradient, whereas complex step differentiation relies on numerical approximations that are valid only when the stepsize ϵ is small. In Figure 4 we show how the calculation of the gradient of the function $\sin(x^2)$ is performed by these two methods. Whereas the second component of the dual number has the exact derivative of the function $\sin(x^2)$ with respect to x , it is not until we take $\epsilon \rightarrow 0$ that we obtain the derivative in the imaginary component for the complex step method. The dependence of the complex step differentiation method on the step size gives it a closer resemblance to finite difference methods than to AD using dual numbers. Further notice the complex step method involves more terms in the calculation, a consequence of the fact that second order terms of the form $i^2 = -1$ are transferred to the real part of the complex number, while for dual numbers the terms associated to $\epsilon^2 = 0$ vanish.¹⁵⁵

3.9.2 AD and symbolic differentiation with sparsity

When sparsity patterns of the gradient/Jacobian to be computed are known, symbolic differentiation can be more efficient than AD.¹³¹ In Section 3.3.3 we discuss how colored AD can be used to efficiently compute a sparse Jacobian. However, colored AD has the limitation that an extremely sparse matrix can have no rows or columns independent of each other. Consider the arrowhead matrix given by

$$J_{\text{arrowhead}} = \begin{bmatrix} \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & & & \\ \bullet & & \bullet & & \\ \bullet & & & \bullet & \\ \bullet & & & & \bullet \end{bmatrix}. \quad (60)$$

In this case, both reverse mode AD and forward mode AD will have to perform n VJPs and JVPs, respectively (for $J_{\text{arrowhead}} \in \mathbb{R}^{n \times n}$ arrowhead matrix), and there is no computational benefit of coloring the matrix here. Instead, symbolic differentiation constructs a symbolic representation of the sparse Jacobian and can fill the

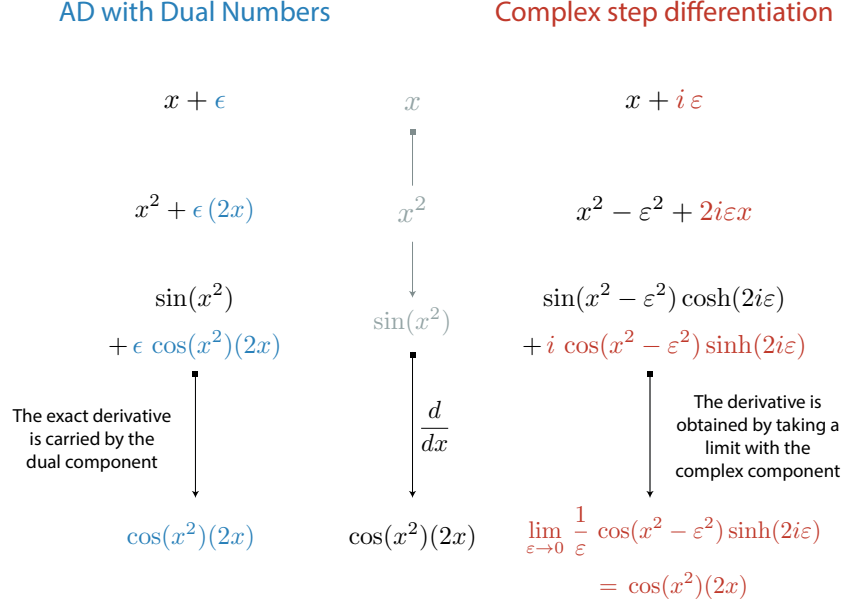


Figure 4: Comparison between AD implemented with dual numbers and complex step differentiation. For the simple case of the function $f(x) = \sin(x^2)$, we can see how each operation is carried in the forward step by the dual component (blue) and the complex component (red). Whereas AD gives the exact gradient at the end of the forward run, in the case of the complex step method we need to take the limit in the imaginary part.

Jacobian with $n + 2 \cdot (n - 1)$ computations, where each computation is significantly cheaper than each VJP or JVP.

Consider the example given in⁹⁵

$$L(x) = \ln \Psi_0 = \ln \sum_{i=1}^d w_i \Psi_i(e^{z_i \sin(x)}) \quad (61)$$

with $0 \geq w_i$ and $z_i \in \mathbb{R}$ some coefficients and $\Psi_i : \mathbb{R} \mapsto \mathbb{R}_{\geq 0}$. Computer algebra systems are able to handle the common sub-expression $\cos(x)/\Psi_0$ in $L'(x)$ to reduce the total number of multiplication involved in the derivative with respect to both forward and reverse AD, which instead will evaluate the sub-product with respect to the term $\cos(x)/\Psi_0$.

3.9.3 Discrete adjoints and reverse AD

Both discrete adjoint methods and reverse AD are classified as reverse and discrete methods (see Figure 1). Furthermore, both methods introduce an intermediate adjoint associated to the partial derivative of the loss function (output variable) with respect to intermediate variables of the forward computation. In the case of reverse AD, this adjoint is defined with the notation \bar{w} (Equation (26)), while in the discrete adjoint method this correspond to each one of the variables $\lambda_1, \lambda_2, \dots, \lambda_N$ (Equation (50)). In this section we will show that both methods are mathematically equivalent,^{140,260} but naive implementations using reverse AD can result in sub-optimal performances compared to the one obtained by directly employing the discrete adjoint method.¹⁰

In order to have a better idea of how this work in the case of a numerical solver, let us consider again the case of a one-step explicit method, non necessarily linear, where the updates u_i satisfy the equation $u_{i+1} = g_i(u_i, \theta)$. Following the same schematics than in Figure 3, we represent the computational graph of the numerical method in Figure using the intermediate variables g_1, g_2, \dots, g_{N-1} . The dual/adjoint variables

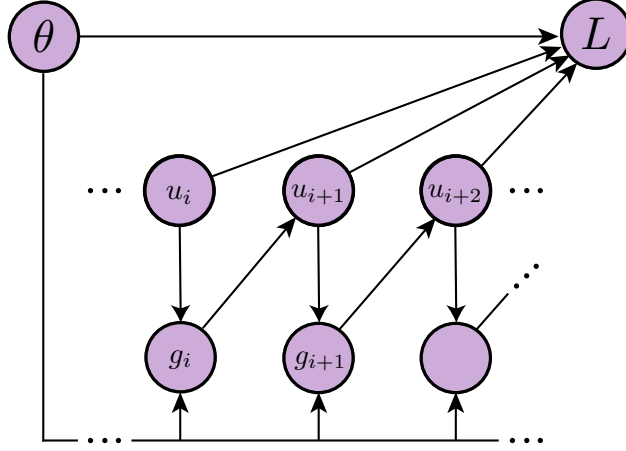


Figure 5: Computational graph associated to the discrete adjoint method. Reverse AD applied on top of the computational graph leads to the update rules for the discrete adjoint. The adjoint variable λ_i in the discrete adjoint method coincides with the adjoint variable \bar{g}_i defined in the backpropagation step.

defined in reverse AD in this computational graph are given by

$$\bar{g}_i = \frac{\partial u_{i+1}}{\partial g_i} \bar{u}_{i+1} = \frac{\partial L}{\partial u_{i+1}} + \frac{\partial g_{i+1}}{\partial u_i} \bar{g}_{i+1}. \quad (62)$$

The updates of \bar{g}_i then mathematically coincide with the updates in reverse mode of the adjoint variable λ_i (see Equation (51)).

Modern numerical solvers use functions g_i that correspond to nested functions, meaning $g_i = g_i^{(k_i)} \circ g_i^{(k_i-1)} \circ \dots \circ g_i^{(1)}$. This is certainly the case for implicit methods when u_{i+1} is the solution of an iterative Newton's method (Equation (5)); or in cases where the numerical solver includes internal iterative sub-routines.¹⁰ If the number of intermediate function is large, reverse AD will result in a large computational graph, potentially leading to excessive memory usage and slow computation.^{10,152} A solution to this problem is to introduce a customized *super node* that directly encapsulates the contribution to the full adjoint in g_i without computing the adjoint for each intermediate function $g_i^{(j)}$. Provided with the value of the Jacobian matrices $\frac{\partial g_i}{\partial u_i}$ and $\frac{\partial g_i}{\partial \theta}$, we can use the implicit function theorem to find the

In both cases, the discrete adjoint method can be implemented directly on top of a reverse AD tool that allows customized adjoint calculation.¹⁹⁰ Furthermore, notice that instead of the full Jacobian, reverse methods only required to compute VJPs.

3.9.4 Consistency: forward AD and sensitivity equations

The sensitivity equations can also be solved in discrete forward mode by numerically discretizing the original ODE and later deriving the discrete sensitivity equations.¹⁴⁹ For most cases, this leads to the same result than in the continuous case.²⁵⁹ We can numerically solve for the sensitivity s by extending the parameter θ to a multidimensional dual number

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \longrightarrow \begin{bmatrix} \theta_1 + \epsilon_1 \\ \theta_2 + \epsilon_2 \\ \vdots \\ \theta_p + \epsilon_p \end{bmatrix} \quad (63)$$

where $\epsilon_i \epsilon_j = 0$ for all pairs of i and j (see Section 3.3.1.1). The dependency of the solution u of the ODE on the parameter θ is now expanded following Equation (23) as

$$u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \rightarrow \begin{bmatrix} u_1 + \sum_{j=1}^p \frac{\partial u_1}{\partial \theta_j} \epsilon_j \\ u_2 + \sum_{j=1}^p \frac{\partial u_2}{\partial \theta_j} \epsilon_j \\ \vdots \\ u_p + \sum_{j=1}^p \frac{\partial u_n}{\partial \theta_j} \epsilon_j \end{bmatrix} = u + s \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_p \end{bmatrix}, \quad (64)$$

that is, the dual component of the vector u corresponds exactly to the sensitivity matrix s . This implies forward AD applied to any multistep linear solver will result in the application of the same solver to the sensitivity equation (Equation (38)). For example, for the forward Euler method this gives

$$\begin{aligned} u^{t+1} + s^{t+1} \epsilon &= u^t + s^t \epsilon + \Delta t f(u^t + s^t \epsilon, \theta + \epsilon, t) \\ &= u^t + f(u^t, \theta, t) + \Delta t \left(\frac{\partial f}{\partial u} s^t + \frac{\partial f}{\partial \theta} \right) \epsilon. \end{aligned} \quad (65)$$

The dual component corresponds to the forward Euler discretization of the sensitivity equation, with s^t the temporal discretization of the sensitivity $s(t)$.

The consistency result for discrete and continuous methods also holds for Runge-Kutta methods.²⁴¹ When introducing dual numbers, the Runge-Kutta scheme in Equation (2) gives the following identities

$$u^{n+1} + s^{n+1} \epsilon = s_n + \Delta t_n \sum_{i=1}^s b_i \dot{k}_i \quad (66)$$

$$k_i + \dot{k}_i \epsilon = f \left(u^n + \sum_{j=1}^s a_{ij} k_j + \left(s^n + \sum_{j=1}^s a_{ij} \dot{k}_j \right) \epsilon, \theta + \epsilon, t_n + c_i \Delta t_n \right) \quad (67)$$

with \dot{k}_i the dual variable associated to k_i . The partial component in Equation (67) carrying the coefficient ϵ gives

$$\begin{aligned} \dot{k}_i &= \frac{\partial f}{\partial u} \left(u^n + \sum_{j=1}^s a_{ij} k_j, \theta, t_n + c_i \Delta t_n \right) \left(s^n + \sum_{j=1}^s a_{ij} \dot{k}_j \right) \\ &+ \frac{\partial f}{\partial \theta} \left(u^n + \sum_{j=1}^s a_{ij} k_j, \theta, t_n + c_i \Delta t_n \right), \end{aligned} \quad (68)$$

which coincides with the Runge-Kutta scheme we will obtain for the original sensitivity equation. This means that forward AD on Runge-Kutta methods leads to solutions for the sensitivity that have the same convergence properties of the forward solver.

3.9.5 Consistency: discrete and continuous adjoints

As previously mentioned, the difference between the discrete and continuous adjoint methods is that the former follows the discretize-then-differentiate approach (also known as finite difference of adjoints²¹⁶). In contrast, continuous adjoint equations are derived by directly operating on the differential equation, without a priori consideration of the numerical scheme used to solve it. In some sense then, we can think of the discrete adjoint $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)$ in Equation (50) as the discretization of the continuous adjoint $\lambda(t)$.

A natural question then is if these two methods effectively compute the same gradient, i.e., if the discrete adjoint consistently approximate its continuous counterpart. Furthermore, since the continuous adjoint method requires to numerically solve the adjoint, we are interested in the relative accuracy of the forward and backwards step. It has been shown that for both explicit and implicit Runge-Kutta methods, as long as the coefficients in the numerical scheme given in Equation (2) satisfy the condition $b_i \neq 0$, $i = 1, 2, \dots, s$, then the ...^{96,205,206,241}

For multilinear solver, ...

Although the discrete adjoint has been shown to be consistent with the continuous adjoint method, the code generated when using the discrete adjoint using AD tools (see Section 3.9.3) can be sub-optimal and manual modification of the differentiation code are required to guarantee correctness.^{9,57}

Furthermore, adjoint methods can fail in chaotic systems.²⁴⁵ In the more general case, both methods can give different computational results.²¹⁶

Some works have shown that continuous adjoints can lead to unstable sensitivities.¹¹⁴

When using Runge-Kutta methods, it is important that both the forward and backwards solver use the same Runge-Kutta coefficients (Equation (2)) in order to both have the same level of accuracy.¹⁰

4 Implementation: A computer science perspective

In this section, we address how these different methods are computationally implemented and how to decide which method to use depending on the scientific task. In order to address this point, it is important to make one further distinction between the methods introduced in Section 3, i.e., between those that apply direct differentiation at the algorithmic level and those that are based on numerical solvers. The former are easier to implement since they are agnostic with respect to the mathematical and numerical properties of the ODE; however, they tend to be either inaccurate, memory-expensive, or at times unfeasible for large models. The latter family of methods that are based on numerical solvers include the sensitivity equations and the adjoint methods, both discrete and continuous; they are more difficult to implement and for real case applications require complex software implementations, but they are also more accurate and adequate. This section is then divided in two parts:

- **Direct methods.** (Section 4.1) Their implementation occurs at a higher hierarchy than the numerical solver software. It includes finite differences, AD, complex step differentiation.
- **Solver-based methods.** Their implementation occurs at the same level of the numerical solver. It includes
 - sensitivity equations (Section 4.2.1)
 - adjoint methods: discrete and continuous (Section 4.2.2)

Despite the fact that these methods can be (in principle) implemented in different programming languages, here we decided to use the Julia programming language for the different examples. Julia is a recent but mature programming language that has already a large tradition in implementing packages aiming to advance differentiable programming,^{23,24} which a strong emphasis in differential equation solvers¹⁸⁸ and sensitivity analysis.¹⁸⁹ Nevertheless, in reviewing existing work, we will also point to applications developed in other programming languages.

The GitHub repository `DiffEqSensitivity-Review` contains both text and code used to generate this manuscript. See Appendix C for a complete description of the scripts provided. The symbol ♣ will be used to reference code generated figures.

4.1 Direct methods

Direct methods are implemented independent of the structure of the differential equation and the numerical solver used to solve it. These include finite differences, complex step differentiation, and both forward and reverse mode AD.

4.1.1 Finite differences

Finite differences are easy to implement manually, do not require much software support, and provide a direct way of approximating a gradient. In Julia, these methods are implemented in `FiniteDiff.jl` and `FiniteDifferences.jl`, which already include subroutines to determine step-sizes. However, finite differences are less accurate and as costly as forward AD⁹³ and complex-step differentiation. Figure 6 illustrates the error in computing the gradient of a simple loss function for both true analytical solution

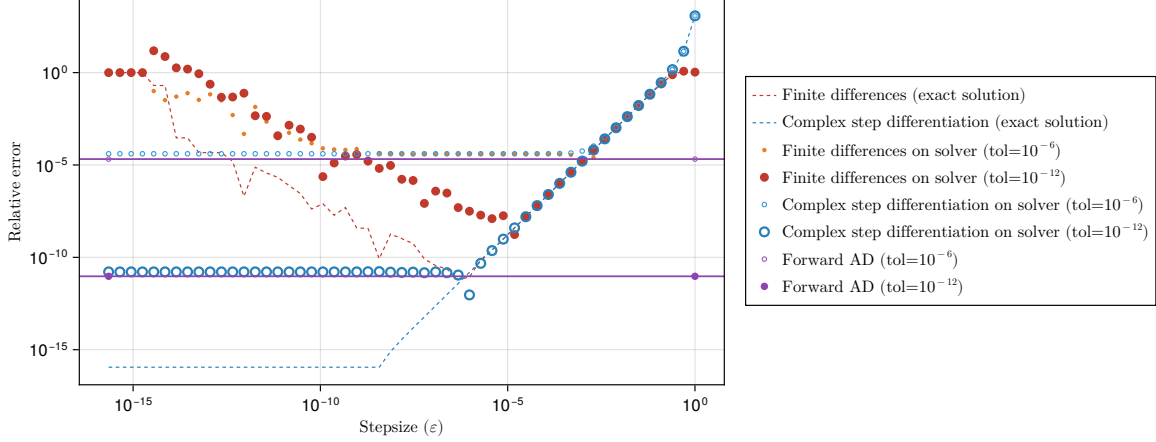


Figure 6: Absolute relative error when computing the gradient of the function $u(t) = \sin(\omega t)/\omega$ with respect to ω at $t = 10.0$ as a function of the stepsize ε . Here, $u(t)$ corresponds to the solution of the differential equation $u'' + \omega^2 u = 0$ with initial condition $u(0) = 0$ and $u'(0) = 1$. The blue dots correspond to the case where the relative error is computed with finite differences. The red and orange lines are for the case where $u(t)$ is numerically computed using the default Tsitouras solver²³² from `OrdinaryDiffEq.jl` using different tolerances. The error when using a numerical solver is larger and it is dependent on the numerical precision of the numerical solver. See ♣₁ in Appendix for code.

and numerical solution of a system of ODEs as a function of the stepsize ε using finite differences. Here we consider the solution of the differential equation $u'' + \omega^2 u = 0$ with initial condition $u(0) = 0$ and $u'(0) = 1$, which has analytical solution $u_{\text{true}}(t) = \sin(\omega t)/\omega$. The numerical solution $u_{\text{num}}(t)$ can be founded by solving the system of ODEs

$$\begin{cases} \frac{du}{dt} = v & u(0) = 0 \\ \frac{dv}{dt} = -\omega^2 u & v(0) = 1. \end{cases} \quad (69)$$

The loss function used to differentiate is given by $L(\theta) = u(10)$. We see that finite differences are inaccurate for computing the derivative of u_{true} with respect to ω (that is, $u'_{\text{true}} = \cos(\omega t) - \sin(\omega t)/\omega^2$) when the stepsize ε is both too small and too large (red dashed line). When the derivative is instead computed using the numerical solution $u_{\text{num}}(t)$ is instead computed with a numerical solver (red circles), the accuracy of the derivative further deteriorates due to approximation errors in the solver. This effect is dependent on the numerical solver tolerance. Notice that in general we do not know the true solution of the differential equation, so the derivative of u_{true} obtained using finite differences just serves as a lower bound of the error we expect to see when performing sensitivity analysis on top of the numerical solver.

4.1.2 Automatic differentiation

The AD algorithms described in Section 3.3 implement algorithmic differentiation using different methods, namely *operator overloading* for AD based on dual numbers, and *source code transformation* for both forward and reverse AD based on the computational graph.¹⁵⁵ In this section we are going to cover how Forward AD is implemented using dual numbers, postponing the implementation of AD based on computational graphs for Reverse AD (Section 4.1.2.2).

4.1.2.1 Forward AD based on dual numbers

Implementing forward AD using dual numbers is usually carried out using *operator overloading*.¹⁷⁰ This means expanding the object associated to a numerical value to include the tangent and extending the definition of atomic algebraic functions. In Julia, this can be done by relying on multiple dispatch.²³ The following example illustrates how to define a dual number and its associated binary addition and multiplication extensions.

```

using Base: @kwdef

@kwdef struct DualNumber{F <: AbstractFloat}
    value::F
    derivative::F
end

# Binary sum
Base.:(+)(a::DualNumber, b::DualNumber) = DualNumber(value = a.value + b.value,
    derivative = a.derivative + b.derivative)

# Binary product
Base.:*(a::DualNumber, b::DualNumber) = DualNumber(value = a.value * b.value,
    derivative = a.value*b.derivative + a.derivative*b.value)

```

We further overload base operations for this new type to extend the definition of standard functions by simply applying the chain rule and storing the derivative in the dual variable following Equation (23):

```

function Base.:(sin)(a::DualNumber)
    value = sin(a.value)
    derivative = a.derivative * cos(a.value)
    return DualNumber(value=value, derivative=derivative)
end

```

In the Julia ecosystem, `ForwardDiff.jl` implements forward mode AD with multidimensional dual numbers¹⁹⁷ and the sensitivity method `ForwardDiffSensitivity` implements forward differentiation inside the numerical solver using dual numbers. Figure 6 shows the result of performing forward AD inside the numerical solver. We can see that for this simple example Forward AD performs as good as the best output of finite differences and complex step differentiation (see Section 4.1.3) when optimizing by the stepsize ε . Further notice that AD is not subject to numerical errors due to floating point arithmetic.⁹⁵

Implementations of forward AD using dual numbers and computational graphs require a number of operations that increases with the number of variables to differentiate, since each computed quantity is accompanied by the corresponding gradient calculations.⁹³ This consideration also applies to the other forward methods, including finite differences and complex-step differentiation, which makes forward models prone to the curse of dimensionality with respect to the number of parameters considered.

Although AD is always algorithmically correct, when combined with a numerical solver *AD can be numerically incorrect* and result in wrong gradient calculations.⁵⁷ To illustrate this point, consider the following example of a simple system of ODEs

$$\begin{cases} \frac{du_1}{dt} = au_1 - u_1u_2 & u_1(0) = 1 \\ \frac{du_2}{dt} = -au_2 + u_1u_2 & u_2(0) = 1 \end{cases} \quad (70)$$

with a the parameter with respect to which we want to differentiate. In the simple case of $a = 1$, the solutions of the ODE are constant functions $u_1(t) \equiv u_2(t) \equiv 1$. As explained in Appendix B, traditional adaptive stepsize solvers used for just solving ODEs are designed to control for numerical errors in the ODE solution but not in its sensitivities when coupled with an internal AD method. This can lead to wrong gradient calculations that propagate through the numerical solver without further warning.

4.1.2.2 Reverse AD based on computational graph

In contrast to finite differences, forward AD, and complex-step differentiation, reverse AD is the only of this family of methods that propagates the gradient in reverse mode by relying on analytical derivatives of primitive functions, which in Julia are available via `ChainRules.jl`. Since this requires the evaluation intermediate variables, reverse AD requires a more delicate protocol of how to store intermediate variables in memory and make them accessible during the backwards pass.

Backwards AD can be implemented via *pullback* functions,¹¹⁰ a method also known as *continuation-passing style*.²⁴² In the backward step, this executes a series of function calls, one for each elementary operation. If one of the nodes in the graph w is the output of an operation involving the nodes v_1, \dots, v_m , where $v_i \rightarrow w$ are all nodes in the graph, then the pullback $\bar{v}_1, \dots, \bar{v}_m = \mathcal{B}_w(\bar{w})$ is a function that accepts

gradients with respect to w (defined as \bar{w}) and returns gradients with respect to each v_i (defined as \bar{v}_i) by applying the chain rule. Consider the example of the multiplicative operation $w = v_1 \times v_2$. Then

$$\bar{v}_1, \bar{v}_2 = v_2 \times \bar{w}, v_1 \times \bar{w} = \mathcal{B}_w(\bar{w}), \quad (71)$$

which is equivalent to using the chain rule as

$$\frac{\partial \ell}{\partial v_1} = \frac{\partial}{\partial v_1}(v_1 \times v_2) \frac{\partial \ell}{\partial w}. \quad (72)$$

A crucial distinction between AD implementations based on computational graphs is between *static* and *dynamical* methods.²⁰ In the case of a static implementations, the computational graph is constructed before any code is executed, which are encoded and optimized for performance within the graph language. For static structures such as neural networks, this is ideal.¹ However, two major drawbacks of static methods are composability with existing code, including support of custom types, and adaptive control flow, which is a common feature of numerical solvers. These issues are addressed in reverse AD implementations using *tracing*, where the program structure is transformed into a list of pullback functions that built the graph dynamically at runtime. Popular libraries in this category are `Tracker.jl` and `ReverseDiff.jl`. There also exist source-to-source AD system that achieve highest performance at the same time they support arbitrary control flow structure. These include `Zygote.jl`,¹¹¹ `Enzyme.jl`,¹⁶⁸ and `Diffractor.jl`. The existence of multiple AD packages lead to the development of `AbstractDifferentiation.jl` which allows to combine different methods.²⁰⁹

4.1.2.3 Checkpointing

In opposition to forward methods, all reverse methods, including backpropagation and adjoint methods, require to access For a numerical solver, the number of memory required can be very large, involving a total of at least $\mathcal{O}(nk)$ terms, with k the number of steps of the numerical solver. Checkpointing is a technique that can be used for all the backwards methods that avoids storing all the intermediate states by balancing storing and recomputation to recover the required state exactly.⁹⁵

Also known as windowing, checkpointing is a technique that sets a trade-off between memory and time by saving intermediate states of the solution in the forward pass and recalculating the solution between intermediate states in the backwards mode.^{95,210}

4.1.3 Complex step differentiation

Modern software already have support for complex number arithmetic, making complex step differentiation very easy to implement. In Julia, complex analysis arithmetic can be easily carried inside the numerical solver. The following example shows how to extend the numerical solver used to solve the ODE in Equation (69) to support complex numbers.

```
function dyn!(du::Array{Complex{Float64}}, u::Array{Complex{Float64}}, p, t)
    ω = p[1]
    du[1] = u[2]
    du[2] = - ω^2 * u[1]
end

tspan = [0.0, 10.0]
du = Array{Complex{Float64}}{([0.0])}
u0 = Array{Complex{Float64}}{([0.0, 1.0])}

function complexstep_differentiation(f::Function, p::Float64, ε::Float64)
    p_complex = p + ε * im
    return imag(f(p_complex)) / ε
end

complexstep_differentiation(x -> solve(ODEProblem(dyn!, u0, tspan, [x]), Tsit5())
    ).u[end][1], 20., 1e-3)
```

Figure 6 further shows the result of performing complex step differentiation using the same example as in Section 4.1.1. We can see from both exact and numerical solution that complex-step differentiation does not

suffer from small values of ε , meaning that ε can be chosen arbitrarily small¹⁵⁵ as long as it does not reach the underflow threshold.⁹⁰ Notice that for large values of the stepsize ε complex step differentiation gives similar results than finite differences, while for small values of ε the performance of complex step differentiation is slightly worse than forward AD. This result emphasize the observation made in Section 3.9.3, complex step differentiation has many aspects in common with finite differences and AD based on dual numbers.

However, the difference between the methods also makes the complex step differentiation sometimes more efficient than both finite differences and AD,¹³¹ an effect that can be counterbalanced by the number of extra unnecessary operations that complex arithmetic requires (see last column in Figure 4).¹⁵⁴

4.2 Solver-based methods

Sensitivity methods based on numerical solvers tend to be better adapted to the structure and properties of the underlying ODE (stiffness, stability, accuracy) but are also more difficult to implement. This difficulty arises from the fact that the sensitivity method needs to deal with some numerical and computational considerations, including i) how to handle matrix/Jacobian-vector products; ii) numerical stability of the forward/reverse solver; and iii) memory-time tradeoff. These factors are further exacerbated by the number of ODEs and parameters in the model. Just a few modern scientific software have the capabilities of solving forward ODE and computing their sensitivities at the same time. These include CVODES within SUNDIALS in C;^{107,213} ODESSA¹³⁷ and FATODE (discrete adjoints)²⁵⁹ both in Fortran; `SciMLSensitivity.jl` in Julia;¹⁸⁹ `Dolfin-adjoint` based on the FEniCS Project;^{64,159} `DENSERKS` in Fortran;⁹ `DASPKADJOINT`;⁴¹ and `DiffRax` in Python.¹²⁵

It is important to remark that the underlying machinery of all solvers relies on solvers for linear systems of equations, which can be solved in dense, band (sparse), and Krylov mode. Another important consideration is that all these methods have subroutines to compute the VJPs involved in the sensitivity and adjoint equations. This calculation is carried out by another sensitivity method, usually finite differences or AD, which plays a central role when analyzing the accuracy and stability of the adjoint method.

4.2.1 Sensitivity equation

For systems of equations with few number of parameters, this method is useful since the system of $n(p + 1)$ equations composed by Equations (1) and (38) can be solved using the same precision for both solution and sensitivity numerical evaluation. Furthermore, this is a forward method and it does not required saving the solution in memory. The simplicity of the sensitivity method makes it available in most software for sensitivity analysis. In Julia, the `ForwardSensitivity` methods implement continuous sensitivity analysis, which performs forward AD on the solver via `ForwardDiff.jl` (see Section 3.9.4).

However, for stiff systems of ODEs the use of the sensitivity equations is unfeasible.¹²⁷ For systems of stiff-ODEs, the cost of numerically stable solvers is cubic in the number of ODEs,²⁴⁸ making the complexity of the sensitivity method $\mathcal{O}(n^3 p^3)$. This complexity makes this method useless for models with many ODEs and/or parameters.

4.2.1.1 Computing VJPs inside the solver

All the solver-based methods has to faced the challenge of how to compute large VJPs. In the case of the sensitivity equation, this correspond to the row/column terms in $\frac{\partial f}{\partial u} s$ in Equation (38). For the adjoint equations, although an efficient trick has been used to remove the computationally expensive VJP, we still need to evaluate the term $\lambda^T \frac{\partial G}{\partial \theta}$ for the discrete adjoint method in Equation (45), and $\lambda^T \frac{\partial f}{\partial \theta}$ for the continuous adjoint method in Equation (59). Therefore, the choice of the specific algorithm to compute VJPs can have significant impact in the overall performance.

In SUNDIALS, the VJPs involved in the sensitivity and adjoint method are handled using finite differences unless specified by the user.¹⁰⁷ In FATODE, these can be computed with finite differences, AD or provides by the user. In the Julia ecosystem, different AD packages are available for this task (see Section 4.1.2.2), including `ForwardDiff.jl`, `ReverseDiff.jl`, `Zygote.jl`,¹¹¹ `Enzyme.jl`,¹⁶⁸ `Tracker.jl`. These will compute the VJPs using some for of AD, which will result in correct calculations but potentially sub-optimal code. For these cases, customized VJPs function can be passed to the sensitivity methods using the `autojacvec`.

4.2.2 Adjoint methods

For complex and large systems, direct methods for computing the gradient on top of the numerical solver can be memory expensive due to the large number of function evaluations required by the solver and the later store of the intermediate states. For these cases, adjoint-based methods allows us to compute the gradients of a loss function by instead computing the adjoint that serves as a bridge between the solution of the ODE and the final gradient. If well the adjoint method offers considerate advantages when working with complex systems, since we are dealing with a new differential equation special care needs to be taken with respect to numerical efficiency and stability.

4.2.2.1 Discrete adjoint method

As we discuss in Section 3.9.3, the discrete adjoint method can be directly been implemented using reverse AD. In the Julia SciML ecosystem, `ReverseDiffAdjoint` performs reverse AD on the numerical solver via `ReverseDiff.jl`; `ZygoteAdjoint` via `Zygote.jl`; and `TrackerAdjoint` via `Tracker.jl`. In all these cases, a custom pullback function needs to be specified that specifies how VJPs are computed thought the numerical solver.¹⁹⁰

4.2.2.2 Continuous adjoint method

	Method	Stability	Stiff Performance	Memory
Discrete	ReverseDiffAdjoint	Best	$\mathcal{O}(n^3 + p)$	High
	ZygoteAdjoint			
	TrackerAdjoint	Best	$\mathcal{O}(n^3 + p)$	High
Continuous	Sensitivity equation	Good	$\mathcal{O}(n^3 p^3)$	$\mathcal{O}(1)$
	Backsolve adjoint \triangleleft	Poor	$\mathcal{O}((n + p)^3)$	$\mathcal{O}(1)$
	Backsolve adjoint \blacktriangleleft	Medium	$\mathcal{O}((n + p)^3)$	$\mathcal{O}(K)$
	Interpolating adjoint \triangleleft	Good	$\mathcal{O}((n + p)^3)$	High
	Interpolating adjoint \blacktriangleleft	Good	$\mathcal{O}((n + p)^3)$	$\mathcal{O}(K)$
	Quadrature adjoint	Good	$\mathcal{O}(n^3 + p)$	High
	Gauss adjoint	...	$\mathcal{O}(n^3 + p)$..

Table 1: Methods that are being checkpointed are indicated with the symbol \blacktriangleleft , with the number K corresponding to the number of checkpoints.

The continuous adjoint methods offers a series of advantages over the discrete method and the rest of the forward methods previously discussed. Just as the discrete adjoint methods and backpropagation, the bottleneck is how to solve for the adjoint $\lambda(t)$ due to its dependency with VJPs involving the state $u(t)$. Effectively, notice that Equation (57) involves the terms $f(u, \theta, t)$ and $\frac{\partial h}{\partial u}$, which are both functions of $u(t)$. In opposition to the discrete adjoint methods, notice that here the full continuous trajectory $u(t)$ is needed, instead of its discrete pointwise evaluation. There are two principal ways of addressing the evaluation of $u(t)$ during the backwards step.

- (i) **Interpolation.** During the forward model, we can store in memory intermediate states of the numerical solution that allow the dense evaluation of the numerical solution at any given time, which is a requirement of continuous methods in opposition to discrete. This can be done using dense output formulas, for example by adding extra stages to the Runge-Kutta scheme (Equation (2)) that allows to define a continuous interpolation, a method known as continuous Runge-Kutta.^{10,248} When using

checkpointing, intermediate variables are saved and the interpolation between them is re-computed on demand.

- (ii) **Backsolving.** Solve again the original ODE together with the adjoint as the solution of the reversed augmented system⁴⁵

$$\frac{d}{dt} \begin{bmatrix} u \\ \lambda \\ \frac{dL}{d\theta} \end{bmatrix} = \begin{bmatrix} -f \\ -\frac{\partial f}{\partial u}^T \lambda - \frac{\partial h}{\partial y}^T \\ -\lambda^T \frac{\partial f}{\partial \theta} - \frac{\partial h}{\partial \theta} \end{bmatrix} \quad \begin{bmatrix} u \\ \lambda \\ \frac{dL}{d\theta} \end{bmatrix} (t_1) = \begin{bmatrix} u(t_1) \\ \frac{\partial L}{\partial u(t_1)} \\ \lambda(t_0)^T s(t_0) \end{bmatrix}. \quad (73)$$

An important problem with this approach is that computing the ODE backwards $\frac{du}{dt} = -f(u, \theta, t)$ can be unstable and lead to large numerical errors.^{127,261} One way of solving this system of equations that ensures stability is by using implicit methods. However, this requires cubic time in the total number of ordinary differential equations, leading to a total complexity of $\mathcal{O}((n+p)^3)$ for the adjoint method.

Both interpolating and backsolve adjoint methods can be implemented along with a checkpointing scheme. This is implemented in `Checkpointing.jl`.²¹⁰

When dealing with stiff differential equations, special considerations need to be taken into account. Two alternatives are proposed in,¹²⁷ the first referred to as *Quadrature Adjoint* produces a high order interpolation of the solution $u(t)$ as we move forward, then solve for λ backwards using an implicit solver and finally integrating $\frac{dL}{d\theta}$ in a forward step. This reduces the complexity to $\mathcal{O}(n^3 + p)$, where the cubic cost in the number of ODEs comes from the fact that we still need to solve the original stiff differential equation in the forward step. A second but similar approach is to use an implicit-explicit (IMEX) solver, where we use the implicit part for the original equation and the explicit for the adjoint. This method also has a complexity of $\mathcal{O}(n^3 + p)$.

4.2.2.3 Solving the quadrature

Another computational consideration is how the integral in Equation (59) is numerically evaluated. Some methods save computation by noticing that the last step in the continuous adjoint method of evaluating $\frac{dL}{d\theta}$ is an integral instead of an ODE, and then can be evaluated as such without the need to include it in the tolerance calculation inside the numerical solver.¹²⁶ Numerical integration, also known as quadrature integration, consists in approximating integrals by finite sums of the form

$$\int_{t_0}^{t_1} F(t)dt \approx \sum_{i=1}^K \omega_i F(\tau_i), \quad (74)$$

where the evaluation of the function occurs in certain knots $t_0 \leq \tau_1 < \dots < \tau_K \leq t_1$, and ω_i are weights. Weights and knots are obtained in order to maximize the order in which polynomials are exactly integrated.²²⁴

Different quadrature methods are based on different choices of the knots and associated weights. Between these methods, the Gaussian quadrature is the faster method to evaluate one-dimensional integrals.¹⁷⁵

5 Recommendations

For sufficient small systems of less than 100 parameters and ODEs, Forward AD is the most efficient method, outperforming sensitivity and adjoint methods.¹⁴⁹

Adjoint methods are computationally more efficient for complicated numerical solvers, but they are also more difficult to implement.

If well direct methods like AD and complex-step differentiation can be costly to use directly on a numerical solver, they can still be used to compute the sensitivities required for the calculation of the more efficient adjoint method (see discussion in Section 4.2.1.1).

As an AD algorithm, we support the use of forward and reverse AD as the method of choice over finite differences, complex step differentiation, and symbolic differentiation. The first two do not really provide an advantage over AD in terms of precision, required some tuning of the stepsize ε , and On the othe side,

if well symbolic differentiation can be more efficient in neasted cases (see ...) or when the spartity patter of the Jacobian is known, in general this advantage is not drasmatic in most real cases.

When using discrete methods (discretize-then-differentiate), it is important to be aware that the differentiation machinery is applied after algorithm to solve the differential equation has been specified, meaning that the computed derivatives are not just with respect to the numerical solution, but also with respect to the algorithm used.⁵⁷ This is certainly the case when the iterative solver have adaptive stepsize controllers that depend on the parameter to differentiate, as we explored in Section 4.1.2.1. Although some solutions has been proposed to solve this in the case of discrete methods,⁵⁷ this is a problem that continuous methods do not have since they apply the differentiation step before the numerical algorithm has been specified.

Furthermore, continuous sensitivity analysis is more efficient while discrete adjoint method is more stable (discussion in appendix of¹⁸⁹)

Generalizations to DAE⁴¹

5.1 Chaotic systems

Both forward and adjoint sensitivity analysis methods from the previous chapters encounter challenges and become less useful when applied to chaotic systems. To illustrate this, let us consider long-time-averaged quantities

$$\langle L(\theta) \rangle_T = \frac{1}{T} \int_0^T L(u(t), \theta) dt, \quad (75)$$

of chaotic systems, where $L(u(t), \theta)$ is the instantaneous objective and $u(t)$ denotes the state of the dynamical system at time t . For ergodic dynamical systems, $\lim_{T \rightarrow \infty} \langle L(\theta) \rangle_T$ depends solely on the governing dynamical system and is independent of the specific choice of trajectory $u(t)$. In particular, $\lim_{T \rightarrow \infty} \langle L(\theta) \rangle_T$ does not depend on the initial condition. Under the assumption of uniform hyperbolic systems, it is possible to derive closed-form expressions and differentiability conditions for $\langle L(\theta) \rangle_T$.^{202,203} However, computing derivatives using numerical methods of statistical quantities of the form (75) with respect to the vector parameter θ in chaotic dynamical systems remains challenging due to the *butterfly effect*, i.e. small changes in the initial state or parameter can result in large differences in a later state. As a consequence, the solutions of the forward and adjoint sensitivity equations blow up (exponentially fast) instead of converging to the actual derivative. To address these issues, various modifications and methods have been proposed, including approaches based on ensemble averages,^{61,134} the Fokker-Planck equation,^{27,230} the fluctuation-dissipation theorem,^{5,6,138} shadowing lemma,^{26,28,172–174,243,244,246} and modifications of Ruelle’s formula.^{44,171}

In Julia this is implemented in the sensitivity method `AdjointLSS` and `NILSS`

6 Conclusions

We also encourage the interested reader to direct their attention to other comprehensive works on automatic differentiation;²⁰ adjoint methods, ...

Appendices

A Lagrangian formulation of the adjoint method

The adjoint equation can be derived directly from the Following the analysis in,⁸⁶ we decided to present both approaches here, although we prefer with the duality viewpoint introduced in the main text since we believe is more commonly used and easy to understand for newcomers.

In this section we are going to derive the adjoint equation for both discrete and continuous methods using the Lagrangian formulation of the adjoint. It is important to mention that this is different than using the Lagrange multipliers approach, a common confusion in the literature.⁸⁷ Conceptually, the method is the same in both discrete and continuous case, with the difference that we manipulate linear algebra objects for the former and continuous operators for the later.

A.1 Discrete adjoint

Following the notation introduced in Section 3.7, we first define the Lagrangian function $I(U, \theta)$ as

$$I(U, \theta) = L(u; \theta) + \lambda^T G(U; \theta), \quad (76)$$

where $\lambda \in \mathbb{R}^{n_k}$ is, in principle, any choice of vector. Given that for every choice of the parameter θ we have $G(U; \theta) = 0$, we have that $I(U, \theta) = L(U, \theta)$. Then, the gradient of L with respect to the vector parameter θ can be computed as

$$\begin{aligned} \frac{dL}{d\theta} &= \frac{dI}{d\theta} = \frac{\partial L}{\partial \theta} + \frac{\partial L}{\partial U} \frac{\partial U}{\partial \theta} + \lambda^T \left(\frac{\partial G}{\partial U} + \frac{\partial G}{\partial U} \frac{\partial U}{\partial \theta} \right) \\ &= \frac{\partial L}{\partial \theta} + \lambda^T \frac{\partial G}{\partial U} + \left(\frac{\partial L}{\partial \theta} + \lambda^T \frac{\partial G}{\partial U} \right) \frac{\partial U}{\partial \theta}. \end{aligned}$$

The important trick in the last term involved grouping all the terms involving the sensitivity $\frac{\partial U}{\partial \theta}$ together. In order to avoid the computation of the sensitivity at all, we can define λ as the vector that makes the last term in the right hand side of Equation (??) which results in the same results we obtained in Equation (44) and (45), that is,

$$\frac{\partial L}{\partial \theta} = -\lambda^T \frac{\partial G}{\partial U}. \quad (77)$$

Finally, the gradient can be computed as

$$\frac{dL}{d\theta} = \frac{\partial L}{\partial \theta} + \lambda^T \frac{\partial G}{\partial U}. \quad (78)$$

A.2 Continuous adjoint

For the continuous adjoint method, we proceed the same way by defining the Lagrangian $I(\theta)$ as

$$I(\theta) = L(\theta) - \int_{t_0}^{t_1} \lambda(t)^T \left(\frac{du}{dt} - f(u, \theta, t) \right) dt \quad (79)$$

where $\lambda(t) \in \mathbb{R}^n$ is any function. It is important to mention that, in principle, there is connection yet between $\lambda(t)$ and the Lagrange multiplier associated to the constraint.⁸⁷ Instead, the condition

$$\int_{t_0}^{t_1} \lambda(t)^T \left(\frac{du}{dt} - f(u, \theta, t) \right) dt = 0 \quad \text{for all function } \lambda : [t_0, t_1] \mapsto \mathbb{R}^n \quad (80)$$

correspond to the weak formulation of the differential equation (1).³⁶ As long as the differential equation is satisfied, we have $I(\theta) = L(\theta)$ for all choices of the vector parameter θ . Now,

$$\frac{dL}{d\theta} = \frac{dI}{d\theta} = \int_{t_0}^{t_1} \left(\frac{\partial h}{\partial \theta} + \frac{\partial h}{\partial u} \frac{\partial u}{\partial \theta} \right) dt - \int_{t_0}^{t_1} \lambda(t)^T \left(\frac{d}{dt} \frac{du}{d\theta} - \frac{\partial f}{\partial u} \frac{du}{d\theta} - \frac{\partial f}{\partial \theta} \right) dt. \quad (81)$$

Notice that the term involved in the second integral is the same we found when deriving the sensitivity equations. We can derive an easier expression for the last term using integration by parts. Using our usual definition of the sensitivity $s = \frac{du}{d\theta}$, and performing integration by parts in the term $\lambda^T \frac{d}{dt} \frac{du}{d\theta}$ we derive

$$\begin{aligned} \frac{dL}{d\theta} = \int_{t_0}^{t_1} \left(\frac{\partial h}{\partial \theta} + \lambda^T \frac{\partial f}{\partial \theta} \right) dt - \int_{t_0}^{t_1} \left(-\frac{d\lambda^T}{dt} - \lambda^T \frac{\partial f}{\partial u} - \frac{\partial h}{\partial u} \right) s(t) dt \\ - \left(\lambda(t_1)^T s(t_1) - \lambda(t_0)^T s(t_0) \right). \end{aligned} \quad (82)$$

Now, we can force some of the terms in the last equation to be zero by solving the following adjoint differential equation for $\lambda(t)^T$ in backwards mode

$$\frac{d\lambda}{d\theta} = - \left(\frac{\partial f}{\partial u} \right)^T \lambda - \left(\frac{\partial h}{\partial u} \right)^T, \quad (83)$$

with final condition $\lambda(t_1) = 0$.

It is easy to see that this derivation is equivalent to solving the Karush-Kuhn-Tucker (KKT) conditions.

A.3 The adjoint from a functional analysis perspective

B When AD is algorithmically correct but numerically wrong

In this section, we are going to consider certain errors that can potentially show up when combining AD with a numerical solver. Numerical solvers for differential equations usually estimate internally a scaled error computed as define an error term computed as^{97,188}

$$\text{Err}_{\text{scaled}}^{n+1} = \left(\frac{1}{n} \sum_{i=1}^n \left(\frac{\text{err}_i^{n+1}}{\text{abstol} + \text{rtol} \times M} \right)^2 \right)^{\frac{1}{2}}, \quad (84)$$

with **abstol** and **rtol** the absolute and relative solver tolerances (customize), respectively; M is the maximum expected value of the numerical solution; and err_i^{n+1} is an estimation of the numerical error at step $n+1$. Common choices for these include $M = \max(u_i^{n+1}, \hat{u}_i^{n+1})$ and $\text{err}_i^{n+1} = u_i^{n+1} - \hat{u}_i^{n+1}$, but these can vary between solvers. Estimations of the local error err_i^{n+1} can be based on two approximation to the solution based on Runge-Kutta pairs,^{97,194} or in theoretical upper bounds provided by the numerical solver. The choice of the norm $\frac{1}{\sqrt{n}} \|\cdot\|_2$ for computing the total error $\text{Err}_{\text{scaled}}$, sometimes known as Hairer norm, has been the standard for a long time¹⁹⁴ and it is based on the assumption that a mere increase in the size of the systems of ODEs (e.g., by simply duplicating the ODE system) should not affect the stepsize choice, but other options can be considered.⁹⁷

The goal of a stepsize controller is to pick Δt_{n+1} as large as possible (so the solver requires less total steps) at the same time that $\text{Err}_{\text{scaled}} \leq 1$. One of the most used methods to archive this is the proportional-integral controller (PIC) that updates the stepsize according to^{194,248}

$$\Delta t_{n+1} = \eta \Delta t_n \quad \eta = w_{n+1}^{\beta_1/q} w_n^{\beta_2/q} w_{n-1}^{\beta_3/q} \quad (85)$$

with $w_{n+1} = 1/\text{Err}_{\text{scaled}}^{n+1}$ the inverse of the scaled error estimates; β_1, β_2 , and β_3 numerical coefficients defined by the controller; and q the order of the numerical solver. If the stepsize Δt_{n+1} proposed in Equation (85) to update from u^{n+1} to u^{n+2} does not satisfy $\text{Err}_{\text{scaled}}^{n+2} \leq 1$, a new smaller stepsize is proposed. When $\eta < 1$ (which is the case for simple controllers with $\beta_2 = \beta_3 = 0$), Equation (85) can be used for the local update. It is also common to restrict $\eta \in [\eta_{\min}, \eta_{\max}]$ so the stepsize does not change abruptly.⁹⁷

When performing forward AD thought numerical solver, the error used in the stepsize controller needs to naturally account for both the errors induced in the numerical solution of the original ODE and the errors in the dual component carrying the value of the sensitivity. This relation between the numerical solver and AD has been made explicit when we presented the relationship between forward AD and the sensitivity equations (Section 3.6, Equation (65)). To illustrate this, let us consider the simple ODE example from Section 4.1.2.1, consisting in the system of equations

$$\begin{cases} \frac{du_1}{dt} = au_1 - u_1u_2 & u_1(0) = 1 \\ \frac{du_2}{dt} = -au_2 + u_1u_2 & u_2(0) = 1. \end{cases} \quad (86)$$

Notice that for $a = 1$ this ODE admits a simple analytical solution $u_1(t) = u_2(t) = 1$ for all times t , making this problem very simple to solve numerically. The following code solves for the derivative with respect to the parameter a using two different methods¹. The second method using forward AD with dual numbers declares the `internalsnorm` argument according to Equation (84).

```
using SciMLSensitivity, OrdinaryDiffEq, Zygote, ForwardDiff

function fiip(du, u, p, t)
    du[1] = p[1] * u[1] - u[1] * u[2]
    du[2] = -p[1] * u[2] + u[1] * u[2]
end

p = [1.]
u0 = [1.0; 1.0]
prob = ODEProblem(fiip, u0, (0.0, 10.0), p);

# Correct gradient computed using
```

¹Full code available at ...


```

grad0 = Zygote.gradient(p->sum(solve(prob, Tsit5(), u0=u0, p=p, sensealg =
    ForwardSensitivity(), saveat = 0.1, abstol=1e-12, reltol=1e-12)), p)
# grad0 = ([212.71042521681443],)

# Original AD with wrong norm
grad1 = Zygote.gradient(p->sum(solve(prob, Tsit5(), u0=u0, p=p, sensealg =
    ForwardDiffSensitivity(), saveat = 0.1, internalnorm = (u,t) -> sum(abs2,u/
    length(u)), abstol=1e-12, reltol=1e-12)), p)
# grad1 = ([6278.15677493293],)

```

The reason why the two methods give different answers is because the error estimation by the stepsize controller is ignoring numerical errors in the dual component. In the later case, the local estimated error is drastically underestimated to $\text{err}_i^{n+1} = 0$, which makes the stepsize Δt_{n+1} to increase by a multiplicative factor at every step. This can be fixed by instead considering a norm that accounts for both the primal and dual components in the forward pass,

$$\text{Err}_{\text{scaled}}^{n+1} = \left[\frac{1}{n(p+1)} \left(\sum_{i=1}^n \left(\frac{u_i^{n+1} - \hat{u}_i^{n+1}}{\text{abstol} + \text{reltol} \max(u_i^{n+1}, \hat{u}_i^{n+1})} \right)^2 + \sum_{i=1}^n \sum_{j=1}^p \left(\frac{s_{ij}^{n+1} - \hat{s}_{ij}^{n+1}}{\text{abstol} + \text{reltol} \max(s_{ij}^{n+1}, \hat{s}_{ij}^{n+1})} \right)^2 \right) \right]^{\frac{1}{2}}, \quad (87)$$

which now gives the right answer

```

sse(x::Number) = x^2
sse(x::ForwardDiff.Dual) = sse(ForwardDiff.value(x)) + sum(sse, ForwardDiff.
    partials(x))

totallength(x::Number) = 1
totallength(x::ForwardDiff.Dual) = totallength(ForwardDiff.value(x)) + sum(
    totallength, ForwardDiff.ppartials(x))
totallength(x::AbstractArray) = sum(totallength,x)

grad3 = Zygote.gradient(p->sum(solve(prob, Tsit5(), u0=u0, p=p, sensealg =
    ForwardDiffSensitivity(), saveat = 0.1, internalnorm = (u,t) -> sqrt(sum(x-
    >sse(x),u) / totallength(u)), abstol=abstol, reltol=reltol)), p)
# grad3 = ([212.71042521681392],)

```

Notice that current implementations of forward AD inside `SciMLSensitivity.jl` already accounts for this and there is no need to specify the internal norm.

C Supplementary code

♣₁ 323

♣₂ 323

References

1. ABADI, M. et al. “TensorFlow: A System for Large-Scale Machine Learning”. In: *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation*. OSDI’16. Savannah, GA, USA: USENIX Association, 2016, pp. 265–283.
2. ABARBANEL, H. D. I., ROZDEBA, P. J., and SHIRMAN, S. “Machine Learning: Deepest Learning as Statistical Data Assimilation Problems”. In: *Neural Computation* 30.8 (Aug. 2018), pp. 2025–2055. doi: 10.1162/neco_a_01094.
3. ABDELHAFEZ, M., SCHUSTER, D. I., and KOCH, J. “Gradient-based optimal control of open quantum systems using quantum trajectories and automatic differentiation”. In: *Phys. Rev. A* 99.5 (2019), p. 052327. doi: 10.1103/PhysRevA.99.052327.
4. ABDELHAFEZ, M. et al. “Universal gates for protected superconducting qubits using optimal control”. In: *Physical Review A* 101.2 (2020), p. 022321. doi: 10.1103/PhysRevA.101.022321.
5. ABRAMOV, R. V. and MAJDA, A. J. “Blended response algorithms for linear fluctuation-dissipation for complex nonlinear dynamical systems”. In: *Nonlinearity* 20.12 (2007), p. 2793. doi: <https://iopscience.iop.org/article/10.1088/0951-7715/20/12/004>.
6. ABRAMOV, R. V. and MAJDA, A. J. “New approximations and tests of linear fluctuation-response for chaotic nonlinear forced-dissipative dynamical systems”. In: *Journal of Nonlinear Science* 18 (2008), pp. 303–341. doi: <https://doi.org/10.1007/s00332-007-9011-9>.
7. AIDE, T. M., CORRADA-BRavo, C., CAMPOS-CERQUEIRA, M., MILAN, C., VEGA, G., and ALVAREZ, R. “Real-Time Bioacoustics Monitoring and Automated Species Identification”. In: *PeerJ* 1.1 (July 16, 2013), e103. doi: 10.7717/peerj.103.
8. ÅKESSON, A., CURTSDOTTER, A., EKLÖF, A., EBENMAN, B., NORBERG, J., and BARABÁS, G. “The Importance of Species Interactions in Eco-Evolutionary Community Dynamics under Climate Change”. In: *Nature Communications* 12.1 (Dec. 6, 2021), p. 4759. doi: 10.1038/s41467-021-24977-x.
9. ALEXE, M. and SANDU, A. “DENSERKS: Fortran sensitivity solvers using continuous, explicit Runge-Kutta schemes”. In: (2007).
10. ALEXE, M. and SANDU, A. “Forward and adjoint sensitivity analysis with continuous explicit Runge-Kutta schemes”. In: *Applied Mathematics and Computation* 208.2 (2009), pp. 328–346. doi: 10.1016/j.amc.2008.11.035.
11. ALLAIRE, G., DAPOGNY, C., and FREY, P. “Shape optimization with a level set based mesh evolution method”. In: *Computer Methods in Applied Mechanics and Engineering* 282 (2014), pp. 22–53.
12. ALSOS, I. G. et al. *Using Ancient Sedimentary DNA to Forecast Ecosystem Trajectories under Climate Change*. preprint. In Review, Nov. 7, 2023. doi: 10.21203/rs.3.rs-3542192/v1.
13. ARRAZOLA, J. M. et al. “Differentiable quantum computational chemistry with PennyLane”. In: *arXiv* (2021). doi: 10.48550/arxiv.2111.09967.
14. ASCHER, U. M. *Numerical methods for evolutionary differential equations*. SIAM, 2008.
15. ASCHER, U. M. and GREIF, C. *A First Course in Numerical Methods*. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2011. doi: 10.1137/9780898719987.
16. BARNOSKY, A. D. et al. “Approaching a State Shift in Earth’s Biosphere”. In: *Nature* 486.7401 (2012), pp. 52–58. doi: 10.1038/nature11018.
17. BARTON, R. R. “Computing Forward Difference Derivatives In Engineering Optimization”. In: *Engineering Optimization* 20.3 (1992), pp. 205–224. doi: 10.1080/03052159208941281.
18. BAUER, F. L. “Computational Graphs and Rounding Error”. In: *SIAM Journal on Numerical Analysis* 11.1 (1974), pp. 87–96. doi: 10.1137/0711010.
19. BAUER, P., THORPE, A., and BRUNET, G. “The quiet revolution of numerical weather prediction”. In: *Nature* 525.7567 (2015), pp. 47–55. doi: 10.1038/nature14956.
20. BAYDIN, A. G., PEARLMUTTER, B. A., RADUL, A. A., and SISKIND, J. M. “Automatic Differentiation in Machine Learning: A Survey”. In: *J. Mach. Learn. Res.* 18.1 (Jan. 2017), pp. 5595–5637.
21. BENNETT, C. H. “Logical Reversibility of Computation”. In: *IBM Journal of Research and Development* 17.6 (1973), pp. 525–532. doi: 10.1147/rd.176.0525.
22. BETANCOURT, M. “A Conceptual Introduction to Hamiltonian Monte Carlo”. In: *arXiv* (2017). doi: 10.48550/arxiv.1701.02434.
23. BEZANSON, J., EDELMAN, A., KARPINSKI, S., and SHAH, V. B. “Julia: A Fresh Approach to Numerical Computing”. In: *SIAM Review* 59.1 (2017), pp. 65–98. doi: 10.1137/141000671.
24. BEZANSON, J., KARPINSKI, S., SHAH, V. B., and EDELMAN, A. “Julia: A Fast Dynamic Language for Technical Computing”. In: *arXiv* (2012). doi: 10.48550/arxiv.1209.5145.
25. BLESSING, S. et al. “Testing variational estimation of process parameters and initial conditions of an earth system model”. In: *Tellus A* 66.0 (2014), p. 22606. doi: 10.3402/tellusa.v66.22606.
26. BLONIGAN, P. J. “Adjoint sensitivity analysis of chaotic dynamical systems with non-intrusive least squares shadowing”. In: *Journal of Computational Physics* 348 (2017), pp. 803–826. doi: <https://doi.org/10.1016/j.jcp.2017.08.002>.
27. BLONIGAN, P. J. and WANG, Q. “Probability density adjoint for sensitivity analysis of the mean of chaos”. In: *Journal of Computational Physics* 270 (2014), pp. 660–686. doi: <https://doi.org/10.1016/j.jcp.2014.04.027>.
28. BLONIGAN, P. J. and WANG, Q. “Multiple shooting shadowing for sensitivity analysis of chaotic dynamical systems”. In: *Journal of Computational Physics* 354 (2018), pp. 447–475. doi: <https://doi.org/10.1016/j.jcp.2017.10.032>.
29. BOCQUET, M., BRAJARD, J., CARRASSI, A., and BERTINO, L. “Data Assimilation as a Learning Tool to Infer Ordinary Differential Equation Representations of Dynamical Models”. In: *Nonlinear Processes in Geophysics* 26.3 (July 10, 2019), pp. 143–162. doi: 10.5194/npg-26-143-2019.
30. BOLIBAR, J., SAPIENZA, F., MAUSSION, F., LGUENSAT, R., WOUTERS, B., and PÉREZ, F. “Universal differential equations for glacier ice flow modelling”. In: *Geoscientific Model Development* 16.22 (2023), pp. 6671–6687. doi: 10.5194/gmd-16-6671-2023.
31. BOUSSANGE, V., BECKER, S., JENTZEN, A., KUCKUCK, B., and PELLISSIER, L. “Deep Learning Approximations for Non-Local Nonlinear PDEs with Neumann Boundary Conditions”. In: *Partial Differential Equations and Applications* 4.6 (Dec. 1, 2023), p. 51. doi: 10.1007/s42985-023-00244-0.
32. BOUSSANGE, V. and PELLISSIER, L. “Eco-Evolutionary Model on Spatial Graphs Reveals How Habitat Structure Affects Phenotypic Differentiation”. In: *Communications Biology* 5.1 (Dec. 6, 2022), p. 668. doi: 10.1038/s42003-022-03595-3.
33. BOUSSANGE, V., VILMELIS, P., and PELLISSIER, L. “Mini-Batching Ecological Data to Improve Ecosystem Models with Machine Learning”. In: (2022).
34. BRADLEY, A. M. *PDE-constrained optimization and the adjoint method*. Tech. rep. Technical Report. Stanford University. https://cs.stanford.edu/textasciitilde_ambrad..., 2013.
35. BRAJARD, J., CARRASSI, A., BOCQUET, M., and BERTINO, L. “Combining Data Assimilation and Machine Learning to Infer Unresolved Scale Parameterization”. In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 379.2194 (2021). doi: 10.1098/rsta.2020.0086.

36. BRÉZIS, H. *Functional analysis, Sobolev spaces and partial differential equations*. Vol. 2. 3. Springer, 2011.
37. BRYSON, A., HO, Y.-C., and SIOURIS, G. "Applied Optimal Control: Optimization, Estimation, and Control". In: *Systems, Man and Cybernetics, IEEE Transactions on* 9 (July 1979), pp. 366–367. doi: 10.1109/TSMC.1979.4310229.
38. BUI-THANH, T., BURSTEDDE, C., GHATTAS, O., MARTIN, J., STADLER, G., and WILCOX, L. C. "Extreme-scale UQ for Bayesian inverse problems governed by PDEs". In: *IEEE Computer Society Press* (2012), p. 3.
39. BUIZZA, R. and PALMER, T. N. "The singular-vector structure of the atmospheric global circulation". In: *Journal of the Atmospheric Sciences* 52.9 (1995), pp. 1434–1456. doi: 10.1175/1520-0469(1995)052<1434:tsvsot>2.0.co;2.
40. BURNHAM, K. P. and ANDERSON, D. R. "Multimodel Inference". In: *Sociological Methods & Research* 33.2 (Nov. 30, 2004), pp. 261–304. doi: 10.1177/0049124104268644.
41. CAO, Y., LI, S., and PETZOLD, L. "Adjoint sensitivity analysis for differential-algebraic equations: algorithms and software". In: *Journal of Computational and Applied Mathematics* 149.1 (2002), pp. 171–191. doi: 10.1016/S0377-0427(02)00528-9.
42. CARRASSI, A., BOCQUET, M., BERTINO, L., and EVENSEN, G. "Data Assimilation in the Geosciences: An Overview of Methods, Issues, and Perspectives". In: *WIREs Climate Change* 9.5 (Sept. 9, 2018), pp. 1–50. doi: 10.1002/wcc.535.
43. CHALMANDRIER, L. et al. "Linking Functional Traits and Demography to Model Species-Rich Communities". In: *Nature Communications* 12.1 (Dec. 2021), p. 2724. doi: 10.1038/s41467-021-22630-1.
44. CHANDRAMOORTHY, N. and WANG, Q. "Efficient computation of linear response of chaotic attractors with one-dimensional unstable manifolds". In: *SIAM Journal on Applied Dynamical Systems* 21.2 (2022), pp. 735–781. doi: <https://doi.org/10.1137/21M1405599>.
45. CHEN, R. T., RUBANOVA, Y., BETTENCOURT, J., and DUVENAUD, D. K. "Neural ordinary differential equations". In: *Advances in neural information processing systems* 31 (2018).
46. CLIFFORD, "Preliminary sketch of biquaternions". In: *Proceedings of the London Mathematical Society* 1.1 (1871), pp. 381–395.
47. COURTIER, P. and TALAGRAND, O. "Variational Assimilation of Meteorological Observations With the Adjoint Vorticity Equation. II: Numerical Results". In: *Quarterly Journal of the Royal Meteorological Society* 113.478 (1987), pp. 1329–1347. doi: 10.1002/qj.49711347813.
48. COVENEY, P. V., DOUGHERTY, E. R., and HIGHFIELD, R. R. "Big data need big theory too". In: *Philosophical Transactions of the Royal Society of London. Series A: Physical and Engineering Sciences* 374.2080 (2016), pp. 20160153–11. doi: 10.1098/rsta.2016.0153.
49. COX, D. R. and EFRON, B. "Statistical thinking for 21st century scientists". In: *Science Advances* 3.6 (2017), e1700768. doi: 10.1126/sciadv.1700768.
50. CURTSDOTTER, A. et al. "Ecosystem Function in Predator–Prey Food Webs—Confronting Dynamic Models with Empirical Data". In: *Journal of Animal Ecology* 88.2 (Feb. 7, 2019). Ed. by D. STOUFFER, pp. 196–210. doi: 10.1111/1365-2656.12892.
51. DANDEKAR, R., RACKAUCKAS, C., and BARBASTATHIS, G. "A Machine Learning-Aided Global Diagnostic and Comparative Tool to Assess Effect of Quarantine Control in COVID-19 Spread". In: *Patterns* 1.9 (2020), p. 100145. doi: 10.1016/j.patter.2020.100145.
52. DIMET, F.-X. L., NAVON, I. M., and DAESCU, D. N. "Second-Order Information in Data Assimilation*". In: *Monthly Weather Review* 130.3 (2002), pp. 629–648. doi: 10.1175/1520-0493(2002)130<0629:soiida>2.0.co;2.
53. DORIGO, T. et al. "Toward the End-to-End Optimization of Particle Physics Instruments with Differentiable Programming: a White Paper". In: *arXiv* (2022). doi: 10.48550/arxiv.2203.13818.
54. DORMANN, C. F. "Promising the Future? Global Change Projections of Species Distributions". In: *Basic and Applied Ecology* 8.5 (Sept. 3, 2007), pp. 387–397. doi: 10.1016/j.baae.2006.11.001.
55. DÜRRBAUM, A., Klier, W., and HAHN, H. "Comparison of Automatic and Symbolic Differentiation in Mathematical Modeling and Computer Simulation of Rigid-Body Systems". In: *Multibody System Dynamics* 7.4 (2002), pp. 331–355. doi: 10.1023/a:1015523018029.
56. DUTKIEWICZ, S., FOLLOWS, M. J., HEIMBACH, P., and MARSHALL, J. "Controls on ocean productivity and air-sea carbon flux: An adjoint model sensitivity study". In: *Geophysical Research Letters* 33.2 (2006), pp. 159–4. doi: 10.1029/2005gl024987.
57. EBERHARD, P. and BISCHOF, C. "Automatic differentiation of numerical integration algorithms". In: *Mathematics of Computation* 68.226 (1996), pp. 717–731. doi: 10.1090/S0025-5718-99-01027-3.
58. ELLIOTT, C. "The simple essence of automatic differentiation". In: *Proceedings of the ACM on Programming Languages* 2.ICFP (2018), p. 70. doi: 10.1145/3236765.
59. ELLIOTT, J. and PERAIRE, J. "Aerodynamic design using unstructured meshes". In: *Fluid Dynamics Conference* (1996). This has an example of the hardcore adjoint method implemented for aerodynamics. It may help to read this to see how the adjoint equations is being solved and the size of the problem. doi: 10.2514/6.1996-1941.
60. ERRICO, R. M. "What Is an Adjoint Model?" In: *Bulletin of the American Meteorological Society* 78.11 (1997), pp. 2577–2591. doi: 10.1175/1520-0477(1997)078<2577:wiaam>2.0.co;2.
61. EYINK, G., HAINE, T., and LEA, D. "Ruelle's linear response formula, ensemble adjoint schemes and Lévy flights". In: *Nonlinearity* 17.5 (2004), p. 1867.
62. FARRELL, B. "Optimal Excitation of Neutral Rossby Waves". In: *Journal of the Atmospheric Sciences* 45.2 (1988), pp. 163–172. doi: 10.1175/1520-0469(1988)045<0163:oeonrw>2.0.co;2.
63. FARRELL, B. F. and IOANNOU, P. J. "Generalized Stability Theory. Part I: Autonomous Operators". In: *Journal of the Atmospheric Sciences* 53.14 (1996), pp. 2025–2040. doi: 10.1175/1520-0469(1996)053<2025:gstpia>2.0.co;2.
64. FARRELL, P. E., HAM, D. A., FUNKE, S. W., and ROGNES, M. E. "Automated Derivation of the Adjoint of High-Level Transient Finite Element Programs". In: *SIAM Journal on Scientific Computing* 35.4 (2013), pp. C369–C393. doi: 10.1137/120873558.
65. FERREIRA, D., MARSHALL, J., and HEIMBACH, P. "Estimating Eddy Stresses by Fitting Dynamics to Observations Using a Residual-Mean Ocean Circulation Model and Its Adjoint". In: *Journal of Physical Oceanography* 35.10 (2005), pp. 1891–1910. doi: 10.1175/jpo2785.1.
66. FIKE, J. A. "Multi-objective optimization using hyper-dual numbers". PhD thesis. 2013.
67. FORGET, G., CAMPIN, J.-M., HEIMBACH, P., HILL, C. N., PONTE, R. M., and WUNSCH, C. "ECCO version 4: an integrated framework for non-linear inverse modeling and global ocean state estimation". In: *Geoscientific Model Development* 8.10 (2015), pp. 3071–3104. doi: 10.5194/gmd-8-3071-2015.
68. FRANK, S. A. "Automatic Differentiation and the Optimization of Differential Equation Models in Biology". In: *Frontiers in Ecology and Evolution* 10 (2022).
69. FRANKLIN, O. et al. "Organizing Principles for Vegetation Dynamics". In: *Nature Plants* 6.5 (2020), pp. 444–453. doi: 10.1038/s41477-020-0655-x.
70. FROLOV, S. et al. "Road Map for the Next Decade of Earth System Reanalysis in the United States". In: *Bulletin of the American Meteorological Society* 104.3 (2023), E706–E714. doi: 10.1175/bams-d-23-0011.1.
71. GÁBOR, A. and BANGA, J. R. "Robust and Efficient Parameter Estimation in Dynamic Models of Biological Systems". In: *BMC Systems Biology* 9.1 (Dec. 29, 2015), p. 74. doi: 10.1186/s12918-015-0219-2.

72. GAIKWAD, S. S., HASCOET, L., NARAYANAN, S. H. K., CURRY-LOGAN, L., GREVE, R., and HEIMBACH, P. "SICOPOLIS-AD v2: tangent linear and adjoint modeling framework for ice sheet modeling enabled by automatic differentiation tool Tapenade". In: *Journal of Open Source Software* 8.83 (2023), p. 4679. doi: 10.21105/joss.04679.
73. GAIKWAD, S. S. et al. "MITgcm-AD v2: Open source tangent linear and adjoint modeling framework for the oceans and atmosphere enabled by the Automatic Differentiation tool Tapenade". In: *arXiv* (2024).
74. GBIF: THE GLOBAL BIODIVERSITY INFORMATION FACILITY. "What Is GBIF?" In: (2022).
75. GEARY, W. L. et al. "A Guide to Ecosystem Models and Their Environmental Applications". In: *Nature Ecology & Evolution* 4.11 (Nov. 14, 2020), pp. 1459–1471. doi: 10.1038/s41559-020-01298-8.
76. GEBREMEDHIN, A. H., MANNE, F., and POTHEN, A. "What color is your Jacobian? Graph coloring for computing derivatives". In: *SIAM review* 47.4 (2005), pp. 629–705.
77. GELBRECHT, M., WHITE, A., BATHIANY, S., and BOERS, N. "Differentiable programming for Earth system modeling". In: *Geoscientific Model Development* 16.11 (2023), pp. 3123–3135. doi: 10.5194/gmd-16-3123-2023.
78. GELMAN, A., CARLIN, J. B., STERN, H. S., DUNSON, D. B., VEHTARI, A., and RUBIN, D. B. *Bayesian data analysis*. CRC press, 2013.
79. GEORGIEVA, N. K., GLAVIC, S., BAKR, M. H., and BANDLER, J. W. "Feasible Adjoint Sensitivity Technique for Em Design Optimization". In: *IEEE Transactions on Microwave Theory and Techniques* 50.12 (2002), p. 2751. doi: 10.1109/tmmt.2002.805131.
80. GHARAMTI, M. et al. "Ensemble Data Assimilation for Ocean Biogeochemical State and Parameter Estimation at Different Sites". In: *Ocean Modelling* 112 (Apr. 2017), pp. 65–89. doi: 10.1016/j.ocemod.2017.02.006.
81. GHATTAS, O. and WILLCOX, K. "Learning physics-based models from data: perspectives from inverse problems and model reduction". In: *Acta Numerica* 30 (2021), pp. 445–554. doi: 10.1017/s0962492921000064.
82. GIERING, R. and KAMINSKI, T. "Recipes for adjoint code construction". In: *ACM Trans Math Softw* 24.4 (1998), pp. 437–474. doi: 10.1145/293686.293695.
83. GIERING, R. and KAMINSKI, T. "Recipes for adjoint code construction". In: *ACM Transactions on Mathematical Software (TOMS)* 24.4 (1998), pp. 437–474. doi: 10.1145/293686.293695.
84. GIERING, R., KAMINSKI, T., TODLING, R., ERRICO, R., GELARO, R., and WINSLOW, N. "Automatic Differentiation: Applications, Theory, and Implementations". In: *Lecture Notes in Computational Science and Engineering* (2006), pp. 275–284. doi: 10.1007/3-540-28438-9_24.
85. GILES, M. B. and PIERCE, N. A. "An introduction to the adjoint approach to design". In: *Flow, Turbulence and Combustion* 65.3 (2000), pp. 393–415.
86. GILES, M. B. and PIERCE, N. A. "An Introduction to the Adjoint Approach to Design". In: *Flow, Turbulence and Combustion* 65.3–4 (2000), pp. 393–415. doi: 10.1023/a:1011430410075.
87. GIVOLI, D. "A tutorial on the adjoint method for inverse problems". In: 380 (2021), p. 113810. doi: 10.1016/j.cma.2021.113810.
88. GODWIN, C. M., CHANG, F.-H., and CARDINALE, B. J. "An Empiricist's Guide to Modern Coexistence Theory for Competitive Communities". In: *Oikos* 129.8 (2020), pp. 1109–1127. doi: 10.1111/oik.06957.
89. GOERZ, M. H., CARRASCO, S. C., and MALINOVSKY, V. S. "Quantum optimal control via semi-automatic differentiation". In: *Quantum* 6 (2022), p. 871. doi: 10.22331/q-2022-12-07-871.
90. GOLDBERG, D. "What every computer scientist should know about floating-point arithmetic". In: *ACM Computing Surveys (CSUR)* 23.1 (1991), pp. 5–48. doi: 10.1145/103162.103163.
91. GOLDBERG, D. and HEIMBACH, P. "Parameter and state estimation with a time-dependent adjoint marine ice sheet model". In: *The Cryosphere* 7.6 (2013), pp. 1659–1678.
92. GORBAN, A. and WUNSCH, D. "The general approximation theorem". In: *1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No.98CH36227)*. Vol. 2. 1998, 1271–1274 vol.2. doi: 10.1109/IJCNN.1998.685957.
93. GRIEWANK, A. "On Automatic Differentiation". In: *Mathematical Programming: Recent Developments and Applications* (1989).
94. GRIEWANK, A. "Who invented the reverse mode of differentiation". In: *Documenta Mathematica, Extra Volume ISMP 389400* (2012).
95. GRIEWANK, A. and WALTHER, A. *Evaluating Derivatives*. 2008. doi: 10.1137/1.9780898717761.
96. HAGER, W. W. "Runge-Kutta methods in optimal control and the transformed adjoint system". In: *Numerische Mathematik* 87.2 (2000), pp. 247–282. doi: 10.1007/s002110000178.
97. HAIRER, E., WANNER, G., and NØRSETT, S. *Solving ordinary differential equations I*. Springer Berlin Heidelberg New York.
98. HARTIG, F. et al. "Connecting Dynamic Vegetation Models to Data - an Inverse Perspective: Dynamic Vegetation Models - an Inverse Perspective". In: *Journal of Biogeography* 39.12 (Dec. 2012), pp. 2240–2252. doi: 10.1111/j.1365-2699.2012.02745.x.
99. HASCOET, L. and PASQUAL, V. "The Tapenade Automatic Differentiation Tool: Principles, Model, and Specification". In: *ACM Transactions on Mathematical Software* 39.3 (2013), pp. 1–43. doi: 10.1145/2450153.2450158.
100. HASCOET, L. and MORLIGHEM, M. "Source-to-source adjoint Algorithmic Differentiation of an ice sheet model written in C". In: *Optimization Methods and Software* 33.4–6 (2018), pp. 829–843.
101. HASTIE, T., TIBSHIRANI, R., FRIEDMAN, J. H., and FRIEDMAN, J. H. *The elements of statistical learning: data mining, inference, and prediction*. Vol. 2. Springer, 2009.
102. HEIMBACH, P., MENEMENLIS, D., LOSCH, M., CAMPIN, J.-M., and HILL, C. "On the formulation of sea-ice models. Part 2: Lessons from multi-year adjoint sea-ice export sensitivities through the Canadian Arctic Archipelago". In: *Ocean Modelling* 33.1–2 (2010), pp. 145–158. doi: 10.1016/j.ocemod.2010.02.002.
103. HEIMBACH, P. and BUGNION, V. "Greenland ice-sheet volume sensitivity to basal, surface and initial conditions derived from an adjoint model". In: *Annals of Glaciology* 50.52 (2009), pp. 67–80.
104. HEIMBACH, P., HILL, C., and GIERING, R. "An efficient exact adjoint of the parallel MIT General Circulation Model, generated via automatic differentiation". In: *Future Generation Computer Systems* 21.8 (2005), pp. 1356–1371. doi: 10.1016/j.future.2004.11.010.
105. HEIMBACH, P. and LOSCH, M. "Adjoint sensitivities of sub-ice-shelf melt rates to ocean circulation under the Pine Island Ice Shelf, West Antarctica". In: *Annals of Glaciology* 53.60 (2012), pp. 59–69. doi: 10.3189/2012/aog60a025.
106. HIGGINS, S. I., SCHEITER, S., and SANKARAN, M. "The Stability of African Savannas: Insights from the Indirect Estimation of the Parameters of a Dynamic Model". In: *Ecology* 91.6 (June 2010), pp. 1682–1692. doi: 10.1890/08-1368.1.
107. HINDMARSH, A. C. et al. "SUNDIALS: Suite of nonlinear and differential/algebraic equation solvers". In: *ACM Transactions on Mathematical Software (TOMS)* 31.3 (2005), pp. 363–396.
108. HU, R. "Supersonic biplane design via adjoint method". PhD thesis. 2010.
109. HUANG, J., SMITH, T. M., HENRY, G. M., and GEJUN, R. A. V. D. "Strassen's Algorithm Reloaded". In: *SC16: International Conference for High Performance Computing, Networking, Storage and Analysis* (2016), pp. 690–701. doi: 10.1109/sc.2016.58.

110. INNES, M. “Don’t Unroll Adjoint: Differentiating SSA-Form Programs”. In: *arXiv* (2018).
111. INNES, M. et al. “A Differentiable Programming System to Bridge Machine Learning and Scientific Computing”. In: *arXiv* (2019). DOI: 10.48550/arxiv.1907.07587.
112. IPSEN, I. C. F. and MEYER, C. D. “The Idea Behind Krylov Methods”. In: *The American Mathematical Monthly* 105.10 (1998), pp. 889–899. DOI: 10.1080/00029890.1998.12004985.
113. JAMESON, A. “Aerodynamic design via control theory”. In: *Journal of Scientific Computing* 3.3 (1988), pp. 233–260. DOI: 10.1007/bf01061285.
114. JENSEN, J. S., NAKSHATRALA, P. B., and TORTORELLI, D. A. “On the consistency of adjoint sensitivity analysis for structural optimization of linear dynamic problems”. In: *Structural and Multidisciplinary Optimization* 49.5 (2014), pp. 831–837. DOI: 10.1007/s00158-013-1024-4.
115. JETZ, W. et al. “Essential Biodiversity Variables for Mapping and Monitoring Species Populations”. In: *Nature Ecology and Evolution* 3.4 (2019), pp. 539–551. DOI: 10.1038/s41559-019-0826-1.
116. JIRARI, H. “Optimal control approach to dynamical suppression of decoherence of a qubit”. In: *Europhysics Letters* 87.4 (2009), p. 40003. DOI: 10.1209/0295-5075/87/40003.
117. JIRARI, H. “From quantum optimal control theory to coherent destruction of tunneling”. In: *The European Physical Journal B* 92 (2019), pp. 1–8. DOI: 10.1140/epjb/e2018-90231-5.
118. JOHNSON, J. B. and OMLAND, K. S. “Model Selection in Ecology and Evolution”. In: *Trends in Ecology & Evolution* 19.2 (Feb. 2004), pp. 101–108. DOI: 10.1016/j.tree.2003.10.013.
119. JOHNSON, S. G. “Notes on Adjoint Methods for 18.335”. In: 2012.
120. JOUVET, G. “Inversion of a Stokes glacier flow model emulated by deep learning”. In: *Journal of Glaciology* 69.273 (2023), pp. 13–26.
121. JOUVET, G., CORDONNIER, G., KIM, B., LÜTHI, M., VIELI, A., and ASCHWANDEN, A. “Deep learning speeds up ice flow modelling by several orders of magnitude”. In: *Journal of Glaciology* (2021), pp. 1–14. DOI: 10.1017/jog.2021.120.
122. JUEDES, D. W. *A taxonomy of automatic differentiation tools*. Tech. rep. Argonne National Lab., IL (United States), 1991.
123. KANTOROVICH, L. V. “On a mathematical symbolism convenient for performing machine calculations”. In: *Dokl. Akad. Nauk SSSR*. Vol. 113. 4. 1957, pp. 738–741.
124. KARCZMARZUK, J. “Functional Differentiation of Computer Programs”. In: *Proceedings of the Third ACM SIGPLAN International Conference on Functional Programming*. ICFP ’98. Baltimore, Maryland, USA: Association for Computing Machinery, 1998, pp. 195–203. DOI: 10.1145/289423.289442.
125. KIDGER, P. “On Neural Differential Equations”. PhD thesis. University of Oxford, 2021.
126. KIDGER, P., CHEN, R. T. Q., and LYONS, T. J. “Hey, that’s not an ODE”: Faster ODE Adjoints via Seminorms”. In: *Proceedings of the 38th International Conference on Machine Learning*. Ed. by M. MEILA and T. ZHANG. Vol. 139. Proceedings of Machine Learning Research. PMLR, 18–24 Jul 2021, pp. 5443–5452.
127. KIM, S., JI, W., DENG, S., MA, Y., and RACKAUCKAS, C. “Stiff neural ordinary differential equations”. en. In: *Chaos: An Interdisciplinary Journal of Nonlinear Science* 31.9 (Sept. 2021), p. 093122. DOI: 10.1063/5.0060697.
128. KINGMA, D. P. and BA, J. “Adam: A Method for Stochastic Optimization”. Dec. 22, 2014.
129. KOCHKOV, D. et al. “Neural General Circulation Models”. In: *arXiv* (2023). DOI: 10.48550/arxiv.2311.07222.
130. LANGLAND, R. H. and BAKER, N. L. “Estimation of observation impact using the NRL atmospheric variational data assimilation adjoint system”. In: *Tellus A: Dynamic Meteorology and Oceanography* 56.3 (2004), pp. 189–201. DOI: 10.3402/tellusa.v56i3.14413.
131. LANTOINE, G., RUSSELL, R. P., and DARGENT, T. “Using Multicomplex Variables for Automatic Computation of High-Order Derivatives”. In: *ACM Transactions on Mathematical Software (TOMS)* 38.3 (2012), p. 16. DOI: 10.1145/2168773.2168774.
132. LAUE, S. *On the Equivalence of Forward Mode Automatic Differentiation and Symbolic Differentiation*. 2019. DOI: 10.48550/ARXIV.1904.02990.
133. LEA, D. J., ALLEN, M. R., and HAINE, T. W. N. “Sensitivity analysis of the climate of a chaotic system”. In: *Tellus A: Dynamic Meteorology and Oceanography* 52.5 (2000), pp. 523–532. DOI: 10.3402/tellusa.v52i5.12283.
134. LEA, D. J., ALLEN, M. R., and HAINE, T. W. “Sensitivity analysis of the climate of a chaotic system”. In: *Tellus A: Dynamic Meteorology and Oceanography* 52.5 (2000), pp. 523–532. DOI: https://doi.org/10.3402/tellusa.v52i5.12283.
135. LEA, D. J., HAINE, T. W. N., ALLEN, M. R., and HANSEN, J. A. “Sensitivity analysis of the climate of a chaotic ocean circulation model”. In: *Quarterly Journal of the Royal Meteorological Society* 128.586 (2002), pp. 2587–2605. DOI: 10.1256/qj.01.180.
136. LECUN, Y., BENGIO, Y., and HINTON, G. “Deep Learning”. In: *Nature* 521.7553 (May 27, 2015), pp. 436–444. DOI: 10.1038/nature14539.
137. LEIS, J. R. and KRAMER, M. A. “Algorithm 658: ODESSA—an Ordinary Differential Equation Solver with Explicit Simultaneous Sensitivity Analysis”. In: *ACM Trans. Math. Softw.* 14.1 (Mar. 1988), pp. 61–67. DOI: 10.1145/42288.214371.
138. LEITH, C. E. “Climate response and fluctuation dissipation”. In: *Journal of Atmospheric Sciences* 32.10 (1975), pp. 2022–2026. DOI: https://doi.org/10.1175/1520-0469(1975)032<2022:CRAFD>2.0.CO;2.
139. LEUNG, N., ABDELHAFEZ, M., KOCH, J., and SCHUSTER, D. “Speedup for quantum optimal control from automatic differentiation based on graphics processing units”. In: *Physical Review A* 95.4 (2017), p. 042318. DOI: 10.1103/PhysRevA.95.042318.
140. LI, D., XU, K., HARRIS, J. M., and DARVE, E. “Coupled time-lapse full-waveform inversion for subsurface flow problems using intrusive automatic differentiation”. In: *Water Resources Research* 56.8 (2020), e2019WR027032.
141. LI, X., WONG, T.-K. L., CHEN, R. T., and DUVERNAUD, D. “Scalable gradients for stochastic differential equations”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2020, pp. 3870–3882.
142. LION, S. “Theoretical Approaches in Evolutionary Ecology: Environmental Feedback as a Unifying Perspective”. In: *American Naturalist* 191.1 (2018), pp. 21–44. DOI: 10.1086/694865.
143. LIONS, J. L. *Optimal control of systems governed by partial differential equations*. Vol. 170. Springer, 1971.
144. LIU, C., KÖHL, A., and STAMMER, D. “Adjoint-Based Estimation of Eddy-Induced Tracer Mixing Parameters in the Global Ocean”. In: *Journal of Physical Oceanography* 42.7 (2012), pp. 1186–1206. DOI: 10.1175/jpo-d-11-0162.1.
145. LOGAN, L. C., NARAYANAN, S. H. K., GREVE, R., and HEIMBACH, P. “SICOPOLIS-AD v1: an open-source adjoint modeling framework for ice sheet simulation enabled by the algorithmic differentiation tool OpenAD”. In: *Geoscientific Model Development* 13.4 (2020), pp. 1845–1864.
146. LYNESS, J. N. “Numerical algorithms based on the theory of complex variable”. In: *Proceedings of the 1967 22nd national conference on -* (1967), pp. 125–133. DOI: 10.1145/800196.805983.
147. LYNESS, J. N. and MOLER, C. B. “Numerical Differentiation of Analytic Functions”. In: *SIAM Journal on Numerical Analysis* 4.2 (1967), pp. 202–210. DOI: 10.1137/0704019.
148. LYU, G., KÖHL, A., MATEI, I., and STAMMER, D. “Adjoint-Based Climate Model Tuning: Application to the Planet Simulator”. In: *Journal of Advances in Modeling Earth Systems* 10.1 (2018), pp. 207–222. DOI: 10.1002/2017ms001194.

149. MA, Y., DIXIT, V., INNES, M. J., GUO, X., and RACKAUCKAS, C. “A comparison of automatic differentiation and continuous sensitivity analysis for derivatives of differential equation solutions”. In: *2021 IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE. 2021, pp. 1–9.
150. MACAYEAL, D. R. “The basal stress distribution of Ice Stream E, Antarctica, inferred by control methods”. In: *Journal of Geophysical Research: Solid Earth* 97.B1 (1992), pp. 595–603.
151. MANZYUK, O., PEARLMUTTER, B. A., RADUL, A. A., RUSH, D. R., and SISKIND, J. M. “Perturbation confusion in forward automatic differentiation of higher-order functions”. In: *Journal of Functional Programming* 29 (2019), e12.
152. MARGOSSIAN, C. C. “A Review of automatic differentiation and its efficient implementation”. In: *arXiv* (2018). DOI: 10.48550/arxiv.1811.05031.
153. MAROTZKE, J., GIERING, R., ZHANG, K. Q., STAMMER, D., HILL, C., and LEE, T. “Construction of the adjoint MIT ocean general circulation model and application to Atlantic heat transport sensitivity”. In: *Journal of Geophysical Research* 104.29 (1999), pp. 529–548. DOI: 10.1029/1999jc900236.
154. MARTINS, J. R. A., STURDZA, P., and ALONSO, J. J. “The complex-step derivative approximation”. In: *ACM Transactions on Mathematical Software (TOMS)* 29 (2003), pp. 245–262. DOI: 10.1145/838250.838251.
155. MARTINS, J., STURDZA, P., and ALONSO, J. “The connection between the complex-step derivative approximation and algorithmic differentiation”. In: *39th Aerospace Sciences Meeting and Exhibit*. 2001, p. 921.
156. MATHUR, R. “An analytical approach to computing step sizes for finite-difference derivatives”. PhD thesis. 2012.
157. MCGREIVY, N., HUDSON, S., and ZHU, C. “Optimized finite-build stellarator coils using automatic differentiation”. In: *Nuclear Fusion* 61 (2021), p. 026020. DOI: 10.1088/1741-4326/abcd76.
158. MEEHL, G. A. et al. “Initialized Earth System prediction from subseasonal to decadal timescales”. In: *Nature Reviews Earth & Environment* (2021), pp. 1–18. DOI: 10.1038/s43017-021-00155-x.
159. MITUSCH, S. K., FUNKE, S. W., and DOKKEN, J. S. “dolfin-adjoint 2018.1: automated adjoints for FEniCS and Firedrake”. In: *Journal of Open Source Software* 4.38 (2019), p. 1292. DOI: 10.21105/joss.01292.
160. MOHAMMADI, B. and PIRONNEAU, O. “Shape optimization in fluid mechanics”. In: *Annual Review of Fluid Mechanics* 36.1 (2004), pp. 255–279. DOI: 10.1146/annurev.fluid.36.050802.121926.
161. MOHAMMADI, B. and PIRONNEAU, O. *Applied shape optimization for fluids*. OUP Oxford, 2009.
162. MOLESKY, S., LIN, Z., PIGGOTT, A. Y., JIN, W., VUCKOVIĆ, J., and RODRIGUEZ, A. W. “Inverse design in nanophotonics”. In: *Nature Photonics* 12.11 (2018), pp. 659–670. DOI: 10.1038/s41566-018-0246-9.
163. MOLNAR, C., CASALICCHIO, G., and BISCHL, B. “Interpretable Machine Learning – A Brief History, State-of-the-Art and Challenges”. In: *arXiv* (2020). DOI: 10.48550/arxiv.2010.09337.
164. MOORE, A. M. and KLEEMAN, R. “The singular vectors of a coupled ocean-atmosphere model of Enso. I: Thermodynamics, energetics and error growth”. In: *Quarterly Journal of the Royal Meteorological Society* 123.540 (1997), pp. 953–981. DOI: 10.1002/qj.49712354009.
165. MOORE, A. M. and KLEEMAN, R. “The singular vectors of a coupled ocean-atmosphere model of Enso. II: Sensitivity studies and dynamical interpretation”. In: *Quarterly Journal of the Royal Meteorological Society* 123.540 (1997), pp. 983–1006. DOI: 10.1002/qj.49712354010.
166. MORLIGHEM, M., SEROUSSI, H., LAROUR, E., and RIGNOT, E. “Inversion of basal friction in Antarctica using exact and incomplete adjoints of a higher-order model”. In: *Journal of Geophysical Research: Earth Surface* 118.3 (2013), pp. 1746–1753.
167. MOSBEUX, C., GILLET-CHAULET, F., and GAGLIARDINI, O. “Comparison of adjoint and nudging methods to initialise ice sheet model basal conditions”. In: *Geoscientific Model Development* 9.7 (2016), pp. 2549–2562.
168. MOSES, W. and CHURAVY, V. “Instead of Rewriting Foreign Code for Machine Learning, Automatically Synthesize Fast Gradients”. In: *Advances in Neural Information Processing Systems*. Ed. by H. LAROCHELLE, M. RANZATO, R. HADSELL, M. F. BALCAN, and H. LIN. Vol. 33. Curran Associates, Inc., 2020, pp. 12472–12485.
169. NAUMANN, U. *The Art of Differentiating Computer Programs*. SIAM, 2011. DOI: 10.1137/1.9781611972078.
170. NEUENHOFEN, M. “Review of theory and implementation of hyper-dual numbers for first and second order automatic differentiation”. In: *arXiv* (2018). DOI: 10.48550/arxiv.1801.03614.
171. NI, A. “Fast linear response algorithm for differentiating chaos”. In: *arXiv preprint arXiv:2009.00595* (2020). DOI: <https://arxiv.org/abs/2009.00595v5>.
172. NI, A. and TALNIKAR, C. “Adjoint sensitivity analysis on chaotic dynamical systems by Non-Intrusive Least Squares Adjoint Shadowing (NILSAS)”. In: *Journal of Computational Physics* 395 (2019), pp. 690–709. DOI: <https://doi.org/10.1016/j.jcp.2019.06.035>.
173. NI, A. and WANG, Q. “Sensitivity analysis on chaotic dynamical systems by Non-Intrusive Least Squares Shadowing (NILSS)”. In: *Journal of Computational Physics* 347 (2017), pp. 56–77. DOI: <https://doi.org/10.1016/j.jcp.2017.06.033>.
174. NI, A., WANG, Q., FERNANDEZ, P., and TALNIKAR, C. “Sensitivity analysis on chaotic dynamical systems by finite difference non-intrusive least squares shadowing (FD-NILSS)”. In: *Journal of Computational Physics* 394 (2019), pp. 615–631. DOI: <https://doi.org/10.1016/j.jcp.2019.06.004>.
175. NORCLIFFE, A. and DEISENROTH, M. P. “Faster Training of Neural ODEs Using Gauß-Legendre Quadrature”. In: *arXiv* (2023). DOI: 10.48550/arxiv.2308.10644.
176. ODEN, J. T., MOSER, R., and GHATTAS, O. “Computer Predictions with Quantified Uncertainty, Part II”. In: *SIAM News* 43.10 (2010), pp. 1–4.
177. OKTAY, D., MCGREIVY, N., ADUOL, J., BEATSON, A., and ADAMS, R. P. “Randomized Automatic Differentiation”. In: *arXiv* (2020). DOI: 10.48550/arxiv.2007.10412.
178. ONKEN, D. and RUTHOTTO, L. “Discretize-Optimize vs. Optimize-Discretize for Time-Series Regression and Continuous Normalizing Flows”. In: *arXiv* (2020). DOI: 10.48550/arxiv.2005.13420.
179. PAGE, K. M. and NOWAK, M. A. “Unifying Evolutionary Dynamics”. In: *Journal of Theoretical Biology* 219.1 (Nov. 2002), pp. 93–98. DOI: 10.1006/jtbi.2002.3112.
180. PALMER, T. N., BUIZZA, R., MOLteni, F., CHEN, Y.-Q., and CORTI, S. “Singular vectors and the predictability of weather and climate”. In: *Philosophical Transactions of the Royal Society of London. Series A: Physical and Engineering Sciences* 348.1688 (1994), pp. 459–475. DOI: 10.1098/rsta.1994.0105.
181. PALMIERI, G., TIBONI, M., and LEGNANI, G. “Analysis of the Upper Limitation of the Most Convenient Cadence Range in Cycling Using an Equivalent Moment Based Cost Function”. In: *Mathematics* 8.11 (2020). DOI: 10.3390/math8111947.
182. PANTEL, J. H. and BECKS, L. “Statistical Methods to Identify Mechanisms in Studies of Eco-Evolutionary Dynamics”. In: *Trends in Ecology & Evolution* (June 2023), S0169534723000800. DOI: 10.1016/j.tree.2023.03.011.
183. PAREDES, J. A., HUPKENS, K., and STOCKER, B. D. *Rsofun: A Model-Data Integration Framework for Simulating Ecosystem Processes*. Nov. 24, 2023. DOI: 10.1101/2023.11.24.568574. URL: <https://www.biorxiv.org/content/10.1101/2023.11.24.568574v1> (visited on 11/28/2023). preprint.

184. PIRONNEAU, O. "Optimal shape design for elliptic systems". In: *System Modeling and Optimization: Proceedings of the 10th IFIP Conference New York City, USA, August 31–September 4, 1981*. Springer. 2005, pp. 42–66.
185. PONTARP, M., BRÄNNSTRÖM, Å., and PETCHEY, O. L. "Inferring Community Assembly Processes from Macroscopic Patterns Using Dynamic Eco-evolutionary Models and Approximate Bayesian Computation (ABC)". In: *Methods in Ecology and Evolution* 10.4 (Apr. 2019). Ed. by T. POISOT, pp. 450–460. doi: 10.1111/2041-210X.13129.
186. RABIER, F., JÄRVINEN, H., KLINKER, E., MAHFOUF, J. F., and SIMMONS, A. "The ECMWF operational implementation of four-dimensional variational assimilation. I: Experimental results with simplified physics". In: *Quarterly Journal of the Royal Meteorological Society* 126.564 (2000), pp. 1143–1170. doi: 10.1002/qj.49712656415.
187. RABIER, F. and COURTIER, P. "Four-Dimensional Assimilation In the Presence of Baroclinic Instability". In: *Quarterly Journal of the Royal Meteorological Society* 118.506 (1992), pp. 649–672. doi: 10.1002/qj.49711850604.
188. RACKAUCKAS, C. and NIE, Q. "DifferentialEquations.jl – A Performant and Feature-Rich Ecosystem for Solving Differential Equations in Julia". In: *Journal of Open Research Software* 5.1 (2016), p. 15. doi: 10.5334/jors.151.
189. RACKAUCKAS, C. et al. "Universal differential equations for scientific machine learning". In: *arXiv preprint arXiv:2001.04385* (2020).
190. RACKAUCKAS, C. et al. "Generalized physics-informed learning through language-wide differentiable programming". In: (2021).
191. RAISSI, M., PERDIKARIS, P., and KARNIADAKIS, G. "Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations". In: *Journal of Computational Physics* 378 (2019), pp. 686–707. doi: 10.1016/j.jcp.2018.10.045.
192. RAMSAY, J. and HOOKER, G. *Dynamic data analysis*. Springer, 2017.
193. RAMSUNDAR, B., KRISHNAMURTHY, D., and VISWANATHAN, V. "Differentiable Physics: A Position Piece". In: *arXiv* (2021). doi: 10.48550/arxiv.2109.07573.
194. RANOGA, H., DALCIN, L., PARSANI, M., and KETCHESON, D. I. "Optimized Runge-Kutta Methods with Automatic Step Size Control for Compressible Computational Fluid Dynamics". In: *Communications on Applied Mathematics and Computation* 4.4 (2022). Paper with the RDPK3Sp35 method, pp. 1191–1228. doi: 10.1007/s42967-021-00159-w.
195. RASP, S., PRITCHARD, M. S., and GENTINE, P. "Deep Learning to Represent Subgrid Processes in Climate Models". In: *Proceedings of the National Academy of Sciences* 115.39 (Sept. 25, 2018), pp. 9684–9689. doi: 10.1073/pnas.1810286115.
196. RAZAVI, S. et al. "The Future of Sensitivity Analysis: An essential discipline for systems modeling and policy support". In: *Environmental Modelling & Software* 137 (2021), p. 104954. doi: 10.1016/j.envsoft.2020.104954.
197. REVELS, J., LUBIN, M., and PAPAMARKOU, T. "Forward-Mode Automatic Differentiation in Julia". In: *arXiv:1607.07892 [cs.MS]* (2016).
198. ROSENBAUM, B., RAATZ, M., WEITHOFF, G., FUSSMANN, G. F., and GAEDKE, U. "Estimating Parameters From Multiple Time Series of Population Dynamics Using Bayesian Inference". In: *Frontiers in Ecology and Evolution* 6 (JAN Jan. 22, 2019). doi: 10.3389/fevo.2018.00234.
199. RÜDE, U., WILLCOX, K., MCINNES, L. C., and STERCK, H. D. "Research and Education in Computational Science and Engineering". In: *SIAM Review* 60.3 (2018), pp. 707–754. doi: 10.1137/16m1096840.
200. RUDER, S. "An overview of gradient descent optimization algorithms". In: *arXiv preprint arXiv:1609.04747* (2016).
201. RUDIN, C., CHEN, C., CHEN, Z., HUANG, H., SEMENOVA, L., and ZHONG, C. "Interpretable machine learning: Fundamental principles and 10 grand challenges". In: *Statistic Surveys* 16.none (2022), pp. 1–85. doi: 10.1214/21-ss133.
202. RUELLE, D. "Differentiation of SRB states". In: *Communications in Mathematical Physics* 187.1 (1997), pp. 227–241. doi: <https://doi.org/10.1007/s002200050134>.
203. RUELLE, D. "A review of linear response theory for general differentiable dynamical systems". In: *Nonlinearity* 22.4 (2009), p. 855. doi: <https://iopscience.iop.org/article/10.1088/0951-7715/22/4/009>.
204. RUPPERT, K. M., KLINE, R. J., and RAHMAN, M. S. "Past, Present, and Future Perspectives of Environmental DNA (eDNA) Metabarcoding: A Systematic Review in Methods, Monitoring, and Applications of Global eDNA". In: *Global Ecology and Conservation* 17 (Jan. 2019), e00547. doi: 10.1016/j.gecco.2019.e00547.
205. SANDU, A. "On the properties of Runge-Kutta discrete adjoints". In: *Computational Science–ICCS 2006: 6th International Conference, Reading, UK, May 28–31, 2006, Proceedings, Part IV* 6. Springer. 2006, pp. 550–557.
206. SANDU, A. "Solution of inverse problems using discrete ODE adjoints". In: *Large-Scale Inverse Problems and Quantification of Uncertainty* (2011), pp. 345–365.
207. SCHÄFER, F., KLOC, M., BRUDER, C., and LÖRCH, N. "A differentiable programming method for quantum control". In: *Machine Learning: Science and Technology* 1.3 (2020), p. 035009. doi: 10.1088/2632-2153/ab9802.
208. SCHÄFER, F., SEKATSKI, P., KOPPEHÖFER, M., BRUDER, C., and KLOC, M. "Control of stochastic quantum dynamics by differentiable programming". In: *Machine Learning: Science and Technology* 2.3 (2021), p. 035004. doi: 10.1088/2632-2153/abec22.
209. SCHÄFER, F., TAREK, M., WHITE, L., and RACKAUCKAS, C. "AbstractDifferentiation.jl: Backend-Agnostic Differentiable Programming in Julia". In: *arXiv* (2021). doi: 10.48550/arxiv.2109.12449.
210. SCHANEN, M., NARAYANAN, S. H. K., WILLIAMSON, S., CHURAVY, V., MOSES, W. S., and PAEHLER, L. "Transparent Checkpointing for Automatic Differentiation of Program Loops Through Expression Transformations". In: (2023). Ed. by J. MIKYŠKA, C. DE MULATIER, M. PASZYNSKI, V. V. KRZHIZHANOVSKAYA, J. J. DONGARRA, and P. M. SLOOT, pp. 483–497.
211. SCHARTAU, M. et al. "Reviews and Syntheses: Parameter Identification in Marine Planktonic Ecosystem Modelling". In: *Biogeosciences* 14.6 (Mar. 29, 2017), pp. 1647–1701. doi: 10.5194/bg-14-1647-2017.
212. SCHNEIDER, T., LAN, S., STUART, A., and TEIXEIRA, J. "Earth System Modeling 2.0: A Blueprint for Models That Learn From Observations and Targeted High-Resolution Simulations". In: *Geophysical Research Letters* 44.24 (Dec. 28, 2017), pp. 12, 396–417. doi: 10.1002/2017GL076101.
213. SERBAN, R. and HINDMARSH, A. C. "CVODES: the sensitivity-enabled ODE solver in SUNDIALS". In: *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. Vol. 47438. 2005, pp. 257–269.
214. SHEN, C. et al. "Differentiable modelling to unify machine learning and physical models for geosciences". In: *Nature Reviews Earth & Environment* (2023), pp. 1–16. doi: 10.1038/s43017-023-00450-9.
215. SILVA, H. D., GUSTAFSON, J. L., and WONG, W.-F. "Making Strassen Matrix Multiplication Safe". In: *2018 IEEE 25th International Conference on High Performance Computing (HiPC)* 00 (2018), pp. 173–182. doi: 10.1109/hipc.2018.00028.
216. SIRKES, Z. and TZIPERMAN, E. "Finite Difference of Adjoint or Adjoint of Finite Difference?" In: *Monthly Weather Review* 125.12 (1997), pp. 3373–3378. doi: 10.1175/1520-0493(1997)125<3373:fdoaoa>2.0.co;2.
217. SISKIND, J. M. and PEARLMUTTER, B. A. "Perturbation confusion and referential transparency: Correct functional implementation of forward-mode AD". In: (2005).
218. SKEELS, A., BACH, W., HAGEN, O., JETZ, W., and PELLISSIER, L. "Temperature-Dependent Evolutionary Speed Shapes the Evolution of Biodiversity Patterns Across Tetrapod Radiations". In: *Systematic Biology* 0.0 (July 9, 2022). Ed. by L. HARMON, pp. 1–16. doi: 10.1093/sysbio/syac048.

219. SQUIRE, W. and TRAPP, G. "Using Complex Variables to Estimate Derivatives of Real Functions". In: 40 (1998), pp. 110–112. doi: 10.1137/s003614459631241x.
220. STAMMER, D., KÖHL, A., VLASENKO, A., MATEI, I., LUNKEIT, F., and SCHUBERT, S. "A Pilot Climate Sensitivity Study Using the CEN Coupled Adjoint Model (CESAM)". In: *Journal of Climate* 31.5 (2018), pp. 2031–2056. doi: 10.1175/jcli-d-17-0183.1.
221. STAMMER, D. et al. "Global ocean circulation during 1992–1997, estimated from ocean observations and a general circulation model". In: *Journal of Geophysical Research: Oceans* 107.C9 (2002), pp. 1–1-27. doi: 10.1029/2001jc000888.
222. STAMMER, D. "Adjusting Internal Model Errors through Ocean State Estimation". In: *Journal of Physical Oceanography* 35.6 (2005), pp. 1143–1153. doi: 10.1175/jpo2733.1.
223. STEIN, E. M. and SHAKARCHI, R. *Complex analysis*. Vol. 2. Princeton University Press, 2010.
224. STOER, J. and BULIRSCH, R. *Introduction to numerical analysis*. Springer, 2002.
225. STROUWEN, A., NICOLAÏ, B. M., and GOOS, P. "Robust Dynamic Experiments for the Precise Estimation of Respiration and Fermentation Parameters of Fruit and Vegetables". In: *PLOS Computational Biology* 18.1 (Jan. 12, 2022). Ed. by P. MENDES, e1009610. doi: 10.1371/journal.pcbi.1009610.
226. TALAGRAND, O. and COURTIER, P. "Variational Assimilation of Meteorological Observations With the Adjoint Vorticity Equation. I: Theory". In: *Quarterly Journal of the Royal Meteorological Society* 113.478 (1987), pp. 1311–1328. doi: 10.1002/qj.49711347812.
227. THACKER, W. C. "Fitting models to inadequate data by enforcing spatial and temporal smoothness". In: *Journal of Geophysical Research: Oceans* (1978–2012) 93.C9 (1988), pp. 10655–10665. doi: 10.1029/jc093ic09p10655.
228. THACKER, W. C. "The role of the Hessian matrix in fitting models to measurements". In: *Journal of Geophysical Research: Oceans* (1978–2012) 94.C5 (1989), pp. 6177–6196. doi: 10.1029/jc094ic05p6177.
229. THACKER, W. C. and LONG, R. B. "Fitting dynamics to data". In: *Journal of Geophysical Research: Oceans* (1978–2012) 93.C2 (1988), pp. 1227–1240. doi: 10.1029/jc093ic02p01227.
230. THUBURN, J. "Climate sensitivities via a Fokker–Planck adjoint approach". In: *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography* 131.605 (2005), pp. 73–92. doi: https://doi.org/10.1256/qj.04.46.
231. TOMS, B. A., BARNES, E. A., and EBERT-UPHOFF, I. "Physically Interpretable Neural Networks for the Geosciences: Applications to Earth System Variability". In: *Journal of Advances in Modeling Earth Systems* 12.9 (2020), pp. 1–20. doi: 10.1029/2019MS002002.
232. TSITOURAS, C. "Runge–Kutta pairs of order 5(4) satisfying only the first column simplifying assumption". In: *Computers & Mathematics with Applications* 62.2 (2011), pp. 770–775. doi: 10.1016/j.camwa.2011.06.002.
233. TZIPERMAN, E., THACKER, W. C., and LONG, R. B. "Oceanic data analysis using a general circulation model. Part I: Simulations". In: *Journal of Physical Oceanography* 22.12 (1992), pp. 1434–1457. doi: 10.1175/1520-0485(1992)022<1434:odauag>2.0.co;2.
234. TZIPERMAN, E., THACKER, W. C., and LONG, R. B. "Oceanic data analysis using a general circulation model. Part II: A North Atlantic model". In: *Journal of Physical Oceanography* 22.12 (1992), pp. 1458–1485. doi: 10.1175/1520-0485(1992)022<1458:odauag>2.0.co;2.
235. TZIPERMAN, E. and THACKER, W. C. "An Optimal-Control/Adjoint-Equations Approach to Studying the Oceanic General Circulation". In: *Journal of Physical Oceanography* 19.10 (1989), pp. 1471–1485. doi: 10.1175/1520-0485(1989)019<1471:aoceat>2.0.co;2.
236. UTKE, J. et al. "OpenAD/F: A Modular Open-Source Tool for Automatic Differentiation of Fortran Codes". In: *ACM Transactions on Mathematical Software (TOMS)* 34.4 (2008), p. 18. doi: 10.1145/1377596.1377598.
237. VALLIS, G. K. "Geophysical fluid dynamics: whence, whither and why?" In: *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* 472.2192 (2016), pp. 20160140–23. doi: 10.1098/rspa.2016.0140.
238. VAN DEN BERG, N. I. et al. "Ecological Modelling Approaches for Predicting Emergent Properties in Microbial Communities". In: *Nature Ecology & Evolution* (May 16, 2022). doi: 10.1038/s41559-022-01746-7.
239. VASWANI, A. et al. "Attention Is All You Need". In: *Advances in Neural Information Processing Systems*. Ed. by I. GUYON et al. Vol. 30. Curran Associates, Inc., 2017.
240. VILLA, C., CHAPLAIN, M. A. J., and LORENZI, T. "Evolutionary Dynamics in Vascularised Tumours under Chemotherapy: Mathematical Modelling, Asymptotic Analysis and Numerical Simulations". In: *Vietnam Journal of Mathematics* 49.1 (Mar. 6, 2021), pp. 143–167. doi: 10.1007/s10013-020-00445-9.
241. WALTHER, A. "Automatic differentiation of explicit Runge–Kutta methods for optimal control". In: *Computational Optimization and Applications* 36.1 (2007), pp. 83–108. doi: 10.1007/s10589-006-0397-3.
242. WANG, F., ZHENG, D., DECKER, J., WU, X., ESSERTEL, G. M., and ROMPF, T. "Backpropagation with Continuation Callbacks: Foundations for Efficient and Expressive Differentiable Programming". In: *Proceedings of the ACM on Programming Languages* 3.ICFP (2019), p. 96. doi: 10.1145/3341700.
243. WANG, Q. "Forward and adjoint sensitivity computation of chaotic dynamical systems". In: *Journal of Computational Physics* 235 (2013), pp. 1–13. doi: https://doi.org/10.1016/j.jcp.2012.09.007.
244. WANG, Q. "Convergence of the least squares shadowing method for computing derivative of ergodic averages". In: *SIAM Journal on Numerical Analysis* 52.1 (2014), pp. 156–170. doi: https://doi.org/10.1137/130917065.
245. WANG, Q., HU, R., and BLONIGAN, P. "Least Squares Shadowing sensitivity analysis of chaotic limit cycle oscillations". In: *Journal of Computational Physics* 267 (2014), pp. 210–224. doi: https://doi.org/10.1016/j.jcp.2014.03.002.
246. WANG, Q., HU, R., and BLONIGAN, P. "Least squares shadowing sensitivity analysis of chaotic limit cycle oscillations". In: *Journal of Computational Physics* 267 (2014), pp. 210–224. doi: https://doi.org/10.1016/j.jcp.2014.03.002.
247. WANG, Y., LAI, C.-Y., and COWEN-BREEN, C. "Discovering the rheology of Antarctic Ice Shelves via physics-informed deep learning". In: (2022).
248. WANNER, G. and HAIRER, E. *Solving ordinary differential equations II*. Springer Berlin Heidelberg New York.
249. WATTS, M. C. "Modelling and the Monitoring of Mesocosm Experiments: Two Case Studies". In: *Journal of Plankton Research* 23.10 (Oct. 1, 2001), pp. 1081–1093. doi: 10.1093/plankt/23.10.1081.
250. WENG, E. S. et al. "Scaling from Individual Trees to Forests in an Earth System Modeling Framework Using a Mathematically Tractable Model of Height-Structured Competition". In: *Biogeosciences* 12.9 (May 7, 2015), pp. 2655–2694. doi: 10.5194/bg-12-2655-2015.
251. WENGERT, R. E. "A simple automatic derivative evaluation program". In: *Communications of the ACM* 7.8 (1964), pp. 463–464. doi: 10.1145/355586.364791.
252. WIGNER, E. P. "The unreasonable effectiveness of mathematics in the natural sciences". In: *Communications on Pure and Applied Mathematics* 13 (1960), pp. 1–14. doi: 10.1002/cpa.3160130102.
253. WOLFE, P. "Checking the Calculation of Gradients". In: *ACM Transactions on Mathematical Software (TOMS)* 8.4 (1982), pp. 337–343. doi: 10.1145/356012.356013.

254. YAZDANI, A., LU, L., RAISSI, M., and KARNIADAKIS, G. E. “Systems Biology Informed Deep Learning for Inferring Parameters and Hidden Dynamics”. In: *PLOS Computational Biology* 16.11 (Nov. 18, 2020). Ed. by V. HATZIMANIKATIS, e1007575. DOI: 10.1371/journal.pcbi.1007575.
255. ZANNA, L., HEIMBACH, P., MOORE, A. M., and TZIPERMAN, E. “Upper-ocean singular vectors of the North Atlantic climate with implications for linear predictability and variability”. In: *Quarterly Journal of the Royal Meteorological Society* 138.663 (2012), pp. 500–513. DOI: 10.1002/qj.937.
256. ZANNA, L., HEIMBACH, P., MOORE, A. M., and TZIPERMAN, E. “Optimal Excitation of Interannual Atlantic Meridional Overturning Circulation Variability”. In: *Journal of Climate* 24.2 (2011), pp. 413–427. DOI: 10.1175/2010jcli3610.1.
257. ZANNA, L., HEIMBACH, P., MOORE, A. M., and TZIPERMAN, E. “The Role of Ocean Dynamics in the Optimal Growth of Tropical SST Anomalies”. In: *Journal of Physical Oceanography* 40.5 (2010), pp. 983–1003. DOI: 10.1175/2009jpo4196.1.
258. ZDEBOROVÁ, L. “Understanding deep learning is also a job for physicists”. en. In: *Nature Physics* (May 2020). DOI: 10.1038/s41567-020-0929-2.
259. ZHANG, H. and SANDU, A. “FATODE: A library for forward, adjoint, and tangent linear integration of ODEs”. In: *SIAM Journal on Scientific Computing* 36.5 (2014), pp. C504–C523.
260. ZHU, W., XU, K., DARVE, E., and BEROZA, G. C. “A general approach to seismic inversion with automatic differentiation”. In: *Computers & Geosciences* 151 (2021), p. 104751. DOI: 10.1016/j.cageo.2021.104751.
261. ZHUANG, J., DVORNEK, N., LI, X., TATIKONDA, S., PAPADEMETRIS, X., and DUNCAN, J. “Adaptive Checkpoint Adjoint Method for Gradient Estimation in Neural ODE.” In: *Proceedings of machine learning research* 119 (2020), pp. 11639–11649.
262. ZIMMERMANN, N. E., EDWARDS JR, T. C., GRAHAM, C. H., PEARMAN, P. B., and SVENNING, J.-C. “New Trends in Species Distribution Modelling”. In: *Ecography* 33.6 (2010), pp. 985–989. DOI: 10.1111/j.1600-0587.2010.06953.x.