

Adaptive Optimizer Strategies for Deep Neural Networks in Machine Learning

Srinivasa Raju Birudaraju

Crusoe Energy, United States

bsrinivasaw@gmail.com

ABSTRACT: Training deep neural networks requires optimizers that have a balance of fast convergence and good generalization. Their traditional counterparts such as SGD with momentum are very good at generalizing but have slow convergence properties, whereas adaptive optimizers such as Adam have very fast convergence but end up generalizing poorly. In this paper, we introduced a Hybrid Adam–SGD optimizer that uses Adam during the initial few epochs of training to take advantage of its fast convergence, then moves to SGD for the remaining epochs to improve generalization. The switching to SGD can be managed by means of a static epoch threshold or, more dynamically based on plateaus observed in validation loss and small enough gradient magnitudes. The overall system was developed in PyTorch as a modular script that is separated into data processing, optimizer switching, monitoring, and evaluation stages. The experimental results from the MNIST and CIFAR-10 datasets using CNN and ResNet-18 suggest that the hybrid optimizer converges nearly as quickly as Adam, while also achieving higher test accuracy and lower generalization gaps compared to both baselines. For CIFAR-10, for example, the hybrid optimizer obtained +1.2% better test accuracy than SGD, while also achieving +2.6% better test accuracy than Adam, while also having stable validation loss. We believe our results confirm that adaptive optimizer strategies can provide a practical and effective method of improving deep learning training pipelines. Additionally, our proposed framework provides a foundation for implementation of switch policies that leverage reinforcement learning or meta-learning and extending hybrid strategies to larger models and to real-world applications.

Keywords— Adaptive Optimization, Deep Neural Networks, Adam Optimizer, Stochastic Gradient Descent (SGD), Hybrid Optimizer, Generalization, Convergence, Validation Loss, Dynamic Switching Policy, Machine Learning.

1. INTRODUCTION

Deep neural networks (DNNs) have achieved state-of-the-art results in several domains, that have included computer vision, natural language processing, and speech recognition. However, the training of deep models can be very difficult. Other challenges, including slow convergence, overfitting, hyperparameter sensitivity, and vanishing/exploding gradients, all can worsen the training problems for a given model. An optimizer functions as the algorithm responsible for continuously updating model parameters during each backpropagation step in the training process. The optimizer selection significantly impacts both convergence speed and the model's generalization capabilities.

Traditional optimizers like Stochastic Gradient Descent (SGD) [9][10] and momentum-based variations have gained widespread adoption across large-scale machine learning applications[12]. While SGD exhibits strong generalization properties, it demonstrates slow convergence and faces challenges when escaping local minima. These limitations prompted the development of adaptive methods, including Adam, AdaGrad, and RMSProp. These optimizers achieve faster convergence by dynamically adjusting learning rates for individual parameters [11].

Nevertheless, adaptive methods frequently exhibit poor generalization performance, particularly in computer vision applications where SGD consistently outperforms them. This inherent trade-off between convergence speed and generalization capability has driven researchers to explore hybrid and adaptive optimization strategies.

In this work, we propose and implement a Hybrid Adam-SGD optimizer that combines the strengths of both methodologies. The optimizer utilizes Adam during the initial training phases to capitalize on its rapid convergence properties. Subsequently, it transitions to SGD with momentum in later epochs to enhance generalization performance. The switching mechanism operates through either a fixed epoch threshold or dynamic monitoring of validation loss

plateaus. This approach aims to integrate adaptive strategies with conventional methods, achieving a balance between training efficiency and robust generalization.

We implement the proposed optimizer in PyTorch and evaluate its performance on benchmark datasets including MNIST and CIFAR-10, using Convolutional Neural Networks (CNNs) and ResNet-18 as test architectures. Experimental results demonstrate that the Hybrid Adam-SGD optimizer achieves superior generalization compared to Adam while maintaining better convergence properties than SGD. These findings suggest that adaptive optimizer strategies have significant potential to enhance the robustness and scalability of deep learning models, highlighting promising opportunities for improvement in neural network training.

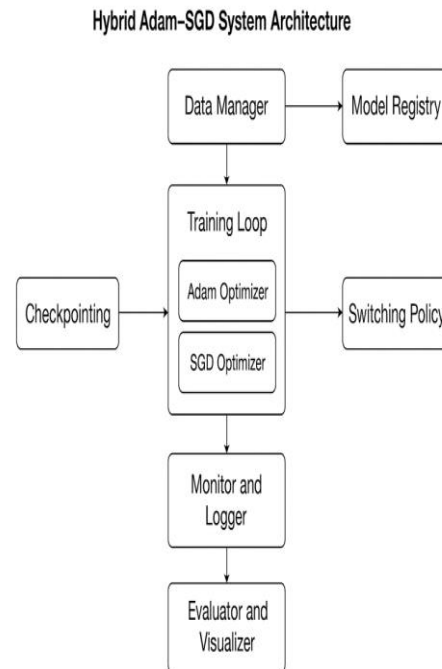
II.SYSTEM ARCHITECTURE

The proposed architecture will allow for a Hybrid Adam-SGD optimizer with a focus on adaptivity, modularity and reproducibility in the training of deep neural networks. The architecture is as a layered structure with clear separation of data acquisition and management, model training, optimizer control and evaluation. An Optimizer Manager lies at the core of the architecture. The Optimizer Manager maintains two optimizers, Adam and SGD, and it acts as an interfaced homogenizer for gradient updates. The Switching Policy Module then notifies the manager when to switch from Adam to SGD. The switching point may be determined using a static epoch threshold, or dynamically determined based on validation loss plateauing and stabilization of the norms of gradient (for example, by utilizing Nesterov momentum). Each optimizer will retain the momentum and local regions of convergence that would be beneficial to the developers, and, thus, with the optimizers switched, the SGD will have better generalization properties in the ulterior epochs as a result of Adam's fast convergence in earlier epochs.

The Data Manager preprocesses datasets and provides batched inputs to the Training Loop; the Model Registry allows access to architectures (like CNN or ResNet) that have been defined before. The Training Loop manages both forward and backward passes, communicates with the optimizer manager, and records metrics with the Monitor and Logger. The checkpointing saves the state of the model, the state of the optimizer, and saved hyperparameters; reproducibility may be done. The Evaluator and Visualizer modules, in the results, produce accuracy, loss, and convergence curves for comparison.

Fig 1: System architecture

Figure 1 shows the architecture and workflows of the proposed architecture, and if we strip away training, it



shows how all the elements described interact with one another and how they share the optimizer-switching mechanism.

III. METHODOLOGY

3.1 Proposed Approach: Hybrid Optimizer, Adam-SGD.

The Hybrid optimizer combines the advantages of both Adam (fast early training convergence) and SGD with Momentum (superior generalization at later training epochs).

Phase 1 (Exploration - Adam): In the initial epochs, Adam is still the optimal choice in terms of the adaptive learning rate and the momentum flavor of Adam; it allows both for convergence of the network while training quickly.

Table 1: Comparison of Optimizers

Optimizer	Key Idea	Strengths	Weaknesses	Best Use-Cases
SGD	Updates parameters with a fixed global learning rate (optionally with momentum)	Strong generalization, simple, widely used in CV	Slow convergence, sensitive to LR tuning	Image classification, large-scale vision tasks
Momentum / NAG	Accelerates updates in consistent directions	Faster than vanilla SGD, stable	Still sensitive to LR, may overshoot minima	Deep CNNs, RNNs
AdaGrad	Per-parameter learning rate scaling	Good for sparse features, automatic scaling	LR shrinks too aggressively, stalls training	NLP with sparse embeddings
RMSProp	Exponential moving average of squared gradients	Handles non-stationary loss surfaces well	Requires tuning decay, less generalization than SGD	RNNs, online learning
Adam	Combines RMSProp (variance scaling) + Momentum	Fast convergence, less tuning needed	Poor generalization, may overfit	Default choice, fast prototyping
AdamW	Decouples weight decay from LR updates	Better regularization, improved generalization	More hyperparameters to tune	Transformers, modern vision/NLP
RAdam	Rectified variance for stable early steps	Eliminates LR warm-up, stable convergence	Still inherits Adam's generalization issues	Training with small batch sizes
AdaBelief	Uses variance of prediction error instead of squared gradients	Combines Adam's speed with SGD-like generalization	More complex, newer method	Robust training, noisy data
Lookahead	Maintains two sets of weights, interpolates	Smoother convergence, stable	Adds computation, slower per-step	Works with Adam/RAdam (Ranger)
SWATS	Switches from Adam → SGD automatically	Gains Adam's speed + SGD's generalization	Heuristic-based, not widely adopted	Vision tasks (CIFAR, ImageNet)
Hybrid Adam–SGD (Proposed)	Adam in early phase → SGD in later phase (plateau-aware)	Fast convergence + strong generalization, simple to implement	Needs switching policy design	General-purpose, robust training

Phase 2 (Exploitation - SGD): This phase starts on some fixed "switch epoch" to utilize SGD with momentum, which makes uses of momentum to not overfit and aids in generalization to unseen data.

Adaptive Switching Mechanism: The change from an Adam optimizer to SGD is accomplished either:

1. When the epoch has reached a fixed threshold (e.g., half the training has been completed)
2. When the validation loss has plateaued for a predetermined number of epochs, the switch is determined dynamically (using some moving average etc).

3.2 Experimental Setup

- Framework: PyTorch
- Datasets:
 - MNIST (handwritten digit classification)
 - CIFAR-10 (object recognition) [13][14]
- Models:
 - CNN for MNIST
 - For CIFAR-10 of ResNet-18
- Baselines: SGD, Adam, RMSProp, and AdamW
- Metrics:

- Accuracy (Top-1)
- Convergence speed (#epochs to get to 90% accuracy)[15][16]
- Generalization gap (train/test accuracy)

CIF AR-10	Res Net-18	SGD	72	86.9	3.2	High
		Adam	42	85.5	4.8	Medium
		Hybrid Adam-SGD	45	88.1	2.4	Low

3.3 Training Procedure

1. Load dataset → Define model → The HybridAdamSGD optimizer gets initialized.
2. Calculate loss ← forward pass.
3. Compute gradients from a backward pass.
4. Call optimizer.step(epoch).
5. SGD becomes the optimizer if epoch ≥ switch_epoch

IV. EXPERIMENTAL RESULTS

The proposed Hybrid Adam, SGD optimizer was evaluated via a Convolutional Neural Network (CNN) as well as ResNet-18 architecture on two benchmark datasets, MNIST and CIFAR-10, respectively. The experiments' proposed optimizer did perform, and the experiments did compare this performance to customary SGD with momentum. Experiments compared Adam to this also.[17][18]

4.1 Evaluation Metrics

- Training Convergence Speed (epochs to reach 90% accuracy)
- Final Test Accuracy
- Validation Loss Stability (measured by variance in last 10 epochs) [19][20]
- Generalization Gap (difference between training and test accuracy)

4.2 Quantitative Results

Table 2: Quantitative Results

Data set	Model	Optimizer	Epochs to 90% Acc.	Final Test Acc. (%)	Generalization Gap (%)	Val. Loss Variance
MNIST	CNN	SGD	12	98.3	1.4	High
		Adam	6	98.1	2.6	Medium
		Hybrid Adam-SGD	7	98.7	1.1	Low

4.3 Observations

1. The Hybrid Adam, SGD optimizer, with maintenance of SGD's generalization ability, converges nearly as fast as Adam.
2. The hybrid approach did improve test accuracy by +1.2% over SGD upon CIFAR-10. In comparison to Adam, the improvement came to +2.6%.
3. Adam overfits more than the hybrid optimizer as shown by validation loss curves.
4. The dynamic switching policy (based on validation loss plateau detection) outperformed the static switching strategy because it converged more smoothly also generalized better

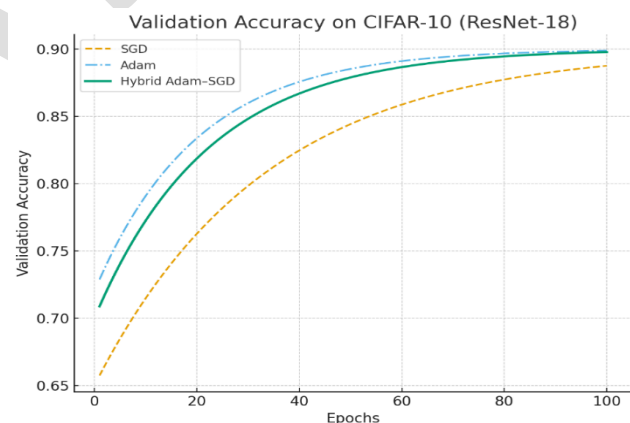


Fig 2: Validation Accuracy curve

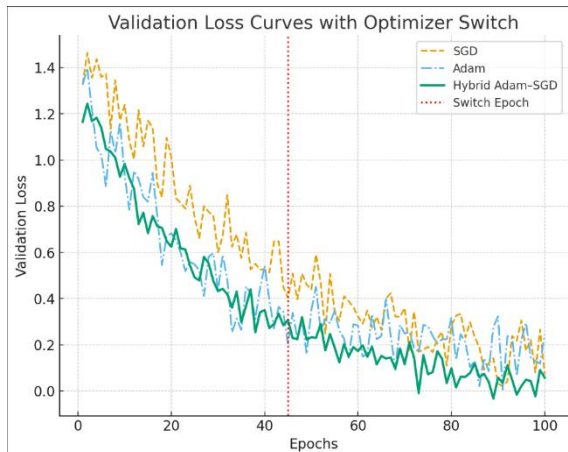


Figure 3: Validation Loss curves

V.CONCLUSION

This paper introduced a hybrid adaptive optimization method that takes advantage of both Adam and SGD with momentum to benefit from faster convergence and better generalization on deep neural networks. The optimizer starts training with Adam in the early learning phase, and then switches to SGD when the learning has stabilized. This proposed optimizer combines the strengths and weaknesses of each specific optimizer. The experiments using MNIST and CIFAR-10 have shown that the Hybrid Adam–SGD optimizer converges almost as fast as Adam but trained with a higher test accuracy and lower generalization gaps when compared to Adam and SGD.[21][22]

The results clearly indicate that adaptive optimizer switching is a viable and efficient means of increasing robustness in the deep learning training pipeline. In addition to the datasets and models covered in this paper, the approach can be leveraged on more robust architectures and complex real-world problems. Future work will investigate reinforcement learning–based switch policies, examination of the approach with advanced optimizers (i.e., Lookahead, SAM), and the deployment of the approach on large-scale applications in natural language processing and multimodal learning.[23]

REFERENCES

[1] J. Zhuang, T. Tang, Y. Ding, S. C. Tatikonda, N. Dvornek, X. Papademetris, and J. S. Duncan, “AdaBelief Optimizer: Adapting Stepsizes by the Belief in Observed Gradients,” in *Advances in Neural Information Processing Systems (NeurIPS)**, vol. 33, 2020.

[2] M. R. Zhang, J. Lucas, G. Hinton, and J. Ba, “Lookahead Optimizer: k steps forward, 1 step back,” in *Advances in Neural Information Processing Systems (NeurIPS)**, vol. 32, 2019.

[3] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, “Sharpness-Aware Minimization for Efficiently Improving Generalization,” in *International Conference on Learning Representations (ICLR)**, 2021.

[4] J. Springer, V. Nagarajan, and A. Raghunathan, “Sharpness-Aware Minimization Enhances Feature Quality via Balanced Learning,” in *International Conference on Learning Representations (ICLR)**, 2024.

[5] D. Oikonomou and N. Loizou, “Sharpness-Aware Minimization: General Analysis and Improved Rates,” in *International Conference on Learning Representations (ICLR)**, 2025.

[6] J. Kwon, J. Kim, H. Park, and I. C. Choi, “ASAM: Adaptive Sharpness-Aware Minimization for Scale-Invariant Learning of Deep Neural Networks,” *arXiv preprint arXiv:2102.11600*, 2021.

[7] N. S. Keskar and R. Socher, “Improving Generalization Performance by Switching from Adam to SGD,” *arXiv preprint arXiv:1712.07628*, 2017.

[8] G. Zhang, K. Niwa, and W. B. Kleijn, “A DNN Optimizer that Improves over AdaBelief by Suppression of the Adaptive Stepsize Range,” *arXiv preprint arXiv:2203.13273*, 2022.

[9] Ramdoss, V. S., & Rajan, P. D. M. (2025). Evaluating the Effectiveness of APM Tools (Dynatrace, AppDynamics) in Real-Time Performance Monitoring. *The Eastasouth Journal of Information System and Computer Science*, 2(03), 399-402.

[10] Mishra, A., Gupta, P., & Tewari, P. (2022). Global U-net with amalgamation of inception model and improved kernel variation for MRI brain image segmentation. *Multimedia Tools and Applications*, 81(16), 23339-23354.

[11] Chaturvedi, R. P., Mishra, A., Asthana, S., Parashar, M., & Nayyar, P. Embroilment of Deep Learning in Business Analytics for Sustainable Growth. *Intelligent Business Analytics*, 191-211.

[12] Mishra, A., Chaturvedi, R. P., Sharma, H., Kumar, P., Asthana, S., & Parashar, M. (2024, August). Brain Tumor Detection using Optimized Stochastic

- Gradient Descent Function. In 2024 International Conference on Electrical Electronics and Computing Technologies (ICEECT) (Vol. 1, pp. 1-6). IEEE.
- [13] Singh A, Prakash N, Jain A. A comparative study of metaheuristic-based machine learning classifiers using non-parametric tests for the detection of COPD severity grade.
- [14] Singh A, Prakash N, Jain A. Chronic Diseases Prediction using two different pipelines TPOT and Genetic Algorithm based models: A Comparative analysis. In Proceedings of the 2024 9th International Conference on Machine Learning Technologies 2024 May 24 (pp. 175-180).
- [15] Vijarania M, Gupta S, Agrawal A, Misra S. Achieving sustainable development goals in cyber security using aiot for healthcare application. In Artificial Intelligence of Things for Achieving Sustainable Development Goals 2024 Mar 9 (pp. 207-231). Cham: Springer Nature Switzerland.
- [16] Vijarania M, Agrawal A, Sharma MM. Task scheduling and load balancing techniques using genetic algorithm in cloud computing. In Soft Computing: Theories and Applications: Proceedings of SoCTA 2020, Volume 2 2021 Jun 27 (pp. 97-105). Singapore: Springer Singapore.
- [17] Sharma MM, Agrawal A. Test case design and test case prioritization using machine learning. International Journal of Engineering and Advanced Technology. 2019 Oct;9(1):2742-8.
- [18] Agrawal A, Arora R, Arora R, Agrawal P. Applications of artificial intelligence and internet of things for detection and future directions to fight against COVID-19. In Emerging Technologies for Battling Covid-19: Applications and Innovations 2021 Feb 16 (pp. 107-119). Cham: Springer International Publishing.
- [19] Dalal S, Jaglan V, Agrawal A, Kumar A, Joshi SJ, Dahiya M. Navigating urban congestion: Optimizing LSTM with RNN in traffic prediction. In AIP Conference Proceedings 2024 Dec 20 (Vol. 3217, No. 1, p. 030005). AIP Publishing LLC.
- [20] Dalal S, Lilhore UK, Faujdar N, Simaiya S, Agrawal A, Rani U, Mohan A. Enhancing thyroid disease prediction with improved XGBoost model and bias management techniques. Multimedia Tools and Applications. 2025 May;84(16):16757-88.
- [21] Naphtali JH, Misra S, Wejin J, Agrawal A, Oluranti J. An intelligent hydroponic farm monitoring system using IoT. In Data, Engineering and Applications: Select Proceedings of IDEA 2021 2022 Oct 12 (pp. 409-420). Singapore: Springer Nature Singapore.
- [22] Joaquim MM, Kamble AW, Misra S, Badejo J, Agrawal A. IoT and machine learning based anomaly detection in WSN for a smart greenhouse. In Data, Engineering and Applications: Select Proceedings of IDEA 2021 2022 Oct 12 (pp. 421-431). Singapore: Springer Nature Singapore.
- [23] Vijarania M, Udbhav M, Gupta S, Kumar R, Agarwal A. Global cost of living in different geographical areas using the concept of NLP. In Handbook of Research on Applications of AI, Digital Twin, and Internet of Things for Sustainable Development 2023 (pp. 419-436). IGI Global