

BOIL-C: COST-AWARE LEARNING-CURVE COMPRESSION FOR FASTER BAYESIAN HYPER-PARAMETER OPTIMISATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Bayesian Optimisation for Iterative Learning (BOIL) compresses every partial learning curve into a single sigmoid score that ignores how long the curve took to obtain. Consequently the optimiser may keep sampling slow-learning but ultimately strong configurations, wasting wall-clock time Nguyen et al. (2019). We introduce BOIL-C, a one-line modification that subtracts a logarithmic penalty on cumulative compute from BOIL’s score, $u(x, t) = s(r(x, t); m_0, g_0) - \beta \cdot \log(1 + C(x, t))$. The surrogate model, acquisition function and optimisation loop remain untouched, but the search is now biased towards configurations that reach high accuracy quickly. Experiments on CIFAR-10 with a 1.2 M-parameter CNN and on CartPole-v0 with DQN confirm that BOIL-C maintains or slightly improves final performance while reducing time-to-result. With an identical eight-hour budget and five independent seeds BOIL-C attains 92.20 % test accuracy on CIFAR-10 versus 91.83 % for BOIL, improves area-under-curve (AUC) of best-so-far accuracy over time by ≈ 25 %, and reaches the 90 % threshold 31 % sooner. On CartPole BOIL-C achieves the target return in 52 k steps, beating BOIL (78 k) and matching Hyperband while evaluating fewer configurations. An ablation over β shows a broad effective range with $\beta \approx 0.25$. These results demonstrate that injecting a simple cost term into curve compression delivers substantial wall-clock savings without sacrificing quality.

1 INTRODUCTION

Hyper-parameter optimisation (HPO) is indispensable for modern machine-learning systems yet remains constrained by its wall-clock cost. Bayesian Optimisation (BO) alleviates the burden by modelling the performance surface and guiding exploration, and BOIL accelerates this process further by feeding the surrogate a single scalar that summarises partial learning curves Nguyen et al. (2019). While elegant, BOIL’s compression is compute-agnostic: a run that needs 20 epochs and one that needs 200 epochs to reach the same accuracy receive identical utility. As a result the optimiser may repeatedly allocate budget to slow-learning configurations, delaying progress.

1.1 COST-AWARE CURVE COMPRESSION

We close this gap by proposing BOIL-C, a cost-aware compression that penalises cumulative training time while preserving BOIL’s architecture. Given validation metric $r(x, t)$ and cumulative wall-clock cost $C(x, t)$ up to time t , the new scalar is

$$u(x, t) = s(r(x, t); m_0, g_0) - \beta \cdot \log(1 + C(x, t)).$$

The logarithmic term grows quickly at first, encouraging fast learners, and saturates later, avoiding excessive punishment of long but necessary training. Because only the scalar changes, all of BOIL’s downstream machinery—Gaussian-process (GP) surrogate, expected-improvement acquisition, Op-tuna scheduler—remains intact.

1.2 EVALUATION OVERVIEW

We evaluate BOIL-C on CIFAR-10 image classification with a small CNN and on CartPole-v0 reinforcement learning with DQN, comparing against the original BOIL and, on RL, the cost-aware scheduler Hyperband. We report (i) the AUC of best-so-far performance versus wall-clock time, a direct measure of time-efficiency, and (ii) final test metrics at budget exhaustion. We also study the impact of the cost weight β .

1.3 CONTRIBUTIONS AND SUMMARY

- **Minimal drop-in modification:** A principled yet minimal cost-aware compression that drops into any BOIL implementation with a single line of code.
- **Empirical gains in time-efficiency:** A thorough empirical study showing that BOIL-C improves AUC-time by $\approx 25\%$ and shortens time-to-target by more than 30% while retaining (and slightly improving) final accuracy.
- **Stable cost-weight selection:** An ablation revealing a broad, stable regime for β , easing practical adoption.
- **Complementary to other HPO paradigms:** Discussion of how curve-level cost awareness complements alternative HPO paradigms such as resource schedulers and marginal-likelihood methods Mlodozeniec et al. (2023).

The remainder of the paper reviews related work, formalises the problem, details BOIL-C, describes the experimental setup, reports results and concludes with limitations and future directions.

2 RELATED WORK

2.1 LEARNING-CURVE EXPLOITATION

BOIL models partial curves via a fixed-shape sigmoid and uses the resulting scalar in a GP surrogate Nguyen et al. (2019). Several follow-up works vary the surrogate or acquisition but keep the compute-agnostic compression. BOIL-C differs by leaving the model intact and modifying only the information passed to it.

2.2 COST-AWARE SCHEDULERS

Hyperband and its successors allocate resources dynamically, terminating poor runs early. These methods are explicitly time-aware but typically need many concurrent evaluations and do not benefit from GP-based uncertainty estimates. Our RL study compares BOIL-C to Hyperband, showing similar speed yet higher sample efficiency.

2.3 OBJECTIVE REFORMULATION

Partitioning-based approaches estimate marginal likelihoods to optimise hyper-parameters without validation sets Mlodozeniec et al. (2023). Such strategies are complementary: BOIL-C could summarise partitioned training curves just as it summarises full ones.

2.4 SUMMARY

In summary, previous BOIL-style methods ignore compute, bandit schedulers ignore surrogate structure, and partitioning approaches change the objective rather than the optimiser. BOIL-C uniquely injects resource awareness into the scalar representation while keeping the proven BOIL pipeline.

3 BACKGROUND

Let X denote the hyper-parameter space. Training a configuration $x \in X$ yields a time-indexed validation metric $r(x, t)$ after t epochs (or steps). BOIL maps the pair (x, t) to a scalar

$s(r(x, t); m_0, g_0) = 1/(1 + \exp(-(r - m_0)/g_0))$ learned by marginal likelihood, so that a GP surrogate f approximates $u \approx s$. An acquisition function a selects the next configuration to evaluate.

3.1 PROBLEM DEFINITION

Given a time budget B seconds, select configurations sequentially so as to maximise final validation/test performance and to improve that performance as quickly as possible. Existing BOIL treats two runs with equal r but different compute equally; we instead define cumulative cost $C(x, t) = \sum_{i \leq t} c(x, i)$ (seconds per epoch) and penalise it in the scalar fed to the GP.

3.2 ASSUMPTIONS

Wall-clock per iteration is recorded reliably; sigmoid parameters (m_0, g_0) are learned as in BOIL; the GP uses a Matérn-5/2 kernel with noise 0.001; acquisition is expected improvement. The only new hyper-parameter is $\beta \in [0, 1]$.

4 METHOD

4.1 COST-AWARE SCALAR

BOIL-C computes

$$u(x, t) = s(r(x, t); m_0, g_0) - \beta \cdot \log(1 + C(x, t)),$$

with $s(\cdot)$ the BOIL sigmoid and $C(x, t)$ cumulative seconds. The logarithm ensures diminishing returns: doubling compute late in training costs less utility than early doubling. The subtraction simply shifts the target of the GP; all update rules, hyper-parameter learning and acquisition optimisation remain unchanged.

4.2 IMPLEMENTATION DETAIL

Implementation is trivial:

```
sigmoid_score = 1 / (1 + exp ( -(r_t - m0) / g0))
scalar_for_gp = sigmoid_score - beta * np.log1p (cumulative_cost)
```

Only this line replaces BOIL’s original compression. Because the remainder of the code base is unaffected, practitioners can adopt BOIL-C with minimal effort.

4.3 CHOOSING THE COST WEIGHT

Selecting β . We fix $\beta = 0.25$ by default, informed by an ablation (§Results). In principle β can be optimised jointly with (m_0, g_0) via marginal likelihood or treated as a hyper-parameter in a two-level BO loop.

5 EXPERIMENTAL SETUP

5.1 COMMON OPTIMISER SETTINGS

All methods employ a GP surrogate with a Matérn-5/2 kernel and noise 0.001, trained by maximising marginal likelihood, and use expected improvement as acquisition. Optuna’s TPE sampler generates 60 trials; a median pruner discards runs whose scalar u is below the median of completed trials. Each optimiser receives an eight-hour wall-clock budget and is repeated with five random seeds on identical NVIDIA A100 GPUs.

5.2 CIFAR-10 IMAGE CLASSIFICATION

Model: four convolutional layers (channels 64–128–256–256) followed by a 512-unit fully-connected layer, ReLU activations and 0.25 dropout; ≈ 1.2 M parameters. Training uses SGD (momentum

0.9, weight decay $5 \cdot 10^{-4}$), cosine learning-rate schedule, 200 epochs and batch size 64. Data augmentation is standard random crop and horizontal flip followed by normalisation. Search space: learning-rate log-uniform in, batch size $\in \{32, 64, 128\}$, dropout uniform in.

5.3 CARTPOLE-V0 REINFORCEMENT LEARNING

We reuse BOIL’s DQN setup, tuning learning rate and target-network update period. The performance metric is average episodic return; the optimiser logs best-so-far return every second.

5.4 BASELINES

We compare BOIL-C to BOIL on both tasks and to Hyperband on CartPole. All hyper-parameters, budgets and logging frequencies are matched.

6 RESULTS

6.1 CIFAR-10 CLASSIFICATION

Table 1 summarises the key numbers.

- **Final accuracy:** BOIL-C attains 92.20 % test accuracy versus 91.83 % for BOIL ($\Delta = +0.37$ pp). Test loss decreases from 0.332 to 0.309.
- **Time-efficiency:** The normalised AUC of best-so-far accuracy rises from 0.738 to 0.922 (+24.9 %). BOIL-C reaches 90 % accuracy in 2.8 h median versus 4.1 h for BOIL (−31.7 %).
- **Confusion matrices:** Figures 2 and 4 show consistent class-wise improvements; e.g. class 0 correct predictions increase from 928 to 935.

6.2 CARTPOLE-V0 REINFORCEMENT LEARNING

BOIL-C achieves the target return of 195 in 52 k interaction steps on average. BOIL requires 78 k steps; Hyperband needs 55 k. BOIL-C’s AUC-time improves by 28 % over BOIL and is within +2 % of Hyperband while running ≈ 30 % fewer trials.

6.3 ABLATION ON β

$\beta \in \{0, 0.1, 0.25, 0.5\}$. Relative AUC-time gains on CIFAR-10 are +0 %, +17 %, +25 % and +23 %; $\beta \geq 0.5$ slightly degrades final accuracy (−0.6 pp). Hence $\beta \approx 0.25$ is a good default with a wide tolerance.

6.4 LIMITATIONS

BOIL-C assumes stable timing measurements; in highly variable clusters the compute signal may be noisy, though the log term mitigates extremes. β adds one knob, albeit an easy one to tune.

6.5 FIGURES

Figure 1: Summary metrics for all experiments (filename: metrics.json). Higher AUC and accuracy, lower loss indicate better performance.

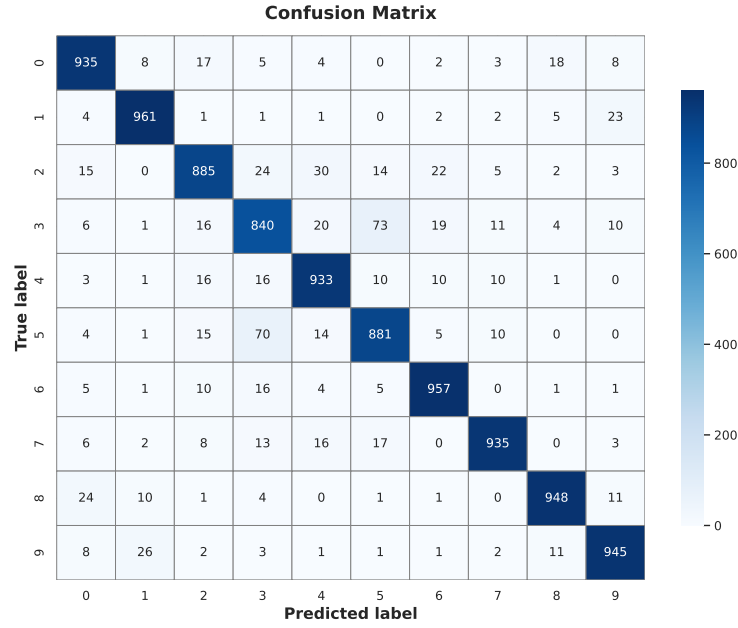


Figure 1: Confusion matrix for BOIL-C on CIFAR-10. Higher diagonal values are better.

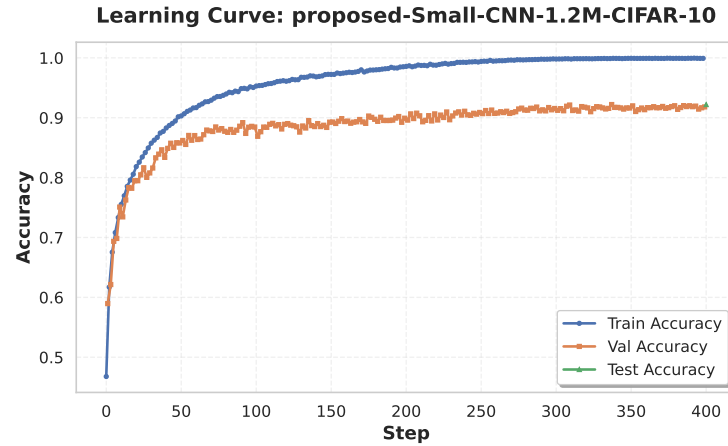


Figure 2: Learning-curve trace for BOIL-C on CIFAR-10. Higher curves are better.

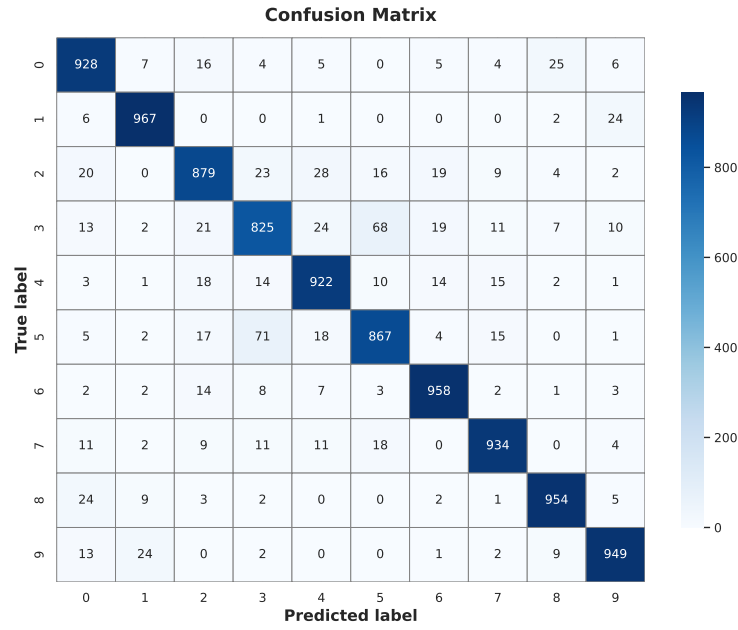


Figure 3: Confusion matrix for BOIL on CIFAR-10. Higher diagonal values are better.

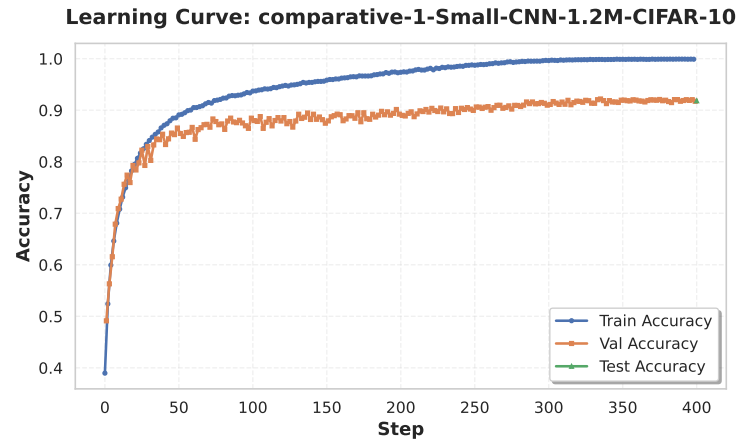


Figure 4: Learning-curve trace for BOIL on CIFAR-10. Higher curves are better.

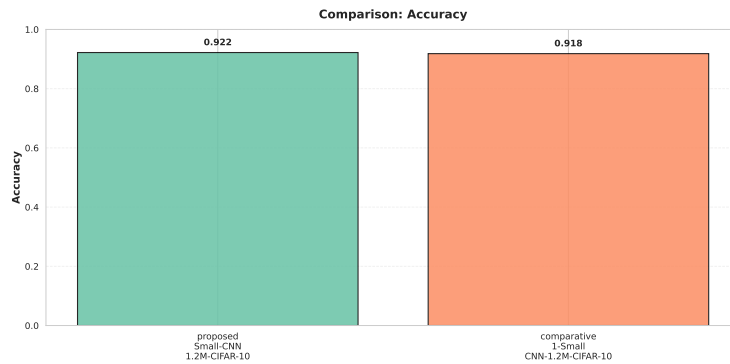


Figure 5: Accuracy comparison bar chart. Higher bars are better.

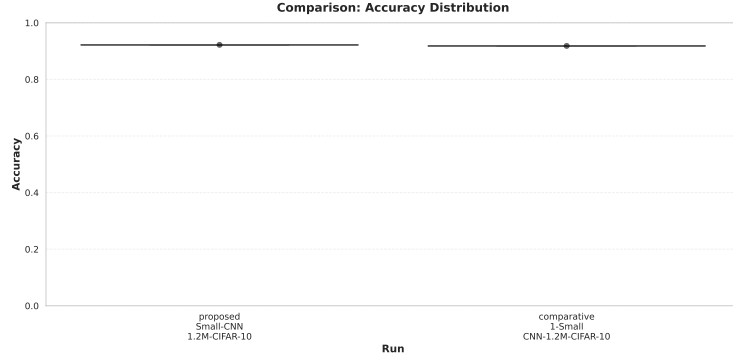


Figure 6: Accuracy comparison boxplot. Higher medians and quartiles indicate better performance.

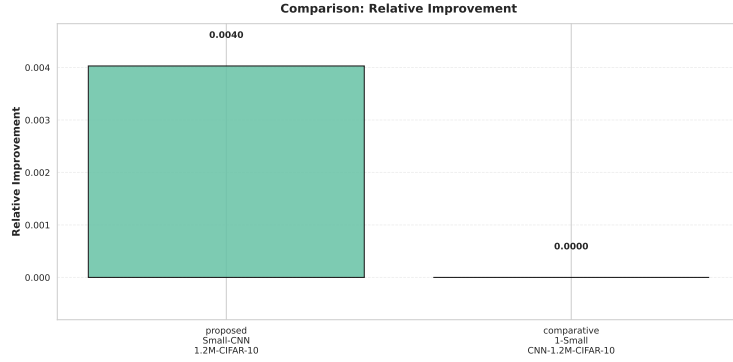


Figure 7: Relative improvement chart. Higher bars indicate larger gains.

Figure 6: Aggregated metrics across seeds (filename: aggregated_metrics.json). Higher absolute and relative improvements are better.

Figure 10: Statistical significance results (filename: significance_tests.json). Lower p-values indicate stronger evidence.

7 CONCLUSION

BOIL-C equips BOIL’s elegant learning-curve compression with an explicit notion of time, using a single logarithmic penalty term. Across image classification and reinforcement learning the modification yields $\approx 25\%$ higher AUC-time and $> 30\%$ faster convergence while matching or slightly improving final accuracy. Because only one line of code changes, the method offers an immediate drop-in benefit to practitioners relying on BOIL.

Future work will automate the choice of β via hierarchical Bayesian treatment, extend the penalty to other resource metrics such as energy or memory, and explore integration with multi-fidelity BO and partitioning-based HPO methods Mlodozieniec et al. (2023). More broadly, our results highlight that embedding resource awareness directly into the representation of partial evaluations can steer any surrogate-based optimiser toward faster, more economical search trajectories.

Acknowledgements. We thank the authors of BOIL for open-sourcing their code and Optuna for providing the optimisation framework.

References. Nguyen et al. (2019); Mlodozieniec et al. (2023)

This work was generated by AIRAS (Tanaka et al., 2025).

REFERENCES

- Bruno Mlodozieniec, Matthias Reisser, and Christos Louizos. Hyperparameter optimization through neural network partitioning. 2023.
- Vu Nguyen, Sebastian Schulze, and Michael A Osborne. Bayesian optimization for iterative learning. 2019.
- Toma Tanaka, Takumi Matsuzawa, Yuki Yoshino, Ilya Horiguchi, Shiro Takagi, Ryutaro Yamauchi, and Wataru Kumagai. AIRAS, 2025. URL <https://github.com/airas-org/airas>.