

The Brittleness of Priors: An Empirical Case for Adaptive Neural Architectures

Dawud Hage¹

¹Independent Researcher

Abstract

Deep and compositional architectures enjoy strong theoretical advantages, yet practical models still hinge on *fixed architectural priors* whose alignment with data is rarely scrutinized. We report a rigorous isoparametric study comparing fully-connected ReLU MLPs, residual MLPs, and periodic-activation networks (SIREN) on canonical 1D function approximation. To isolate architecture from capacity, we employ a strict isoparametric control, matching total parameter counts across all models; to avoid memorization, we inject Gaussian noise, hold out validation data, apply early stopping, and evaluate both interpolation and extrapolation error. Our results reveal a fundamental *brittleness of priors*: no single static architecture dominates. Shallow-and-wide MLPs excel on smooth, non-compositional targets; deeper MLPs win on compositional signals; periodic activations are decisively superior on oscillatory targets; and residual connections, despite optimization benefits, underperform when their identity-preserving bias misaligns with the target family. These empirical trends cohere with theory on depth-enabled expressivity and spectral bias. We conclude that maximal generalization arises only when architectural bias matches data structure, motivating a shift from static design towards heterogeneous modular systems and, in the longer term, *Architecturally Plastic Networks* (APNs).

1 Introduction

The expressive power of depth and composition is well established: deep networks can approximate certain function classes exponentially more efficiently than shallow ones [1, 2]. At the same time, empirical studies reveal a *spectral bias*, showing that standard networks fit low-frequency structure first and struggle with high-frequency content [3]. Specialized priors—for example sinusoidal activations (SIREN) and Fourier features—can invert this behavior by giving the network a bias toward high-frequency structure [4, 5]. We term the strong dependence of generalization on this match between architectural bias and data structure the *brittleness of priors*. This work contributes a controlled, foundational study demonstrating that when model capacity is

held fixed, a misalignment between an architecture’s prior and the data’s underlying structure severely degrades generalization.

2 Isoparametric Experimental Design

To fairly evaluate architectural priors, we designed a rigorous experimental protocol. We compare three architectural families: (i) plain ReLU MLPs of varying depth, (ii) residual MLPs with matching depth, and (iii) periodic-activation networks (SIREN). These models are trained to approximate three canonical 1D function classes: (a) smooth non-compositional polynomials (e.g., x^2), (b) purely oscillatory signals (e.g., $\sin x$), and (c) mixed-compositional functions (e.g., $\sin x + x^2$).

A core principle of our study is a strict *isoparametric control*: matching the total number of trainable parameters across all models—to ensure that performance differences are due to architectural priors, not model capacity. To ensure models learn the underlying function rather than memorizing noisy data, we employ a robust training procedure. We sample points on a fixed interval, add Gaussian noise to the training set, hold out a clean validation split, and train with Adam, using early stopping on the validation MSE to select the best model. We report mean \pm std error across multiple random seeds, evaluating both in-domain (interpolation) and out-of-domain (extrapolation) performance against the true, noise-free function.

3 Results and Key Findings

Our experiments yield three consistent findings. First, *no single architecture is universally optimal*. Different architectural motifs excel on different targets: shallow-and-wide MLPs dominate on smooth polynomials; deeper MLPs are superior for compositional or high-variation signals; and SIREN is decisively best on oscillatory targets. This directly confirms the brittleness of static priors.

Second, *residual bias can misalign*. Despite their well-known optimization advantages, residual MLPs consistently underperform their plain MLP counterparts when learning these functions from scratch.

089 This suggests that ease of optimization does not com-
090 pensate for a fundamentally mismatched inductive
091 bias when the target function is not an identity-like
092 mapping.

093 Third, *extrapolation rankings are consistent*. The
094 performance advantages observed during interpola-
095 tion largely carry over to out-of-domain evaluation.
096 For instance, periodic priors maintain their supe-
097 riority on oscillatory targets beyond the training
098 interval. These empirical patterns align with es-
099 tablished theory on depth and composition [1, 2],
100 spectral-bias analyses [3], and network expressivity
101 characterizations [6, 7].

102 4 Discussion

103 Our evidence that generalization hinges on prior-
104 data alignment has direct implications for current
105 SOTA paradigms. Mixture-of-Experts (MoE) sys-
106 tems, for instance, typically deploy dozens of archi-
107 tecturally homogeneous experts [8]; our results argue
108 compellingly for using a *heterogeneous* toolkit of ex-
109 perts that can cover a more diverse set of inductive
110 biases.

111 This finding also informs the trajectory of Neural
112 Architecture Search (NAS). The field has moved
113 beyond simple design-time searches, with innova-
114 tions like one-shot supernet [9] and hypernetwork-
115 based weight sharing [10] making the search pro-
116 cess more efficient. The frontier has pushed fur-
117 ther into dynamic, training-time paradigms, with so-
118 phisticated methods like Bayesian Population-Based
119 Training that co-evolve an entire population of mod-
120 els, jointly optimizing their weights and hyperpa-
121 rameters throughout a single training run [11]. Our
122 work provides the empirical justification to take
123 the next logical step: moving from these advanced
124 training-time methods to a fully dynamic, *inference-*
125 *time architectural adaptivity*. The pronounced per-
126 formance gaps we observed suggest that a model
127 capable of selecting the correct architectural motif
128 for a given input could achieve superior performance
129 and efficiency compared to any single, static model,
130 however well-optimized.

131 This motivates a clear path forward, as laid out
132 in established roadmaps for the field [12]. We pro-
133 pose a "Chimera" benchmark—a synthetic regime-
134 switching time series—as a concrete testbed to eval-
135 uate and drive progress on this front.

136 5 Conclusion

137 A controlled, isoparametric comparison reveals that
138 architectural priors are powerful but brittle. Max-
139 imal generalization arises only when these priors
140 align with data structure. This foundational finding
141 serves as a robust proof-of-concept, justifying a shift

in focus from designing single, static architectures 142
towards creating adaptive systems that can leverage 143
architectural diversity. The long-term goal this work 144
motivates is the development of *Architecturally Plas-* 145
tic Networks (APNs): models that learn to adapt 146
not only their weights but also their own structure 147
in response to data. 148

References 149

- [1] M. Telgarsky. "Benefits of depth in neural 150
networks". In: *Conference on learning theory*. 151
PMLR. 2016, pp. 1517–1539. 152
- [2] T. Poggio, A. Banburski, and Q. Liao. "The- 153
oretical issues in deep networks". In: *Pro- 154
ceedings of the National Academy of Sciences* 155
117.48 (2020), pp. 30039–30045. 156
- [3] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, 157
M. Lin, F. Hamprecht, Y. Bengio, and A. 158
Courville. "On the spectral bias of neural net- 159
works". In: *International conference on ma- 160
chine learning*. PMLR. 2019, pp. 5301–5310. 161
- [4] V. Sitzmann, J. Martel, A. Bergman, D. Lind- 162
dell, and G. Wetzstein. "Implicit neural repre- 163
sentations with periodic activation functions". 164
In: *Advances in neural information processing 165
systems* 33 (2020), pp. 7462–7473. 166
- [5] M. Tancik, P. Srinivasan, B. Mildenhall, S. 167
Fridovich-Keil, N. Raghavan, U. Singhal, R. 168
Ramamoorthi, J. Barron, and R. Ng. "Fourier 169
features let networks learn high frequency func- 170
tions in low dimensional domains". In: *Ad- 171
vances in neural information processing sys- 172
tems* 33 (2020), pp. 7537–7547. 173
- [6] G. Montúfar, R. Pascanu, K. Cho, and Y. Ben- 174
gio. "On the number of linear regions of deep 175
neural networks". In: *Advances in neural in- 176
formation processing systems* 27 (2014). 177
- [7] M. Raghu, B. Poole, J. Kleinberg, S. Gan- 178
guli, and J. Sohl-Dickstein. "On the expressive 179
power of deep neural networks". In: *internat- 180
ional conference on machine learning*. PMLR. 181
2017, pp. 2847–2854. 182
- [8] N. Shazeer, A. Mirhoseini, K. Maziarz, A. 183
Davis, Q. Le, G. Hinton, and J. Dean. "Out- 184
rageously large neural networks: The sparsely- 185
gated mixture-of-experts layer". In: *arXiv 186
preprint arXiv:1701.06538* (2017). 187
- [9] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. 188
Han. "Once for All: Train One Network and 189
Specialize it for Efficient Deployment". In: *In- 190
ternational Conference on Learning Represen- 191
tations*. 2020. 192

- 193 [10] A. Brock, T. Lim, J. M. Ritchie, and N. We-
194 ston. “Smash: one-shot model architecture
195 search through hypernetworks”. In: *arXiv*
196 *preprint arXiv:1708.05344* (2017).
- 197 [11] X. Wan, C. Lu, J. Parker-Holder, P. J. Ball, V.
198 Nguyen, B. Ru, and M. A. Osborne. “Bayesian
199 Generational Population-Based Training”. In:
200 *Proceedings of the First Conference on Auto-*
201 *mated Machine Learning*. Vol. 188. Proceed-
202 ings of Machine Learning Research. PMLR,
203 2022, pp. 1–27.
- 204 [12] F. Hutter, L. Kotthoff, and J. Vanschoren,
205 eds. *Automated Machine Learning: Methods,*
206 *Systems, Challenges*. Chapter 3: Neural Archi-
207 tecture Search. Springer, 2019.