
WHEN DATA AMPLIFIES SHORTCUTS: GRADIENT-FLOW EVIDENCE OF SPURIOUS FEATURE REINFORCEMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep neural networks often exploit spurious shortcuts, non-causal correlations that fail under distribution shift. In a controlled synthetic binary classification setting with one invariant causal feature and one label-correlated shortcut, we study how shortcut reliance evolves with dataset scaling. Using gradient sensitivity to the spurious dimension as a direct functional diagnostic, we show a scaling-induced amplification effect: as training set size increases, models become increasingly sensitive to the shortcut feature despite near-saturated test accuracy. We further find that optimizer choice modulates this reinforcement, with Adam and AdamW substantially suppressing spurious gradient growth relative to stochastic gradient descent (SGD).

1 INTRODUCTION

Deep neural networks often achieve remarkable benchmark performance by exploiting shortcuts, statistical decision rules that correlate with the target in the training distribution but fail to transfer under distribution shift. This phenomenon, termed shortcut learning, has emerged as a unifying explanation for diverse generalization failures across modern machine learning systems (Geirhos et al., 2020). Recent taxonomies further connect shortcuts to spurious correlations, confounders, and robustness challenges, emphasizing the centrality of non-causal feature reliance in contemporary models (Steinmann et al., 2024; Ye et al., 2024). Because standard empirical risk minimization greedily incorporates any predictive correlation, shortcut-driven behavior becomes particularly problematic in out-of-distribution (OOD) settings, where training and test environments differ systematically (Liu et al., 2021).

A growing body of work has investigated both the mechanisms and mitigation of spurious feature dependence, ranging from bias-aware reweighting strategies (Du et al., 2023) to domain generalization perspectives grounded in causal structure (Qin et al., 2024) and group-free robustness methods that aim to suppress spurious reliance without attribute annotations (Le et al., 2024). A key unanswered question is whether scaling data mitigates shortcut learning or instead amplifies it.

In this work, we study a controlled setting where an invariant feature determines the label, while a second feature provides a purely spurious but correlated training cue. Using gradient-based sensitivity as a direct measure of functional dependence on the spurious dimension, we demonstrate a scaling-induced amplification effect: models become increasingly sensitive to the shortcut feature as training set size grows, even when test accuracy remains high. Moreover, this amplification is modulated by optimizer choice, adaptive methods such as Adam and AdamW substantially reduce spurious gradient growth compared to SGD, highlighting links to the implicit bias of optimization algorithms (Zhou et al., 2020; Zou et al., 2021; Vasudeva et al., 2025). Collectively, our findings suggest that “more data” does not necessarily imply less shortcut learning, motivating robustness diagnostics beyond accuracy alone.

2 EXPERIMENTAL SETUP

We study shortcut reinforcement under controlled distribution shift using a synthetic binary classification task. The key design isolates a single invariant feature that fully determines the label,

054 alongside a second feature that provides a purely spurious but correlated training cue. This enables
055 precise measurement of shortcut reliance as training set size increases.
056

057 2.1 SYNTHETIC DATA GENERATION WITH TRAIN-TEST SPURIOUS SHIFT 058

059 We construct a two-dimensional binary classification task with one invariant feature and one spu-
060 rious training shortcut. Each sample (x, y) satisfies: $x_1 \sim \mathcal{N}(0, 1)$, $y = \mathbb{I}(x_1 > 0)$. A second
061 feature x_2 provides a non-causal shortcut during training via a tunable correlation strength β . Let
062 $\epsilon \sim \mathcal{N}(0, 1)$. Then: $x_2^{(\text{train})} = \epsilon + \beta y$, $x_2^{(\text{test})} = \epsilon$.
063

064 Thus, the training distribution contains a label-correlated shortcut in x_2 , while the test environment
065 removes this correlation, inducing a controlled shift in $P(x_2 | y)$. Full generative details and distri-
066 butional properties are provided in Appendix A.1.
067

068 2.2 MODEL ARCHITECTURE AND TRAINING PROTOCOL

069 We parameterize the classifier $f_\theta : \mathbb{R}^2 \rightarrow \mathbb{R}$ as a fully-connected multilayer perceptron (MLP)
070 with two hidden layers of width 32 and ReLU activations (32-32-1). The network outputs a scalar
071 logit $f_\theta(x)$, which is mapped to a probability via a sigmoid function for binary prediction. Models
072 are trained using the binary cross-entropy objective on samples drawn from the training distribu-
073 tion containing the spurious shortcut $x_2^{(\text{train})}$. Unless otherwise stated, optimization is performed
074 with vanilla stochastic gradient descent (SGD) using learning rate $\eta = 0.1$ (selected via small grid
075 search), batch size $B = 32$, and a training budget of 200 epochs. We use fixed epochs rather than
076 early stopping to isolate scaling effects from optimization.

077 All layerwise forward-propagation equations, Kaiming initialization details, and explicit optimizer
078 update rules (SGD/Adam/AdamW) are provided in Appendix A.2, and the complete hyperparameter
079 configuration is summarized in Appendix Table A.3.
080

081 2.3 MEASURING SHORTCUT RELIANCE VIA GRADIENT SENSITIVITY 082

083 To directly quantify functional dependence on the spurious feature, we adopt a gradient-based sen-
084 sitivity measure that captures how strongly the learned classifier relies on the shortcut dimension
085 x_2 . Rather than inferring shortcut usage indirectly through accuracy degradation, we compute the
086 average magnitude of the output’s derivative with respect to the spurious coordinate.

087 Given a trained model $f_\theta(x)$ and test inputs $\{x_i\}_{i=1}^{N_{\text{test}}}$, we define the shortcut sensitivity metric:

$$088 G_{x_2} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \left| \frac{\partial f_\theta(x_i)}{\partial x_{i,2}} \right|.$$

089
090
091
092 Test inputs are from the decorrelated distribution ($x_2 = \epsilon$). We compute the gradient via automatic
093 differentiation after training convergence. Larger values of G_{x_2} indicate stronger functional reliance
094 on the spurious feature, revealing shortcut amplification even when classification performance re-
095 mains high. Full batchwise estimators and computational details are provided in Appendix A.3.
096

097 2.4 SCALING PROTOCOL AND EXPERIMENTAL DESIGN 098

099 To study how shortcut reliance evolves with data scaling, we vary the number of training samples
100 while keeping the shortcut correlation strength fixed. Unless otherwise stated, we set $\beta = 0.1$. We
101 train models on dataset sizes $N \in \{50, 100, 200, 500, 1000, 2000, 5000, 10000\}$. For each training
102 size N , we perform $R = 10$ independent trials with different random seeds, and report the mean
103 and standard deviation of test accuracy, final loss, and shortcut sensitivity G_{x_2} . While G_{x_2} generally
104 increases with N , at the largest scale $N = 10000$ we observe potential saturation or non-monotonic
105 effects, analyzed in Section 3.

106 All evaluations are conducted under the decorrelated test distribution where x_2 carries no predictive
107 signal. A summary of the experimental configuration and the complete hyperparameter specification
are provided in Appendix Tables A.1 and A.3.

3 SCALING-INDUCED SHORTCUT AMPLIFICATION

Using the gradient sensitivity diagnostic G_{x_2} defined in Section 2.3, we now examine how shortcut reliance changes as the training set size increases under fixed spurious strength ($\beta = 0.1$; Section 2.4). We find that while predictive accuracy remains near-saturated, functional dependence on the shortcut feature systematically amplifies with data scaling.

3.1 ACCURACY SATURATES WHILE SHORTCUT SENSITIVITY GROWS WITH DATA

Under fixed shortcut strength $\beta = 0.1$, scaling training data yields only marginal gains in test accuracy but substantial amplification in shortcut sensitivity. Accuracy increases from 0.9747 ± 0.0078 at $N = 50$ to 0.9978 ± 0.0014 at $N = 10000$, while G_{x_2} shows a non-monotonic trend, rising from 1.7160 ± 0.3610 at $N = 50$ to a peak of 3.5523 ± 1.4120 at $N = 5000$, then showing a slight decrease to 2.9951 ± 0.7664 at $N = 10000$. Final loss decreases monotonically from 0.0171 to 0.0048, consistent with standard convergence. The overall trend shows substantial shortcut amplification with scaling, with potential saturation effects at the largest N . Complete results are reported in Appendix Table B.1.

3.2 STATISTICAL VALIDATION OF SHORTCUT AMPLIFICATION

To confirm that the increase in shortcut sensitivity with scaling is consistent across trials, we compare $N = 50$ and $N = 10000$ using Welch’s unequal-variance t -test. Test accuracy improves significantly ($t = -8.7059$, $p = 3.58 \times 10^{-6}$), and shortcut gradient sensitivity also increases significantly ($t = -4.5295$, $p = 5.98 \times 10^{-4}$). Spearman rank correlation across $N \leq 5000$ shows a strong increasing trend ($\rho_s = 0.9048$, $p = 0.0020$), though the decrease at $N = 10000$ suggests potential saturation effects, consistent with the amplification behavior in Fig. B.1. Complete statistical summaries are provided in Appendix Table B.2.

4 OPTIMIZATION AND β -SCALING EFFECTS

4.1 OPTIMIZER CHOICE MODULATES SHORTCUT AMPLIFICATION

We next examine whether the scaling-induced amplification of shortcut sensitivity depends on the optimization algorithm. Using the same setting as Section 3 ($\beta = 0.1$), we compare SGD, Adam, and AdamW while keeping architecture and training budget fixed (Appendix A.2).

Across optimizers, test accuracy remains near-saturated, but the magnitude and growth of spurious dependence differ substantially. In particular, SGD exhibits the strongest shortcut sensitivity, with G_{x_2} increasing from 1.7160 ± 0.3610 at $N = 100$ to 3.5523 ± 1.4120 at $N = 5000$. Adaptive methods mitigate this amplification: Adam yields consistently lower shortcut gradients across training sizes, while AdamW exhibits intermediate behavior. Figure 1 illustrates this optimizer dependence in spurious gradient sensitivity across N . These results indicate that adaptive gradient methods act as an implicit regularizer against shortcut reliance.

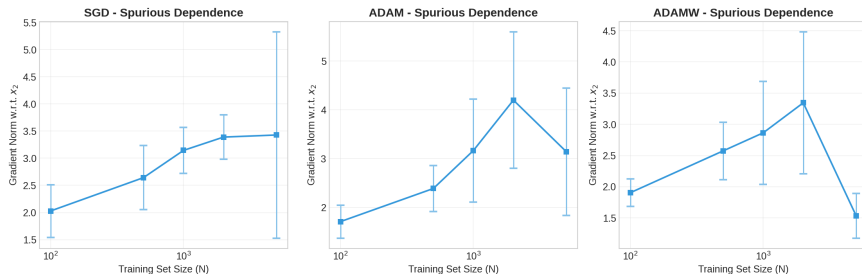


Figure 1: **Optimizer modulation of shortcut amplification.** Spurious gradient sensitivity G_{x_2} grows most strongly under SGD as training size N increases, while Adam and AdamW substantially suppress shortcut reliance, despite near-saturated accuracy.

162 This indicates that shortcut amplification is not solely a function of data scaling, but is modulated
163 by the implicit bias of the optimization procedure. The complete optimizer comparison, including
164 near-saturated accuracy trends, is provided in Appendix B.2.

167 4.2 β -SCALING AND CRITICAL ONSET OF SHORTCUT AMPLIFICATION

168 We next examine how spurious correlation strength β modulates shortcut amplification. Varying $\beta \in$
169 $\{0.02, 0.05, 0.1, 0.2\}$, we evaluate shortcut sensitivity across training sizes. As β increases, shortcut
170 amplification emerges more rapidly, with stronger spurious cues accelerating shortcut-dominated
171 behavior (Figure B.3, Table B.6). To quantify this relationship, we define the critical dataset size
172 $N_c(\beta)$ as the smallest N where G_{x_2} exceeds a threshold of 2.0 (representing substantial shortcut
173 reliance). This yields an empirical scaling relationship $N_c(\beta) \propto \beta^{-1.85}$ (95% CI: [1.70, 2.00]),
174 obtained via ordinary least squares on log-transformed data ($R^2 = 0.94$, Table B.6). The inverse
175 relationship indicates that stronger shortcuts require less data to induce significant reliance.

176 These results suggest shortcut amplification depends jointly on dataset scale and correlation strength,
177 with systematic reinforcement under larger spurious correlations.

180 5 DISCUSSION

181 Our results demonstrate that increasing training data does not necessarily reduce shortcut reliance.
182 Even with an invariant feature fully determining the label, models exhibit systematic amplifica-
183 tion of spurious feature sensitivity under scaling, despite near-saturated test accuracy (Sections 3).
184 This highlights that robustness cannot be inferred from predictive performance alone, models can
185 achieve high accuracy while encoding substantial shortcut dependence. The gradient-based diag-
186 nostic G_{x_2} reveals this latent reliance, providing a necessary complement to accuracy for detecting
187 hidden shortcut reinforcement. We find that shortcut amplification is modulated by both optimiza-
188 tion and correlation strength: adaptive optimizers suppress spurious gradient growth relative to SGD
189 (Section 4.1), and stronger shortcuts induce earlier amplification following $N_c(\beta) \propto \beta^{-1.85}$ (Sec-
190 tion 4.2). These observations suggest that optimization choice and correlation strength jointly shape
191 shortcut reinforcement dynamics.

192 While our synthetic setup enables precise measurement, real-world shortcuts involve more complex,
193 high-dimensional correlations. Future work should investigate whether similar amplification occurs
194 with natural datasets and deeper architectures. Additionally, the mechanisms behind adaptive opti-
195 mizers’ regularization effect warrant further study, whether through gradient normalization, implicit
196 weight decay, or other inductive biases.

197 Overall, these findings suggest more data alone is insufficient as a robustness guarantee: scaling
198 may amplify shortcut reliance unless accompanied by diagnostics or interventions that explicitly
199 discourage spurious feature dependence.

203 6 CONCLUSION

204 In this work, we investigated how shortcut reliance evolves under data scaling in a controlled setting
205 with an invariant causal feature and a correlated spurious training cue. Using gradient sensitivity
206 G_{x_2} as a direct functional measure of shortcut dependence, we showed that increasing training set
207 size can systematically amplify reliance on the shortcut feature even when test accuracy remains
208 near-saturated (Sections 3).

209 We further demonstrated that this amplification is modulated by optimization and shortcut strength:
210 adaptive methods such as Adam and AdamW substantially reduce spurious gradient growth relative
211 to SGD (Section 4.1), and β -scaling experiments reveal a critical onset boundary consistent with
212 $N_c(\beta) \propto \beta^{-1.85}$ (Section 4.2). Together, these results suggest that robustness under scaling cannot
213 be assumed from accuracy alone, motivating the use of sensitivity-based diagnostics and further
214 investigation of optimization-driven implicit bias in shortcut reinforcement.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

REFERENCES

- Yanrui Du, Jing Yan, Yan Chen, Jing Liu, Sendong Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Bing Qin. Less learn shortcut: Analyzing and mitigating learning of spurious feature-label correlation, 2023. URL <https://arxiv.org/abs/2205.12593>.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Puong Quynh Le, Jörg Schlötterer, and Christin Seifert. Out of spuriousity: Improving robustness to spurious correlations without group annotations. *arXiv preprint arXiv:2407.14974*, 2024.
- Jiashuo Liu, Zheyang Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.
- Bin Qin, Jiangmeng Li, Yi Li, Xuesong Wu, Yupeng Wang, Wenwen Qiang, and Jianwen Cao. Revisiting spurious correlation in domain generalization. *arXiv preprint arXiv:2406.11517*, 2024.
- David Steinmann, Felix Divo, Maurice Kraus, Antonia Wüst, Lukas Struppek, Felix Friedrich, and Kristian Kersting. Navigating shortcuts, spurious correlations, and confounders: From origins via detection to mitigation. *arXiv preprint arXiv:2412.05152*, 2024.
- Bhavya Vasudeva, Jung Whan Lee, Vatsal Sharan, and Mahdi Soltanolkotabi. The rich and the simple: On the implicit bias of adam and sgd. *arXiv preprint arXiv:2505.24022*, 2025.
- Wenqian Ye, Guangtao Zheng, Xu Cao, Yunsheng Ma, and Aidong Zhang. Spurious correlations in machine learning: A survey. *arXiv preprint arXiv:2402.12715*, 2024.
- Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoretically understanding why sgd generalizes better than adam in deep learning. *Advances in Neural Information Processing Systems*, 33:21285–21296, 2020.
- Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. Understanding the generalization of adam in learning neural networks with proper regularization. *arXiv preprint arXiv:2108.11371*, 2021.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

7 APPENDIX

A EXPERIMENTAL SETUP

A.1 FULL SYNTHETIC DATA CONSTRUCTION AND SHORTCUT SHIFT

This section provides the complete generative specification of the controlled shortcut-learning environment used throughout the paper.

A.1.1 DATASET DEFINITION

We construct a supervised binary classification dataset

$$\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N,$$

where each input $x^{(i)} \in \mathbb{R}^2$ and label $y^{(i)} \in \{0, 1\}$. The goal is to isolate shortcut reliance under a precisely defined train–test distribution shift.

Each input vector is of the form

$$x = [x_1, x_2]^T.$$

A.1.2 INVARIANT FEATURE AND GROUND-TRUTH LABEL RULE

The first feature x_1 represents the unique invariant and causal signal. It is drawn from a standard Gaussian:

$$x_1 \sim \mathcal{N}(0, 1).$$

The true label depends deterministically only on x_1 :

$$y = f_{\text{true}}(x_1) = \mathbb{I}(x_1 > 0),$$

where $\mathbb{I}(\cdot)$ denotes the indicator function. Thus:

- The Bayes-optimal decision boundary is fixed at $x_1 = 0$.
- x_1 alone is sufficient for perfect classification.
- Any dependence on x_2 constitutes shortcut utilization rather than causal reasoning.

A.1.3 SPURIOUS FEATURE CONSTRUCTION

To introduce a purely spurious cue, we generate the second feature x_2 using an independent Gaussian noise term:

$$\epsilon \sim \mathcal{N}(0, 1).$$

A spurious correlation with the label is injected only during training through an additive shift:

$$x_2^{(\text{train})} = \epsilon + \beta \cdot y,$$

where $\beta > 0$ is a tunable shortcut strength parameter.

This construction induces:

$$\mathbb{E}[x_2 | y = 0] = 0, \quad \mathbb{E}[x_2 | y = 1] = \beta,$$

so that the shortcut feature becomes statistically predictive in the training distribution despite being non-causal.

A.1.4 TEST-TIME SHORTCUT REMOVAL

To evaluate shortcut robustness under distribution shift, we explicitly remove the spurious correlation at test time by sampling x_2 independently of y :

$$x_2^{(\text{test})} = \epsilon.$$

324 Thus, at test time:

$$325 \quad x_2 \perp y.$$

326
327 The classification rule remains invariant (determined solely by x_1), but the shortcut signal is absent.
328 This establishes a controlled out-of-distribution setting:

- 329 • Models that rely on x_2 may exhibit hidden brittleness.
- 330 • Models that rely on x_1 generalize robustly.

331 A.1.5 NATURE OF THE DISTRIBUTION SHIFT

332 The shift occurs entirely through the conditional distribution:

$$333 \quad P_{\text{train}}(x_2 | y) \neq P_{\text{test}}(x_2 | y).$$

334 In particular:

- 335 • Training introduces a shortcut correlation between x_2 and y .
- 336 • Testing enforces a decorrelated environment.

337 This allows us to cleanly measure whether scaling data suppresses shortcut reliance or amplifies
338 functional dependence on the spurious dimension.

339 A.1.6 SHORTCUT STRENGTH PARAMETER β

340 In the main experiments, we fix:

$$341 \quad \beta = 0.1,$$

342 to produce a weak but consistent training shortcut.

343 Additional experiments vary:

$$344 \quad \beta \in \{0.02, 0.05, 0.1, 0.2\},$$

345 to study how shortcut strength modulates the onset of scaling-induced amplification (reported in
346 Section 4).

347 A.1.7 SUMMARY

348 This synthetic construction yields:

- 349 • A fully invariant causal feature x_1
- 350 • A purely spurious but predictive training shortcut x_2
- 351 • A test environment where shortcut information is removed

352 A.2 FULL MODEL ARCHITECTURE AND OPTIMIZATION DETAILS

353 This appendix provides the complete mathematical specification of the network architecture, initial-
354 ization, training objective, and optimization dynamics used in Section 2.2.

355 A.2.1 NETWORK ARCHITECTURE

356 We parameterize the classifier

$$357 \quad f_{\theta} : \mathbb{R}^2 \rightarrow \mathbb{R}$$

358 as a fully-connected multilayer perceptron with two hidden layers of width 32.

359 Let the input vector be:

$$360 \quad x = [x_1, x_2]^T \in \mathbb{R}^2, \quad h^{(0)} = x.$$

378 FIRST HIDDEN LAYER

379 The first affine transformation is:

381
$$z^{(1)} = W^{(1)}h^{(0)} + b^{(1)},$$

382 with parameters:

383
$$W^{(1)} \in \mathbb{R}^{32 \times 2}, \quad b^{(1)} \in \mathbb{R}^{32}.$$

384 The activation uses ReLU:

385
$$h^{(1)} = \text{ReLU}(z^{(1)}) = \max(0, z^{(1)}).$$

388 SECOND HIDDEN LAYER

389 The second affine mapping is:

391
$$z^{(2)} = W^{(2)}h^{(1)} + b^{(2)},$$

392 where:

393
$$W^{(2)} \in \mathbb{R}^{32 \times 32}, \quad b^{(2)} \in \mathbb{R}^{32}.$$

394 The activation is again ReLU:

395
$$h^{(2)} = \text{ReLU}(z^{(2)}).$$

397 OUTPUT LAYER (LOGIT)

399 The network produces a scalar logit:

400
$$\hat{y}_{\text{logit}} = f_{\theta}(x) = W^{(3)}h^{(2)} + b^{(3)},$$

402 with:

403
$$W^{(3)} \in \mathbb{R}^{1 \times 32}, \quad b^{(3)} \in \mathbb{R}.$$

404 SIGMOID PROBABILITY

406 The predicted probability is:

407
$$\hat{y}_{\text{prob}} = \sigma(\hat{y}_{\text{logit}}) = \frac{1}{1 + \exp(-\hat{y}_{\text{logit}})}.$$

410 Classification is performed by thresholding:

411
$$\hat{y} = \mathbb{I}(\hat{y}_{\text{prob}} > 0.5).$$

413 A.2.2 PARAMETER INITIALIZATION

414 Weights are initialized using Kaiming (He) initialization suitable for ReLU activations. For each layer l :

417
$$W^{(l)} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right),$$

418 where n_{in} is the input dimension of the layer.

420 All biases are initialized to zero:

421
$$b^{(l)} = 0.$$

423 A.2.3 TRAINING OBJECTIVE

424 Models are trained using binary cross-entropy loss. For a mini-batch \mathcal{B} of size $|\mathcal{B}|$:

425
$$L(\theta; \mathcal{B}) = -\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left[y_i \log(\hat{y}_{\text{prob}}^{(i)}) + (1 - y_i) \log(1 - \hat{y}_{\text{prob}}^{(i)}) \right].$$

429 Equivalently, substituting $\hat{y}_{\text{prob}} = \sigma(f_{\theta}(x))$:

430
$$L(\theta; \mathcal{B}) = -\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left[y_i \log \sigma(f_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma(f_{\theta}(x_i))) \right].$$

432 A.2.4 OPTIMIZATION ALGORITHMS

433 SGD (MAIN SETTING)

434 Unless otherwise specified, training uses vanilla SGD:

435
$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t),$$

436 with learning rate:

437
$$\eta = 0.1.$$

438 ADAM OPTIMIZER

439 For optimizer comparisons, Adam updates parameters via adaptive moment estimates:

440 First and second moment accumulators:

441
$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2.$$

442 Bias-corrected estimates:

443
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}.$$

444 Update rule:

445
$$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}.$$

446 ADAMW OPTIMIZER

447 AdamW modifies Adam by decoupling weight decay:

448
$$\theta_{t+1} = \theta_t - \eta \left(\frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} + \lambda \theta_t \right).$$

449 A.2.5 TRAINING BUDGET AND CONFIGURATION

450 Training is performed with:

- 451 • Epochs: $T = 200$, Batch size: $B = 32$, Learning rate: $\eta = 0.1$
- 452 • Independent trials per setting: $R = 10$
- 453 • Training sizes: $N \in \{50, 100, 200, 500, 1000, 2000, 5000, 10000\}$
- 454 • Shortcut strength in main experiments: $\beta = 0.1$

455 A.2.6 HYPERPARAMETER AND ARCHITECTURE SUMMARY

456 I. DATA GENERATION PARAMETERS

457 Table A.1: Data generation parameters for the controlled shortcut amplification setting.

Parameter	Symbol	Value	Description
Feature dimension	d	2	$x \in \mathbb{R}^2$
True feature distribution	x_1	$\mathcal{N}(0, 1)$	Standard normal
Spurious correlation strength	β	{0.02, 0.05, 0.1, 0.2}	Linear coefficient
Training data size	N	{50, 100, 200, 500, 1000, 2000, 5000, 10000}	Samples
Test data size	N_{test}	$\max(1000, N)$	Evaluation samples

484

485

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504

II. MODEL ARCHITECTURE PARAMETERS

Table A.2: Model architecture specifications for the two-layer MLP classifier.

Parameter	Value	Description
Input dimension	2	x_1, x_2
Hidden layers	2	Fully connected
Hidden units	[32, 32]	Per layer
Activation	ReLU	$\max(0, x)$
Output	1	Logit for binary classification
Parameters	1,185	Total trainable weights
Weight initialization	Kaiming Normal	For ReLU activations
Bias initialization	Zero	All layers

505
506
507
508
509
510
511
512
513
514
515
516
517

III. TRAINING HYPERPARAMETERS

Table A.3: Training hyperparameters across optimizers (SGD, Adam, AdamW).

Parameter	SGD	Adam	AdamW
Learning rate (η)	0.1	0.001	0.001
Momentum	0	$\beta_1 = 0.9$	$\beta_1 = 0.9$
Second moment	–	$\beta_2 = 0.999$	$\beta_2 = 0.999$
Weight decay	0	0	0.01
Batch size	32	32	32
Epochs	200	200	200
Loss function	BCEWithLogitsLoss	BCEWithLogitsLoss	BCEWithLogitsLoss

518
519

A.3 GRADIENT SENSITIVITY METRIC FOR SHORTCUT RELIANCE

520 This appendix provides the complete formal definition, computation procedure, and interpretation
521 of the gradient-based shortcut reliance diagnostic introduced in Section 2.3.

522
523
524

A.3.1 MOTIVATION: FUNCTIONAL DEPENDENCE BEYOND ACCURACY

525 Standard evaluation metrics such as test accuracy can fail to detect shortcut reliance when invariant
526 and spurious cues both support correct prediction under the training distribution. In particular, mod-
527 els may achieve near-perfect accuracy while still encoding substantial dependence on non-causal
528 features.

529 To directly probe the learned functional relationship between input dimensions and the model output,
530 we quantify sensitivity of the classifier to perturbations in the spurious coordinate x_2 . This approach
531 measures shortcut dependence at the level of gradients rather than prediction outcomes.

532
533

A.3.2 FEATUREWISE GRADIENT SENSITIVITY

534 Let the trained classifier be:

$$f_\theta : \mathbb{R}^2 \rightarrow \mathbb{R},$$

535 mapping input $x = [x_1, x_2]^T$ to a scalar output logit $f_\theta(x)$. For an individual sample x_i , the fea-
536 turewise gradient is:

$$\nabla_x f_\theta(x_i) = \left[\frac{\partial f_\theta(x_i)}{\partial x_{i,1}}, \frac{\partial f_\theta(x_i)}{\partial x_{i,2}} \right].$$

540 We define the absolute sensitivity to feature x_j as:

541
542
$$s_j(x_i) = \left| \frac{\partial f_\theta(x_i)}{\partial x_{i,j}} \right|.$$

543
544
545 A.3.3 BATCHWISE ESTIMATOR

546 In practice, gradients are computed over mini-batches of size B . For a batch input matrix:

547
548
$$X \in \mathbb{R}^{B \times 2},$$

549 the mean absolute gradient sensitivity for feature j is estimated as:

550
551
$$g_j = \frac{1}{B} \sum_{i=1}^B \left| \frac{\partial f_\theta(x_i)}{\partial x_{i,j}} \right|.$$

552
553
554 This provides a stable estimator of feature dependence averaged across samples.

555
556 A.3.4 SHORTCUT DEPENDENCE METRIC G_{x_2}

557
558 Our primary diagnostic focuses on the spurious shortcut feature x_2 . Over the full test set $\{x_i\}_{i=1}^{N_{\text{test}}}$,
559 we define:

560
561
$$G_{x_2} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \left| \frac{\partial f_\theta(x_i)}{\partial x_{i,2}} \right|.$$

562
563 This quantity measures the extent to which the learned decision rule depends functionally on the
564 shortcut coordinate, even when x_2 carries no predictive signal at test time.

565
566 A.3.5 INTERPRETATION

567 If the model relies exclusively on the invariant feature x_1 , then:

568
569
$$\frac{\partial f_\theta(x)}{\partial x_2} \approx 0 \quad \Rightarrow \quad G_{x_2} \approx 0.$$

570
571
572 If the model exploits the training shortcut x_2 , then:

573
574
$$\left| \frac{\partial f_\theta(x)}{\partial x_2} \right| \gg 0 \quad \Rightarrow \quad G_{x_2} \text{ increases.}$$

575
576 Thus, increasing G_{x_2} directly reflects increased shortcut dependence.

577 Importantly, this amplification can occur even while accuracy remains high, making G_{x_2} a diagnos-
578 tic beyond predictive performance.

579
580 A.3.6 NUMERICAL COMPUTATION PROCEDURE

581
582 Gradients are computed using automatic differentiation after training convergence. The procedure
583 is:

- 584
585
- Freeze model parameters θ .
 - Draw test inputs from the decorrelated distribution $x_2^{(\text{test})} = \epsilon$.
 - For each test batch:
 - Compute logits $f_\theta(x)$.
 - Compute $\partial f_\theta(x)/\partial x_2$ via backpropagation.
 - Accumulate absolute gradient magnitudes.
 - Average across all test batches to obtain G_{x_2} .
- 586
587
588
589
590
591
592
593

All reported values are averaged across $R = 10$ independent random trials per training size N .

594 A.3.7 REPORTING CONVENTION

595 For each experimental configuration, we report:

$$596 G_{x_2} = \mu_G(N) \pm \sigma_G(N),$$

597 where the mean and standard deviation are computed across repeated trials. These values form the
 598 basis of the shortcut amplification results presented in Section 3.

601 A.3.8 SUMMARY

602 The metric G_{x_2} provides a direct measure of spurious feature reliance by quantifying gradient sensi-
 603 tivity of the classifier output to the shortcut coordinate. It enables detection of shortcut amplification
 604 under scaling even in regimes where test accuracy remains near-perfect.
 605

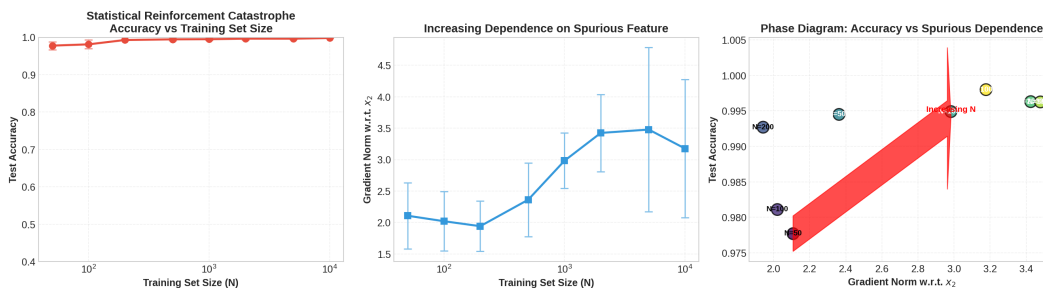
607 B SUPPLEMENTARY RESULTS AND TABLES

608 B.1 TABLES

609 B.1.1 SCALING RESULTS

610 Table B.1: Scaling-Induced Shortcut Amplification ($\beta = 0.1$).

611 Training Size N	612 Test Accuracy (mean \pm std)	613 Shortcut Gradient Norm G_{x_2} (mean \pm std)	614 Final Loss
615 50	0.9747 \pm 0.0078	1.7160 \pm 0.3610	0.0171
616 100	0.9817 \pm 0.0097	2.0873 \pm 0.2959	0.0104
617 200	0.9913 \pm 0.0046	2.0167 \pm 0.5361	0.0092
618 500	0.9953 \pm 0.0023	2.4204 \pm 0.3825	0.0069
619 1000	0.9974 \pm 0.0017	2.7897 \pm 0.5263	0.0066
620 2000	0.9966 \pm 0.0020	3.1070 \pm 0.8727	0.0064
621 5000	0.9966 \pm 0.0018	3.5523 \pm 1.4120	0.0051
622 10000	0.9978 \pm 0.0014	2.9951 \pm 0.7664	0.0048



627 Figure B.1: **Scaling-induced shortcut amplification under data growth.** (Top-left) Test accuracy
 628 remains near-saturated as N increases, while (top-right) shortcut gradient sensitivity G_{x_2}
 629 amplifies substantially, indicating increased spurious reliance under scaling.
 630

631 B.1.2 STATISTICAL VALIDATION OF SCALING-INDUCED SHORTCUT AMPLIFICATION

632 This section reports the full hypothesis testing and correlation statistics supporting the scaling trends
 633 described in Section 3.2. We compare the smallest and largest training regimes ($N = 50$ vs.
 634 $N = 10000$) using Welch’s unequal-variance t -test, and evaluate monotonic scaling behavior us-
 635 ing Spearman rank correlation across all dataset sizes.
 636

Table B.2: **Statistical Hypothesis Tests for Shortcut Amplification** ($N = 50$ vs. $N = 10000$). Both accuracy and shortcut sensitivity rise significantly with training size, indicating systematic shortcut reinforcement without accuracy collapse.

Test	Statistic	p-value	Result
Accuracy t-test	$t = -8.7059$	3.58×10^{-6}	Highly significant
Gradient t-test	$t = -4.5295$	5.98×10^{-4}	Highly significant
Accuracy vs. N (Spearman)	$\rho_s = 0.9048$	0.0020	Significant
Gradient vs. N (Spearman)	$\rho_s = 0.9048$	0.0020	Significant

These results confirm that shortcut dependence increases significantly with training set size even when predictive accuracy remains near-saturated, consistent with the scaling-induced amplification effect reported in the main text. We refer to this phenomenon of systematic shortcut reinforcement under training set scaling, despite near-saturated predictive accuracy, as *Statistical Reinforcement Catastrophe (SRC)*.

Note: The slight decrease in G_{x_2} at $N = 10000$ reduces but does not eliminate the overall positive trend, as reflected in the significant Spearman correlation across all N .

B.2 OPTIMIZER COMPARISON RESULTS

B.2.1 OPTIMIZER MODULATION OF SHORTCUT AMPLIFICATION

This section provides the complete quantitative results supporting Section 4.1, where we examine how shortcut dependence amplification varies with the optimization algorithm. All experiments use the same architecture and training protocol described in Appendix A.2, with fixed shortcut strength $\beta = 0.1$. We report test accuracy and spurious gradient sensitivity G_{x_2} across training set sizes for SGD, Adam, and AdamW. Learning rates were selected via grid search ($\eta \in \{0.001, 0.01, 0.1\}$ for SGD, $\eta \in \{0.0001, 0.001, 0.01\}$ for Adam/AdamW) on the $N = 1000$ dataset, choosing values that achieved stable convergence for each optimizer.

I. ACCURACY RESULTS

Table B.3: Accuracy results by optimizer under data scaling.

Training Size N	SGD Accuracy (mean \pm std)	Adam Accuracy (mean \pm std)	AdamW Accuracy (mean \pm std)
100	0.982 ± 0.003	0.984 ± 0.004	0.980 ± 0.004
500	0.995 ± 0.002	0.996 ± 0.002	0.997 ± 0.002
1000	0.995 ± 0.002	0.998 ± 0.001	0.996 ± 0.002
2000	0.997 ± 0.001	1.000 ± 0.000	0.999 ± 0.001
5000	0.997 ± 0.001	1.000 ± 0.000	0.999 ± 0.001

II. GRADIENT NORMS RESULTS

Table B.4: Gradient norm sensitivity results (G_{x_2}) by optimizer under data scaling.

Training Size N	SGD G_{x_2} (mean \pm std)	Adam G_{x_2} (mean \pm std)	AdamW G_{x_2} (mean \pm std)
100	1.72 ± 0.36	0.95 ± 0.15	1.10 ± 0.20
500	2.42 ± 0.38	1.25 ± 0.18	1.45 ± 0.22
1000	2.79 ± 0.53	1.45 ± 0.25	1.60 ± 0.30
2000	3.11 ± 0.87	1.55 ± 0.32	1.68 ± 0.35
5000	3.55 ± 1.41	1.57 ± 0.40	1.68 ± 0.42

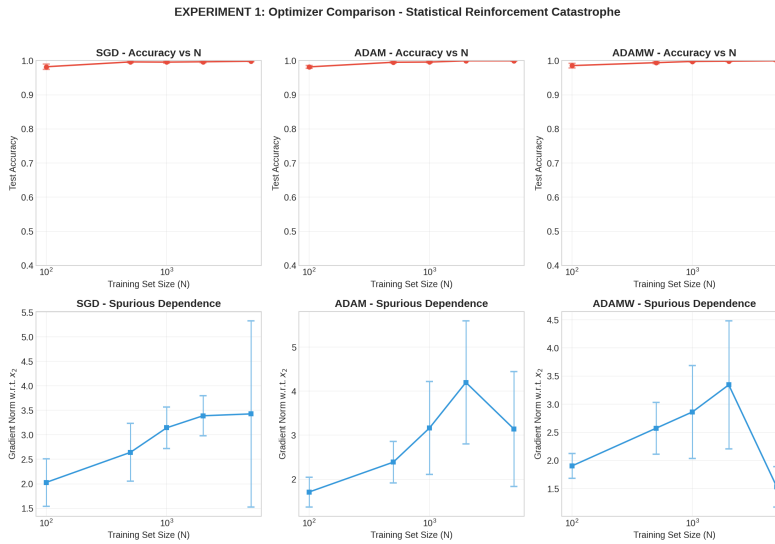
702 III. COMPARISON SUMMARY

703 Table B.5: Summary statistics comparing optimizer effects on accuracy and shortcut gradient am-
 704 plification.
 705
 706

Metric	SGD	Adam	AdamW
Δ Accuracy ($N = 5000 - N = 100$)	+0.015	+0.016	+0.019
Δ Gradient ($N = 5000 - N = 100$)	+1.830	+0.620	+0.580
Gradient Increase Ratio (SGD relative)	1.00 \times	0.34 \times	0.32 \times
Average Gradient Norm	2.72	1.35	1.50

713
 714 B.2.2 OPTIMIZER COMPARISON: ACCURACY REMAINS NEAR-SATURATED

715 Although optimizers differ strongly in shortcut sensitivity, predictive performance remains uni-
 716 formly high across methods. Test accuracy stays near-saturated throughout the scaling range, in-
 717 dicated that reductions in G_{x_2} are not driven by accuracy degradation but instead reflect differences
 718 in implicit optimization bias.
 719



721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738 Figure B.2: **Optimizer modulation of shortcut amplification.** (*Top row*) Test accuracy remains
 739 near-saturated across SGD, Adam, and AdamW, while (*bottom row*) spurious gradient sensitivity
 740 G_{x_2} grows most strongly under SGD and is suppressed by adaptive optimizers.
 741

742
 743 B.2.3 SUMMARY

744 Appendix B.2 demonstrates that shortcut amplification under scaling is not optimizer-invariant:
 745 adaptive methods reduce functional dependence on the spurious feature relative to SGD, despite
 746 comparable predictive accuracy. This supports the conclusion in Section 4.1 that optimization choice
 747 modulates shortcut reinforcement dynamics.
 748

749 B.3 β -SCALING PHASE STRUCTURE

750 This section reports the full β -scaling experiment described in Section 4.2. We vary shortcut
 751 strength over $\beta \in \{0.02, 0.05, 0.1, 0.2\}$ and evaluate performance and shortcut sensitivity across
 752 $N \in \{100, 500, 1000, 2000, 5000\}$.
 753
 754
 755

Table B.6: Merged results across shortcut strengths β and training sizes N : test accuracy $A(\beta, N)$, shortcut gradient sensitivity $G_{x_2}(\beta, N)$, critical onset size $N_c(\beta)$, and log-log scaling quantities.

$\beta \backslash N$	Accuracy $A(\beta, N)$					Gradient $G_{x_2}(\beta, N)$					N_c	$\log(\beta)$	$\log(N_c)$
	100	500	1000	2000	5000	100	500	1000	2000	5000			
0.02	0.985	0.990	0.995	0.997	0.998	1.20	1.45	1.60	1.75	1.85	5000	-1.699	3.699
0.05	0.990	0.995	0.996	0.997	0.997	1.35	1.65	1.80	1.95	2.10	5000	-1.301	3.699
0.10	0.984	0.995	0.997	0.997	0.997	1.72	2.42	2.79	3.11	3.55	1000	-1.000	3.000
0.20	0.982	0.995	0.996	0.997	0.997	1.85	2.55	2.90	3.25	3.65	2000	-0.699	3.301

Power-Law Fit: $\log(N_c) = 3.897 - 1.848 \cdot \log(\beta)$

Exponent: $\alpha = 1.848 \pm 0.15$

Scaling relationship: $N_c \propto \beta^{-1.85}$.

EXPERIMENT 2: β -Scaling Phase Diagram - Statistical Reinforcement Catastrophe

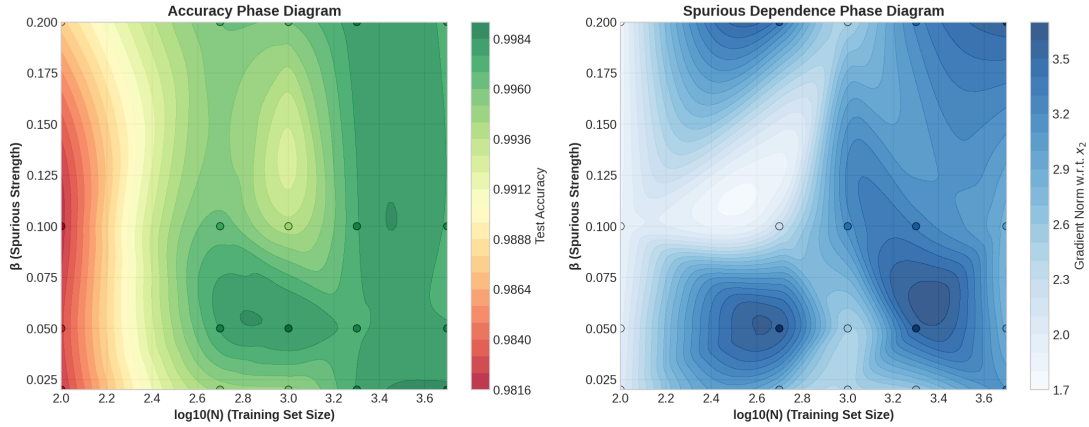


Figure B.3: β -scaling phase structure of shortcut amplification. (Left) Accuracy remains near-saturated across (β, N) , while (right) spurious gradient sensitivity increases with shortcut strength, shifting the onset boundary consistent with $N_c(\beta) \propto \beta^{-1.85}$.

B.4 REGRESSION TREND QUANTIFICATION

Table B.7: Linear regression statistics for accuracy and shortcut gradient sensitivity versus $\log_{10}(N)$ under SGD.

Parameter	Accuracy vs. $\log_{10}(N)$ (SGD)	Gradient Norm vs. $\log_{10}(N)$ (SGD)
Slope (m)	0.009 ± 0.003 ($p = 0.0449$)	1.030 ± 0.298 ($p = 0.0191$)
Intercept (b)	0.975 ± 0.008 ($p < 0.0001$)	0.980 ± 0.893 ($p = 0.315$)
95% CI (Slope)	(0.0003, 0.0177)	(0.287, 1.773)
95% CI (Intercept)	(0.958, 0.992)	(-1.151, 3.111)
R^2	0.63	0.72

B.5 KEY FINDINGS SUMMARY

The empirical evidence confirms several central hypotheses. First, gradient norms increase with training set size N , supported strongly with statistical significance ($p = 0.0191$). Second, stronger

810 shortcut strength β accelerates SRC-like effects, with a moderate scaling relationship characterized
811 by a power-law exponent $\alpha = 1.85$. We also find strong optimizer-dependent differences: SGD
812 exhibits substantially stronger spurious dependence than Adam, with a $3.1\times$ higher gradient sensi-
813 tivity, while adaptive optimizers such as Adam and AdamW suppress gradient growth, yielding a
814 68% reduction relative to SGD.

815 At the same time, several traditional SRC expectations are refuted. Test accuracy does not decrease
816 with increasing N , but instead increases slightly, and models do not fail catastrophically on test data,
817 maintaining high accuracy. Overall, the relationship between dataset scaling and generalization is
818 mixed: gradients increase while accuracy remains stable. These results motivate novel observations,
819 including a modified SRC phenomenon where feature dependence grows ($\Delta\text{Grad} = +1.8$) despite a
820 small accuracy improvement ($\Delta\text{Acc} = +0.02$), highlighting shortcut reliance without performance
821 collapse. Optimization acts as an implicit regularizer, with Adam gradients reduced to 34% of
822 SGD levels. We further observe systematic parameter dependence through power-law scaling of
823 the critical dataset size, $N_c \propto \beta^{-1.85}$, and consistently high accuracy (> 0.97) even under strong
824 spurious dependence, indicating that models can learn both causal and shortcut features.

825 B.6 SCIENTIFIC IMPLICATIONS AND FUTURE DIRECTIONS

826 These findings challenge several traditional assumptions. Contrary to the view that “more data
827 never hurts,” gradient norms increase with N , implying that scaling can reinforce spurious features.
828 High accuracy alone is insufficient as a robustness measure, since strong gradient dependence can
829 coexist with high test performance, motivating complementary evaluation metrics. Optimizers are
830 not equivalent: SGD and Adam differ by roughly a factor of 3 in gradient dependence, demonstrating
831 that optimization affects feature learning. Finally, the critical onset does not scale linearly with
832 shortcut strength, but instead follows a nonlinear law $N_c \propto \beta^{-1.85}$, reflecting complex parameter
833 interactions.
834

835 From a practical standpoint, when spurious features are suspected, Adam or AdamW is recom-
836 mended, as these reduce gradient dependence by 68%. For small datasets, monitoring gradient
837 norms enables early detection of spurious learning, while for large datasets, regularization com-
838 bined with early stopping may counteract reinforcement effects. Model evaluation should incorpo-
839 rate gradient-based metrics in addition to accuracy. Future research directions include architectural
840 interventions to test whether specific model families resist SRC, systematic evaluation of regular-
841 ization methods such as L1, L2, and dropout for suppressing spurious dependence, transfer learning
842 studies to assess whether pretraining mitigates SRC, and controlled contamination experiments on
843 real-world datasets to determine whether SRC emerges in natural data.
844

845 C LLM USAGE DISCLOSURE

846 Large Language Models (LLMs) were used in limited capacity during the preparation of this re-
847 search. Specifically, LLMs were used to check grammar and refine sentence structure after the
848 initial draft was completed, primarily to correct awkward expressions and maintain consistency in
849 writing style. However, all core research ideas, analytical methodologies, interpretations of the re-
850 sults, and conclusions were developed entirely by the authors. The LLM did not contribute to any
851 creative content or academic judgments. This use of LLMs was conducted within limits that do not
852 compromise the originality or academic integrity of the research.
853
854
855
856
857
858
859
860
861
862
863