

Objective discovery of dominant dynamical processes with machine learning

Bryan Kaiser (✉ bkaiser@lanl.gov)

Los Alamos National Laboratory <https://orcid.org/0000-0002-4652-6935>

Juan Saenz

Los Alamos National Laboratory

Maike Sonnewald

Princeton University

Daniel Livescu

Los Alamos National Laboratory

Physical Sciences - Article

Keywords: Dynamical Regime Identification, Verification Criterion, Unsupervised Learning Framework, ad hoc Conventional Analyses

Posted Date: July 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-745356/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Objective discovery of dominant dynamical processes
2 with machine learning

3 LA-UR-21-27162

4 Bryan E. Kaiser^{1*}, Juan A. Saenz¹, Maike Sonnewald^{2,3,4}, and Daniel
5 Livescu⁵

6 ¹Los Alamos National Laboratory, X-Computational Physics Division XCP-4,
7 Los Alamos, NM 87545, USA

8 ²Princeton University, Program in Atmospheric and Oceanic Sciences, Princeton, NJ 08540, USA

9 ³NOAA/OAR Geophysical Fluid Dynamics Laboratory, Ocean and Cryosphere Division,
10 Princeton, NJ 08540, USA

11 ⁴University of Washington, School of Oceanography, Seattle, WA, USA

12 ⁵Los Alamos National Laboratory, Computer Computational and Statistical Physics Division
13 CCS-2, Los Alamos, NM 87545

14 * corresponding author e-mail: bkaiser@lanl.gov

15 Significant advances in the understanding and modeling of dynamical sys-
16 tems has been enabled by the identification of processes that locally and approx-
17 imately dominate system behavior,¹ or dynamical regimes. The conventional
18 regime identification method involves tedious and *ad hoc* parsing of data to ju-
19 diciously obtain scales to ascertain which governing equation terms are dominant
20 in each regime. Surprisingly, no objective and universally applicable criterion
21 exists to robustly identify dynamical regimes in an unbiased manner, neither
22 for conventional nor for machine learning-based methods of analysis. Here, we
23 formally define dynamical regime identification as an optimization problem by
24 using a verification criterion, and we show that an unsupervised learning frame-
25 work can automatically and credibly identify regimes. This eliminates reliance
26 upon conventional analyses, with vast potential to accelerate discovery.² Our
27 verification criterion also enables unbiased comparison of regimes identified by
28 different methods. In addition to diagnostic applications, the verification crite-
29 rion and learning framework are immediately useful for data-driven dynamical
30 process modeling,^{3,4,5,6,7} and are relevant to researchers interested in the de-
31 velopment of inherently interpretable methods⁸ for scientific machine learning.
32 Automation of this kind of approximate mechanistic analysis is necessary for
33 scientists to gain new dynamical insights from increasingly large data streams.

34 Observations of dynamical systems often exhibit patterns of spatial and/or temporal spar-
35 sity within the terms of the relevant governing equations. A *dominant balance*⁵⁰ is a subset
36 of governing equation terms which, on average, dominate the remaining equation terms by
37 at least an order of magnitude. Observed dominant balances are often *non-asymptotic*, with
38 no obvious parameter permitting a series expansion of the equation terms and guarantee-
39 ing that the error associated with each neglected term vanishes uniformly in the parameter
40 limits. A *dynamical regime* is a dominant balance together with boundaries in space and/or
41 time. Crucially, dominant balances are dominant *relative* to the magnitude of the negligible

42 terms within the same regime.

43 The importance of the relative nature of regime identification is best illustrated by
44 d’Alembert’s zero drag paradox,¹⁰ which took over 150 years to be resolved by Prandtl.¹¹
45 The paradox emerged from the assumption of a global, rather than relative, threshold for the
46 importance of frictional forces in fluid flow. Its resolution by Prandtl, who pointed out that
47 frictional forces cannot be ignored in surface boundary layers, informs contemporary knowl-
48 edge of aerodynamic stall and thus partially enables the high safety standards of commercial
49 air travel we enjoy today. Prandtl used the conventional method for dynamical regime identi-
50 fication, which involves forming regime hypotheses from visual cues and dynamical intuition
51 and then testing the hypotheses by comparison with the observed equation term magnitudes
52 for each regime. The tediousness and *ad hoc* manner of this method of regime identifica-
53 tion implies that a formal method is required to automatically identify regimes within large
54 data sets, such as can be found in fields as diverse as nonlinear waves, plasma dynamics,
55 earthquake dynamics, general relativity, quantum field theory, biochemical reaction-diffusion
56 dynamics, fibrillation dynamics, epilepsy, and turbulent flows, fiber optics, biofilm dynamics,
57 weather, and climate dynamics.^{12, 13, 14, 15, 16}

58 In this Article, we propose an objective and robust method for identifying non-asymptotic
59 regimes. We 1) formalize the regime identification problem as an optimization problem, 2)
60 propose a verification criterion to quantify optimal solutions to this problem, 3) propose a
61 custom dominant balance hypothesis selection algorithm, and 4) propose an unsupervised
62 learning framework^{17, 18} for solving the problem, drawing upon recent successes of unsuper-
63 vised learning as a tool for partitioning regimes^{45, 50} and upon the success of dimensional-
64 ity reduction algorithms²⁰ at selecting dominant balance hypotheses.⁵⁰ The automation of
65 regime identification is a necessary first step towards the development of a learning agent
66 capable of developing parameterizations of chaotic dynamics by deploying the same logic
67 system as human scientists: the scientific method. The goal of this approach is to develop
68 scientific machine learning models that are as intelligible as they are predictive, attributes

69 that are not necessarily inversely proportional to one another despite the widely held beliefs
70 to the contrary.⁸

71 **Problem formulation**

72 Given the array of data $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_N]$, consisting of N observations of the D -dimensioned
73 vector of equation terms \mathbf{e}_n , we seek to label each observation with a D -dimensioned hy-
74 pothesis vector \mathbf{h}_n , where $h_{ni} \in \{0, 1\}$ for each n^{th} observation of the i^{th} equation term. We
75 assume that the equation is closed, $\sum_{i=1}^D e_{ni} = 0$, for all observations. The entire array of
76 data is labeled by $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]$, and zeros in each hypothesis vector \mathbf{h}_n indicate equation
77 terms in \mathbf{e}_n that are neglected. We choose a verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H})$, such that the
78 optimal fit hypotheses, \mathbf{H}_{opt} , can be obtained by varying the hypotheses \mathbf{H} to find

$$\mathbf{H}_{\text{opt}} = \begin{cases} \underset{\mathbf{H}}{\text{argmax}} \mathcal{V}(\mathbf{E}, \mathbf{H}) & \text{if } \max \mathcal{V}(\mathbf{E}, \mathbf{H}) > \mathcal{V}(\mathbf{E}, \mathbf{1}) \\ \mathbf{1} & \text{if } \max \mathcal{V}(\mathbf{E}, \mathbf{H}) \leq \mathcal{V}(\mathbf{E}, \mathbf{1}) \end{cases}, \quad (1)$$

79 where $\mathbf{1}$ is an array of ones indicating all equation terms are retained for the entire data array.
80 We use the notation conventions of Bishop,²¹ where scalars are italicized, lower case bold
81 represents one dimensional arrays, and upper case bold represents two or higher dimensional
82 arrays.

83 We propose Equation 1 as a definition of the dynamical regime identification problem,
84 in which one seeks to partition the observations \mathbf{E} into different regimes with different dom-
85 inant balances, as labeled by \mathbf{H}_{opt} . The dominant balances within \mathbf{H}_{opt} can be assigned by
86 conventional methods,²² or they can be assigned by using clustering algorithms to partition
87 data into regimes and subsequently by using dimensionality reduction algorithms to select
88 dominant balances for each regime,⁵⁰ as we shall describe below. The verification criterion
89 $\mathcal{V}(\mathbf{E}, \mathbf{H})$ defines the optimal regime hypotheses *and* it permits objective comparison of the
90 optimal regime hypotheses from different methods.

91 **Verification criterion**

92 We seek an intelligible *verification criterion*, $\mathcal{V}(\mathbf{E}, \mathbf{H})$, where we define intelligible as both
 93 interpretable⁸ *and* congruent with domain knowledge.²³ We propose a definition of optimal
 94 dominant balances as balances that satisfy two conditions for each regime,

- 95 1. the magnitude difference between the selected dominant terms and the negligible terms
 96 must be maximized;
- 97 2. the magnitude difference between the terms within the selected dominant set must be
 98 minimized.

99 If the first condition is not satisfied, then all equation terms should be retained, i.e., they are
 100 all equally dominant. We propose a verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H})$ that is the sample-space
 101 weighted average of the *local magnitude score* for the n^{th} sample,

$$\mathcal{M}_n(\mathbf{e}_n, \mathbf{h}_n) = \frac{\Gamma_n}{1 + \Omega_n} \in [0, 1], \quad (2)$$

102 where Γ_n is the normalized difference of the log of the smallest magnitude equation term of
 103 the set of terms that are considered dominant and the largest magnitude equation term of
 104 the set of terms that are considered negligible. We thus refer to Γ_n as the gap in magnitude
 105 and \mathcal{M}_n as the local magnitude score. Ω_n is a penalty imposed by the difference between
 106 maximum and minimum magnitudes within the set of terms that are considered dominant.
 107 Definitions of Γ_n and Ω_n are provided in *Methods*. The score measures the consistence of
 108 local truncations of the equation with the observed magnitudes of equation terms.

109 We propose the weighted average of the score $\mathcal{M}_n(\mathbf{e}_n, \mathbf{h}_n)$, when averaged over N samples,
 110 as the verification criterion,

$$\mathcal{V}(\mathbf{E}, \mathbf{H}) = \frac{\sum_{n=1}^N w_n \cdot \mathcal{M}_n(\mathbf{e}_n, \mathbf{h}_n)}{\sum_{n=1}^N w_n}, \quad (3)$$

111 where the array of weights $\mathbf{w} = [w_1, \dots, w_N]$ are the discrete differentials of the observed

112 domain, e.g. space and/or time differentials. For example, if the N observations of data
113 set \mathbf{E} are of equation terms distributed across a two-dimensional space, then the verification
114 criterion is the area-weighted average of all scores for each observation. Optimal verification
115 criteria approach unity, $\mathcal{V}(\mathbf{E}, \mathbf{H}) \rightarrow 1$, to favor the identification of regimes that contain
116 dominant balances that reflect the two conditions listed above.

117 **Unsupervised learning framework**

118 We propose an unsupervised machine learning framework^{17,18} that automatically discovers
119 regimes by using the verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H})$ (Equation 3) to solve the problem defined
120 by Equation 1. The framework is depicted in Figure 1, in which the regime identification
121 problem is broken into partitioning, hypothesis selection, and hypothesis testing tasks. The
122 left column outlines the conventional *ad hoc* method of regime identification, and the right
123 column depicts our framework. Our framework intentionally emulates the scientific method:
124 the hypothesized regimes \mathbf{H} are tested by evaluating their fit to the equation data \mathbf{E} by
125 using the verification criterion.

126 The first task shown in Figure 1, row 1, is to partition \mathbf{E} into different regimes. For
127 humans, this task is often the mere act of visually recognizing the difference in dynamics
128 from one sampled region to another. Sonnewald *et al.*⁴⁵ first suggested that the heuristic
129 act of recognizing different regimes can be formulated as a partitioning problem that can be
130 credibly solved using clustering algorithms. Clustering algorithms, a class of unsupervised
131 machine learning algorithms that yield a finite set of categories according to similarities
132 or relationships among its objects,^{24,25} reveal underlying patterns of sparsity in the data.
133 However, resulting clusters are sensitive to the choice of algorithm parameters,⁴⁴ and there
134 is no definition of a cluster that is universal to all clustering algorithms.²⁷

135 The second task, shown in the second row of Figure 1, is to select hypotheses \mathbf{H} for
136 all samples. Humans typically perform this task by estimating characteristic scales from
137 observations and choosing a threshold for each regime by which some terms are deemed

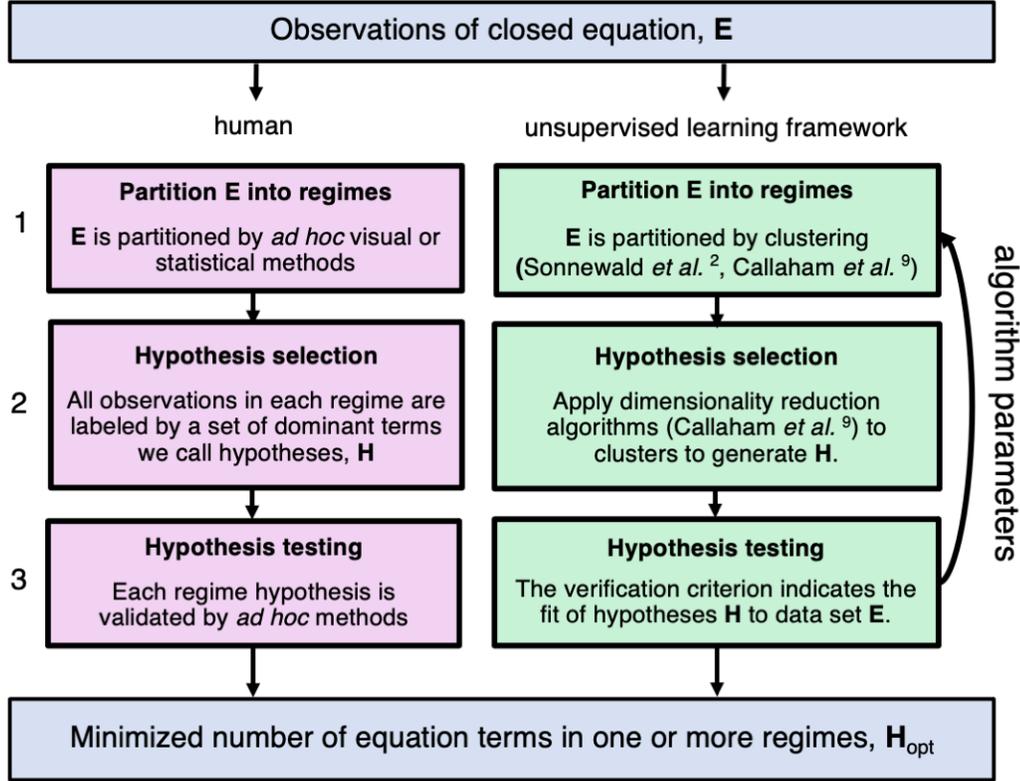


Figure 1: **Diagram of the dynamical regime problem.** Partitioning and empirical scaling analysis performed by a human (left column), and algorithms capable of performing said tasks (right column). The loop over algorithm parameters illustrates the procedure for obtaining \mathbf{H}_{opt} in Equation 1.

138 negligible⁴² for all samples within a regime. Callaham *et al.*⁵⁰ proposed sparse principal
 139 component analysis²⁹ (SPCA) for hypothesis selection because it labels features with small
 140 variances as negligible by performing a least absolute shrinkage and selection operator³⁰
 141 (LASSO) regression on the principal axes from principal component analysis. This applica-
 142 tion of SPCA, or any other dimensionality reduction technique that pertains to convex data,²⁰
 143 is geometrically and statistically consistent with EM clustering algorithms (e.g. K -means,
 144 Gaussian Mixture Models) because both algorithms assume convex, uni-modal, zero-skew
 145 data. We propose a custom algorithm for hypothesis selection, combinatorial hypothesis
 146 selection (CHS, see Methods), in which dominant balances are selected by calculating the
 147 magnitude score (Equation 2) for every possible dominant balance applied to the cluster-
 148 average equation terms and then selecting balance associated with the highest score.

149 The final task shown in Figure 1 is to measure the fit of hypotheses \mathbf{H} to the data \mathbf{E} .
 150 This task was conventionally performed indirectly through *post hoc* validation of models con-
 151 structed using relevant identified regimes. Crucially, the framework applies to any choice of
 152 clustering and hypothesis selection algorithms and, therefore, allows for objective evaluation
 153 and comparisons of different algorithms. We have formalized direct verification of hypothe-
 154 ses by defining the regime identification problem in Equation 1 and proposing a verification
 155 criterion $\mathcal{V}(\mathbf{E}, \mathbf{H})$.

156 The computational complexity of the framework depends on the algorithms chosen for
 157 clustering and hypothesis selection because the complexity of verification criterion is $\mathcal{O}(N)$.
 158 In *Methods* we demonstrate that the computation time of a single pass through the framework
 159 scales polynomially with sample size N for all combinations of a non-parametric and a
 160 parametric clustering algorithm paired with SPCA hypothesis selection and CHS. However,
 161 practical application of the framework requires that the user search a subset of the potentially
 162 infinite range of possible algorithm parameters. Thus, familiarity with the chosen algorithms
 163 and the statistical properties of the data set will reduce the overall computational complexity
 164 and expedite dynamical regime discovery.

165 Examples

166 **Global ocean barotropic vorticity** Sonnewald *et al.*⁴⁵ used K -means clustering, man-
 167 ual hypothesis selection, and algorithm- and problem-specific verification criteria to discover
 168 new and canonical oceanic dynamical regimes. They used the time-mean vertically integrated
 169 barotropic vorticity equation,

$$\overbrace{\nabla \cdot (f\mathbf{U})}^{\text{advection of planetary vorticity}} = \underbrace{\frac{\nabla p_b \times \nabla H}{\rho}}_{\text{bottom pressure torque}} + \underbrace{\frac{\nabla \times \boldsymbol{\tau}}{\rho}}_{\text{wind \& bottom stress curl}} + \underbrace{\nabla \times \mathbf{A}}_{\text{nonlinear torque}} + \underbrace{\nabla \times \mathbf{B}}_{\text{diffusive torque}}, \quad (4)$$

170 which describes the balance of processes that control the rate of solid body rotation of a
 171 column of seawater (see *Methods* for more details).

172 Figures 2a) and 2b) show the optimal dominant balances for each regime and their spa-
 173 tial distributions, respectively, for K -means clustering and CHS, which are quantitatively
 174 similar and qualitatively consistent with the results of Sonnewald *et al.*⁴⁵ The optimal veri-
 175 fication criterion, $\mathcal{V}(\mathbf{E}, \mathbf{H}) = 0.90$, was evaluated at $K = 49$. This result is consistent with
 176 the range of prescribed clusters chosen by Sonnewald *et al.*⁴⁵ using information theoretic
 177 and a custom geographic convergence verification criteria. Figures 2c) and 2d) show the
 178 optimal results for HDBSCAN clustering and CHS, corresponding to a verification criterion
 179 of $\mathcal{V}(\mathbf{V}, \mathbf{H}) = 0.87$. While the the K -means or HDBSCAN clustering results identify similar
 180 mid-latitude balances, the K -means results score higher and include nonlinear balances in
 181 expected locations such as the Gulf Stream on the United States eastern seaboard.

Tumor angiogenesis reaction-diffusion Anderson^{46,32} calculated numerical solutions
 with different permutations of terms eliminated to identify dominant processes in tumor
 angiogenesis (the process by which tumors develop blood flow). We demonstrate that our
 framework directly identifies which terms are dominant without the need for multiple simu-
 lations. The tumor-induced angiogenesis model of Anderson *et al.*⁴⁶ is composed of conser-
 vation laws of three continuous variables (see *Methods*), where the endothelial-cell density
 per unit area (cells that rearrange and migrate from preexisting vasculature to form new
 capillaries), n , is governed by

$$\frac{\partial n}{\partial t} = \overbrace{d\nabla^2 n}^{\text{random motility}} - \overbrace{\nabla \cdot (\chi n \nabla c)}^{\text{chemotaxis}} - \overbrace{\nabla \cdot (\rho n \nabla f)}^{\text{haptotaxis}}, \quad (5)$$

$$= d\nabla^2 n - \chi n \nabla^2 c - \chi \nabla n \cdot \nabla c - n \nabla \chi \cdot \nabla c - \rho n \nabla^2 f - \rho \nabla n \cdot \nabla f, \quad (6)$$

182 where $\chi(c) = \chi_0/(1 + \alpha_0 c)$. Figure 2e) shows the absolute time rate of change of cells as
 183 endothelial cell growth propagates towards the tumor. Figures 2f) and Figures 2g) show

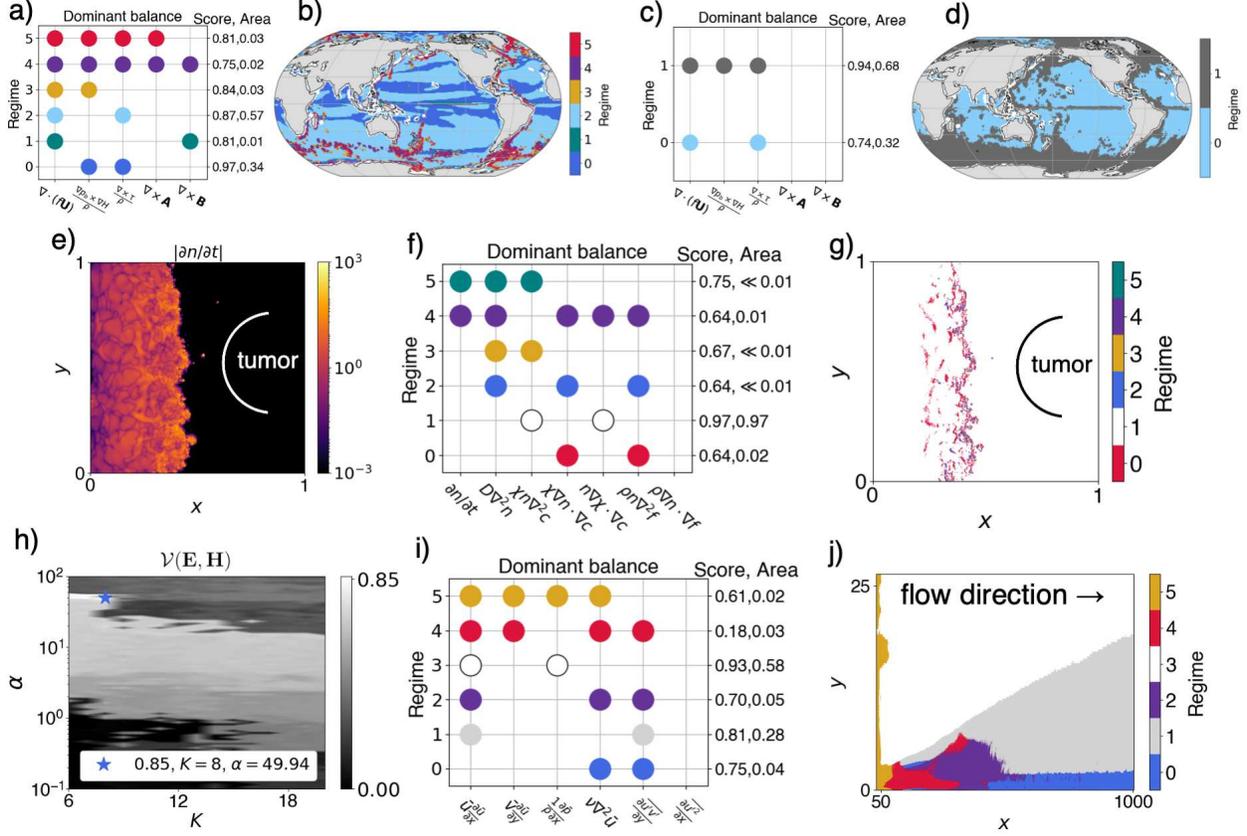


Figure 2: **Examples of framework-identified regimes in two-dimensional spatial domains.** a) to d) show the dominant balances and spatial distributions of regimes in ocean vorticity determined by K -means (a,b) and HDBSCAN clustering (c,d) with CHS. e) to g) show endothelial cell growth rates, dominant balances, and spatial distributions of regimes, respectively, determined by K -means clustering with CHS. h) to j) show the optimization over the verification criterion, dominant balances, and spatial distributions of regimes, respectively, determined by GMM clustering with SPCA hypothesis selection for a spatially developing turbulent boundary layer.

184 the dominant balances and spatial distributions of the optimal regimes identified by using
 185 K -means clustering and CHS. The optimal verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H}) = 0.96$ occurred
 186 at $K = 9$. The results from only one simulation suggest that the fastest cell growth is a
 187 residual of a dominant chemotactic-haptotactic balance, $\chi \nabla n \cdot \nabla c \sim \rho n \nabla^2 f$, (cluster 0, red)
 188 in the regions of tissue.

189 **Turbulent boundary layers** Canonical turbulent boundary layer dynamical regimes,²²
 190 previously identified by Callahan *et al.*⁵⁰ using Gaussian Mixture Model (GMM) clustering

191 with SPCA hypothesis selection and no quantitative verification criteria, can be identified
 192 automatically. Turbulent boundary layers (TBLs) develop as a high-speed flow blows over
 193 non-deformable surfaces. The equation that governs the velocity in the direction of the mean
 194 flow, u , is

$$\overbrace{\overline{u} \frac{\partial \overline{u}}{\partial x} + \overline{v} \frac{\partial \overline{u}}{\partial y}}^{\text{mean momentum flux divergence}} = \underbrace{-\frac{1}{\rho} \frac{\partial \overline{p}}{\partial x}}_{\text{mean pressure gradient}} + \underbrace{\nu \nabla^2 \overline{u}}_{\text{mean momentum diffusion}} - \underbrace{\frac{\partial \overline{u'v'}}{\partial y} - \frac{\partial \overline{u'^2}}{\partial x}}_{\text{turbulent momentum flux divergence}}, \quad (7)$$

195 where the velocity and pressure fields (u, v, p) have been decomposed into mean and fluctu-
 196 ating components denoted by overbars and primes, respectively. The x direction points in
 197 the downwind direction, and the y direction points in the direction normal to the surface.

198 Figures 2h) shows the framework optimization over LASSO regression coefficient α and
 199 prescribed number of clusters K with the optimal verification criterion of $\mathcal{V}(\mathbf{E}, \mathbf{H}) = 0.85$
 200 for $K = 8$ and $\alpha = 49.94$. The optimal regimes are shown in Figures 2i) and 2j). The
 201 regimes are consistent with the results of Callaham *et al.*⁵⁰ and with domain knowledge,³³
 202 but notably our framework required no fluid dynamical knowledge.

203 Conclusions

204 We, for the first time, have formalized the dynamical regime identification problem by defin-
 205 ing it as the maximization of a verification criterion (Equations 1-3). Our formalism is
 206 independent of the method by which the optimization problem is solved, thus transforming
 207 a previously *ad hoc* method of dynamical analysis into a method for objective regime identi-
 208 fication. This formalism is the foundation of our unsupervised learning framework, in which
 209 equation data is partitioned by clustering algorithms,^{50,45} labeled as dominant balances by
 210 dimensionality reduction algorithms,⁵⁰ and the fit of the labels to the data is evaluated by
 211 the verification criterion. This framework is repeated over a user-specified range of algorithm
 212 parameters to find the optimal regimes as defined by the highest verification criterion. We
 213 show that our framework yields results consistent with domain knowledge and previous stud-

ies,^{50,45} and we emphasize that the framework is broadly applicable to chaotic systems. We anticipate that this work could dramatically expedite the discovery of unknown regimes in new data and accelerate efforts in data-driven dynamical process modeling.^{3,4,5,6,7,34,?} This work can be seen as a first step in a new paradigm in the development of machine learning for scientific applications in which the algorithms are inherently intelligible because they are constructed in a manner that explicitly incorporates the scientific method.

Main references

¹ Barenblatt, G. I. *Scaling, Self-similarity, and Intermediate Asymptotics: Dimensional Analysis and Intermediate Asymptotics*. 14 (1996).

² Sonnewald, M. *et al.* Bridging observation, theory and numerical simulation of the ocean using machine learning. *arXiv preprint, arXiv:2104.12506* (2021).

³ Rudy, S. H., Brunton, S. L., Proctor, J. L. & Kutz, J. N. Data-driven discovery of partial differential equations. *Science Advances* **3**, e1602614 (2017).

⁴ Raissi, M. Deep hidden physics models: Deep learning of nonlinear partial differential equations. *The Journal of Machine Learning Research* **19**, 932–955 (2018).

⁵ Rackauckas, C. *et al.* Universal differential equations for scientific machine learning. *arXiv preprint arXiv:2001.04385* (2020).

⁶ Reichstein, M. *et al.* Deep learning and process understanding for data-driven earth system science. *Nature* **566**, 195–204 (2019).

⁷ Schneider, T., Lan, S., Stuart, A. & Teixeira, J. Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations. *Geophysical Research Letters* **44**, 12–396 (2017).

- 236 ⁸ Rudin, C. Stop explaining black box machine learning models for high stakes decisions
237 and use interpretable models instead. *Nature Machine Intelligence* **1**, 206–215 (2019).
- 238 ⁹ Callaham, J. L., Koch, J. V., Brunton, B. W., Kutz, J. N. & Brunton, S. L. Learning
239 dominant physical processes with data-driven balance models. *Nature communications* **12**,
240 1–10 (2021).
- 241 ¹⁰ d’Alembert, J. L. R. *Essai d’une Nouvelle Théorie de la Rrésistance des Fluides* (1752).
- 242 ¹¹ Prandtl, L. Über flussigkeitsbewegung bei sehr kleiner reibung. *Verhandl. III, Internat.*
243 *Math.-Kong., Heidelberg, Teubner, Leipzig, 1904* 484–491 (1904).
- 244 ¹² Strogatz, S. H. *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology,*
245 *Chemistry, and Engineering* (1994).
- 246 ¹³ Blow, K. J. & Wood, D. Theoretical description of transient stimulated raman scattering
247 in optical fibers. *IEEE Journal of Quantum Electronics* **25**, 2665–2673 (1989).
- 248 ¹⁴ Seminara, A. *et al.* Osmotic spreading of bacillus subtilis biofilms driven by an extracellular
249 matrix. *Proceedings of the National Academy of Sciences* **109**, 1116–1121 (2012).
- 250 ¹⁵ Vallis, G. K. *Atmospheric and Oceanic Fluid Dynamics* (2017).
- 251 ¹⁶ Peixoto, J. P. & Oort, A. H. *Physics of climate* (1992).
- 252 ¹⁷ Kohavi, R. & John, G. H. Wrappers for feature subset selection. *Artificial intelligence* **97**,
253 273–324 (1997).
- 254 ¹⁸ Dy, J. G. & Brodley, C. E. Feature selection for unsupervised learning. *Journal of machine*
255 *learning research* **5**, 845–889 (2004).
- 256 ¹⁹ Sonnewald, M., Wunsch, C. & Heimbach, P. Unsupervised learning reveals geography of
257 global ocean dynamical regions. *Earth and Space Science* **6**, 784–794 (2019).

- 258 ²⁰ Van Der Maaten, L., Postma, E. & Van den Herik, J. Dimensionality reduction: a com-
259 parative. *J Mach Learn Res* **10**, 13 (2009).
- 260 ²¹ Bishop, C. M. *Pattern Recognition and Machine Learning* (2006).
- 261 ²² Tennekes, H. & Lumley, J. *A First Course in Turbulence* (1972).
- 262 ²³ Wang, J., Oh, J., Wang, H. & Wiens, J. Learning credible models. In *Proceedings of the*
263 *24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,
264 2417–2426 (2018).
- 265 ²⁴ MacQueen, J. *et al.* Some methods for classification and analysis of multivariate obser-
266 vations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and*
267 *probability*, vol. 1, 281–297 (Oakland, CA, USA, 1967).
- 268 ²⁵ Hartigan, J. A. *Clustering algorithms* (John Wiley & Sons, Inc., 1975).
- 269 ²⁶ Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
270 *Research* **12**, 2825–2830 (2011).
- 271 ²⁷ Estivill-Castro, V. Why so many clustering algorithms: a position paper. *ACM SIGKDD*
272 *explorations newsletter* **4**, 65–75 (2002).
- 273 ²⁸ Zohuri, B. *Dimensional Analysis beyond the Pi Theorem* (2017).
- 274 ²⁹ Zou, H., Hastie, T. & Tibshirani, R. Sparse principal component analysis. *Journal of*
275 *computational and graphical statistics* **15**, 265–286 (2006).
- 276 ³⁰ Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal*
277 *Statistical Society: Series B (Methodological)* **58**, 267–288 (1996).
- 278 ³¹ Anderson, A. R. & Chaplain, M. A. J. Continuous and discrete mathematical models of
279 tumor-induced angiogenesis. *Bulletin of mathematical biology* **60**, 857–899 (1998).

280 ³² Anderson, A. R., Chaplain, M. A., Newman, E. L., Steele, R. J. & Thompson, A. M. Math-
281 ematical modelling of tumour invasion and metastasis. *Computational and mathematical*
282 *methods in medicine* **2**, 129–154 (2000).

283 ³³ Schetz, J. A. & Bowersox, R. D. *Boundary Layer Analysis* (2011).

284 ³⁴ Schmidt, M. & Lipson, H. Distilling free-form natural laws from experimental data. *science*
285 **324**, 81–85 (2009).

286 **Methods**

287 **Data Availability Statement**

288 The tumor angiogenesis reaction-diffusion and synthetic datasets generated during and anal-
289 ysed during the current study are available from the corresponding author on reasonable
290 request. The Estimating the Circulation and Climate of the Ocean (ECCO) ocean state es-
291 timate (<https://data.nas.nasa.gov/ecco/>), version 4 release 2,^{35,36,37,38} was used to estimate
292 the global ocean vorticity budget. The spatially-developing turbulent boundary layer data is
293 available in the Johns Hopkins University turbulence database (<http://turbulence.pha.jhu.edu>)
294 in the developing boundary layer repository.³⁹

295 **Code Availability Statement**

296 The tumor angiogenesis reaction-diffusion code is publicly available at
297 github.com/bekaiser/tumor_dynamics. Pending permission from Los Alamos National Lab-
298 oratory, the unsupervised learning framework and synthetic data generation codes will be
299 publicly available at github.com/bekaiser.

300 **An intelligible verification criterion**

301 The local order-of-magnitude score, $\mathcal{M}_n(\mathbf{e}_n, \mathbf{h}_n)$, hereafter the local magnitude score, per-
302 taining to a single observation, measures the magnitude gap between dominant terms $\mathbf{h}_n \cdot \mathbf{e}_n$
303 and negligible terms $|\mathbf{h}_n - \mathbf{1}| \cdot \mathbf{e}_n$ in a single observation \mathbf{e}_n (recall that the terms that are
304 selected as dominant are labeled by $h_{ni} = 1$, and the neglected terms are labeled by $h_{ni} = 0$).
305 Define $F = \{1, \dots, D\}$ as the *index set*⁴⁰ of the indices of the full set equation terms in vector
306 \mathbf{e}_n , such that

$$\mathbf{e}_n = \bigcup_{i \in F} e_{ni}, \quad (8)$$

307 for observation n such that $1 \leq n \leq N$.

308 We refer to the binary sets that represent the dominant terms as *hypotheses* because
 309 they represent informal equation truncations that are not guaranteed to have asymptotic
 310 properties. The hypotheses for the entire data set \mathbf{E} form an array, \mathbf{H} , which has the same
 311 dimensions as \mathbf{E} , [number of samples \times number of equation terms]. The hypothesis vectors
 312 for each observation can be expressed as

$$\mathbf{h}_n = \bigcup_{i \in \mathbf{F}} h_{ni}, \quad (9)$$

313 where \mathbf{h}_n is an indicator function⁴¹ that consists entirely of ones and zeros, which represent
 314 selected dominant terms and negligible terms, respectively.

315 The indices of elements in \mathbf{e}_n that are selected as dominant terms by the hypothesis \mathbf{h}_n
 316 form the selection index set S_n , where

$$S_n \subseteq \mathbf{F}. \quad (10)$$

317 The number of selected elements may vary for each observation n , and if $S_n = \mathbf{F}$ then $\mathbf{h}_n = \mathbf{1}$
 318 and no equation terms are neglected. It follows that the remainder index set R_n for the n^{th}
 319 observation is defined by set subtraction

$$R_n = \mathbf{F} - S_n, \quad (11)$$

320 and, therefore, the remainder index set and selected index set are non-overlapping,

$$R_n \cap S_n = \emptyset. \quad (12)$$

Thus the cardinality, or size, of the selected index set and remainder index set are $2 \leq \text{card}(S_n) \leq D$ and $0 \leq \text{card}(R_n) \leq D - 2$, respectively. The lower bound of two selected terms is not necessary nor required; we impose it because a dominant balance of just one

term is conceptually ambiguous. Let the arrays of selected and remainder equation terms from \mathbf{e}_n be \mathbf{s}_n and \mathbf{r}_n , respectively. \mathbf{s}_n and \mathbf{r}_n are normalized by the smallest element of \mathbf{e}_n and defined as

$$\mathbf{s}_n = \frac{\bigcup_{i \in S_n} |e_{ni}|}{\min(\bigcup_{i \in F} |e_{ni}|)}, \quad (13)$$

$$\mathbf{r}_n = \frac{\bigcup_{i \in R_n} |e_{ni}|}{\min(\bigcup_{i \in F} |e_{ni}|)}, \quad (14)$$

321 respectively. If $\min(\bigcup_{i \in F} |e_{ni}|) = 0$, then the minimum non-zero absolute valued element
 322 of \mathbf{e}_n replaces the denominators in Equations 13 and 14. Let the relative magnitude gap
 323 between the normalized subsets, Γ , be defined as a scalar for each n^{th} observation:

$$\Gamma_n = \begin{cases} \frac{\log_{10}(\min(\mathbf{s}_n) - \max(\mathbf{r}_n))}{\log_{10}(\min(\mathbf{s}_n) + \max(\mathbf{r}_n))} & \text{if } \min(\mathbf{s}_n) > \max(\mathbf{r}_n) \\ 0 & \text{if } \min(\mathbf{s}_n) \leq \max(\mathbf{r}_n) \end{cases}. \quad (15)$$

324 The magnitude gap Γ is normalized such that $\Gamma \in [0, 1]$, by imposing the floor condition *if*
 325 $\Gamma < 0$ *then* $\Gamma = 0$ to correct for spurious large amplitude negative values of Γ that arise as
 326 $\min(\mathbf{s}_n) \rightarrow \max(\mathbf{r}_n)$. $\Gamma \rightarrow 1$ as the number of orders of magnitude between the element with
 327 the minimum absolute value of the selected subset and element with the maximum absolute
 328 value of the remainder subset approaches infinity, and if $\Gamma = 1$ then the hypothesis that
 329 the terms in the selected subset dominate the terms in the remainder subset is exact (for
 330 any numerical implementation the optimal is limited by machine precision, so $\Gamma \approx 1$). If
 331 the magnitude difference between the two subsets vanishes, then $\Gamma \rightarrow 0$, and if a remainder
 332 subset term exceeds the absolute magnitude of the selected subset, then the hypothesis does
 333 not represent a dominant balance, and $\Gamma = 0$ is prescribed.

334 Since the goal is to choose the selected subset, \mathbf{s}_n , such that it corresponds to the domi-
 335 nant terms, the feature magnitudes of the selected subset should be approximately the same.
 336 Otherwise, the smallest magnitude term(s) in the selected subset should be removed from

337 that subset and added to the remainder subset. To penalize large absolute magnitude dif-
 338 ferences within the selected subset, we introduce the scalar penalty for the n^{th} observation,
 339

$$\Omega_n = \log_{10}(\max(\mathbf{s}_n)) - \log_{10}(\min(\mathbf{s}_n)) \in [0, \infty). \quad (16)$$

340 A base 10 logarithm is chosen for the penalty because it corresponds most directly to the
 341 notion of orders of magnitude. If $\Omega_n \rightarrow 0$, the absolute magnitudes of the selected subset
 342 terms approach uniformity.

343 **Scaling and the local magnitude score**

344 The local magnitude score is a) invariant to the magnitude of the feature vector \mathbf{x}_n and b)
 345 invariant to the sign of the elements of the feature vector,

$$\mathcal{M}_n(\mathbf{e}_n, \mathbf{h}_n) = \mathcal{M}_n(\pm c\mathbf{e}_n, \mathbf{h}_n), \quad (17)$$

346 where c is a positive scalar constant. Therefore, the score is invariant to the choice of di-
 347 mensional or non-dimensional equations, and, equivalently, it can be applied to Buckingham
 348 Π theorem to identify dominant Π groups. Readers unfamiliar with the Buckingham Π the-
 349 orem and how it pertains to the non-dimensionalization of partial differential equations are
 350 referred to Zohuri.⁴²

351 **Combinatorial hypothesis selection**

352 We propose a simple hypothesis selection algorithm that we will refer to as the combinatorial
 353 hypothesis selection (CHS) algorithm. Since the number of all possible hypotheses for an
 354 equation is a permutation of two types (0 or 1) with repetition allowed, the number of possible
 355 hypotheses is $\mathcal{O}(2^D)$. If the number of equation terms, D , is not large, then hypotheses
 356 can be feasibly generated by calculating the magnitude score (Equation 2) for all possible
 357 hypotheses and then selecting the hypothesis that is awarded the highest score. Equation 2

358 can be applied to a single data sample or to an average of samples. The exponential time
 359 complexity limits the feasibility of computing CHS to equations with relatively few terms.

360 Synthetic data for complexity analyses

Consider a two-dimensional array of data with an even number of equation terms, where half of the terms are two orders of magnitude larger in one half of the domain and vice versa, with no variability in the x direction. Figure 3a), in *Extended data* shows the synthetic data e_{ni} consisting of $D = 8$ equation terms, featuring two regimes in which dominant terms have magnitudes of $\mathcal{O}(10)$ and negligible terms have magnitudes of $\mathcal{O}(10^{-1})$. The regimes are separated by a discontinuity at $y = 0.5$. Multiplicative sinusoidal noise is added to give the two regions variance that is proportional to 10% of the signal amplitude in each region. The two regions can be considered dynamical regimes: in each, half of the terms dominate the other equation terms by two orders of magnitude. The regions of dominant terms are prescribed by the Heaviside step function \mathcal{H} , such that:

$$\mathbf{e}_i(x, y) = (-1)^i \eta(y) (\lambda \mathcal{H}(\phi) + \beta), \quad (18)$$

$$\eta(y) = \eta_0 \sin(\omega y), \quad (19)$$

$$\phi = \begin{cases} y - 0.5 & \forall i < D/2 \\ 0.5 - y & \forall i \geq D/2 \end{cases}, \quad (20)$$

361 where x and y are spatial coordinates. The equation closes exactly for all N samples,
 362 $\sum_{i=1}^D e_{ni} = 0$, and the prescribed coefficients are $\lambda = 10^1$, $\beta = 10^{-1}$, $\eta_0 = 10^{-1}$, and $\omega = 10\pi$,
 363 Once again, e_{ni} is the n^{th} observation of the i^{th} feature.

364 Figures 3b), 3c), and 3d) show the results using K -means clustering and SPCA hypoth-
 365 esis selection. Figure 3b) shows the variation of the verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H})$ with α ,
 366 the LASSO regression coefficient for SPCA, and K , the prescribed number of clusters for
 367 K -means clustering. The optimal is marked with the blue star, $\mathcal{V}(\mathbf{E}, \mathbf{H}) = 0.996$, though

368 much of the white band in Figure 3b) corresponds to equivalently optimal results. Figure 3c)
369 shows the dominant balances of the optimal regimes and Figure 3d) shows the spatial distri-
370 bution of the optimal regimes. Identical optimal regimes were identified by using K -means
371 clustering with CHS, by using Hierarchical Density-Based Scan (HDBSCAN) clustering and
372 SPCA hypothesis selection, and HDBSCAN and CHS. The optimal regimes are robust be-
373 cause the magnitude separation between dominant and negligible terms is at least two orders
374 of magnitude everywhere and spatial boundaries of the regimes are discontinuous.

375 While comprehensive complexity analyses are beyond the scope of this Article, we can
376 infer some general properties of the framework’s time complexity. Exhaustive searches over
377 algorithm parameters may very well be NP-hard. The search over K , the prescribed number
378 of clusters for K -means, to minimize the sum of the square of the Euclidean distance of
379 each data point to its nearest center is NP-hard even for just two equation terms,⁴³ $D = 2$.
380 CHS is prohibitively complex at large numbers of equation terms D because its complexity
381 scales with the number of possible dominant balances, $\mathcal{O}(2^D)$. However, SPCA hypothesis
382 selection adds an additional parameter for optimization; therefore, we recommend CHS for
383 equations with fewer terms than 10.

384 The average wall time elapsed for the framework computations over algorithm parameter
385 ranges are shown in Figure 3. The algorithm parameter ranges were specified as follows: for
386 K -means clustering the number of prescribed clusters K was specified as $K = \{2, \dots, 10\}$
387 and the other hyperparameters were the default choices as provided by SciKit Learn.⁴⁴
388 For HDBSCAN clustering the prescribed minimum number of samples for a cluster was
389 specified as 100 samples, and the minimum cluster size was varying from 2000 samples to
390 3000 samples. For hypothesis selection by SPCA, the LASSO regression coefficient was varied
391 between 10^{-2} and 10^2 . Figures 3e) and 3f) show the time complexity with respect to number
392 of samples N and number of equation terms D , respectively. Each point represents the
393 average wall time for one pass through the framework (Figure 1). Figure 3e) shows that the
394 computation time scales polynomially with sample size N for all algorithm choices. Figure

395 3f) shows that CHS becomes prohibitively complex with increasing number of equation terms
 396 because its complexity scales with $\mathcal{O}(2^D)$. However, SPCA adds an additional parameter
 397 for optimization; therefore, we recommend choosing CHS if the number of equation terms is
 398 less than 10.

399 **Global ocean barotropic vorticity**

400 We use the vorticity data processing method of Sonnewald *et al.*,⁴⁵ who computed a 20-
 401 year mean of the ECCO data at 1° resolution to calculate terms of the vertically integrated
 402 barotropic vorticity equation. The variables in Equation 4 are: f is the Coriolis parameter, \mathbf{U}
 403 is the vertically integrated horizontal velocity, p_b is the bottom pressure, H is the depth, ρ is a
 404 reference density, $\boldsymbol{\tau}$ represents surface stress, ∇ is applied only to the horizontal coordinates,
 405 \mathbf{A} contains nonlinear horizontal momentum fluxes, and \mathbf{B} contains linear horizontal diffusive
 406 fluxes.

407 **Tumor angiogenesis reaction-diffusion**

We solve the non-dimensional, tumor-induced angiogenesis governing equations of Anderson
*et al.*⁴⁶ The tumor angiogenic factor concentration, c (chemicals secreted by the tumor
 that promote angiogenesis), and the fibronectin concentration, f (macromolecules that are
 secreted by n and stimulate the directional migration of n), are governed by

$$\frac{\partial f}{\partial t} = \beta n - \gamma n f, \quad (21)$$

$$\frac{\partial c}{\partial t} = -\eta c n. \quad (22)$$

408 Endothelial cell migration up the fibronectin concentration gradient is termed haptotaxis,^{47,48}
 409 while endothelial cell migration up the gradient of tumor angiogenic factor concentration is
 410 termed chemotaxis.⁴⁹

411 In Figure 2c), the tumor is located at $x, y = 1, 0.5$, and the endothelial cell growth

412 is propagating in the positive x direction towards the tumor. Figures 2d) and 2d) show
 413 the regime distributions and dominant balances, respectively, for the optimal results for
 414 K -means clustering and CHS hypothesis selection.

415 We numerically solve the same problem as Anderson *et al.*,⁴⁶ with the exception that 1%
 416 amplitude red noise was added to the initial c and f fields to provide additional variability
 417 for illustrative purposes. A second-order accurate finite difference code was used to calculate
 418 each term in the expanded form of the endothelial cell density equation, such that \mathbf{E} is
 419 composed of observations of the terms in Equation 6). We employ the same boundary
 420 conditions, initial conditions, and constant coefficients (d , α_0 , χ_0 , ρ , β , γ , and η) as⁴⁶ at
 421 double the resolution. Second-order finite differences were employed for spatial derivatives,
 422 and 4th-order adaptive Runge-Kutta was employed for the temporal evolution. No flux
 423 boundary conditions were applied to all four boundaries of the square domain:

$$\mathbf{n} \cdot (d\nabla n - \chi(c)n\nabla c - \rho n\nabla f) = 0, \quad (23)$$

424 where \mathbf{n} is the unit normal vector to the boundaries. The initial conditions, for a circular
 425 tumor (TAF distribution) some distance from three clusters of endothelial cells, are:

$$c(x, y, 0) = \begin{cases} 1, & 0 \leq r \leq 0.1 \\ \frac{(\nu-r)^2}{\nu-r_0}, & 0.1 < r \leq 1 \end{cases}, \quad (24)$$

426 where $r = \sqrt{(x - x_0)^2 + (y - y_0)^2}$.

$$f(x, y, 0) = ke^{-\frac{x^2}{\epsilon_1}}, \quad (25)$$

427

$$n(x, y, 0) = e^{-\frac{x^2}{\epsilon_2}} \sin^2(6\pi y), \quad (26)$$

428 where $\nu = (\sqrt{5} - 0.1)/(\sqrt{5} - 1)$, $r_0 = 0.1$, $x_0 = 1$, $y_0 = 1/2$, $k = 0.75$, $\epsilon_1 = 0.45$, $\epsilon_2 = 0.001$.

429 The constant coefficients were specified as $d = 0.00035$, $\alpha_0 = 0.6$, $\chi_0 = 0.38$, $\rho = 0.34$,
430 $\beta = 0.05$, $\gamma = 0.1$, and $\eta = 0.1$.

431 **Turbulent boundary layers**

432 We use the same turbulence simulation data set and data processing method as Callaham *et*
433 *al.*,⁵⁰ where ρ and ν are constants that represent the fluid density and kinematic viscosity,
434 respectively. The overbar averaging operator represents averaging over the spanwise direction
435 as well as averaging over time, and the diffusion operator is defined as $\nabla^2 = \partial^2/\partial x^2 + \partial^2/\partial y^2$.

436 **Methods references**

437 ³⁵ Forget, G. *et al.* Ecco version 4: An integrated framework for non-linear inverse modeling
438 and global ocean state estimation (2015).

439 ³⁶ Wunsch, C. & Heimbach, P. Dynamically and kinematically consistent global ocean cir-
440 culation and ice state estimates. In *International Geophysics*, vol. 103, 553–579 (Elsevier,
441 2013).

442 ³⁷ ECCO Consortium. A twenty-year dynamical oceanic climatology: 1994-2013. part 1: Ac-
443 tive scalar fields: Temperature, salinity, dynamic topography, mixed-layer depth, bottom
444 pressure. (2017).

445 ³⁸ ECCO Consortium. A twenty-year dynamical oceanic climatology: 1994-2013. part 2:
446 Velocities, property transports, meteorological variables, mixing coefficients (2017).

447 ³⁹ Zaki, T. A. From streaks to spots and on to turbulence: exploring the dynamics of
448 boundary layer transition. *Flow, turbulence and combustion* **91**, 451–473 (2013).

449 ⁴⁰ Munkres, J. R. *Topology* (2000).

- 450 ⁴¹ Cormen, T. H., Leiserson, C. E., Rivest, R. L. & Stein, C. *Introduction to Algorithms*
451 (2009).
- 452 ⁴² Zohuri, B. *Dimensional Analysis beyond the Pi Theorem* (2017).
- 453 ⁴³ Mahajan, M., Nimbhorkar, P. & Varadarajan, K. The planar k-means problem is np-hard.
454 *Theoretical Computer Science* **442**, 13–21 (2012).
- 455 ⁴⁴ Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
456 *Research* **12**, 2825–2830 (2011).
- 457 ⁴⁵ Sonnewald, M., Wunsch, C. & Heimbach, P. Unsupervised learning reveals geography of
458 global ocean dynamical regions. *Earth and Space Science* **6**, 784–794 (2019).
- 459 ⁴⁶ Anderson, A. R. & Chaplain, M. A. J. Continuous and discrete mathematical models of
460 tumor-induced angiogenesis. *Bulletin of mathematical biology* **60**, 857–899 (1998).
- 461 ⁴⁷ Carter, S. B. Principles of cell motility: the direction of cell movement and cancer invasion.
462 *Nature* **208**, 1183–1187 (1965).
- 463 ⁴⁸ Carter, S. B. Haptotaxis and the mechanism of cell motility. *Nature* **213**, 256–260 (1967).
- 464 ⁴⁹ Sholley, M., Ferguson, G., Seibel, H., Montour, J. & Wilson, J. Mechanisms of neovascular-
465 ization. vascular sprouting can occur without proliferation of endothelial cells. *Laboratory*
466 *investigation; a journal of technical methods and pathology* **51**, 624–634 (1984).
- 467 ⁵⁰ Callaham, J. L., Koch, J. V., Brunton, B. W., Kutz, J. N. & Brunton, S. L. Learning
468 dominant physical processes with data-driven balance models. *Nature communications* **12**,
469 1–10 (2021).

470 **Acknowledgements**

471 This work was performed under the auspices of DOE. Financial support comes partly from
472 Los Alamos National Laboratory (LANL), Laboratory Directed Research and Development
473 (LDRD) project "Machine Learning for Turbulence," 20180059DR. LANL, an affirmative
474 action/equal opportunity employer, is managed by Triad National Security, LLC, for the
475 National Nuclear Security Administration of the U.S. Department of Energy under contract
476 89233218CNA000001. Computational resources were provided by the Institutional Comput-
477 ing (IC) program at LANL.

478 MS acknowledges funding from Cooperative Institute for Modeling the Earth System,
479 Princeton University, under Award NA18OAR4320123 from the National Oceanic and At-
480 mospheric Administration, U.S. Department of Commerce. The statements, findings, con-
481 clusions, and recommendations are those of the authors and do not necessarily reflect the
482 views of Princeton University, the National Oceanic and Atmospheric Administration, or the
483 U.S. Department of Commerce.

484 **Author contributions**

485 B.K. drafted the manuscript, contributed to the design of the framework, wrote all software,
486 and performed all example analyses. J.S. substantively revised the manuscript, contributed
487 to the design of the framework, and contributed to the problem formulation. M.S. substan-
488 tively revised the manuscript, contributed to the software implementation of the framework,
489 and contributed the vorticity data set. D.L. revised the manuscript, contributed to the
490 design of the framework, and substantively contributed to the problem formulation.

491 **Competing interest declaration**

492 The authors declare no competing interests.

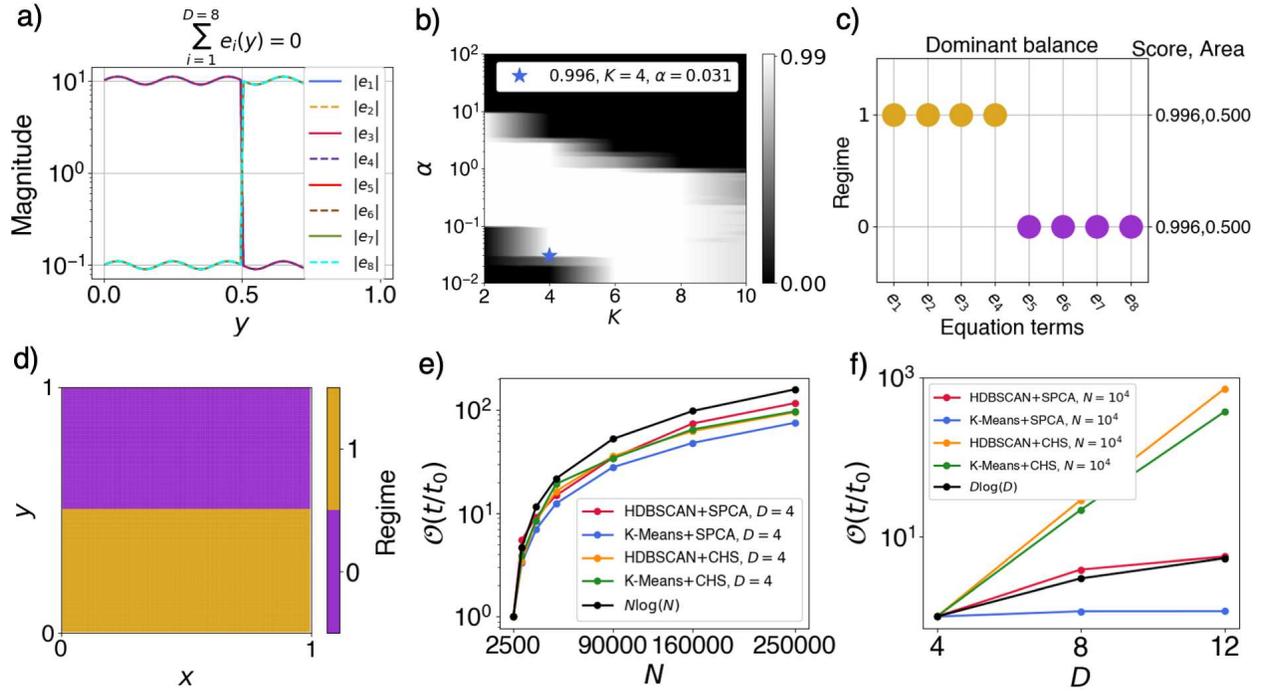


Figure 3: **Synthetic data and computational complexity.** The synthetic data for all y at constant x is shown in a). b) shows the variation global magnitude score as the number of prescribed clusters for K -means and the LASSO regression coefficient α for SPCA are varied. c) and d) show the dominant balances and their spatial distribution for the optimal results that occur in the white band in b) which correspond to the verification criterion $\mathcal{V}(\mathbf{E}, \mathbf{H}) = 0.996$. e) and f) show the variation in wall time as a function of sample size and number of equation terms, respectively, for difference algorithm choices.