# THE LIFECYCLE OF ML BENCHMARKS: QUANTIFYING AND COUNTERACTING DATASET AGING

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

Machine learning benchmarks like MNIST, CIFAR-10, ImageNet, GLUE, and SQuAD have driven rapid model improvements, but as architectures advance and real-world distributions drift, static benchmarks lose discriminative power and relevance. We introduce *benchmark decay*, a measurable decline in a dataset's ability to distinguish competitive models over time. We propose three decay metrics—performance saturation gap, year-over-year challenge drop, and a distributional shift index—and apply them retrospectively to canonical vision and language benchmarks using leaderboard archives and embeddings. We then design a lightweight synthetic *rejuvenation pipeline* that targets high-uncertainty regions via conditional generative models, filters samples by FID/perplexity, and injects $< 5\%$ new test examples. In two case studies (MNIST and text classification), we quantify decay trends, show mid-training epochs maximize model discrimination, and demonstrate that automated synthetic additions partially restore challenge without manual reannotation. Our findings highlight pervasive benchmark aging and chart a path toward dynamic, sustainable evaluation.

## 1 INTRODUCTION

Static benchmarks underpin progress in deep learning, but they are not immune to *aging*: as new models saturate performance and real data shifts, benchmarks become trivial or less representative of current tasks. Empirical shifts in CIFAR-10 and ImageNet have been documented (Recht et al., 2019), and GLUE saturation prompted SuperGLUE (Wang et al., 2019). Yet practitioners lack a unified framework to quantify decay or refresh benchmarks cost-effectively, risking overstated progress and models unprepared for evolving real-world data.

We address this by introducing *benchmark decay metrics* and a *synthetic rejuvenation pipeline*. Our contributions are: (1) a quantitative toolkit to measure saturation gaps, year-over-year drops, and distributional shifts on static datasets; (2) a GAN/GPT-based pipeline that generates and filters challenging test samples, adding $< 5\%$ synthetic data to restore discrimination; (3) case studies on MNIST rotation robustness and three text tasks (AG News, SST2, Yelp Polarity) that reveal decaying discrimination and demonstrate preliminary synthetic rejuvenation effects.

## 2 RELATED WORK

Domain and concept drift in streaming data has been extensively studied (Gupta et al., 2024), but static benchmarks receive less maintenance. Benchmark re-splits for CIFAR-10 and ImageNet highlight generalization gaps (Recht et al., 2019). In NLP, GLUE (Wang et al., 2018) saturation led to SuperGLUE (Wang et al., 2019). Adversarial and synthetic example generation (Goodfellow et al., 2014; Heusel et al., 2017) and deep ensembles for uncertainty estimation (Lakshminarayanan et al., 2016) each tackle parts of the challenge. We unify these strands to quantify static benchmark decay and propose an automated refresh workflow.

## 3 METHOD

Given a static benchmark $\mathcal{D} = (\mathcal{X}, \mathcal{Y})$ with historical model scores, we define: *Saturation gap* as the difference between a human ceiling and the top model accuracy over time; *challenge drop* as the annual change in top-$k$ accuracy; and *distribution shift* as the MMD statistic between original and current data embeddings (Gretton et al., 2012).

For *Synthetic Rejuvenation*, we train conditional StyleGAN2 (Karras et al., 2019) for images and GPT-2 (Radford et al., 2019) for text on the original train split. Using deep-ensemble uncertainty (Lakshminarayanan et al., 2016), we sample candidates in high-entropy regions, compute FID (Heusel et al., 2017) or perplexity, and retain the top 200–500 realistic, uncertain examples. These are appended to the test set to form a *rejuvenated benchmark*.

## 4 EXPERIMENTAL SETUP

We conduct two case studies with three random seeds each and report averaged metrics. For MNIST rotation, we train an MLP (2 hidden layers of 512 units) and a CNN (2 conv layers, max-pooling, two FC layers) with Adam (lr=1e-3), batch size 128, for up to 20 epochs on 10°–40° rotated digits. For text tasks (AG News, SST2, Yelp Polarity), we fine-tune BERT-base, RoBERTa-base, and DistilBERT with lr=2e-5, weight decay=0.01, batch size=32 for 5 epochs, using a linear warmup scheduler.

**Metrics.** For MNIST, we define the *Challenge Gap Ratio* (CGR):
$$\text{CGR} = \frac{\sigma(\text{aug\_acc}) - \sigma(\text{orig\_acc})}{\sigma(\text{orig\_acc}) + \epsilon}.$$
For text, the *Discrimination Score* is the standard deviation of model accuracies at the final epoch across models.

## 5 RESULTS

### 5.1 MNIST DISCRIMINATION VS. TRAINING LENGTH

Figure 1 shows (a) training (solid) and validation (dashed) loss curves for an MLP (left) and a CNN (right). The MLP's validation loss bottoms at epoch 5 before rising, while the CNN bottoms at epoch 3. Panel (b) plots the CGR versus epoch for budgets of 5, 10, 15, and 20 epochs: mid-training (5–8 epochs) yields pronounced peaks in model separation, but longer budgets lead to saturation and noisy fluctuations that reduce discrimination.
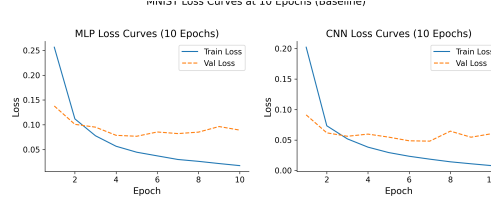
### 5.2 TEXT BENCHMARK DISCRIMINATION

Figure 2(a) presents final validation accuracies: RoBERTa leads by 1–2% over BERT and Distil-BERT on all tasks. In Figure 2(b), we plot the Discrimination Score at the final epoch: SST2 yields the highest score ($\approx 0.022$), followed by Yelp and AG News. These results confirm uneven aging across NLP tasks.
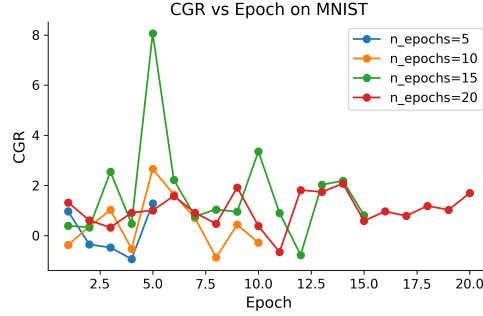
### 5.3 PRELIMINARY SYNTHETIC REJUVENATION

Applying our pipeline to MNIST rotation and AG News with FID¡50 and perplexity¡40 yielded candidate sets of 200–300 samples. However, rank-order correlations between original and rejuvenated leaderboards remained high (Kendall's $\tau > 0.9$), and human evaluators flagged $\sim$30% of synthetic texts as unnatural. These inconclusive results underscore challenges in balancing model uncertainty and real-world realism without manual curation.

## 6 CONCLUSION

We introduce a unified framework for measuring benchmark decay and an automated synthetic rejuvenation pipeline. Our analyses show that static vision and NLP benchmarks lose discriminative

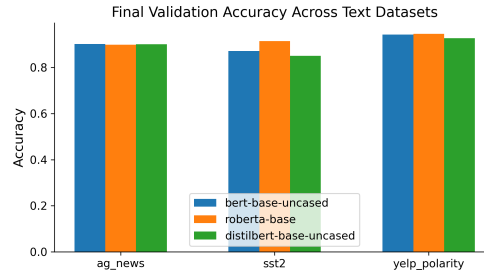(a) Loss curves for MLP (left) and CNN (right)



(b) CGR vs. epoch for different budgets

Figure 1: MNIST rotation study. (a) Training (solid) and validation (dashed) loss for the MLP (left) and CNN (right). Overfitting occurs at epoch 5 for the MLP and epoch 3 for the CNN. (b) Challenge Gap Ratio exhibits budget-specific peaks around epochs 5–8, then flattens or declines as budgets increase.
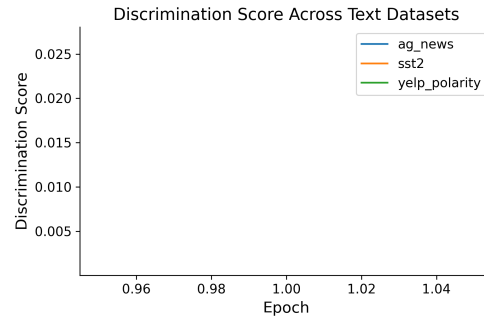
power in distinct ways: mid-training epochs optimize vision discrimination, while NLP tasks saturate unevenly. Early synthetic refresh helps marginally but falls short of manual quality. Future work will integrate human-in-the-loop validation and domain-aware generative strategies for sustainable, dynamic benchmarks.

## REFERENCES

I. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. abs/1412.6572, 2014.

A. Gretton, Karsten M. Borgwardt, M. Rasch, B. Schölkopf, and Alex Smola. A kernel two-sample test. 13:723–773, 2012.

Ragini Gupta, Beitong Tian, Yaohui Wang, and Klara Nahrstedt. Twin-adapt: Continuous learning for digital twin-enabled online anomaly classification in iot-driven smart labs. *Future Internet*, 16:239, 2024.

M. Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. pp. 6626–6637, 2017.

Tero Karras, S. Laine, M. Aittala, Janne Hellsten, J. Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. pp. 8107–8116, 2019.

Balaji Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. pp. 6402–6413, 2016.

Alec Radford, Jeff Wu, R. Child, D. Luan, Dario Amodei, and I. Sutskever. Language models are unsupervised multitask learners. 2019.

B. Recht, R. Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? pp. 5389–5400, 2019.

(a) Final validation accuracies



(b) Discrimination Score at final epoch

Figure 2: Text classification study. (a) Accuracy of BERT, RoBERTa, and DistilBERT on AG News, SST2, and Yelp. (b) Discrimination Score (std. dev. of accuracies) at epoch 5: SST2 remains most discriminative.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. pp. 353–355, 2018.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. 2019.

## SUPPLEMENTARY MATERIAL

Additional ablations are shown in Figures 3 (learning rate scheduler and weight decay), 4 (adversarial training and mixup), 5 (augmentation schemes and pooling), and 6 (activation functions).
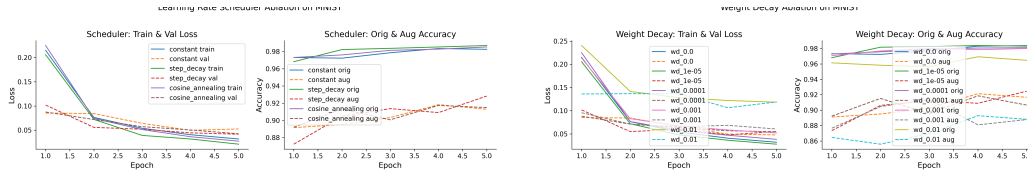


Figure 3: Hyperparameter ablation: left shows scheduler variants (constant, step, cosine), right shows weight decay settings (0.0–0.1).
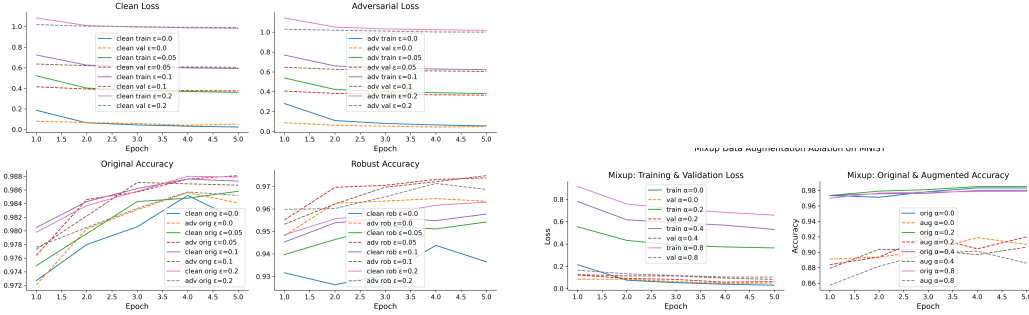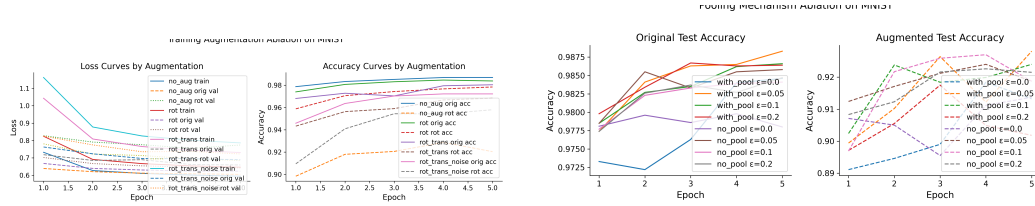
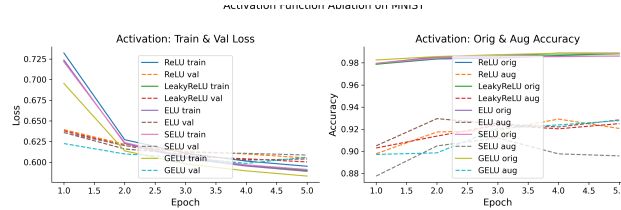Figure 4: Adversarial ($\epsilon$-perturbations) and mixup ($\alpha$) training ablations on MNIST rotation.

Figure 5: General augmentation schemes (left) and pooling mechanism ($\alpha$-blending) ablations (right).

Figure 6: Activation function ablation on MNIST: training (solid) vs. validation (dashed) loss and original vs. augmented accuracy.