
Clinically-Guided Counterfactuals (C^3): Physics and Pathology-Aware Augmentation and Evaluation for Robust Medical Imaging Models

Anonymous Author(s)

Abstract

1 Clinical deployment of imaging AI remains fragile: routine distribution
2 shifts—scanner vendor and reconstruction kernel, MRI protocol updates, dose
3 and slice profile changes, patient positioning and demographics, and device op-
4 tics—can degrade performance in ways that standard leaderboards and generic
5 augmentations fail to predict. We ask whether robustness and calibration can
6 be improved, without compromising clinical validity, by training and evaluating
7 models against *label-preserving, clinically grounded counterfactuals*. We intro-
8 duce *Clinically-Guided Counterfactuals (C^3)*, a framework that (i) unifies physics-
9 informed acquisition perturbations with tightly constrained, pathology-preserving
10 semantic edits; (ii) screens all counterfactuals through a conservative validity gate;
11 and (iii) reports *shift-stable utility*, a worst-case case-level score complementary to
12 AUROC, Dice, ECE, and Brier. Across chest X-ray (CheXpert→MIMIC-CXR),
13 MS brain MRI segmentation (multi-site→held-out site), and diabetic retinopathy
14 grading (EyePACS→Messidor-2), C^3 delivers consistent OOD gains (e.g.,
15 macro-AUROC +0.035 on CXR; lesion-wise Dice +0.044 on MRI; DR AUROC
16 +0.036), tighter calibration, reduced prediction volatility under realistic shifts, and
17 interpretable robustness diagnostics suitable for deployment checks.

18 1 Introduction

19 Deep models for medical imaging often underperform when faced with routine distribution shifts (4),
20 from reconstruction kernels and protocol changes to illumination and sensor variation (3). Standard
21 data augmentation and aggregate metrics (e.g., AUROC, Dice) offer limited visibility into per-case
22 worst-case behavior, which is crucial for safe deployment.

23 We propose to reframe robustness around each image’s *clinically admissible neighborhood*:
24 acquisition-realistic and pathology-preserving variants for which the diagnostic label or lesion
25 mask should remain invariant. Concretely, for input x with label y (classification) or mask m
26 (segmentation), we define transformations $\{t_\phi\}_{\phi \in \Phi}$ and the counterfactual neighborhood $\mathcal{N}(x) =$
27 $\{t_\phi(x) : \phi \in \Phi, \text{label}(t_\phi(x)) = y \text{ or mask}(t_\phi(x)) = m\}$. Within this neighborhood, models
28 should maintain consistent predictions and stable calibration.

29 Main Contributions:

- 30 • *Clinically-Guided Counterfactuals (C^3)*: a framework combining (i) physics-aware acqui-
31 sition perturbations, (ii) a pathology-preserving editor for small but clinically plausible
32 semantic/context edits, and (iii) a *validity gate* that filters counterfactuals via conservative
33 equivalence tests and sparse human anchors.
- 34 • *Training objective and evaluation*: a regularized objective that enforces prediction con-
35 sistency across $\mathcal{N}(x)$, and a deployment-facing *shift-stable utility* metric summarizing

36 worst-case calibrated performance per exam, reported alongside AUROC, Dice, ECE, and
37 Brier.

- 38 • *Evidence across three modalities/tasks:* CXR classification, MS MRI lesion segmentation,
39 and DR grading show consistent OOD improvements, better calibration, and reduced volatility
40 under scanner/protocol/camera shifts, with ablations demonstrating the complementary
41 roles of physics transforms, the editor, and the validity gate.

42 2 Methodology

43 2.1 Counterfactual neighborhoods and consistency objective

44 For classification with model f_θ , let $p_\theta(\cdot | x)$ denote predictive probabilities. For segmentation, let
45 $\hat{m}_\theta(x)$ be the predicted mask. We train with

$$\mathcal{L} = \mathcal{L}_{\text{sup}}(f_\theta(x), y) + \lambda \mathbb{E}_{x' \sim \mathcal{N}(x)} [D(f_\theta(x), f_\theta(x'))], \quad (1)$$

46 where $D = D_{\text{KL}}(p_\theta(\cdot | x) \| p_\theta(\cdot | x'))$ for classification and

$$D = 1 - \frac{2\langle \hat{m}_\theta(x), \hat{m}_\theta(x') \rangle + \epsilon}{\|\hat{m}_\theta(x)\|_1 + \|\hat{m}_\theta(x')\|_1 + \epsilon} \quad (2)$$

47 for segmentation.

48 For evaluation, we compute per-case *shift-stable utility*

$$U(x) = \min_{x' \in \mathcal{N}(x)} s(f_\theta(x'), y) \quad \text{or} \quad U(x) = \min_{x' \in \mathcal{N}(x)} \text{Dice}(\hat{m}_\theta(x'), m), \quad (3)$$

49 which lower-bounds performance under admissible clinical variation.

50 2.2 Physics-aware acquisition perturbations

51 We implement modality-specific operators that approximate routine acquisition changes while pre-
52 serving labels/masks:

- 53 • *CT/CXR:* Poisson thinning before FBP for low dose; kernel sharpness variation (soft \leftrightarrow sharp);
54 blur consistent with thicker slices/partial volume; mild beam hardening / scatter shifts of
55 HU distributions.
- 56 • *MRI:* B_0/B_1 bias fields; sequence-appropriate Rician/non-central- χ noise; contrast modula-
57 tion via Bloch-informed lookups (TE/TR/flip); slice-profile broadening.
- 58 • *Fundus:* Illumination geometry, vignetting, and sensor-pattern perturbations matched to
59 camera response, preserving microaneurysms and exudates.

60 2.3 Pathology-preserving editor

61 A diffusion backbone is fine-tuned with weak supervision (reports/labels) to produce low-amplitude,
62 clinically plausible edits (e.g., rib-shadow contrast, subtle effusion haze, projection geometry; illu-
63 mination and vessel-contrast tweaks in fundus). Two soft constraints keep edits near the clinical
64 manifold: (i) a lesion-mask consistency penalty discouraging changes to annotated pathology, and
65 (ii) a text-image agreement term over a curated pathology vocabulary to stabilize the global clinical
66 description.

67 2.4 Validity gate

68 Before inclusion in training/evaluation, counterfactuals must satisfy conservative tests:

- 69 • *Classification:* $\|p_\theta(\cdot | x) - p_\theta(\cdot | x')\|_\infty \leq \delta$ and matching predicted class argmax.
- 70 • *Segmentation:* teacher/consensus masks must meet $\text{IoU}(m^*(x), m^*(x')) \geq \tau$.

71 Thresholds (δ, τ) are set using radiologist-audited anchors.

72 **2.5 Experimental settings**

73 We hold architectures/optimizers/schedules fixed to isolate C^3 :

74 1. *CXR classification*: DenseNet-121 on CheXpert (1), OOD on MIMIC-CXR (2); 5 findings
75 (Atelectasis, Cardiomegaly, Consolidation, Edema, Pleural Effusion).

76 2. *MS MRI segmentation*: 3D U-Net on a multi-site cohort; OOD evaluation on a held-out site;
77 lesion- and volume-wise Dice.

78 3. *Fundus DR grading*: EfficientNet-B3 on EyePACS; OOD on Messidor-2; AUROC for
79 referable DR and ECE.

80 Baselines: standard augmentations, RandAugment/AugMix variants, and TTA.

81 **3 Results and Discussion**

82 **Chest X-ray (CheXpert→MIMIC-CXR).** C^3 improves macro-AUROC from 0.864 (STDAUG)
83 and 0.872 (AUGMIX) to 0.907, with consistent per-pathology gains (e.g., Edema 0.903 → 0.935,
84 Consolidation 0.842 → 0.887, Effusion 0.918 → 0.944). Calibration improves (ECE 5.7% → 2.9%;
85 Brier –13.4% relative). Neighborhood agreement rises 0.73 → 0.86, and shift-stable utility increases
86 0.782 → 0.846, indicating stronger worst-case performance under clinically realistic perturbations.
87 Ablations show physics transforms and the editor contribute additively, while removing the validity
88 gate superficially boosts ID AUROC yet harms worst-case utility and calibration, exposing hidden
89 label drift.

90 **MS MRI segmentation (multi-site→held-out).** Mean lesion-wise Dice increases from 0.598
91 (STDAUG) and 0.624 (AUGMIX) to 0.668 with C^3 ; small-lesion recall ($< 10 \text{ mm}^3$) improves 0.521 →
92 0.603 with 17.8% fewer false negatives. Volume-wise Dice rises 0.706 → 0.744, and volume
93 calibration tightens (slope 0.81 → 0.93). Under simulated slice-thickness increase (1.0 mm → 3.0 mm),
94 Hausdorff-95 decreases from 8.9 mm to 7.3 mm, whereas non- C^3 baselines exceed 9.5 mm. Physics-
95 only (Dice 0.653) and editor-only (0.639) trail the full model (0.668).

96 **Fundus DR (EyePACS→Messidor-2).** AUROC improves 0.842 → 0.878 over STDAUG (and
97 0.855 → 0.878 vs. TTA), while ECE halves (4.2% → 2.0%). Prediction flips under admissible
98 illumination/camera shifts drop from 14.7% to 6.8%. The 10th percentile of per-patient shift-stable
99 utility increases 0.763 → 0.820, and OOD AUROC CIs shrink by 24%, reflecting reduced variance
100 across nuisance factors.

101 **Interpretability and deployment diagnostics.** C^3 attenuates failure modes aligned with radiologist
102 intuition: e.g., CXR reliance on rib-shadow/illumination artifacts; MRI small-lesion under-
103 segmentation exacerbated by bias fields or thicker slices; fundus sensitivity to vascular glare. Per-exam
104 robustness profiles provide actionable diagnostics for model cards and site readiness checks.

105 **Takeaways.** (1) Physics-aware operators anchor robustness to acquisition realities; (2) small,
106 clinically plausible edits broaden coverage of nuisance semantics without drifting labels; (3) a
107 strict validity gate is essential to avoid training on mislabeled counterfactuals; (4) worst-case, case-
108 level reporting (*shift-stable utility*) complements aggregate metrics and reveals deployment-relevant
109 volatility.

110 **4 Conclusion**

111 C^3 reframes robustness around clinically grounded counterfactual neighborhoods, aligning both
112 training and evaluation with how images vary across scanners, protocols, and devices. The approach
113 consistently improves OOD accuracy, calibration, and worst-case stability across CXR, MRI, and
114 fundus tasks, while surfacing interpretable diagnostics for deployment. Because C^3 composes with
115 standard datasets and models, it can be adopted without architectural changes.

116 **5 Future Directions**

- 117 • *Prospective and site-onboarding studies*: Evaluate C³ in prospective multi-site rollouts and
118 during scanner/protocol onboarding to quantify reductions in drift-related incidents.
- 119 • *Expanded modalities and tasks*: Extend physics operators and editors to ultrasound, mam-
120 mography, and pathology WSIs; explore detection and multi-label settings.
- 121 • *Human-in-the-loop validity*: Incorporate lightweight radiologist spot-audits and active
122 sampling to calibrate (δ, τ) and prioritize hard neighborhoods.
- 123 • *Fairness & subpopulation robustness*: Construct neighborhoods that target demo-
124 graphic/device subgroups, pairing shift-stable utility with stratified fairness metrics.
- 125 • *Uncertainty and calibration*: Combine C³ with deep ensembles/temperature scaling under
126 neighborhood perturbations; track case-level reliability diagrams.
- 127 • *Efficiency*: Distill neighborhood training via curriculum or importance sampling; cache
128 reusable physics transforms to limit compute.
- 129 • *Operational tooling*: Package per-exam robustness profiles for model cards and site-
130 readiness checklists; add *pre-deployment* synthetic probes matching local scanner settings.

131 **Potential Negative Societal Impact**

132 C³ could increase reliance on synthetic data; if misused without validity checks, this risks over-
133 confidence and deployment to populations not represented in the original datasets. The framework
134 should not replace external validation on real multi-site cohorts. Conservative validity gates and
135 radiologist spot audits are recommended, and the use of synthetic counterfactuals should be disclosed
136 in documentation and model cards to avoid accidental data leakage or privacy concer

137 **References**

- 138 [1] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute,
139 Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A.
140 Mong, Safwan Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz,
141 Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. CheXpert: A Large Chest Radiograph
142 Dataset with Uncertainty Labels and Expert Comparison. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI 2019)*, pages 590–597. AAAI Press, 2019. URL
143 <https://arxiv.org/abs/1901.07031>.
- 145 [2] Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-
146 ying Deng, Roger G. Mark, Seth J. Berkowitz, and Steven Horng. MIMIC-CXR-JPG, a large
147 publicly available database of labeled chest radiographs. *Scientific Data*, 6:317, 2019. Nature
148 Publishing Group. URL <https://www.nature.com/articles/s41597-019-0322-0>.
- 149 [3] Varun Gulshan, Lily Peng, Marc Coram, Martin C. Stumpe, Derek Wu, Arunachalam
150 Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros,
151 Ramasamy Kim, Rajiv Raman, Phil Nelson, Jessica L. Mega, and Dale R. Webster. Develop-
152 ment and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in
153 Retinal Fundus Photographs. *JAMA*, 316(22):2402–2410. American Medical Association, 2016.
154 URL <https://jamanetwork.com/journals/jama/fullarticle/2588763>.
- 155 [4] Etienne Decencière, Xiwei Zhang, Guénolé Cazuguel, Bruno Lay, Béatrice Cochener, Caroline
156 Trone, Philippe Gain, Raphaël Ordóñez, Pascale Massin, Ali Erginay, Béatrice Charton, and
157 Jean-Claude Klein. Feedback on a publicly distributed image database: the Messidor database.
158 *Image Analysis & Stereology*, 33(3):231–234. Image Analysis & Stereology Society, 2014. URL
159 <http://www.ias-iss.org/ojs/IAS/article/view/1155>.