
Evaluating Foundation Models for the Brain: A Dynamical Systems Perspective

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Foundation models promise to transform neuroscience and brain–computer in-
2 terfaces (BCIs), but their evaluation remains fragmented and often misleading.
3 Standard benchmarks that emphasize in-distribution accuracy fail to capture what
4 truly matters in dynamical domains: the ability to generalize across conditions that
5 differ from those seen during training. In this perspective, we propose a unified
6 framework for evaluating brain foundation models through the lens of dynamical
7 systems theory. We introduce a generalization spectrum—a hierarchy of distribu-
8 tion shifts spanning system identity, parameter regimes, attractor structure, initial
9 conditions, and observation noise—that clarifies what kinds of robustness should
10 be expected from models claiming to be “foundational.” We then map this spectrum
11 onto a brain-specific taxonomy of distribution shifts—surface (hardware and noise),
12 functional (state and task), and structural (subject and species)—to ground the
13 framework in neuroscientific practice. Building on this foundation, we outline suc-
14 cess criteria for brain foundation models: strong out-of-distribution generalization
15 with minimal data, benchmark-validated performance across diverse tasks, and
16 scalable power-law improvements with model and dataset size. This framework
17 provides a common language for machine learning, neuroscience, and clinical
18 research, and offers a roadmap for building evaluation cultures that distinguish
19 narrow task solvers from truly foundational models.

20

1 Introduction

21 Foundation models [1] have begun to reshape research in neuroscience and brain–computer interfaces
22 (BCIs). By training at scale across diverse datasets [3], they promise reusable representations that
23 can accelerate clinical translation, enable more robust neural decoding, and provide a scientific
24 lens onto brain dynamics. Yet, despite rapid progress, the evaluation of brain foundation models
25 (BFMs) remains fragmented and inconsistent. Most current studies rely on narrow downstream
26 benchmarks—such as classification accuracy on specific EEG datasets—that reward interpolation
27 within the training distribution. Such metrics obscure a central scientific and practical question: *can*
28 *these models generalize across the distribution shifts that inevitably arise in neural systems and their*
29 *measurement?*

30 Traditional evaluation frameworks in machine learning, which emphasize performance on held-out
31 i.i.d. test sets, are ill-suited for dynamical systems [8]. Brains, like other complex dynamical systems,
32 exhibit multiple operating regimes, attractors, and sensitivities to initial conditions. Perturbations
33 in hardware, subject state, or physiology generate test conditions that differ qualitatively from
34 training data. Evaluating BFMs therefore requires moving beyond binary distinctions between “in-
35 distribution” and “out-of-distribution.” Instead, it demands a structured understanding of the *spectrum*
36 of generalization challenges that models must confront.

37 In this paper, we propose a framework for evaluating BFM_s that unifies theory, neuroscience practice,
38 and benchmarking protocols. At its core is a **generalization spectrum** derived from dynamical
39 systems theory, which organizes distribution shifts into a principled hierarchy. From highest-to-lowest
40 level of abstraction, they are: system identity, parameter regime, attractor structure, initial conditions,
41 and observational noise. We map this spectrum onto a brain-specific taxonomy of shifts—**surface**
42 (hardware and noise), **functional** (state and task), and **structural** (subject and species)—to clarify
43 how evaluation should reflect neuroscientific realities. Finally, we outline concrete success criteria
44 and benchmarking platforms that can instantiate these ideas, providing a roadmap for community
45 standards. Our aim is to move evaluation from an ad hoc collection of tasks toward a rigorous
46 scientific probe that distinguishes narrow solutions from truly foundational models.

47 2 The Generalization Spectrum for Dynamical Systems

48 At the heart of evaluating brain foundation models lies the question of how well a model generalizes
49 to conditions beyond those represented in its training data. In static machine learning settings,
50 this is often framed as a binary distinction: a model either performs well on a held-out i.i.d. test
51 set (“in-distribution”) or fails on qualitatively different data (“out-of-distribution”). However, this
52 dichotomy collapses in the context of dynamical systems such as the brain. Here, distribution shifts
53 occur along a continuum, ranging from mild perturbations to fundamental changes in the underlying
54 system. Capturing this continuum requires a framework that recognizes multiple, nested levels of
55 generalization.

56 We propose a **generalization spectrum**, which organizes distribution shifts in dynamical systems into
57 five levels of increasing granularity (see Appendix A for a more rigorous mathematical formulation):

- 58 1. **System-level generalization.** Models are trained on one class of dynamical system but
59 tested on a fundamentally different class. For example, a model trained on Lorenz dynamics
60 but tested on Rössler dynamics, or in neuroscience, a model trained on human EEG but
61 applied to rodent electrophysiology. This represents the most extreme form of generalization
62 [27].
- 63 2. **Regime-level generalization.** The governing equations remain the same, but system pa-
64 rameters change, potentially crossing bifurcation points that alter qualitative behavior. In
65 neuroscience, this corresponds to parameter shifts induced by pharmacological manipu-
66 lations or long-term plasticity [26].
- 67 3. **Attractor-level generalization.** Within a fixed system and parameter setting, the dynamics
68 may admit multiple attractors or operating modes. Generalization here requires predicting
69 trajectories that lie in a different basin of attraction than those seen during training. For
70 neural data, this is analogous to capturing transitions between resting state, task-engaged
71 state, or pathological rhythms [9, 12].
- 72 4. **Initial-condition generalization.** Even within a single attractor, trajectories initialized at
73 different states can evolve in ways unseen during training, especially in chaotic systems
74 where small perturbations can amplify over time. For brain data, this includes variability
75 across trials within the same subject and task [18, 12].
- 76 5. **Noise-level generalization.** The system dynamics remain fixed, but the observation channel
77 is perturbed—for example, changes in sensor noise, hardware differences, or preprocessing
78 pipelines. This is the minimal but most practically pervasive form of distribution shift in
79 EEG and other neural recording modalities [13, 10].

80 This spectrum reframes generalization as a structured hierarchy rather than a binary label. It highlights
81 that failures at different levels carry different implications: a model that fails at the noise level may
82 be unsuitable for deployment, while one that succeeds up to attractor-level generalization but not at
83 system-level still has strong scientific value. By explicitly situating evaluation within this hierarchy,
84 we can diagnose the strengths and limits of brain foundation models with greater precision, and align
85 expectations with both scientific and practical goals.

86 **3 Mapping the Spectrum to Brain-Specific Shifts**

87 While the generalization spectrum provides a theoretical backbone grounded in dynamical systems, its
88 value lies in how it maps onto the practical distribution shifts encountered in neuroscience and brain-
89 computer interface (BCI) research. To bridge theory and practice, we introduce a complementary
90 taxonomy of shifts that arise in neural data: **surface**, **functional**, and **structural**. Each corresponds to
91 a subset of levels in the generalization spectrum, anchoring abstract categories of dynamical change
92 in concrete neuroscientific settings.

93 **3.1 Surface Shifts: Noise and Hardware Variability**

94 Surface shifts are the most immediate and practically pervasive. They occur when the observation
95 channel changes while the underlying neural dynamics remain fixed. Examples include variation
96 across EEG caps, differences in amplifier quality, electrode impedance changes, or the presence of
97 movement artifacts. Such shifts map primarily onto the *noise-level* and, in some cases, the *initial-
98 condition* levels of the spectrum. Success at this level indicates robustness to the day-to-day realities
99 of neural data collection, and is a minimal requirement for BFM s intended for deployment [17, 10, 4].

100 **3.2 Functional Shifts: Brain State and Task Dynamics**

101 Functional shifts reflect changes in the operating mode of the brain without altering its structural
102 identity. They include transitions between alertness and drowsiness, resting state and task-engaged
103 state, or healthy and pathological rhythms. These correspond to *attractor-level* and *regime-level*
104 generalization: the system identity is preserved, but parameters or initial conditions shift the brain
105 into qualitatively different modes of operation. Evaluating generalization at this level probes whether
106 a model captures the global structure of neural dynamics beyond the specific states represented in
107 training data [19, 20].

108 **3.3 Structural Shifts: Subject and Species Differences**

109 Structural shifts are the most demanding, arising when the identity of the system itself changes.
110 This includes variation across subjects, recording modalities, or species. Such shifts map directly
111 to *system-level* generalization in the spectrum. For example, a model trained on human EEG that
112 transfers to nonhuman primate recordings must abstract beyond superficial statistics to capture
113 dynamical principles shared across brains. Success at this level is critical for building BFM s that
114 support cross-subject BCIs or comparative neuroscience [25].

115 **3.4 Unifying View**

116 Together, the surface–functional–structural taxonomy provides an interpretable bridge between
117 the abstract dynamical systems spectrum and the practical realities of brain research. It clarifies
118 which generalization challenges a given benchmark actually tests, and highlights where progress
119 is most urgently needed. Surface robustness ensures reliability, functional generalization supports
120 adaptability, and structural generalization enables scalability across individuals and contexts. By
121 situating evaluation within this taxonomy, the neuroscience and ML communities can converge on a
122 common language for diagnosing and comparing the capabilities of brain foundation models.

123 **4 Success Criteria for Brain/EEG Foundation Models**

124 If foundation models are to fulfill their promise for neuroscience and BCIs, their evaluation must
125 go beyond narrow task accuracy. We propose three ranked criteria that define what it means for a
126 brain foundation model (BFM) to be truly *foundational*. Each criterion isolates a distinct property:
127 *transferability*, *breadth*, and *scaling potential*. Together, they offer a principled way to assess whether
128 models capture reusable structure in neural dynamics.

129 **1. Out-of-Distribution Generalization with Data Efficiency (Transferability)**

130 The first and most fundamental requirement is that BFM^s enable rapid transfer to new tasks with
131 minimal labeled data. This criterion probes whether a model’s internal representations encode
132 reusable neural structure, rather than task-specific correlations. Success means that, when faced with
133 an unfamiliar decoding or forecasting problem, the BFM can be adapted with orders-of-magnitude
134 less supervision than a model trained from scratch [27]. Transferability directly reduces the cost of
135 deploying models across subjects, tasks, and clinical settings.

136 **2. Benchmark-Validated Performance (Breadth)**

137 The second requirement is breadth: a BFM must demonstrate consistent superiority across diverse
138 and standardized evaluation suites. While transferability measures how well the model adapts to new
139 situations, benchmark breadth measures how reliably it performs across a spectrum of existing tasks,
140 brain states, modalities, and noise conditions. Strong results on established baselines [21, 15, 20]
141 indicate that the benefits of large-scale pretraining are realized not just in isolated cases, but across
142 the field. Breadth ensures that a model is not a niche tool, but a robust platform others can trust and
143 build upon.

144 **3. Scalable Power-Law Behavior (Scaling Potential)**

145 The third requirement concerns the trajectory of progress: foundational models must improve
146 predictably as more data and parameters are added. The criterion of power-law scaling is the hallmark
147 of foundation models in other domains [14, 11]. Scaling potential provides the strategic justification
148 for long-term investment: if gains follow a reliable scaling curve, then expanding datasets and
149 architectures will continue to pay off. It also offers a common yardstick to compare very different
150 modeling approaches, since scaling laws expose whether progress is driven by structure learning or
151 by brittle overfitting.

152 **4.1 Interpreting Scaling Failures**

153 Consider the case of building an EEG foundation model. When scaling behavior departs from
154 power-law trends, the implications differ depending on context:

- 155 • **Plateau below threshold:** Performance saturates before reaching levels required for practical
156 applications, suggesting that EEG alone may be insufficient unless paired with other
157 modalities.
- 158 • **Plateau above threshold:** Performance saturates at levels well above application requirements,
159 implying that the task may be too simple to justify large-scale foundation modeling.
- 160 • **Model underperforming, data valuable:** If power-law scaling exists but is weak (small
161 exponent), it’s possible that the data contain reusable structure, but current architectures fail
162 to unlock it. Investment and efforts should shift toward advancing ML techniques. Signals
163 of architectural underfitting can be corroborated by intrinsic-dimension estimates of the
164 objective landscape [16].
- 165 • **No power law (noise-dominated regime):** If model/data scaling yields no performance
166 improvements, this suggests that the signal-to-noise ratio of EEG may be too low to support
167 general-purpose foundation models, restricting utility to narrow, task-specific applications.

168 Together, these success criteria establish a rigorous checklist for BFM^s. A model that demonstrates
169 strong OOD transfer, validates across diverse benchmarks, and scales predictably with resources can
170 legitimately be called foundational. Conversely, models that fail to meet these criteria should be
171 recognized as narrow tools, valuable in context but not general platforms for neuroscience.

172 **5 Benchmarking Protocols and Platforms**

173 Establishing rigorous success criteria requires equally rigorous protocols for measurement. In natural
174 language processing and computer vision, shared benchmarks such as GLUE [24, 23] and ImageNet
175 [6] accelerated progress by standardizing evaluation and making comparisons transparent. Analogous

176 platforms in control and RL (e.g., OpenAI Gym [2], DeepMind Control Suite [22], D4RL [7])
177 demonstrate how shared tasks and standard metrics shape progress. For brain foundation models
178 (BFMs), we argue that the community must similarly converge on shared evaluation protocols that
179 probe robustness, efficiency, and scaling across the generalization spectrum. Again, taking EEG
180 foundation models as a case study, we offer EEG-specific diagnostics.

181 **5.1 EEG-Specific Diagnostics**

182 Electroencephalography provides a practical testbed for BFM evaluation due to its accessibility,
183 prevalence, and inherent challenges of noise and variability. We propose a suite of diagnostic tasks
184 that capture distinct dimensions of robustness and generalization:

- 185 1. **Task-specific benchmarking:** Compare foundation models to baselines trained from scratch
186 with matched data.
- 187 2. **OOD task generalization:** Evaluate zero- or few-shot transfer to tasks not included in
188 pretraining.
- 189 3. **Scaling law validation:** Plot loss against dataset and model sizes to test for power-law
190 behavior.
- 191 4. **Embedding richness:** Compare learned embeddings to raw time-series and handcrafted
192 EEG features using a fixed decoder.
- 193 5. **Decoder vs. embedding contribution:** Vary decoder complexity while freezing embeddings
194 to separate representational power from decoding capacity.
- 195 6. **Channel degradation:** Measure performance as EEG channels are progressively removed
196 to simulate hardware constraints.
- 197 7. **Real-world signal robustness:** Quantify performance under common noise sources such as
198 blinks, muscle artifacts, and impedance variability [13, 5].
- 199 8. **Cross-hardware generalization:** Evaluate transfer across medical-grade and consumer
200 EEG systems.

201 This diagnostic suite provides a multi-dimensional profile of a BFM, making it possible to identify
202 strengths, weaknesses, and suitability for practical use cases.

203 **6 Outlook and Call to Action**

204 Brain foundation models promise to accelerate discovery, translation, and application across neuro-
205 science and brain–computer interfaces. Yet without principled evaluation, it is impossible to know
206 whether these models capture reusable structure in brain dynamics or merely interpolate within narrow
207 datasets. We have outlined a framework for rigorous evaluation built on three pillars: a dynamical-
208 systems-inspired *generalization spectrum*, a neuroscience-grounded *surface–functional–structural*
209 *taxonomy*, and concrete *benchmarking protocols*.

210 Progress in this area requires not only better models but also better evaluation cultures. Community
211 benchmark suites—analogous to ImageNet for vision [6] and GLUE for language [24, 23]—are
212 needed to anchor claims of generality, ensure reproducibility, and accelerate progress through shared
213 goals. By organizing evaluation around structured distribution shifts, researchers can identify both
214 the strengths and limitations of a model and draw clearer implications for scientific and clinical use.

215 Looking ahead, we call for collaboration across machine learning, neuroscience, and clinical com-
216 munities to adopt such frameworks. Surface robustness ensures reliability, functional generalization
217 supports adaptability, and structural generalization enables cross-subject and cross-species transfer.
218 Together, these capabilities would make brain foundation models genuinely foundational.

219 In short, foundation models will only be as transformative as the standards by which they are judged.
220 By uniting theory, taxonomy, and practice, evaluation itself can become a scientific probe that
221 distinguishes narrow tools from models capturing the deeper principles of brain dynamics.

222 **References**

223 [1] Bommasani, R. (2021). On the opportunities and risks of foundation models. *arXiv preprint*
224 *arXiv:2108.07258*.

225 [2] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W.
226 (2016). Openai gym.

227 [3] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A.,
228 Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances*
229 *in neural information processing systems*, 33:1877–1901.

230 [4] Casson, A. J., Yates, D. C., Smith, S. J., Duncan, J. S., and Rodriguez-Villegas, E. (2010).
231 Wearable electroencephalography. *IEEE engineering in medicine and biology magazine*, 29(3):44–
232 56.

233 [5] Delorme, A., Sejnowski, T., and Makeig, S. (2007). Enhanced detection of artifacts in eeg data
234 using higher-order statistics and independent component analysis. *Neuroimage*, 34(4):1443–1449.

235 [6] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-
236 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*
237 *Recognition*, pages 248–255.

238 [7] Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. (2021). D4rl: Datasets for deep
239 data-driven reinforcement learning.

240 [8] Gilpin, W. (2021). Chaos as an interpretable benchmark for forecasting and data-driven modelling.
241 *arXiv preprint arXiv:2110.05266*.

242 [9] Göring, N., Hess, F., Brenner, M., Monfared, Z., and Durstewitz, D. (2024). Out-of-domain
243 generalization in dynamical systems reconstruction. *arXiv preprint arXiv:2402.18377*.

244 [10] Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Parkko-
245 nen, L., and Hämäläinen, M. S. (2014). Mne software for processing meg and eeg data. *neuroimage*,
246 86:446–460.

247 [11] Hestness, J., Narang, S., Ardalani, N., Diamos, G. F., Jun, H., Kianinejad, H., Patwary, M.
248 M. A., Yang, Y., and Zhou, Y. (2017). Deep learning scaling is predictable, empirically. *ArXiv*,
249 [abs/1712.00409](https://arxiv.org/abs/1712.00409).

250 [12] Hu, J., Hu, Y., Chen, W., Jin, M., Pan, S., Wen, Q., and Liang, Y. (2024). Attractor memory
251 for long-term time series forecasting: A chaos perspective. *Advances in Neural Information*
252 *Processing Systems*, 37:20786–20818.

253 [13] Jung, T.-P., Makeig, S., Humphries, C., Lee, T.-W., McKeown, M. J., Iragui, V., and Sejnowski,
254 T. J. (2000). Removing electroencephalographic artifacts by blind source separation. *Psychophysiology*,
255 37(2):163–178.

256 [14] Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford,
257 A., Wu, J., and Amodei, D. (2020). Scaling laws for neural language models. *arXiv preprint*
258 *arXiv:2001.08361*.

259 [15] Lawhern, V. J., Solon, A. J., Waytowich, N. R., Gordon, S. M., Hung, C. P., and Lance, B. J.
260 (2018). Eegnet: a compact convolutional neural network for eeg-based brain–computer interfaces.
261 *Journal of neural engineering*, 15(5):056013.

262 [16] Li, C., Farkhoor, H., Liu, R., and Yosinski, J. (2018). Measuring the intrinsic dimension of
263 objective landscapes. In *International Conference on Learning Representations*.

264 [17] Lopez-Calderon, J. and Luck, S. J. (2014). Erplab: an open-source toolbox for the analysis of
265 event-related potentials. *Frontiers in human neuroscience*, 8:213.

266 [18] Lorenz, E. N. (1969). The predictability of a flow which possesses many scales of motion.
267 *Tellus*, 21(3):289–307.

268 [19] Makeig, S., Debener, S., Onton, J., and Delorme, A. (2004). Mining event-related brain
269 dynamics. *Trends in cognitive sciences*, 8(5):204–210.

270 [20] Roy, Y., Banville, H., Albuquerque, I., Gramfort, A., Falk, T. H., and Faubert, J. (2019).
271 Deep learning-based electroencephalography analysis: a systematic review. *Journal of neural
272 engineering*, 16(5):051001.

273 [21] Schirrmeister, R. T., Springenberg, J. T., Fiederer, L. D. J., Glasstetter, M., Eggensperger, K.,
274 Tangermann, M., Hutter, F., Burgard, W., and Ball, T. (2017). Deep learning with convolutional
275 neural networks for eeg decoding and visualization. *Human brain mapping*, 38(11):5391–5420.

276 [22] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., de Las Casas, D., Budden, D., Abdolmaleki,
277 A., Merel, J., Lefrancq, A., Lillicrap, T., and Riedmiller, M. (2018). Deepmind control suite.

278 [23] Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman,
279 S. R. (2020). Superglue: A stickier benchmark for general-purpose language understanding
280 systems.

281 [24] Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. (2019). Glue: A
282 multi-task benchmark and analysis platform for natural language understanding.

283 [25] Wu, D., Xu, Y., and Lu, B.-L. (2020). Transfer learning for eeg-based brain–computer interfaces:
284 A review of progress made since 2016. *IEEE Transactions on Cognitive and Developmental
285 Systems*, 14(1):4–19.

286 [26] Yin, Y., Ayed, I., de Bézenac, E., Baskiotis, N., and Gallinari, P. (2021). Leads: Learning dy-
287 namical systems that generalize across environments. *Advances in Neural Information Processing
288 Systems*, 34:7561–7573.

289 [27] Zhang, Y. and Gilpin, W. (2024). Zero-shot forecasting of chaotic systems. *arXiv preprint
290 arXiv:2409.15771*.

291 **A Rigorous Definitions of the Generalization Spectrum**

292 **Setup and notation.** Let $(\mathcal{X}, \|\cdot\|)$ be a smooth n -dimensional state space and let a (possibly
293 parameterized) dynamical system be given by a smooth vector field $f_{\Sigma, \theta} : \mathcal{X} \rightarrow \mathbb{R}^n$, where $\Sigma \in \mathcal{S}$
294 indexes the *system identity* (e.g., family of governing equations) and $\theta \in \Theta \subset \mathbb{R}^p$ collects continuous
295 parameters. Denote by $\phi_{\Sigma, \theta}^t : \mathcal{X} \rightarrow \mathcal{X}$ the associated flow (or time- t map for discrete time). Let
296 $\mathcal{A}(f_{\Sigma, \theta})$ be the set of attractors of $f_{\Sigma, \theta}$, and for $A \in \mathcal{A}(f_{\Sigma, \theta})$ let μ_A be an invariant probability
297 measure supported on A (e.g., an SRB measure in chaotic regimes). Let $B(A) \subset \mathcal{X}$ denote the basin
298 of attraction of A .

299 We assume an *observation channel* given by a measurable map $h : \mathcal{X} \rightarrow \mathcal{Y}$ (deterministic sensing)
300 and a family of noise laws $\{\mathsf{K}_\eta(\cdot | y)\}_{\eta \in \mathcal{H}}$ on \mathcal{Y} (stochastic sensing), so that observed data Y_t is
301 generated from

$$X_t = \phi_{\Sigma, \theta}^t(X_0), \quad Y_t \sim \mathsf{K}_\eta(\cdot | h(X_t)).$$

302 A training configuration is a tuple

$$\mathbf{M}_{\text{train}} = (\Sigma, \theta, A, \rho, \eta, h),$$

303 where ρ is a distribution of initial conditions supported on $B(A)$ (often $\rho = \mu_A$ for on-attractor
304 sampling). A test configuration \mathbf{M}_{test} is defined analogously, possibly with altered components.

305 A learning algorithm produces a predictor g (e.g., a forecaster, decoder, or controller). Let $\ell :
306 \mathcal{Y}^{\mathbb{N}} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ be a bounded loss comparing a prediction (which may depend on a context window)
307 to the next observation. For a configuration \mathbf{M} , define the *asymptotic pathwise loss*

$$\mathcal{L}(g; \mathbf{M}) = \limsup_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \ell(g(Y_{0:t-1}), Y_t) \right], \quad (1)$$

308 where the expectation is over $X_0 \sim \rho$ and the observation noise. Under ergodicity of μ_A and
309 stationarity of the observation channel, (1) coincides with the μ_A -expectation of a one-step loss by
310 Birkhoff's theorem.

311 We measure *generalization* by comparing the same learned g across configurations via the *generalization gap*

$$\Gamma(g; \mathbf{M}_{\text{train}} \rightarrow \mathbf{M}_{\text{test}}) = \mathcal{L}(g; \mathbf{M}_{\text{test}}) - \mathcal{L}(g; \mathbf{M}_{\text{train}}).$$

313 Below we define five nested families of allowable test configurations; moving outward increases
314 dynamical novelty.

315 **Level 1: System generalization (new governing equations).** Here the test system may belong
316 to a *different* identity class. Formally, let $\Sigma^* \in \mathcal{S}$ satisfy $\Sigma^* \notin \text{supp}(\mathcal{D}_{\text{train}}^\Sigma)$, where $\mathcal{D}_{\text{train}}^\Sigma$ is
317 the distribution over system identities used to generate training data. Allow arbitrary $\theta^* \in \Theta$,
318 $A^* \in \mathcal{A}(f_{\Sigma^*, \theta^*})$, $\rho^* \ll B(A^*)$, while keeping the observation family (h, η) fixed (unless otherwise
319 specified). A *Level 1 test* admits any $\mathbf{M}_{\text{test}} = (\Sigma^*, \theta^*, A^*, \rho^*, \eta, h)$ with Σ^* disjoint from the training
320 support. Success at Level 1 requires small Γ uniformly over such \mathbf{M}_{test} (e.g., in expectation over a
321 held-out $\mathcal{D}_{\text{test}}^{\Sigma^*}$).

322 **Level 2: Regime generalization (parameter/bifurcation shift).** Fix the identity Σ but shift
323 parameters beyond the training support. Let $\Theta_{\text{train}}, \Theta_{\text{test}} \subset \Theta$ be disjoint (up to negligible overlap).
324 A *Level 2 test* fixes Σ and admits any $\mathbf{M}_{\text{test}} = (\Sigma, \theta^*, A^*, \rho^*, \eta, h)$ with $\theta^* \in \Theta_{\text{test}}$, $A^* \in \mathcal{A}(f_{\Sigma, \theta^*})$,
325 $\rho^* \ll B(A^*)$. Note this explicitly includes *bifurcation* events: qualitative changes in $\mathcal{A}(f_{\Sigma, \theta})$ as θ
326 crosses critical sets.

327 **Level 3: Attractor generalization (new invariant set, same system/regime).** Fix (Σ, θ) . Let
328 $A \neq A^*$ be distinct attractors of $f_{\Sigma, \theta}$ with basins $B(A)$ and $B(A^*)$. A *Level 3 test* keeps (Σ, θ) fixed
329 and admits any $\mathbf{M}_{\text{test}} = (\Sigma, \theta, A^*, \rho^*, \eta, h)$ with $\rho^* \ll B(A^*)$. This probes whether g captures
330 global structure beyond the specific invariant set seen in training.

331 **Level 4: Initial-condition generalization (local within-basin shifts).** Fix (Σ, θ, A) and its basin
332 $B(A)$. Let ρ be the training initial distribution (often μ_A or a neighborhood thereof). A *Level 4*
333 *test* perturbs initial conditions within $B(A)$ while keeping all other components fixed: $\mathbf{M}_{\text{test}} =$

334 $(\Sigma, \theta, A, \rho^*, \eta, h)$ with $\text{supp}(\rho^*) \subset B(A)$ and ρ^* *not* absolutely continuous w.r.t. ρ in general (to
 335 allow targeted shifts).

336 For fine-grained control, decompose local perturbations using the Oseledets splitting (or, locally,
 337 the eigenspaces of the Jacobian $J_f(x)$): for $x \in A$, let $\mathcal{E}_s(x)$ and $\mathcal{E}_u(x)$ denote stable and unstable
 338 subspaces. Define *stable-direction tests* by pushing ρ via maps $x \mapsto x + \delta$ with $\delta \in \mathcal{E}_s(x)$ and small
 339 $\|\delta\|$; and *unstable-direction tests* analogously with $\delta \in \mathcal{E}_u(x)$. Unstable-direction tests typically
 340 induce rapid divergence in trajectory prefixes while remaining within $B(A)$.

341 **Level 5: Observation/noise generalization (surface/channel shift).** Fix $(\Sigma, \theta, A, \rho)$. Let the
 342 observation channel vary within a specified family. Write the stochastic sensing as a Markov
 343 kernel $K_\eta(\cdot | h(x))$ on \mathcal{Y} . A *Level 5 test* admits any $\mathbf{M}_{\text{test}} = (\Sigma, \theta, A, \rho, \eta^*, h^*)$ with (h^*, η^*) in a
 344 prescribed perturbation class. For example, one may constrain observation shifts by

$$\sup_{x \in \mathcal{X}} D(K_{\eta^*}(\cdot | h^*(x)) \| K_\eta(\cdot | h(x))) \leq \varepsilon,$$

345 for a divergence D (e.g., total variation, Wasserstein, or KL where defined), or by structural constraints
 346 (e.g., channel dropouts, downsampling, additive noise with altered covariance). This level formalizes
 347 hardware changes, artifacts, and measurement noise mismatches.

348 **Partial order and severity.** The spectrum induces a natural partial order of novelty:

$$\text{Level 5} \preceq \text{Level 4} \preceq \text{Level 3} \preceq \text{Level 2} \preceq \text{Level 1},$$

349 since each outer level relaxes constraints of the inner ones. A *severity index* can be attached to a test
 350 by metrizing each coordinate: (i) a discrete metric on identities Σ ; (ii) a parameter metric $d_\Theta(\theta, \theta^*)$;
 351 (iii) an invariant-set distance $d_{\mathcal{P}}(\mu_A, \mu_{A^*})$ (e.g., Wasserstein on invariant measures); (iv) a local shift
 352 size $\|\delta\|$ and alignment with $\mathcal{E}_{u/s}$ (e.g., via local Lyapunov exponents); (v) an observation-channel
 353 distance $\sup_x D(K_{\eta^*}(\cdot | h^*(x)) \| K_\eta(\cdot | h(x)))$. These compose into a vector-valued difficulty label for
 354 each \mathbf{M}_{test} .

355 **Evaluation primitives.** Given a trained predictor g from $\mathbf{M}_{\text{train}}$, we report: (i) the asymptotic loss
 356 $\mathcal{L}(g; \mathbf{M}_{\text{test}})$ at each level; (ii) the generalization gap $\Gamma(g; \mathbf{M}_{\text{train}} \rightarrow \mathbf{M}_{\text{test}})$; and (iii) *valid prediction*
 357 τ_ϵ (the largest horizon for which mean error stays below ϵ) to separate short-horizon tracking
 358 from long-horizon invariant adherence. For Level 3, invariant-set fidelity can be scored by comparing
 359 empirical measures of forecasts to μ_{A^*} (e.g., via $d_{\mathcal{P}}$) to capture attractor-shape preservation even
 360 when pointwise errors grow.

361 **Remarks.** (1) The framework extends to random or controlled dynamical systems by letting $f_{\Sigma, \theta}$
 362 define a cocycle over a driving process and by augmenting the configuration with control policies; def-
 363 initions are unchanged with flows replaced by skew-product flows. (2) In neuroscience applications,
 364 the Surface–Functional–Structural taxonomy is captured by: Level 5 (surface: sensing/noise), Lev-
 365 els 4–3 (functional: state/mode within fixed (Σ, θ)), and Levels 2–1 (structural: parameter or identity
 366 changes across subjects/species). (3) The dynamical viewpoint avoids i.i.d. assumptions: losses are
 367 time averages under invariant measures, OOD is expressed as explicit changes in $(\Sigma, \theta, A, \rho, h, \eta)$,
 368 and “difficulty” is measured by dynamical distances rather than solely by sample-distribution diver-
 369 gences.