

UDC 681.327.12(045)

DOI:10.18372/1990-5548.88.20966

Anatoly Kot

MINIMUM FIDELITY FOR RELIABLE ARCHITECTURE RANKING IN BAYESIAN NAS FOR OBJECT DETECTION

Educational and Research Institute for Applied System Analysis, National Technical University of Ukraine “Ihor Sikorsky Kyiv Polytechnic Institute” Kyiv, Ukraine
E-mail: anatoly.kot@gmail.com, ORCID 0000-0002-7490-8834

Abstract—This paper addresses the problem of reducing computational costs in Neural Architecture Search for object detection. The key question in low-fidelity approaches is: after what minimum number of epochs does the early training signal already provide acceptable architecture ranking? This paper presents a methodology for determining minimum fidelity based on out-of-sample rank correlation: trials are split into calibration (70%) and test (30%) sets, a composite proxy is built on the first set and evaluated on the second. An empirical study was conducted on 60 architectures trained for 25 epochs on an object detection dataset (6,772 images, 6 classes). The results show that a rank ensemble of three training metrics (*val_loss*, *train_accuracy*, *val_accuracy*) achieves Spearman $\rho = 0.877$ out-of-sample at epoch 3, providing 88% computational savings. A more complex 7-component metric underperforms the simple ensemble due to overfitting on a small sample ($\rho_{\text{test}} = 0.70$ vs. 0.88). The results are locally valid within the studied search space and the specific dataset.

Keywords—Neural architecture search, bayesian optimization, minimum fidelity, proxy metrics, learning curve, TPE, Optuna.

I. INTRODUCTION

Neural Architecture Search (NAS) has become widely used in computer vision tasks, including object detection. The classic approach – fully training each candidate architecture – is computationally impractical for large search spaces ($\sim 10^9$ configurations), since CNN architectures for detection have a much larger hyperparameter space compared to classifiers, and depth, width, activation functions, and optimizers interact in a non-linear way.

Low-fidelity evaluation – training each trial for only a few epochs – reduces costs by orders of magnitude, but introduces a key uncertainty: does the early signal (metrics after k epochs) reliably reflect the final quality of the architecture? If the answer is “yes, already after 3 epochs”, the potential computational savings are over 80%; if “only after 15 epochs”, the benefit is much smaller. Most studies use a fixed low-fidelity evaluation without empirical justification for the choice of the number of epochs [1]. At the same time, the standard approach – evaluating the proxy metric in-sample (on the same trials used for calibration) – systematically overestimates correlation, since the metric is selected and evaluated on the same data [2].

In this paper, we systematically answer the question: after what minimum number of epochs does a composite proxy reliably rank architectures in out-of-sample conditions? To do this, we conducted a large-scale experiment (60 trials \times 25 epochs),

applied strict out-of-sample evaluation, and analyzed the structure of the metric space using PCA.

II. RELATED WORK

Neural Architecture Search methods fall into three categories: evolutionary [3], gradient-based (DARTS [4]), and Bayesian Optimization-based. DARTS and similar approaches require a continuous relaxation of the search space, which limits their applicability to arbitrary architectural spaces. Bayesian Optimization with TPE [5] places no restrictions on parameter types and works effectively with discrete and categorical spaces – which is critical for tasks with mixed architectural and training hyperparameters.

Zero-cost proxies [6] evaluate an architecture without training – using grad norms, Jacobian covariance, and similar metrics. They are instantaneous but show unstable correlation across different tasks [7]; in particular, their effectiveness on detection tasks is significantly lower than on ImageNet. Learning Curve Extrapolation [8] approximates the learning curve with parametric functions and predicts final performance, but requires enough epochs to make reliable extrapolations. Multi-fidelity methods – Hyperband [9], ASHA – adaptively allocate budget, pruning weak configurations in early rounds, but do not test the hypothesis about the reliability of the early signal out-of-sample and do not identify the minimum sufficient fidelity.

The key methodological problem: if a proxy metric is selected and evaluated on the same trials, the correlation is systematically overestimated. [2] demonstrates that in-sample proxy evaluation systematically inflates rank correlation with final rankings compared to independent out-of-sample testing. [10] shows that surrogate models trained on a small fixed space do not transfer to a larger realistic space without significant loss of ranking quality. However, most low-fidelity NAS studies still evaluate proxies in-sample, which limits the practical value of their findings.

In the literature we reviewed, we found no work that provides a quantitative answer to the question: after what minimum number of epochs does a composite proxy reach an acceptable ranking level under strict out-of-sample conditions for object detection with a mixed architectural space. This study is an empirical step in that direction within a specifically defined search space.

III. PROBLEM STATEMENT

Let \mathcal{A} be the space of architectural configurations and $f(\alpha)$ be the final performance of architecture $\alpha \in \mathcal{A}$ after full training (25 epochs). In low-fidelity evaluation, instead of $f(\alpha)$, a proxy $\hat{f}_k(\alpha)$ is used – a metric computed after $k < K$ epochs. The quality of the proxy is measured by rank correlation:

$$\rho_k = \text{Spearman}(\hat{f}_k(\alpha_1), \dots, \hat{f}_k(\alpha_n), f(\alpha_1), \dots, f(\alpha_n)) \quad (1)$$

The key requirement is out-of-sample evaluation: trials are split into \mathcal{D}_{train} and \mathcal{D}_{test} ; proxy \hat{f}_k is calibrated on \mathcal{D}_{train} , and ρ_k is computed on \mathcal{D}_{test} . Minimum fidelity is defined as:

$$k^* = \min \{k : |\rho_k^{test}| \geq \rho_{thr}, \text{ stably} \}, \quad (2)$$

where $\rho_{thr} = 0.80$ is an operational threshold for practically acceptable ranking; this is a working boundary chosen by the authors to make the problem concrete, not a universal industry standard.

The subject of the study is a military-technical object detection dataset (aircraft, tanks, helicopters, cars, ships, missile systems): 6,772 images, 6 classes, YOLO format, 320×320 pixels, train/val split = 5.417/1.355 (80/20). The 320×320 resolution was chosen as a trade-off between detection speed and quality under a limited GPU budget. In this paper, val_accuracy is defined as the fraction of images in the val set for which the model correctly predicted the class of at least one object with IoU ≥ 0.5 and confidence ≥ 0.5 . This metric is deliberately

simplified compared to mAP: it is cheap to compute, numerically stable from the first epochs (when mAP is still close to zero for weak architectures), and sufficient for comparative ranking within a single search space – exactly what is needed from a proxy. The DynamicDetector search space covers $\sim 10^9$ configurations: number of conv blocks (2–5), filter sizes ($\{16, 32, 64, 128\}$), kernel sizes ($\{3, 5\}$), activation functions ($\{\text{ReLU}, \text{LeakyReLU}, \text{GELU}\}$), optimizers ($\{\text{Adam}, \text{AdamW}, \text{SGD}\}$), learning rate ($\{0.0001, 0.001, 0.01\}$), dropout ($\{0.3, 0.5, 0.7\}$), batch size ($\{16, 32, 64\}$). At each of the 25 epochs of each trial, 9 metrics were collected: train_loss, train_accuracy, val_loss, val_accuracy, gap (val_loss – train_loss), loss_cv (std/mean of losses), grad_norm_mean, grad_cv, output_entropy ($-\sum p_i \log p_i$).

IV. APPROACH

The proposed methodology consists of four steps: (1) calibration run – train N architectures for the full budget of K epochs with per-epoch metric collection; (2) out-of-sample split – divide trials into \mathcal{D}_{train} (70%, 42 trials) and \mathcal{D}_{test} (30%, 18 trials); (3) proxy synthesis – on \mathcal{D}_{train} , select metrics and define the weights of the composite proxy; (4) fidelity analysis – for each k , compute ρ_k^{test} and find k^* .

The composite proxy is built as a rank ensemble: each metric x_m is transformed into a rank z-score $\text{ranknorm}(x_m) = (r(x_m) - \bar{r}) / \sigma_r$, and then a weighted sum is computed:

$$\text{Ensemble}(\alpha, k) = \sum_{m \in \mathcal{M}} w_m \cdot \text{sign}(\rho_m) \cdot \text{ranknorm}(x_m^{(k)}), \quad (3)$$

where ρ_m is the Spearman correlation of metric m with $f(\alpha)$ on \mathcal{D}_{train} , and weights $w_m = |\rho_m| / \sum_j |\rho_j|$. Metrics with $|\rho_m| > 0.5$, $p < 0.5$, are included. After calibration, the weights are fixed and applied to \mathcal{D}_{test} without modification.

The ensemble construction procedure consists of two separate steps: first, a correlation analysis is performed on \mathcal{D}_{train} and metrics with the highest and most stable correlation with final accuracy are selected; then the resulting ensemble is independently evaluated on \mathcal{D}_{test} . For comparison, the following are considered: top3_ensemble {val_loss, train_accuracy, val_accuracy}, single baselines (neg_val_loss, train_accuracy, val_accuracy), PCA_PC1, and the full DSS (7 components). The composition of top3_ensemble was determined based on this correlation analysis; cross-validation

over multiple splits was not performed, which is a limitation of the study.

The structure of the metric space is analyzed through PCA on standardized metrics (8 metrics, excluding gap): $\mathbf{Z} = \text{StandardScaler}(\mathbf{X})$, followed by $\text{PCA} \rightarrow [\text{PC}_1, \text{PC}_2, \dots]$. The PC1 loadings show what structure the ensemble metric approximates. The residual informativeness of each metric beyond val_loss is determined through rank-residuals: $\epsilon_m = \text{ranknorm}(x_m) - \hat{\beta}_m \cdot \text{ranknorm}(x_{\text{val_loss}})$; a metric is considered uniquely informative if

$$\begin{aligned} \epsilon_m &= \text{ranknorm}(x_m) \\ -\hat{\beta}_m \cdot \text{ranknorm}(x_{\text{val_loss}} | \text{Spearman}(\epsilon_m, f)) &> 0.25, \\ p &< 0.05. \end{aligned}$$

TABLE I. SPEARMAN ρ AND KENDALL τ OF METRICS (EPOCH 5) VS. FINAL ACCURACY

Metric	ρ (Spearman)	τ (Kendall)	p-value	Strength
val_loss	-0.8522	-0.7099	$<10^{-6}$	STRONG
train_accuracy	+0.8331	+0.6693	$<10^{-6}$	STRONG
val_accuracy	+0.7869	+0.6317	$<10^{-6}$	STRONG
train_loss	-0.7729	-0.6052	$<10^{-6}$	STRONG
output_entropy	-0.6712	-0.5142	$<10^{-6}$	MODERATE
loss_cv	+0.6339	+0.4960	$<10^{-6}$	MODERATE
grad_cv	-0.5162	-0.3527	$<10^{-4}$	MODERATE
grad_norm_mean	+0.3403	+0.2389	0.0078	WEAK
gap	+0.0005	+0.0034	0.9971	NONE

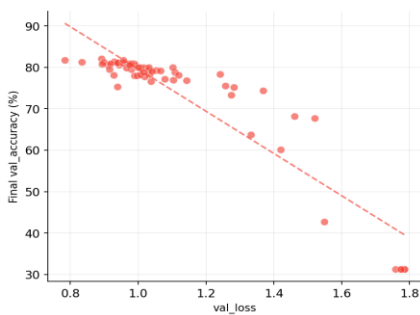


Fig. 1. val_loss at epoch 5 vs. final accuracy (60 trials). Spearman $\rho = -0.8522$

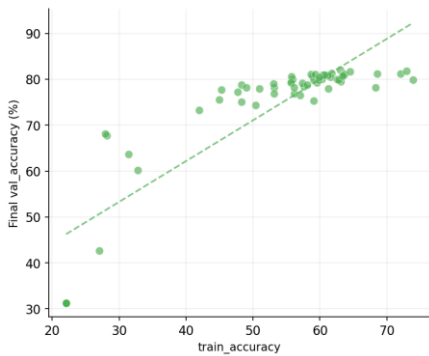


Fig. 2. train_accuracy at epoch 5 vs. final accuracy (60 trials). Spearman $\rho = +0.8331$

Implementation: Optuna [11] (TPE, seed = 42, n_startup_trials = 10, n_trials = 60), PyTorch 2.0+, CUDA 12.1, 25 epochs, SciPy (spearmanr, kendalltau), scikit-learn (PCA, StandardScaler).

V. RESULTS

Correlation analysis at epoch 5 (60 trials) showed that the strongest predictors of final accuracy are val_loss ($\rho = -0.852$, Fig. 1) and train_accuracy ($\rho = +0.833$, Fig. 2); output_entropy shows moderate correlation ($\rho = -0.671$, Fig. 3); gap is excluded due to near-zero correlation ($\rho \approx 0.001$) (Table I). The pairs train_loss/train_accuracy (mutual $\rho = -0.975$) and val_loss/val_accuracy ($\rho = -0.94$) are nearly mirror images and carry the same information (Fig. 4).

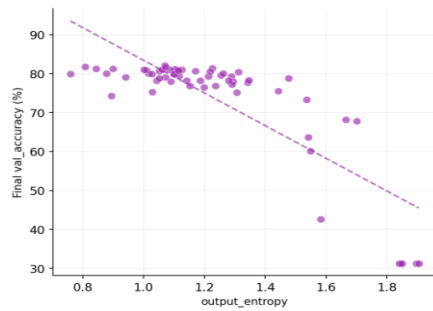


Fig. 3. Output_entropy at epoch 5 vs. final accuracy (60 trials). Spearman $\rho = -0.6712$

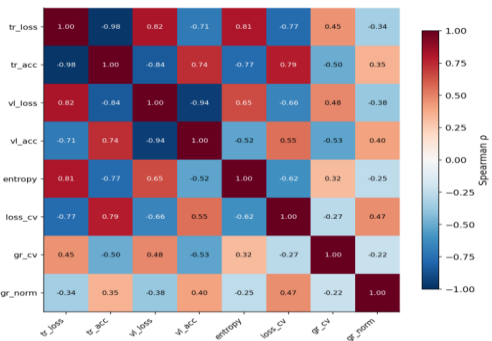


Fig. 4. Cross-correlation matrix of metrics at epoch 5. Two clusters of near-mirror pairs: {train_loss, train_accuracy} and {val_loss, val_accuracy}

Out-of-sample comparison at epoch 5 reveals a significant gap between in-sample and out-of-sample performance: the complex 7-component DSS with in-sample $\rho = 0.862$ drops to $\rho_{\text{test}} = 0.700$ – 23% lower than the best baseline `neg_val_loss` ($\rho_{\text{test}} = 0.797$). This is a sign of overfitting on a small n (Table III).

Analysis of $\rho(k)$ curves across all 25 epochs showed: `top3_ensemble`, `val_accuracy`, and `neg_val_loss` reach the threshold $\rho \geq 0.80$ at epoch 3 ($\rho = 0.877, 0.834, 0.826$ respectively), while

`PCA_PC1` only reaches it at epoch 20 — despite PC1 explaining 75% of metric variance. The advantage of `top3_ensemble` is not in reaching the threshold first, but in combining early signal appearance with a smaller drop amplitude: among methods that reach $\rho \geq 0.80$, it has the highest minimum ρ across epochs (0.652 at epoch 10), while `neg_val_loss` drops to 0.454 and `val_accuracy` to 0.473 (both at epoch 10). `train_accuracy` shows a smaller minimum (0.701) but never crosses the 0.80 threshold within 25 epochs (Tables IV and V, Fig. 5).

TABLE II. RESIDUAL INFORMATIVENESS OF METRICS BEYOND VAL_LOSS

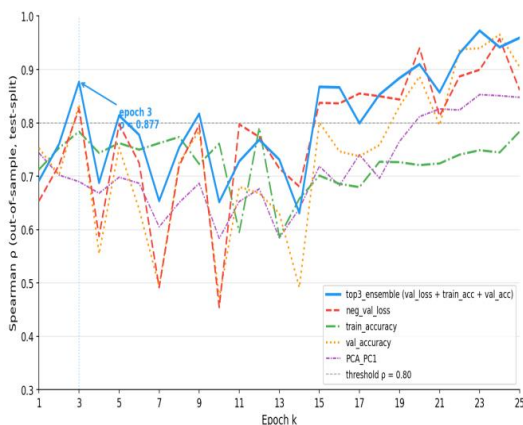
Metric	$\rho(\text{metric, final})$	$\rho(\text{residual, final})$	Unique?
<code>train_accuracy</code>	+0.8331	+0.4307	yes
<code>loss_cv</code>	+0.6339	+0.3346	yes
<code>output_entropy</code>	-0.6712	-0.3333	yes
<code>train_loss</code>	-0.7729	-0.3007	yes
<code>val_accuracy</code>	+0.7869	-0.1031	no (duplicates <code>val_loss</code>)
<code>gap</code>	+0.0005	+0.1827	no
<code>grad_cv</code>	-0.5162	-0.1817	no
<code>grad_norm_mean</code>	+0.3403	+0.1373	no

TABLE III. OUT-OF-SAMPLE PROXY COMPARISON (TEST SPLIT = 18 TRIALS, EPOCH 5)

Method	ρ_{train}	ρ_{test}	τ_{test}	Regret
DSS (7 components)	0.862	0.700	0.564	6.05
<code>train_accuracy</code>	0.824	0.762	0.577	1.85
<code>val_accuracy</code>	0.796	0.753	0.612	0.00
<code>neg_val_loss</code>	0.867	0.797	0.656	0.74

TABLE IV. OUT-OF-SAMPLE $\rho(k)$ BY EPOCH (SELECTED)

Epoch	<code>top3_ensemble</code>	<code>neg_val_loss</code>	<code>train_accuracy</code>	<code>val_accuracy</code>
1	+ 0.692	- 0.654	+ 0.712	+ 0.753
2	+ 0.760	- 0.721	+ 0.753	+ 0.702
3	+ 0.877	- 0.826	+ 0.784	+ 0.834
5	+ 0.814	- 0.797	+ 0.762	+ 0.753
10	+ 0.652	- 0.454	+ 0.761	+ 0.473
15	+ 0.867	- 0.837	+ 0.701	+ 0.800
25	+ 0.959	- 0.861	+ 0.784	+ 0.905

Fig. 5. $\rho(k)$ curves (out-of-sample)

`top3_ensemble` reaches $\rho \geq 0.80$ at epoch 3 and has the highest minimum ρ among methods that cross this threshold (0.652 vs. 0.454 for `neg_val_loss` and 0.473 for `val_accuracy`). `PCA_PC1` reaches the threshold only at epoch 20.

TABLE V. FIRST EPOCH WHERE $|\rho_{\text{TEST}}| \geq 0.80$

Method	First epoch ≥ 0.80	Max ρ	@ epoch
<code>top3_ensemble</code>	epoch 3	0.973	23
<code>val_accuracy</code>	epoch 3	0.965	24
<code>neg_val_loss</code>	epoch 3	0.957	24
<code>PCA_PC1</code>	epoch 20	0.853	23
<code>train_accuracy</code>	never reaches	0.791	12

PCA analysis at epoch 5 confirms the structural nature of the ensemble: PC1 explains 75.0% of variance and has nearly equal loadings (~0.38–0.40) on all main metrics – this is a “general learning quality factor”. PC2 (12.7%) is driven by gradient metrics (grad_cv, loss_cv), which do not correlate with final accuracy in early epochs. Despite the structural similarity between PC1 and top3_ensemble, PCA_PC1 as a predictor underperforms the ensemble (reaching $\rho \geq 0.80$ only at epoch 20), because PCA optimizes the reconstruction of metric variance rather than correlation with the target variable (Figs. 6 and 7).

The best architecture (Trial #12, 82.07%): 5 conv blocks [64→128→32→32→128], kernels [5×5, 3×3, 5×5, 5×5, 3×3], GELU, AdamW, LR = 0.0001, BS=32. All top-5 architectures share: 5 conv blocks, GELU, Adam/AdamW, LR = 0.0001. Worst trial (#30): 31.29% – a gap of 50.78 p.p. (Tables VI, Figs. 8 and 9).



Fig. 6. PC1 (75% of variance) – a uniform ensemble of main metrics (epoch 5, train split)

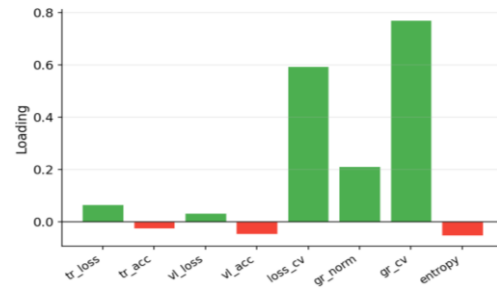


Fig. 7. PC2 (12.7%) – “gradient stability”: dominated by grad_cv and loss_cv

TABLE VI. TOP-5 ARCHITECTURES (60 TRIALS)

Rank	Trial	Best Val Acc	Val Acc @ epoch 3	Architecture	Optimizer
1	12	82.07%	65.54%	5 blocks [64,128,32,32,128], GELU	AdamW, 0.0001
2	45	81.77%	-	5 blocks, GELU	Adam, 0.0001
3	58	81.70%	-	5 blocks, GELU	Adam, 0.0001
4	22	81.33%	-	5 blocks, GELU	AdamW, 0.0001
5	25	81.25%	-	5 blocks, GELU	AdamW, 0.0001

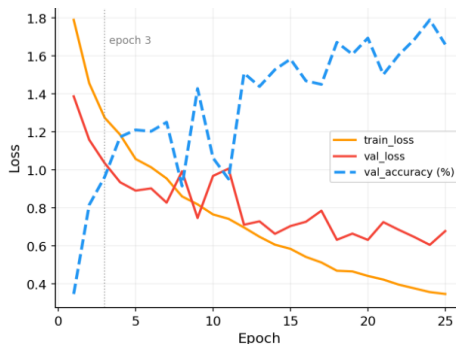


Fig. 8. Learning curves of the best trial (Trial #12, 82.07%). Already at epoch 3, val_loss ≈ 1.04

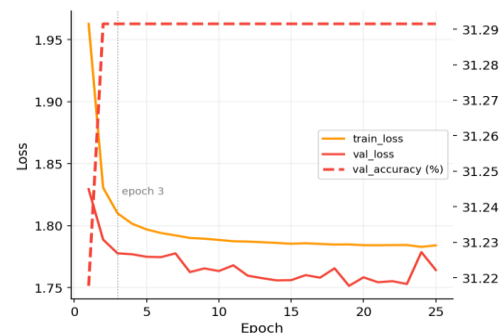


Fig. 9. Learning curves of the worst trial (Trial #30, 31.29%). At epoch 3, val_loss ≈ 1.78 – the gap is clearly visible

TABLE VII. COMPARISON OF EVALUATION STRATEGIES (60 TRIALS)

Strategy	Epochs	Time	ρ out-of-sample	Savings
Full training	25	770 min	oracle	0%
top3_ensemble @ epoch 3	3	~92 min	0.877	88%
top3_ensemble @ epoch 5	5	~154 min	0.814	80%
neg_val_loss @ epoch 5	5	~154 min	0.797	80%
DSS (7 comp.) @ epoch 5	5	~154 min	0.700	80%

VI. CONCLUSIONS

This paper studied the question of minimum fidelity – how many training epochs are sufficient for the early signal to reliably rank architectures in Bayesian NAS for object detection. The central result: a rank ensemble of three training metrics (val_loss, train_accuracy, val_accuracy) reaches the operational threshold of $\rho \geq 0.80$ out-of-sample as early as epoch 3, corresponding to an 88% reduction in computational cost compared to full training. This allows significantly more architectures to be explored for the same budget, or a complete NAS run to be performed on a single GPU.

The key methodological finding is that out-of-sample evaluation is fundamentally important: the 7-component DSS metric achieves in-sample $\rho = 0.86$, but degrades to $\rho = 0.70$ on test trials. This is not an anomaly but a pattern – with small n , a more complex model overfits to the calibration data. The simpler ensemble generalizes better precisely because it does not try to fit the noise.

PCA confirms that top3_ensemble structurally corresponds to the first principal component of the metric space: both are a “general learning quality factor”. However, the direct PCA predictor underperforms the ensemble, because PCA optimizes variance reconstruction rather than correlation with the target variable — a distinction that has practical significance when choosing an approach.

Practical recommendation: for tasks similar to the one studied – a mixed CNN search space for single-domain detection – it is advisable to use the rank ensemble {val_loss, train_accuracy, val_accuracy} at epoch 3 as the objective function for TPE after an initial local calibration on at least 40 trials. Before applying the approach to a new search space or dataset, it is recommended to repeat the calibration experiment to confirm the minimum fidelity.

Scope of validity. The results were obtained within the studied search space (DynamicDetector) and the specific detection dataset. The methodology – out-of-sample proxy evaluation through trial splitting – is transferable; the specific minimum fidelity values (epoch 3) need to be verified on other tasks. Out-of-sample evaluation was performed on a single split (18 test trials), so numerical differences between methods should be interpreted as indicative.

The threshold $\rho_{thr} = 0.80$ is an operational choice made by the authors.

Directions for future work: cross-dataset validation (COCO, Pascal VOC), repeated split stability analysis, adaptive multi-fidelity strategy, meta-learning of ensemble weights for new tasks, extension to transformer-based detectors.

REFERENCES

- [1] T. Elsken, J. H. Metzen, and F. Hutter, “Neural architecture search: A survey,” *JMLR*, 2019. https://doi.org/10.1007/978-3-030-05318-5_11
- [2] C. Sciuto, K. Yu, M. Jaggi, C. Musat, and M. Salzmann, “Evaluating the search phase of neural architecture search,” *ICLR*, 2020.
- [3] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, “Regularized evolution for image classifier architecture search,” *AAAI*, 2019. <https://doi.org/10.1609/aaai.v33i01.33014780>
- [4] H. Liu, K. Simonyan, and Y. Yang, “DARTS: Differentiable architecture search,” *ICLR*, 2019.
- [5] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, “Algorithms for hyper-parameter optimization,” *NeurIPS*, 2011.
- [6] M. S. Abdelfattah, A. Mehrotra, Ł. Dudziak, and N. D. Lane, “Zero-cost proxies for lightweight NAS,” *ICLR*, 2021.
- [7] J. Mellor, J. Turner, A. Storkey, and E. J. Crowley, “Neural architecture search without training,” *ICML*, 2021.
- [8] T. Domhan, J. T. Springenberg, and F. Hutter, “Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves,” *IJCAI*, 2015.
- [9] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, “Hyperband: A novel bandit-based approach to hyperparameter optimization,” *JMLR*, 2017.
- [10] J. Siems, L. Zimmer, A. Zela, J. Lukasik, M. Keuper, and F. Hutter, “NAS-Bench-301 and the case for surrogate benchmarks for NAS,” *NeurIPS (Workshop)*, 2020.
- [11] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” *KDD*, 2019. <https://doi.org/10.1145/3292500.3330701>

Received: February 16, 2026

Accepted: March 14, 2026

Published: April 18, 2026

Kot Anatoly. ORCID 0000-0002-7490-8834. Senior Lecturer.

Artificial Intelligence Department, Educational and Research Institute for Applied System Analysis, National Technical University of Ukraine "Ihor Sikorsky Kyiv Polytechnic Institute," Kyiv, Ukraine.

Education: National Technical University of Ukraine "Ihor Sikorsky Kyiv Polytechnic Institute", (2025).

Research area: artificial Intelligence.

Publications: 12 publications.

E-mail: anatoly.kot@gmail.com

A. Т. Кот. Мінімальна точність для надійного рейтингу архітектури в баєсівському NAS для виявлення об'єктів

Роботу присвячено проблемі зниження обчислювальних витрат у задачі автоматичного синтезу нейронних мереж (Neural Architecture Search, NAS) для детекції об'єктів. Ключовою невизначеністю low-fidelity підходів є питання: після якої мінімальної кількості епох ранній навчальний сигнал вже прийнятно ранжує архітектури? В роботі розроблено методологію визначення мінімальної фіделіті на основі out-of-sample оцінки rank correlation: trials розбиваються на калібраційний (70%) та тестовий (30%) набори, composite proxy синтезується на першому і перевіряється на другому. Проведено емпіричне дослідження на 60 архітектурах, навчених по 25 епох на датасеті детекції об'єктів (6,772 зображень, 6 класів). В результаті встановлено, що rank-ансамбль трьох навчальних метрик (val_loss, train_accuracy, val_accuracy) досягає Spearman $\rho = 0.877$ out-of-sample на епох 3, забезпечуючи 88% економію обчислень. Складна 7-компонентна метрика поступається простому ансамблю через overfitting при малій вибірці ($\rho_{\text{test}} = 0.70$ проти 0.88). Результати є локально валідними у досліджуваному пошуковому просторі та на конкретному датасеті.

Ключові слова: пошук нейронної архітектури, баєсівська оптимізація, мінімальна точність, проксі-метрики, крива навчання, TPE, Optuna.

Кот Анатолій Тарасович. ORCID 0000-0002-7490-8834. Старший викладач.

Кафедра штучного інтелекту, Навчально-науковий інститут прикладного системного аналізу, Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна.

Освіта: Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна, (2025).

Напрямок наукової діяльності: штучні нейронні мережі, програмування.

Кількість публікацій: 12 публікацій.

E-mail: anatoly.kot@gmail.com