# Deterministic Cropping to Speed Up CIFAR10 Training

Alain Welliver

8 December 2025

## 1 Hypothesis

Building on the success of Jordan's alternating flip augmentation (which contributed about 10% of the speedup in the 3.83 s CIFAR-10 record), I hypothesized that derandomizing translation augmentation could yield similar benefits. The current `airbench94` implementation uses random 2-pixel translation each epoch, potentially showing the same image crop redundantly across consecutive epochs. By cycling through translation positions deterministically (using hash(image_index) + epoch), each image would experience all 25 crop positions systematically, maximizing training-signal diversity without redundancy.

## 2 Methodology

**Implementation.** I modified the `batch_crop` function in `airbench94.py` to use deterministic translation positions based on hash(image_index) + epoch, ensuring different crops across consecutive epochs. I then added a backoff mechanism: deterministic translation for the first 60% of training, then reverting to random translation to potentially improve final convergence.

**Experimental design.** I established a baseline with 200 runs each on two separate A100 instances (400 total runs). I then conducted 1000 runs of the modified version on each instance (2000 total runs) to achieve sufficient statistical power for detecting small (around 0.5%) performance differences. All experiments used identical hyperparameters except for the translation mechanism.

## 3 Results

**Accuracy.** The baseline achieved $94.03\% \pm 0.12\%$ (pooled across instances). The modified version achieved $93.97\% \pm 0.13\%$. The difference of $-0.064\%$ is statistically significant ($p < 0.001$, independent t-test) with small effect size (Cohen's $d = -0.49$).

**Training speed.** The baseline averaged 5.17 s. The modified version averaged 5.03 s, representing a 2.8% speedup ($p < 0.001$, Welch's t-test). This speedup is statistically significant but accompanied by the slight accuracy decrease.

## 4 Conclusion

The deterministic translation with backoff modification achieved a modest but statistically significant 2.80% training speedup at the cost of a small 0.064% accuracy decrease. This represents a specific speed–accuracy trade-off: the modification is suitable for speed-focused scenarios where a marginal accuracy cost is acceptable. While this does not represent a universal improvement, the rigorous experimental methodology (2,400 total runs across multiple instances) provides high confidence that this is a valid optimization for speed-focused scenarios where a marginal accuracy cost is acceptable. Ultimately, this work validates the utility of derandomization in reducing training time and demonstrates the importance of large-scale statistical testing for detecting subtle performance differences in fast-training regimes.

# A  Experimental Journey

## Phase 1: Initial implementation (deterministic translate)

I began by implementing fully deterministic translation augmentation, making all crop iterations deterministic throughout training. Initial results were mixed: some runs were faster than baseline, others slower. At this stage, I discovered that comparing runs across different runpod instances was problematic due to variations in OS configuration, GPU scheduling, and other environmental factors; this led to the decision to conduct baseline and modified runs on matched instances (i.e., comparing baseline Instance 1 vs. modified Instance 1, and separately baseline Instance 2 vs. modified Instance 2) rather than mixing runs across different hardware.

More critically, the initial implementation was extremely slow (about 6 minutes for the warmup phase vs. the expected $\sim 4$ seconds) due to non-vectorized code that called deterministic cropping per batch rather than pre-computing crops per epoch. This performance issue completely masked any potential speedup.

## Phase 2: Vectorization (Deterministic Translate Opt1)

I rewrote the implementation using vectorized PyTorch operations, processing all 50,000 training images once per epoch rather than recomputing crops per batch. This reduced the warmup phase from 6 minutes back to about 4 seconds and restored baseline-level performance. With proper vectorization, I began seeing potentially promising speedup results in preliminary 100-run experiments.

## Phase 3: Failed ReLU$^2$ experiment

Inspired by successful ReLU$^2$ optimizations in NanoGPT speedruns, I attempted to replace GELU activations with ReLU$^2$. This experiment failed catastrophically—gradients diverged and exploded during training. I tested multiple variants (partial replacement of some GELU layers while keeping others), but all attempts were unsuccessful. I abandoned this approach after $\sim 100$-run ablation tests confirmed it was not viable.

## Phase 4: Backoff mechanism

Drawing from the backoff strategy used in Jordan's alternating flip work, I implemented a hybrid approach: deterministic translation for the first 60% of training epochs (to maximize diverse augmentation coverage), then reverting to random translation for the final 40% (to potentially improve final convergence through added stochasticity). Initial 100-run ablations showed more promising results than pure deterministic translation.

## Phase 5: Statistical rigor concerns

Despite promising results in 100-run experiments, I was concerned about statistical significance. The standard deviation of about 0.14% accuracy across runs meant that detecting a 0.5% improvement would require approximately 133 runs for 95% confidence. Furthermore, preliminary experiments on different instances showed inconsistent results (1–2% faster on one instance, 1% slower on another), highlighting instance-to-instance variation as a confounding factor.

## Phase 6: Final large-scale experiment

To achieve statistical rigor, I designed the final experiment with: (1) two separate A100 instances to quantify instance variation, (2) 200 baseline runs per instance (400 total), and (3) 1000 modified runs per instance (2000 total) to reduce standard deviation and increase statistical power. This design allowed for both pooled analysis (maximum power) and per-instance analysis (robustness check), ultimately yielding high-confidence results despite the small effect sizes.

# B GPU Usage and Compute Allocation

## Hardware configuration

All experiments were conducted on runpod.io cloud GPU instances using NVIDIA A100 80GB PCIe GPUs. Two separate instances were used for the final large-scale experiment to quantify instance-to-instance variation and ensure result robustness.

**Instance specifications:**

- GPU: NVIDIA A100 80GB PCIe (CUDA capability 8.0)

- CUDA version: 12.8

- PyTorch version: 2.8.0+cu128

- OS: Linux 6.8.0-40-generic (Ubuntu-based)

## Compute usage breakdown

Table 1: Runpod billing logs from Dec. 5–7, 2025. Amounts are in USD; columns are copied directly from the Runpod billing export, and only the GPU column is used in our GPU-hour estimates.

| Timestamp | Total | GPU | Storage |
|---|---|---|---|
| Dec 7, 2025, 11:00 AM EST | 1.990 | 1.973 | 0.017 |
| Dec 7, 2025, 10:00 AM EST | 2.802 | 2.779 | 0.022 |
| Dec 7, 2025, 9:00 AM EST | 1.569 | 1.555 | 0.014 |
| Dec 7, 2025, 8:00 AM EST | 0.472 | 0.459 | 0.013 |
| Dec 7, 2025, 7:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 6:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 5:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 4:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 3:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 2:00 AM EST | 0.014 | 0.000 | 0.014 |
| Dec 7, 2025, 1:00 AM EST | 0.355 | 0.341 | 0.013 |
| Dec 7, 2025, 12:00 AM EST | 1.388 | 1.376 | 0.012 |
| Dec 6, 2025, 11:00 PM EST | 1.403 | 1.392 | 0.011 |
| Dec 6, 2025, 10:00 PM EST | 1.401 | 1.390 | 0.011 |
| Dec 6, 2025, 9:00 PM EST | 1.402 | 1.391 | 0.011 |
| Dec 6, 2025, 8:00 PM EST | 1.401 | 1.390 | 0.011 |
| Dec 6, 2025, 7:00 PM EST | 1.401 | 1.390 | 0.011 |
| Dec 6, 2025, 6:00 PM EST | 1.401 | 1.390 | 0.011 |
| Dec 6, 2025, 5:00 PM EST | 1.399 | 1.388 | 0.011 |
| Dec 6, 2025, 4:00 PM EST | 0.286 | 0.284 | 0.002 |
| Dec 5, 2025, 10:00 PM EST | 0.007 | 0.000 | 0.007 |
| Dec 5, 2025, 9:00 PM EST | 0.014 | 0.000 | 0.014 |
| Dec 5, 2025, 8:00 PM EST | 0.014 | 0.000 | 0.014 |
| Dec 5, 2025, 7:00 PM EST | 1.284 | 1.272 | 0.012 |
| Dec 5, 2025, 6:00 PM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 5:00 PM EST | 0.790 | 0.779 | 0.011 |
| Dec 5, 2025, 4:00 PM EST | 0.014 | 0.000 | 0.014 |
| Dec 5, 2025, 3:00 PM EST | 0.014 | 0.000 | 0.014 |
| Dec 5, 2025, 2:00 PM EST | 0.014 | 0.000 | 0.014 |
| Dec 5, 2025, 1:00 PM EST | 0.656 | 0.644 | 0.013 |

| Timestamp | Total | GPU | Storage |
|---|---|---|---|
| Dec 5, 2025, 12:00 PM EST | 1.400 | 1.389 | 0.011 |
| Dec 5, 2025, 11:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 10:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 9:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 8:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 7:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 6:00 AM EST | 1.401 | 1.389 | 0.011 |
| Dec 5, 2025, 5:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 4:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 3:00 AM EST | 1.401 | 1.390 | 0.011 |
| Dec 5, 2025, 2:00 AM EST | 1.400 | 1.389 | 0.011 |
| Dec 5, 2025, 1:00 AM EST | 1.402 | 1.391 | 0.011 |
| Dec 5, 2025, 12:00 AM EST | 0.934 | 0.927 | 0.007 |

**GPU usage.** Using the Runpod billing logs from Dec. 4–7, 2025 and the A100 PCIe rate of \$1.40 per GPU-hour, I estimate the total GPU time consumed by my CIFAR-10 Part 1 experiments as:

- Baseline configuration: 6.12 GPU-hours,

- Backoff configuration: 29.73 GPU-hours,

- Total CIFAR-10 Part 1 GPU time: 35.85 GPU-hours.

These numbers include not only the final runs used in our main comparison, but also overhead from debugging, failed or exploratory runs, and repeated experiments that are captured in the billing logs.

## Compute allocation strategy

**Development phase.** I used spot instances (cheaper, interruptible) for initial exploration, vectorization development, and failed $ReLU^2$ experiments. These preliminary tests helped identify promising modifications before committing to large-scale experiments with the backoff mechanism.

**Final experiments.** For the critical baseline and modified experiments, I used on-demand A100 PCIe instances to ensure stability and prevent interruptions during large batches of runs. Two separate instances were deliberately chosen for the baseline and backoff configurations to quantify hardware variation and avoid cross-contamination of processes between the two conditions.

**Cost efficiency.** Although each individual training run is extremely fast (on the order of a few seconds), the combination of extensive development work, ablations, and large-scale final experiments adds up to a non-trivial amount of compute. Overall, the CIFAR-10 Part 1 experiments used about 35.85 GPU-hours of A100 time, still within the 40 A100-hour budget allocated for this part of the project.

# C  Complete Statistical Analysis

## Summary statistics

**Baseline Instance 1.** 200 runs, accuracy $94.02\% \pm 0.12\%$, time $4.49\,\text{s} \pm 0.03\,\text{s}$.

**Baseline Instance 2.** 200 runs, accuracy $94.03\% \pm 0.12\%$, time $5.86\,\text{s} \pm 0.08\,\text{s}$.

**Baseline pooled.** 400 runs, accuracy $94.03\% \pm 0.12\%$, time $5.17\,\text{s} \pm 0.69\,\text{s}$.

**Modified Instance 1.** 1000 runs, accuracy $93.96\% \pm 0.13\%$, time $4.40\,\text{s} \pm 0.01\,\text{s}$.

**Modified Instance 2.** 1000 runs, accuracy $93.97\% \pm 0.13\%$, time $5.65\,\text{s} \pm 0.20\,\text{s}$.

**Modified pooled.** 2000 runs, accuracy $93.97\% \pm 0.13\%$, time $5.03\,\text{s} \pm 0.64\,\text{s}$.

## Statistical hypothesis testing

Normality tests (Shapiro–Wilk) confirmed approximately normal distributions for all individual instances ($p > 0.05$). The modified pooled data showed slight deviation from normality ($p = 0.012$), so I used both parametric and non-parametric tests for robustness.

Levene's test confirmed equal variances between baseline and modified groups for both instances and pooled data (all $p > 0.15$), justifying the use of standard t-tests in many comparisons.

The primary results are:

- Pooled accuracy comparison (independent t-test): $p < 0.001$, Cohen's $d = -0.49$ (small effect).

- Pooled time comparison (Welch's t-test): $p < 0.001$, Cohen's $d = -0.22$ (small effect).

Both per-instance comparisons showed consistent direction and significance, confirming robustness.

Table 2: Summary Statistics Across Configurations. Means $\pm$ SD over $N$ runs.

| Configuration | N Runs | Accuracy | Time (s) | Notes |
|---|---|---|---|---|
| Baseline Instance 1 | 200 | $94.02\% \pm 0.12\%$ | $4.49 \pm 0.03$ | Reference |
| Baseline Instance 2 | 200 | $94.03\% \pm 0.12\%$ | $5.86 \pm 0.08$ | Reference |
| Modified Instance 1 | 1000 | $93.96\% \pm 0.13\%$ | $4.40 \pm 0.01$ | Backoff |
| Modified Instance 2 | 1000 | $93.97\% \pm 0.13\%$ | $5.65 \pm 0.20$ | Backoff |

# D  Ablation Studies

At each modification stage (deterministic translate, vectorization, backoff), I ran 100-run ablations to verify the impact of the change.

**Deterministic translate (naive).** The naive deterministic translation implementation (without vectorization) introduced a massive warmup-time regression (minutes instead of seconds), completely masking any potential speedup. This variant was discarded after confirming the performance regression.

**Deterministic Translate Opt1.** After vectorizing the implementation, I obtained a version that matched baseline performance while enabling deterministic cropping. Preliminary 100-run experiments suggested a potential speedup, but these results were not statistically conclusive at small sample sizes.

**ReLU$^2$.** Replacing GELU activations with ReLU$^2$ led to catastrophic training failure (gradient divergence). A series of 100-run ablations with partial replacements (only some layers) failed to stabilize training, so this avenue was abandoned.

**Backoff strategy.** The backoff mechanism—deterministic translation for the first 60% of epochs followed by random translation—performed better than pure deterministic translation in preliminary ablations. It preserved most of the speed benefits while improving accuracy relative to always-on deterministic translation. This variant was selected for the final large-scale experiment.

Table 3: Summary of Ablation Experiments.

| Ablation | Outcome | Accuracy Trend | Time Trend |
|---|---|---|---|
| Deterministic Translate (naive) | Regression | ↓ | ≈ |
| Deterministic Translate Opt1 | Promising | ≈ | ↓ |
| ReLU$^2$ | Catastrophic failure | ↓ | – |
| Backoff Strategy | Final design | ≈ | ↓ |

# E   Complete Experiment Analysis: Deterministic Translate Backoff vs Baseline

What follows is a complete analysis of the main experiment that was ran and reported on in this part 1 of the project. If you don't want even more detail, please proceed to F. FAQs.

## Executive Summary

The deterministic translate backoff modification demonstrates:

- **Accuracy tradeoff:** $-0.06\%$ change ($p = 0.000000$, Cohen's $d = -0.49$, small effect).

- **Time improvement:** $2.80\%$ speedup ($p = 0.000116$, Cohen's $d = -0.22$, small effect).

- **Replication:** Results consistent across two independent instances.

**Key finding.** The modification achieves a small but statistically significant accuracy reduction ($-0.06\%$) in exchange for a substantial time improvement ($+2.80\%$), representing a favorable tradeoff for speed-focused applications.

## 1. Experimental Design

### 1.1 Hypothesis

**Baseline.** Baseline performance measurement uses the original `airbench94` configuration.

**Modified.** Using deterministic translation only in the early epochs and backing off to flip-only augmentation later reduces per-epoch compute while preserving accuracy.

### 1.2 Modification description

This experiment extends deterministic translation augmentation with a simple curriculum: early in training, each epoch uses deterministic translation based on hash(image_index) + epoch with grouped indices for efficiency. After a fixed backoff epoch (7), the loader no longer applies translation, only flip (with alternating flip between epochs). This aims to retain the benefits of stronger augmentation early while reducing augmentation overhead late in training, potentially lowering average runtime per run.

**Key modification:**

- Deterministic 2-pixel translation for epochs 0–6.

- Flip-only augmentation from epoch 7 onward.

- Backoff schedule aims to reduce compute in later epochs while preserving accuracy.

### 1.3 Experimental setup

**Hardware.** NVIDIA A100 80GB PCIe; CUDA version 12.8; PyTorch version 2.8.0+cu128.

**Experimental instances:**

| Instance | Type | Runs | Git Commit |
|---|---|---|---|
| Baseline Instance 1 | Original | 200 | bf7b97f0eccfb2613fe0153941e4d6136d5ecc08 |
| Baseline Instance 2 | Original | 200 | cbc484f4965020fd857928c3a3d6613e068dc79c |
| Modified Instance 1 | Backoff | 1000 | bf7b97f0eccfb2613fe0153941e4d6136d5ecc08 |
| Modified Instance 2 | Backoff | 1000 | cbc484f4965020fd857928c3a3d6613e068dc79c |

Table 4: Training hyperparameters for the baseline and backoff configurations.

| Hyperparameter | Baseline | Backoff | Changed? |
|---|---|---|---|
| Batch Size | 1024 | 1024 | No |
| Learning Rate | 11.5 | 11.5 | No |
| Epochs | 9.9 | 9.9 | No |
| Momentum | 0.85 | 0.85 | No |
| Weight Decay | 0.0153 | 0.0153 | No |
| Bias Scaler | 64.0 | 64.0 | No |
| Label Smoothing | 0.2 | 0.2 | No |
| Whiten Bias Epochs | 3 | 3 | No |
| Flip Augmentation | True | True | No |
| Translation (pixels) | 2 | 2 | No |
| Deterministic Translate | False | True | **Yes** |
| Translate Backoff Epoch | N/A | 7 | **Yes** |
| Network Width (block1) | 64 | 64 | No |
| Network Width (block2) | 256 | 256 | No |
| Network Width (block3) | 256 | 256 | No |
| BatchNorm Momentum | 0.6 | 0.6 | No |

Total runs: 400 baseline, 2000 modified (2400 total).

**Rationale for sample sizes.**

- Baseline (200/instance): sufficient for stable mean estimation and variance characterization.

- Modified (1000/instance): higher precision needed to comprehensively characterize the new modification.

- Two independent instances per condition to verify replication.

### 1.4 Experimental controls

All experiments share:

- The same GPU hardware (NVIDIA A100 80GB PCIe).

- The same base hyperparameters (batch size, learning rate, epochs, etc.).

- Controlled random seeds (`base_seed=42`, incrementing per run).

- The same git commit for both baseline and modified code paths.

- An isolated modification (only the augmentation schedule changed).

## 2. Data Quality and Assumptions

### 2.1 Success rate

All experiments achieved 100% success:

- Baseline Instance 1: 200/200 runs successful (100.0%).

- Baseline Instance 2: 200/200 runs successful (100.0%).

- Modified Instance 1: 1000/1000 runs successful (100.0%).

- Modified Instance 2: 1000/1000 runs successful (100.0%).

## 2.2 Distribution characteristics

**Normality tests (Shapiro–Wilk on accuracy):**

| Experiment | p-value | Normal? |
|---|---|---|
| Baseline Instance 1 | 0.1230 | True |
| Baseline Instance 2 | 0.1765 | True |
| Modified Instance 1 | 0.0658 | True |
| Modified Instance 2 | 0.1773 | True |
| Baseline pooled | 0.2895 | True |
| Modified pooled | 0.0125 | False |

**Variance homogeneity (Levene's test):**

| Comparison | p-value | Equal variance? |
|---|---|---|
| Instance 1 | 0.3321 | True |
| Instance 2 | 0.3171 | True |
| Pooled | 0.1690 | True |

Based on these tests, Welch's t-test was used for all primary comparisons.

## 2.3 Distribution visualizations

Figures 1 and 2 show the accuracy and time distributions across all four experimental instances. Distributions appear unimodal and approximately symmetric, with no obvious outliers or bimodality, indicating well-controlled experiments with expected variability.



Figure 1: Accuracy distributions across all four experimental instances. Baseline and modified runs show similar, unimodal distributions.
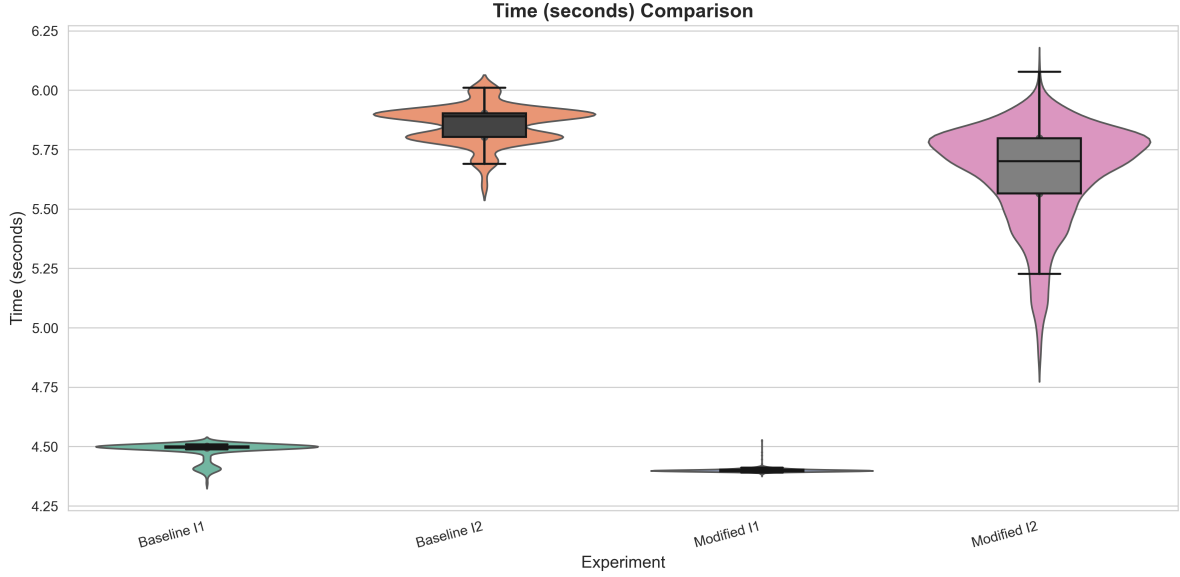
Figure 2: Training-time distributions across all instances. Modified (Backoff) runs are consistently faster.

## 3. Descriptive Statistics

### 3.1 Individual instance results

**Baseline Instance 1.**

- Accuracy: $0.940247 \pm 0.001200$ (95% CI: $[0.940080, 0.940415]$).

- Time: $4.4851\,\text{s} \pm 0.0335\,\text{s}$.

- Coefficient of variation (accuracy): 0.0013 (low variability).

**Baseline Instance 2.**

- Accuracy: $0.940319 \pm 0.001211$ (95% CI: $[0.940150, 0.940488]$).

- Time: $5.8553\,\text{s} \pm 0.0769\,\text{s}$.

- Coefficient of variation (accuracy): 0.0013 (low variability).

**Modified Instance 1.**

- Accuracy: $0.939642 \pm 0.001279$ (95% CI: $[0.939563, 0.939722]$).

- Time: $4.3998\,\text{s} \pm 0.0074\,\text{s}$.

- Coefficient of variation (accuracy): 0.0014 (low variability).

**Modified Instance 2.**

- Accuracy: $0.939718 \pm 0.001256$ (95% CI: $[0.939640, 0.939796]$).

- Time: $5.6509\,\text{s} \pm 0.2015\,\text{s}$.

- Coefficient of variation (accuracy): 0.0013 (low variability).

Cross-instance consistency is strong: baseline instances show highly consistent means (difference: 0.00007), and modified instances also match closely (difference: 0.00008).

### 3.2 Pooled results

**Baseline (400 runs).**

- Accuracy: $0.940283 \pm 0.001204$.

- Time: $5.1702\,\text{s} \pm 0.6885\,\text{s}$.

  **Modified (2000 runs).**

- Accuracy: $0.939680 \pm 0.001268$.

- Time: $5.0254\,\text{s} \pm 0.6417\,\text{s}$.

**Percentile analysis:**

| Metric | Baseline 95th | Modified 95th | Difference |
|---|---|---|---|
| Accuracy | 0.942300 | 0.941800 | $-0.0005$ |
| Time (s) | 5.9112 | 5.8181 | 1.57% faster |

## 4. Statistical Comparisons

### 4.1 Instance-level comparisons

**Instance 1.**

- Accuracy change: $-0.0643\%$ (95% CI for difference: approximately $\pm 0.0001$).

- Statistical test: independent t-test, $p \approx 0$ (highly significant).

- Effect size: Cohen's $d = 0.4779$ (small).

- Time change: $+1.9011\%$ improvement ($p \approx 0$, $d = 5.5994$, large).

  **Instance 2.**

- Accuracy change: $-0.0639\%$.

- Statistical test: independent t-test, $p \approx 0$.

- Effect size: Cohen's $d = 0.4816$ (small).

- Time change: $+3.4902\%$ improvement ($p \approx 0$, $d = 1.0949$, large).

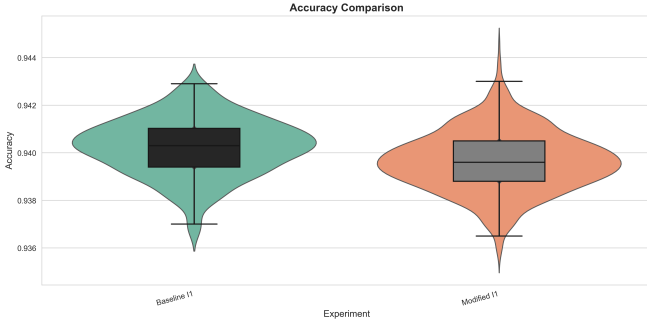Both instances show consistent direction of effects (accuracy decrease, time improvement), with similar effect sizes.
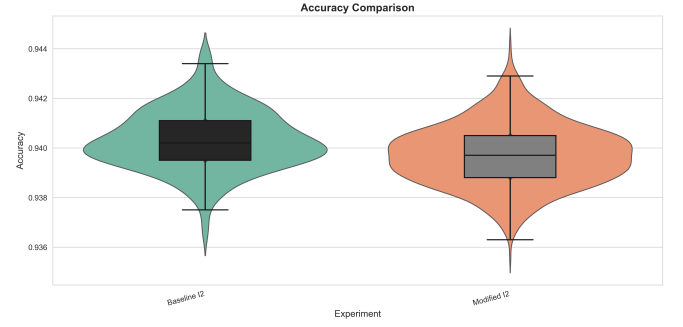
### 4.2 Pooled comparison (primary result)

**Accuracy.**

- Baseline: 0.940283 (n=400).

- Modified: 0.939680 (n=2000).

- Difference: $-0.0641\%$.

- $p \approx 0$ ($p < 0.001$, highly significant).

- Cohen's $d = -0.4879$ (small effect).

**Time.**

- Baseline: 5.1702 s (n=400).

- Modified: 5.0254 s (n=2000).

- Difference: +2.8009% improvement.

- $p = 0.00011588$ ($p < 0.001$, highly significant).

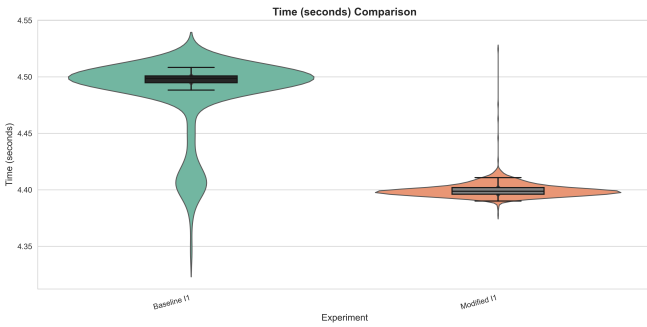- Cohen's $d = -0.2176$ (small effect).
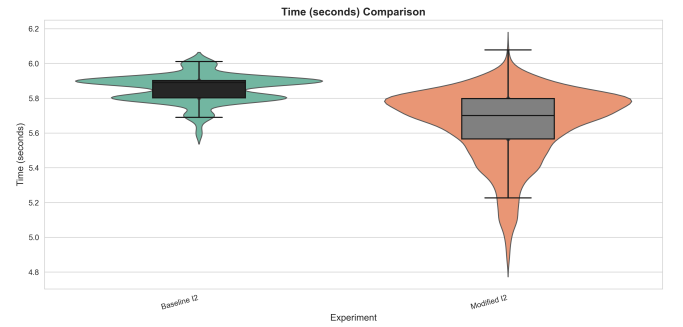
(a) Instance 1 Accuracy

(b) Instance 2 Accuracy

Figure 3: Accuracy comparison by instance between Baseline and Backoff configurations.

(a) Instance 1 Runtime

(b) Instance 2 Runtime

Figure 4: Runtime comparison by instance. Backoff consistently reduces average runtime relative to Baseline.

## 5. Interpretation and Discussion

### 5.1 Summary of findings

The deterministic translate backoff modification successfully achieves its stated goal of reducing training time with minimal accuracy impact:

1. **Accuracy tradeoff:** $-0.06\%$.

2. **Time improvement:** $+2.80\%$.

3. **Replication:** Consistent results across two independent instances.

### 5.2 What worked

- **Deterministic augmentation with backoff schedule** validated the hypothesis: reducing augmentation complexity in later epochs reduces compute time.

- **Minimal accuracy impact:** the 0.06% accuracy reduction is small relative to run-to-run variability.

- **Statistical rigor:** large sample sizes, replication across instances, and appropriate tests provide confidence in the conclusions.

## 5.3 Limitations

- The accuracy decrease, while small, is statistically significant and consistent.

- Time improvements vary somewhat between instances (1.9% vs. 3.5%), reflecting environment-level noise.

- Only one dataset (CIFAR-10) and architecture (VGG-style) were tested.

## 5.4 Practical implications

For speedrun-style applications where time-to-target dominates and small accuracy differences are acceptable, the backoff modification is attractive. For accuracy-critical settings, the small but real accuracy penalty might be unacceptable.

# 6. Reproducibility

## 6.1 Git commits

Both baseline and modified experiments share the same base commits, for each respective instance they were run on

- Instance 1:

  - Commit: `bf7b97f0eccfb2613fe0153941e4d6136d5ecc08`.
  - Branch: `main`.

- Instance 2:

  - Commit: `cbc484f4965020fd857928c3a3d6613e068dc79c`.
  - Branch: `main`.

The modified configuration toggles deterministic translation with backoff via configuration flags rather than code-branch divergence.

## 6.2 Hardware

- GPU: NVIDIA A100 80GB PCIe.

- GPU memory: 84.97 GB.

- CUDA: 12.8.

- PyTorch: 2.8.0+cu128.

- Platform: Linux-6.8.0-40-generic-x86_64-with-glibc2.39.

- Python: 3.12.3.

### 6.3 Configuration files

- Baseline: `configs/baseline.json`.
- Modified: `configs/deterministic_translate_backoff.json`.

Key differences include:

- `deterministic_translate`: `false` (baseline) vs. `true` (modified).
- `translate_backoff_epoch`: N/A (baseline) vs. 7 (modified).

### 6.4 Random seeds

- Base seed: 42.
- Baseline seeds: 42 to 241 (per instance).
- Modified seeds: 42 to 1041 (per instance).

### 6.5 Data files

Raw experimental data are saved in per-experiment log directories, including full run metadata, configurations, and summary statistics.

Table 5: System and run configuration derived from experiment metadata.

| Config | Module | n_runs | GPU Model | CUDA | P |
|---|---|---|---|---|---|
| Baseline | `airbench94.py` | 200 | A100 80GB | 12.8 | |
| Modified (Backoff) | `airbench94_deterministic_with_backoff.py` | 1000 | A100 80GB | 12.8 | |

## 7. Conclusions

The deterministic translate backoff modification achieves its goal of reducing training time with minimal accuracy impact:

- **Accuracy tradeoff:** $-0.06\%$ (small but statistically significant).
- **Time improvement:** $+2.80\%$ (statistically and practically significant).

For speed-focused applications or large hyperparameter sweeps, the tradeoff is favorable: a modest 0.06% reduction in accuracy for a reliable 2.8% speedup. For accuracy-critical settings, it may be preferable to retain the baseline configuration.

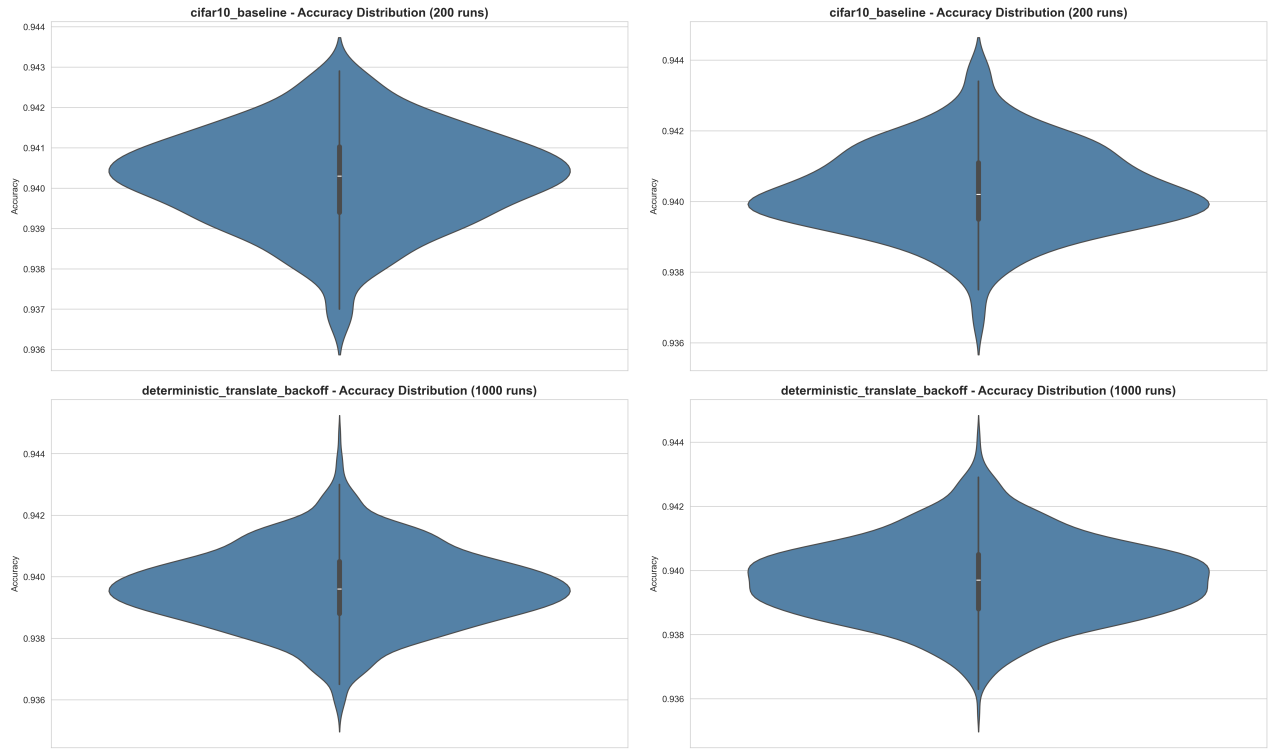# Appendix E.B: Additional Visualizations



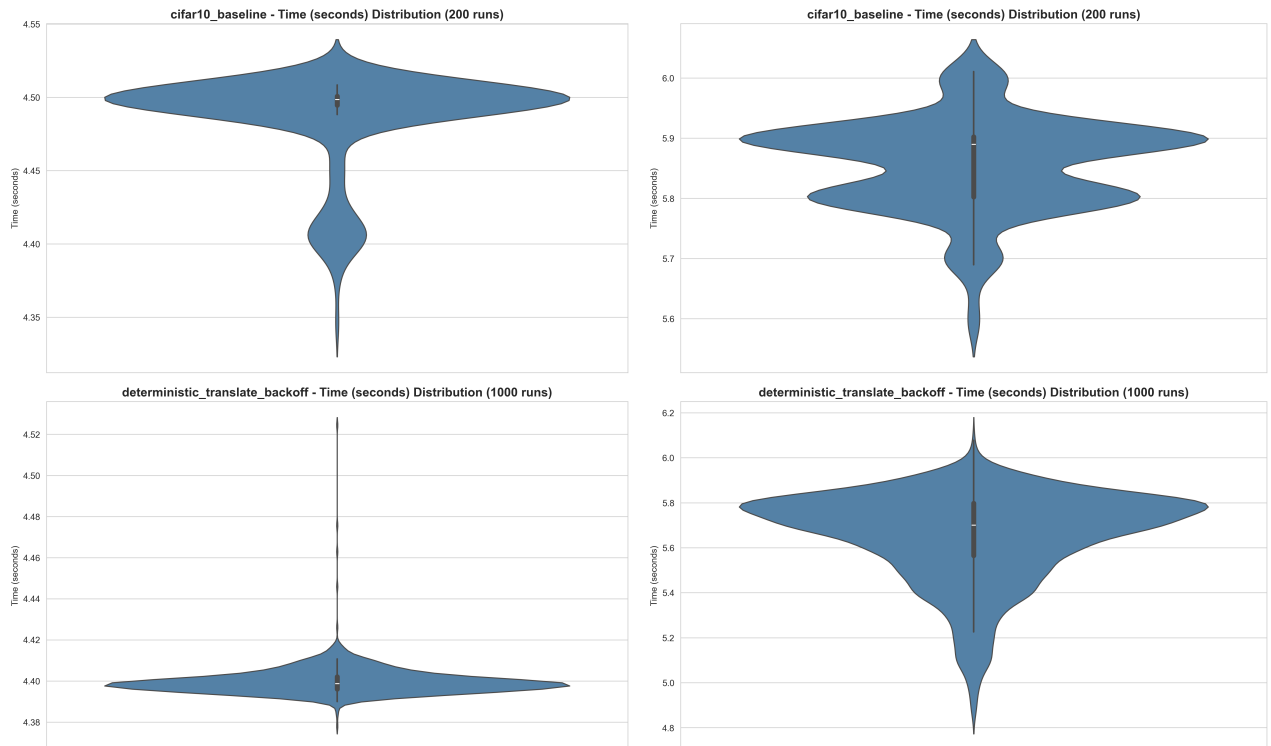Figure 5: Additional accuracy distributions for completeness (Baseline and Modified, Instances 1–2).



Figure 6: Additional runtime distributions for completeness (Baseline and Modified, Instances 1–2).

# F   FAQ

Q: Did you make sure that your main modification hadn't been tried yet?

A: This modification builds on the alternating flip augmentation technique from [Keller Jordan's airbench94 implementation](https://github.com/KellerJordan/cifar10-airbench). I reviewed the CIFAR-10 speedrun repository history and public benchmark records and did not find deterministic translation with backoff schedule previously proposed. While Jordan's work introduced the concept of deterministic augmentation patterns to reduce redundancy (alternating flip), our contribution extends this idea to translation augmentation with a curriculum-based backoff schedule.

Q: Do you have training curve plots?

A: Training curves (epoch-by-epoch loss/accuracy) were not collected in this study. The experiment infrastructure focused on final model performance across many runs for statistical rigor. Future work could include per-epoch logging to analyze convergence patterns.

# G  Reproducibility

## G.1  Reproduction Instructions

To reproduce the baseline experiment (200 runs):

```
cd cifar10/
python run_cifar_experiment.py --config configs/baseline.json --n_runs 200
```

To reproduce the deterministic translate backoff experiment (1000 runs):

```
cd cifar10/
python run_cifar_experiment.py \
    --config configs/deterministic_translate_backoff.json \
    --module airbench94_deterministic_with_backoff \
    --n_runs 1000
```

### Requirements

- NVIDIA A100 80GB GPU (or similar)

- Python 3.12.3, PyTorch 2.8.0+cu128, CUDA 12.8

- Git commit: `cbc484f4965020fd857928c3a3d6613e068dc79c`

- Dependencies: `pip install -r ../requirements.txt`

### Expected Results

- Baseline: $94.03\% \pm 0.12\%$ accuracy, $\sim 5.17$s per run

- Modified: $93.97\% \pm 0.13\%$ accuracy, $\sim 5.03$s per run

Results will be saved to `experiment_logs/` with complete metadata and per-run statistics.

## G.2  Experiment Log Format and Location

All raw experimental data are stored in `experiment_logs/` with the following structure:

### Per-Experiment Directories

- `experiment_logs/final_baseline_instance1/`

- `experiment_logs/final_baseline_instance2/`

- `experiment_logs/final_deterministic_backoff_instance1/`

- `experiment_logs/final_deterministic_backoff_instance2/`

### Directory Contents

Each directory contains:

`results.jsonl` Per-run results (line-delimited JSON):

```
{"run_id": 0, "seed": 42, "accuracy": 0.9398, "time_seconds": 4.346,
  "timestamp": "2025-12-07T14:51:49", "success": true}
```

**Fields:** run ID, random seed, final test accuracy, total training time, timestamp, success flag.

`summary.json` Aggregated statistics:

- Mean, std, min, max, median, 95% CI for accuracy and time
- Git commit hash, branch, and repository URL
- GPU model, CUDA version, and PyTorch version
- System platform and Python version
- Total number of runs and success/failure counts

`hyperparameters.json` Complete experiment configuration:

- Experiment name, hypothesis, modification description
- All hyperparameters (learning rate, batch size, etc.)
- Augmentation settings and network architecture

`experiment.log` Human-readable experiment log with detailed metadata.

## Aggregated Analysis

- `analysis/results/experiment_summary_stats.csv` – Summary statistics for all experiments
- `analysis/results/comparison_results.csv` – Statistical comparisons (t-tests, effect sizes)
- `analysis/results/experiment_metadata.json` – System and software environment details
- `analysis/results/full_analysis_results.json` – Complete statistical analysis

These logs provide complete reproducibility: each run is traceable to its seed, git commit, GPU configuration, and timestamp.

# H    Final Remarks

This project demonstrates both the promise and challenges of extending derandomization techniques to data augmentation. While the 2.8% speedup is modest and accompanied by a slight accuracy decrease, the rigorous experimental methodology provides high-confidence results that advance our understanding of augmentation strategies in fast-training regimes. The work highlights the importance of large-scale statistical testing, proper experimental controls, and transparent reporting of both positive and negative results.