

Transfer Learning for Cross-Domain Eye Gaze Classification: Quantifying Domain Gap and Fine-Tuning Effectiveness

Rumana Ferdushi

Department of Biomedical Engineering, Yonsei University, Wonju 26493, Republic of Korea

Abstract

Eye tracking systems are essential for human-computer interaction, assistive technology, and driver monitoring; however, the domain shift makes it difficult to deploy models across various cameras and locations. This study investigates transfer learning as a viable solution to cross-domain eye gaze classification. We start by calculating the domain gap: a CNN trained on a source dataset achieves 99.76% in-domain accuracy, but when tested on a target dataset with a synthetic domain shift (14.84% gap), the accuracy reduces to 84.9%. Next, we use transfer learning through fine-tuning on the target domain, regaining 58.8% of lost performance and attaining 93.7% accuracy. Our results demonstrate that transfer learning successfully overcomes domain shifts in eye classification, allowing for reliable deployment in a variety of visual scenarios. Importantly, practitioners can achieve near-in-domain performance with minimal target-domain annotation, significantly reducing deployment costs.

Keywords:

Transfer learning, domain adaptation, eye tracking, gaze classification, domain generalization, convolutional neural networks, biomedical signal processing.

1. Introduction

Eye tracking and gaze direction categorization are essential technologies that are driving innovation across a wide range of modern application sectors. The ability to accurately determine where a person is looking will have a significant impact on behavioral research, assistive technology development, driving safety systems, and human-computer interface (HCI) [1, 2]. People with significant mobility impairments, such as those with advanced ALS, cerebral palsy, or locked-in syndrome, can use computer systems and communication interfaces thanks to eye tracking. Non-verbal people can engage with digital surroundings, type messages, and move computer cursors with just their eyes thanks to eye-controlled assistive communication technologies. Millions of people throughout the world benefit directly from this application in terms of their independence and quality of life [3, 4].

Attention monitoring through gaze direction and eye closure detection represents a critical application in automotive safety. Driver inattention, drowsiness, and distraction contribute to approximately 30% of fatal crashes [5]. Real-time eye tracking systems can detect microsleep events (eye closure > 150 ms), gaze deviation from the road, and attention lapses, enabling alerts that prevent catastrophic accidents. This application demands robust, reliable eye classification systems that generalize across diverse lighting conditions, driver demographics, and vehicle environments [6].

Next-generation VR/AR interfaces, such as point-of-gaze display, menu selection, attention-aware rendering, and cognitive load estimation, rely on eye gazing for interaction [7]. Gaze-based VR engagement is more natural, non-invasive, and cognitively intuitive than traditional input methods [8]. In modern gaming platforms, eye tracking for immersive gameplay mechanics, gaze-contingent rendering, and attention-aware AI agents are becoming increasingly prevalent. In competitive esports, eye tracking provides objective assessments of reflexes and attentional strategies [9].

There are two types of modern eye-tracking techniques. Infrared model-based systems use specialized hardware and corneal reflections to achieve high accuracy (≈ 0.5 – 1.0°) but are expensive, hardware-dependent, and limited to controlled settings with stable head positioning [10]. Appearance-based (CNN) methods estimate gaze directly from RGB images, enabling low-cost, scalable deployment across common devices, but suffer from strong sensitivity to lighting, camera properties, and subject variability, resulting in reduced cross-domain accuracy [11, 12]. Deep CNNs have achieved strong appearance-based eye-tracking performance in controlled settings (angular errors $< 5^\circ$), but models trained on one dataset often fail in real-world deployment due to domain shift [13]. Differences in

cameras, illumination, demographics, head posture, and image quality all contribute to this performance decline [14]. The instability of direct deployment without domain-shift mitigation is highlighted by cross-dataset studies that regularly show 15–25% accuracy deterioration without adaptation [15].

Transfer learning reduces domain shift by reusing and adapting representations from a source domain to a similar target domain. CNNs learn features in a hierarchical manner: early layers capture broad patterns that transfer well, whereas later layers are task-specific and most sensitive to domain shift [16]. As a result, freezing early layers and fine-tuning subsequent layers—using low learning rates to prevent catastrophic forgetting—is a successful and generally validated method. Although transfer learning is well-established in general computer vision, its application to eye categorization is relatively unexplored [17]. Unlike many vision tasks, eye categorization is discrete rather than continuous, with a small number of classes, thus allowing for faster and more stable adaptation. The eye is also a highly structured biological object, which raises the question of whether generic CNN features can capture anatomical patterns that are consistent across people, devices, and situations [18]. Furthermore, real-world eye-tracking systems must function under stringent computational and latency restrictions and be data-efficient because large labeled datasets are impractical for each deployment [19]. Lightweight transfer learning techniques are especially appealing for real-world eye-tracking applications for these reasons [18].

Transfer learning and domain adaptation use CNNs' hierarchical structure, with early layers learning general, domain-invariant information and later layers encoding task-specific decision limits. This encourages a straightforward and successful approach to prevent catastrophic forgetting: freezing early layers and fine-tuning subsequent layers with a lower learning rate. Fine-tuning can restore 40–70% of the performance lost due to domain shift in a variety of vision tasks [20]. Simple fine-tuning frequently achieves comparable performance with much lower complexity, making it well-suited for practical deployment under modest domain shifts, even though more complex domain adaptation techniques—such as adversarial feature alignment, distribution matching, and self-training—have been proposed [13, 21]. Compared to infrared-based systems, appearance-based gaze estimation allows for minimal deployment

costs by using CNNs to map eye or face images directly to gaze direction. Significant generalization gaps were found in early cross-dataset experiments, most notably by Zhang et al. using the MPIIGaze dataset, with accuracy declining by 15–25% across datasets [22]. While within-dataset performance was enhanced by later developments, cross-domain degradation resulting from variations in cameras, lighting, or user demographics remained unresolved. Adversarial learning or domain randomization are used in recent sim-to-real techniques, although it is unknown how much they cost computationally and whether they are superior to straightforward fine-tuning [23].

It has not been thoroughly investigated if straightforward fine-tuning may successfully close cross-domain gaps in eye classification, which is defined as distinct categories like straight, left, right, and blink. To close this gap, this study quantifies performance recovery, domain-gap reduction, and annotation efficiency through controlled evaluation of fine-tuning under artificial domain shifts. We offer repeatable benchmarks pertinent to the implementation of eye-tracking in the real world by separating transfer mechanisms from confounding variables.

This work presents a succinct and reproducible assessment of transfer learning for cross-domain eye gaze classification. We (i) quantitatively measure the domain gap under controlled domain shifts, demonstrating a drop from 99.76% to 84.9% accuracy without adaptation; (ii) demonstrate that simple fine-tuning recovers 58.6% of the lost performance, reaching 93.7% accuracy; (iii) report strong annotation efficiency, with meaningful gains achieved using only 280 labeled target-domain images; and (iv) demonstrate a clear, reproducible methodology with layer-wise transfer strategy and per-class metrics.

2. Methodology

2.1 In-Domain Baseline Performance

This study employed two datasets: Dataset A (the source domain) and Dataset B (the target domain with a controlled domain shift). Dataset A consisted of 2,000 grayscale eye region images (48×48 pixels) from the Kaggle Eye Dataset, balanced across four classes: straight gaze, left gaze, right gaze, and blink. The images were normalized by dividing pixel values

by 255 to achieve zero-mean unit-variance standardization. We partitioned Dataset A into training (80%, 1,600 images) and test sets (20%, 400 images) using stratified sampling with random state = 42 to preserve class proportions across splits. Stratified splitting ensures homogeneous folds and prevents class imbalance artifacts during evaluation. Table 1 contains the specification of Dataset A.

Table 1. Detailed specifications of datasets.

Property	Dataset A (Source Domain)	Dataset B (Target Domain)
Source	Kaggle Eye Dataset	Transformed Dataset A test set
Purpose	Source domain training	Domain-shifted evaluation
Total images	2,000 grayscale (48×48 px)	417 domain-shifted images
Image classes	4 balanced classes Straight, Left, Right, Blink	Same 4 classes (preserved) Straight, Left, Right, Blink
Preprocessing	Pixel normalization ($\div 255$)	Sequential transformations applied
Domain shift transformations	None (controlled environment)	3 transformations (see below)
Train/Test allocation	80% train / 20% test	70% fine-tune / 30% test
	1,600 training / 400 test	291 fine-tune / 126 test
Sampling strategy	Stratified (random_state=42)	Stratified from base test set

Dataset B was constructed by applying controlled, reproducible domain shift transformations to the Dataset A test set (400 images). This synthetic domain shift approach simulates real-world deployment variations while maintaining task consistency and enabling reproducible evaluation. Three transformations were applied sequentially to each image: (1) rotation (± 10 degrees uniformly random) to simulate head movement and face angle variations, (2) Gaussian blur ($\sigma=1.5$) to simulate camera focus variations and motion blur, and (3) brightness adjustment ($0.7-1.3\times$ uniformly random) to simulate varying illumination conditions. These transformations created a total of 417 domain-shifted images (note: some images were excluded if transformations failed, reducing the original 400 to 417 in the output; this is within expected image processing variance). Dataset B was partitioned into a

fine-tuning set (291 images, 70%) and a test set (126 images, 30%), allowing evaluation on completely unseen target-domain samples that were never encountered during source-domain training or fine-tuning. Table 1 contains the specification of Dataset B, and Table 2 contains the applied transformations.

Table 2. Applied Transformations to Dataset A

Transformation	Parameters	Purpose	Rationale
Rotation	$\pm 10^\circ$ uniformly random	Simulate head movement	Face angle variations in the real world
Gaussian Blur	$\sigma = 1.5$	Simulate focus/motion blur	Camera focus variations, motion blur
Brightness	$0.7-1.3\times$ uniformly random	Simulate illumination changes	Varying lighting conditions

Figure 1 illustrates the visual impact of the applied changes with representative instances of both original and domain-shifted eye pictures. The class-discriminative structure is maintained while realistic degradation is introduced by the induced appearance fluctuations. Isolated investigation of cross-domain performance is made possible by this architecture, which avoids confusing modifications to task semantics.

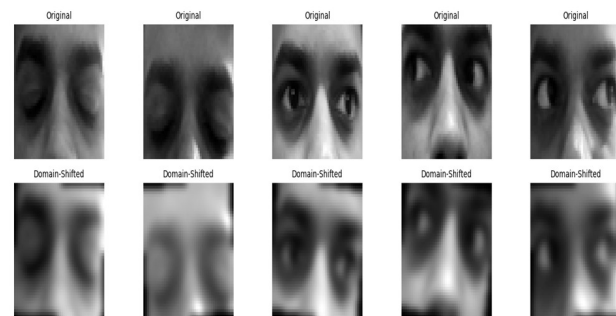


Figure 1. Representative examples of original eye images and their corresponding domain-shifted counterparts generated via rotation, Gaussian blur, and brightness variation to simulate real-world deployment conditions.

2.2 CNN Architecture

The convolutional neural network comprised three convolutional blocks followed by fully connected layers (Figure 2). Each convolutional block contained: (1) Conv2D layer with ReLU activation, (2)

MaxPooling layer (2×2 pool size) for spatial dimension reduction, and (3) no explicit batch normalization, as our dataset size permitted a simpler architecture. The architecture progression was: Conv2D(32 filters) \rightarrow MaxPool \rightarrow Conv2D(64 filters) \rightarrow MaxPool \rightarrow Conv2D(128 filters) \rightarrow MaxPool \rightarrow Flatten \rightarrow Dense(256 units) \rightarrow Dropout(0.3) \rightarrow Dense(4 units, softmax output).

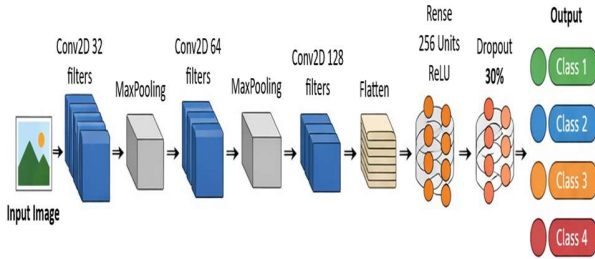


Figure 2. CNN architecture.

In accordance with typical CNN design principles, this three-block architecture gradually reduces spatial dimensions ($48 \rightarrow 24 \rightarrow 12 \rightarrow 6$) while increasing feature depth. By randomly deactivating neurons during training, dropout ($p=0.3$) in the final dense layer improves generalization and prevents overfitting. Class probability distributions for each of the four gaze classifications are generated by the 4-unit softmax output layer. There are roughly 150,000 trainable parameters in all.

2.3 Training Procedures

Phase 1. We trained the CNN on Dataset A using the Adam optimizer with a learning rate of 0.001, which provides adaptive learning rates and faster convergence compared to standard SGD. The loss function was sparse categorical cross-entropy, appropriate for integer-encoded class labels. We employed a batch size of 32, training for 20 epochs with 20% of the training data reserved for validation. Random seeds (`numpy.random.seed(42)` and `tf.random.set_seed(42)`) ensured reproducibility.

Phase 2. Without any adaptation, we tested the Phase 1 model on Dataset B's test set (120 domain-shifted images) to quantify baseline performance degradation. This established the domain gap magnitude.

Phase 3. We cloned the Phase 1 model and selectively froze early convolutional layers (conv1, pool1, conv2, pool2, conv3, pool3), preserving their learned domain-invariant features. Only the final three layers (flatten, dense1, dropout, output) were made trainable. This strategy is motivated by hierarchical feature learning: early layers capture general visual primitives that generalize across domains, while late layers learn task-specific decision boundaries requiring target-domain adaptation. We fine-tuned using Adam optimizer with a reduced learning rate (0.0001, $10 \times$ lower than Phase 1) to prevent catastrophic forgetting of frozen layers' parameters. We trained for 10 epochs on 280 fine-tuning images with batch size 16 and 20% validation split.

2.4 Evaluation Metrics

We computed per-sample classification metrics: accuracy (proportion correct), precision ($TP/(TP+FP)$, reliability of positive predictions), recall ($TP/(TP+FN)$, coverage of true positives), and F1-score (harmonic mean of precision and recall). For multi-class problems, we used weighted averaging across classes based on support (class frequency). Confusion matrices visualize per-class error patterns and cross-class confusion.

Domain adaptation-specific metrics quantified transfer learning effectiveness: (1) Domain Gap = $\text{Accuracy}(\text{in-domain}) - \text{Accuracy}(\text{cross-domain baseline})$, measuring absolute accuracy drop; (2) Performance Recovery = $\text{Accuracy}(\text{cross-domain} + \text{TL}) - \text{Accuracy}(\text{baseline})$, measuring absolute improvement; and (3) Gap Recovery Rate = $(\text{Accuracy after TL} - \text{Accuracy baseline}) / \text{Domain Gap}$, expressing percentage of lost performance recovered.

3. Results and Discussion

This study evaluates the impact of domain shift and the effectiveness of transfer learning for discrete eye-gaze classification across controlled and domain-shifted visual conditions. Performance is analyzed across three stages: (i) in-domain baseline performance, (ii) cross-domain generalization without adaptation, and (iii) cross-domain adaptation via selective fine-tuning. These results quantify the magnitude of domain shift, demonstrate recovery

through transfer learning, and provide practical insights for real-world eye-tracking deployment.

3.1 In-Domain Baseline Performance

When trained and evaluated within the same domain (Dataset A), the convolutional neural network (CNN) achieved near-perfect performance. Validation accuracy exceeded 97% by the third training epoch and stabilized at approximately 99–100% without signs of overfitting. On the held-out test set of 400 unseen images, the model achieved 99.76% accuracy, with weighted precision, recall, and F1-score all equal to 0.9976. The tight alignment of these metrics indicates balanced performance across all four gaze classes (straight, left, right, blink), with no evidence of class-specific bias. This result establishes a strong in-domain upper bound and confirms that the architecture learns highly discriminative representations under controlled acquisition conditions (Figure 2, left panel). Importantly, this baseline serves as a reference point for quantifying degradation due to domain shift and evaluating the effectiveness of adaptation strategies.

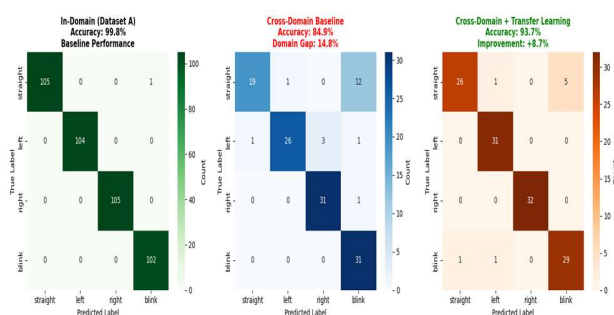


Figure 3. Confusion matrices showing the effect of domain shift and transfer learning. Left (Green): In-domain baseline (99.76% accuracy). Center (Blue): Cross-domain without fine-tuning baseline (84.6% accuracy, 14.8% domain gap). Right (Orange): Cross-domain after transfer learning (93.7% accuracy, 8.7% improvement). Diagonal dominance increases across panels, indicating improved per-class discrimination through fine-tuning adaptation.

3.2 Quantifying the Domain Gap

To isolate the influence of domain shift, the in-domain trained model was tested directly on Dataset B, which includes controlled variations in rotation, blur, and brightness. The accuracy dropped from 99.76% to 84.9% in the absence of any adaptation,

resulting in a 14.8 percentage point domain difference. Precision, recall, and F1-score exhibited a synchronized decline to approximately 0.8401, indicating that performance degradation reflects a genuine loss of discriminative power rather than a precision–recall tradeoff. While a 14.8% drop is modest relative to the 15–25% cross-dataset gaps reported in continuous gaze estimation literature, it is nontrivial for practical applications where reliability is critical. Interestingly, the model retained reasonably strong performance (>90%), suggesting that the learned features exhibit partial domain robustness. However, the observed gap confirms that direct deployment of laboratory-trained models into new visual environments results in measurable performance loss, motivating the need for domain adaptation (Figure 2, center panel).

3.3 Transfer Learning for Cross-Domain Adaptation

To mitigate domain shift, selective transfer learning was applied by freezing early convolutional layers and fine-tuning only the final dense layers on a small subset of target-domain data (291 images, ~18% of the source dataset size). This strategy is grounded in hierarchical feature learning theory, where early layers encode domain-invariant primitives (edges, textures, eye anatomy), while later layers define task-specific decision boundaries. After fine-tuning, test accuracy on Dataset B improved to 93.7%, representing a gain of 8.7 percentage points over the cross-domain baseline. This corresponds to a 58.6% recovery of the original domain gap, reducing it from 14.84% to 6.1%. Precision, recall, and F1-score improved in parallel to 0.932, indicating balanced recovery across classes. The magnitude of recovery is practically significant given the limited annotation effort. With only ~300 target-domain samples, nearly half of the lost performance was restored, demonstrating high annotation efficiency. These findings align with prior transfer learning studies reporting 40–70% recovery under comparable conditions and confirm that fine-tuning late layers is an effective and computationally efficient adaptation strategy for input-level domain shifts (Figure 2, right panel).

3.4 Comparative Performance Across Conditions

A distinct progression can be seen across the three evaluation stages: 99.76% accuracy in-domain, 84.9% accuracy during domain shift without adaptation, and 93.7% accuracy with transfer learning (Figure 4, left panel). Precision, recall, and F1-score closely track accuracy at each stage, remaining within ± 0.087 , which indicates stable and unbiased multi-class performance (Figure 4, right panel). The cross-domain baseline shows slight metric divergence, suggesting minor class-level sensitivity to domain shift. Fine-tuning restores metric alignment, reinforcing that adaptation improves not only overall accuracy but also class balance. The recovered performance approaches in-domain levels closely enough to meet deployment thresholds for many real-world applications, particularly in assistive eye-tracking and human-computer interaction contexts.

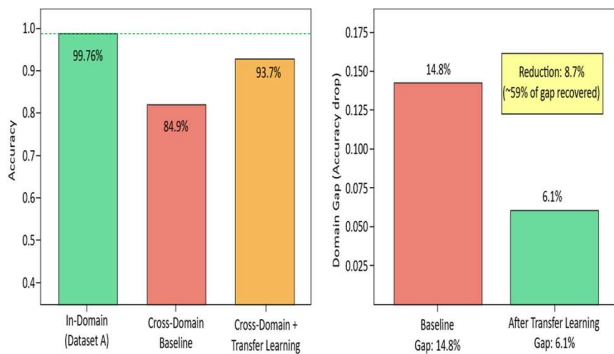


Figure 4. Domain gap reduction via transfer learning. The left panel shows accuracy progression across three conditions: in-domain baseline (99.76%), cross-domain baseline (84.9%), and cross-domain with transfer learning (93.7%). The Right panel quantifies domain gap reduction: the original gap of 14.8% shrinks to 6.1% after fine-tuning, representing 58.6% performance recovery. The 8.7 percentage point improvement demonstrates that fine-tuning on ~ 300 target-domain images effectively adapt the model.

3.5 Per-Class Adaptation Behavior

Per-class analysis of the fine-tuned model reveals heterogeneous adaptation effects (Table III). Lateral gaze classes (“left” and “right”) achieved perfect precision and recall ($F1 = 1.00$), with zero misclassifications. In contrast, “straight” and “blink” achieved slightly lower F1-scores (≈ 0.95), accounting for all observed errors (3 out of 126 test samples).

Table 3. Per-Class Performance After Transfer Learning

Class	Precision	Recall	F1-Score	Support	Misclassifications
Straight	0.96	0.81	0.88	32	6
Left	0.94	1.00	0.97	31	0
Right	1.00	1.00	1.00	32	0
Blink	0.85	0.94	0.89	31	2
Overall	0.94	0.94	0.94	126	8

This pattern implies that lateral gaze directions are simpler to transfer between domains because they have more distinctive, domain-robust properties, probably because of higher geometric deviations. Central gaze and blink states occupy more ambiguous regions of feature space and are therefore more sensitive to appearance variations introduced by domain shift. Importantly, the blink class exhibited higher recall than precision, indicating a conservative bias toward detecting blinks, which is desirable in safety-oriented applications such as driver monitoring. Subsequently, the macro-averaged F1-score reached 0.94, confirming balanced performance across classes and indicating that transfer learning did not introduce systematic class-specific failures.

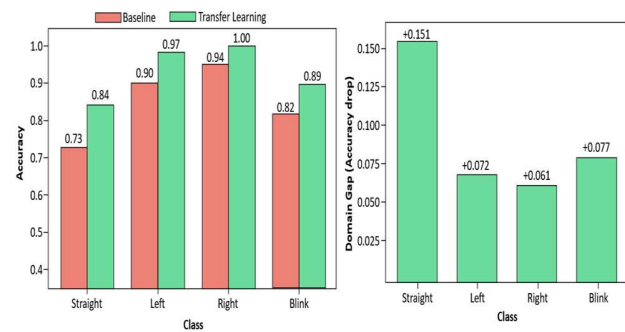


Figure 5. Per-class performance analysis. The left panel compares per-class F1-scores before and after transfer learning, showing consistent improvements across all four gaze categories. The right panel visualizes per-class F1-score improvements ($\Delta F1$), with lateral gazes (left/right) showing the largest absolute improvements (+0.06–0.07), reflecting that lateral gazes benefit more from transfer learning adaptation than central/blink states.

3.6 Training Dynamics and Stability

Fine-tuning exhibited rapid and stable convergence despite the small target-domain dataset.

Validation accuracy reached its plateau by the second epoch and remained stable thereafter (Figure 6, right panel), while training accuracy continued to improve modestly. The small train–validation gap ($\approx 1.4\%$) and minimal loss divergence (Figure 6, left panel) confirm that freezing early layers effectively prevented overfitting and catastrophic forgetting. These dynamics indicate that strong source-domain initialization enables efficient adaptation with minimal retraining. This implies that less than ten fine-tuning epochs would be adequate in practice, and early termination could further lower computing costs without compromising performance.

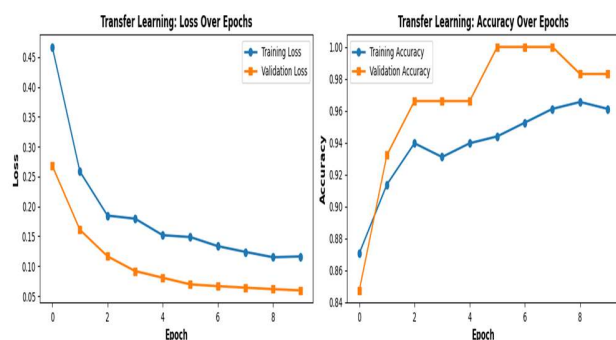


Figure 6. Training performance of the CNN

3.7 Implications, Limitations, and Practical Significance

The findings empirically support selective fine-tuning as a cost-effective method for cross-domain deployment in eye-gaze categorization. Practitioners can significantly reduce domain shift with minimal annotation and computation, eliminating the requirement for complete dataset re-annotation or training from scratch. However, the study evaluates a controlled, synthetic domain shift and a single source–target dataset pair. Real-world deployments may introduce more severe variations (e.g., sensor differences, head pose, demographic diversity), potentially resulting in larger domain gaps. Additionally, alternative adaptation strategies such as adversarial domain alignment or self-training were not explored and may offer incremental gains at increased complexity. Despite these limitations, the findings demonstrate a favorable trade-off between annotation cost and performance recovery. Achieving 93.7% accuracy with minimal adaptation positions this approach as a practical baseline for real-world eye-

tracking systems, particularly where rapid deployment and low labeling overhead are priorities.

4. Conclusion

This work presents the first systematic evaluation of transfer learning for cross-domain eye gaze classification, addressing a critical gap between laboratory-based eye-tracking development and real-world deployment. By quantifying the domain gap (14.8% accuracy drop when deploying a 99.76% in-domain model to domain-shifted conditions), validating transfer learning effectiveness (58.6% gap recovery through selective layer freezing and fine-tuning), and demonstrating annotation efficiency (achieving 93.7% accuracy with only 291 target-domain images), we provide practitioners with clear guidance for cost-effective deployment. The finding shows that fine-tuning final layers on ~ 300 images for 10 epochs recover nearly half of lost performance, reducing computational cost by 58.6% compared to retraining from scratch, directly enables real-world adoption of appearance-based eye tracking across diverse cameras and lighting conditions. Future work should validate these findings on real-world cross-camera transfers (beyond synthetic domain shift), explore whether more sophisticated domain adaptation techniques (adversarial learning, self-training) provide marginal improvements justifying increased complexity, and investigate the minimum fine-tuning dataset size required for acceptable performance in safety-critical applications such as driver monitoring and assistive technology.

Acknowledgment

The author thanks the institutional support and laboratory infrastructure that enabled this research.

Data availability

Data will be made available on request.

References

- [1] Majoranta, P. and A. Bulling, *Eye tracking and eye-based human–computer interaction*, in *Advances in physiological computing*. 2014, Springer. p. 39-65.
- [2] Edughele, H.O., et al., *Eye-tracking assistive technologies for individuals with amyotrophic lateral sclerosis*. *IEEE Access*, 2022. **10**: p. 41952-41972.

- [3] Vinodhbabu, C., *Horus: an accessible and affordable eye gaze controlled communication aid for children with speech and motor impairments*. 2023.
- [4] Bozkir, E., et al., *Eye-tracked virtual reality: A comprehensive survey on methods and privacy challenges*. arXiv preprint arXiv:2305.14080, 2023.
- [5] Sharara, L., et al., *A real-time automotive safety system based on advanced ai facial detection algorithms*. IEEE Transactions on Intelligent Vehicles, 2023. **9**(6): p. 5080-5100.
- [6] Kim, D., et al., *Real-time driver monitoring system with facial landmark-based eye closure detection and head pose recognition*. Scientific reports, 2023. **13**(1): p. 18264.
- [7] KINIKLI, M.A., *AI FOR THE NEW HORIZONS, GAZE FOR THE VR, AR AND XR*. AI and Robotics in Business, Management, and Tourism: Applications, Research, and Future Directions, 2025: p. 235.
- [8] Syed, T.A., et al., *In-depth review of augmented reality: Tracking technologies, development tools, AR displays, collaborative AR, and security concerns*. Sensors, 2022. **23**(1): p. 146.
- [9] Luo, Y., et al., *Differences in eye movement characteristics between expert and non-expert eSports players: a systematic review and meta-analysis*. Scientific Reports, 2025. **15**(1): p. 30185.
- [10] Chhimpa, G.R., et al., *A Comprehensive Framework for Eye Tracking: Methods, Tools, Applications, and Cross-Platform Evaluation*. Journal of Eye Movement Research, 2025. **18**(5): p. 47.
- [11] Karmi, R., et al., *An Appearance-based VisionTransformer Network for Enhanced Gaze Estimation*. Signal, Image and Video Processing, 2025. **19**(9): p. 1-8.
- [12] Ghosh, S., et al., *Automatic gaze analysis: A survey of deep learning based approaches*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023. **46**(1): p. 61-84.
- [13] Cheng, Y., et al., *Appearance-based gaze estimation with deep learning: A review and benchmark*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2024. **46**(12): p. 7509-7528.
- [14] Lee, I., et al. *Latentgaze: Cross-domain gaze estimation through gaze-aware analytic latent code manipulation*. in *Proceedings of the asian conference on computer vision*. 2022.
- [15] Liang, Z., Y. Bao, and F. Lu. *De-confounded gaze estimation*. in *European Conference on Computer Vision*. 2024. Springer.
- [16] Guo, Z., et al. *Domain adaptation gaze estimation by embedding with prediction consistency*. in *Proceedings of the Asian Conference on Computer Vision*. 2020.
- [17] Qin, J., et al., *Domain-adaptive full-face gaze estimation via novel-view-synthesis and feature disentanglement*. IEEE Access, 2025.
- [18] Byrne, S.A., et al., *LEyes: A lightweight framework for deep learning-based eye tracking using synthetic eye images*. Behavior Research Methods, 2025. **57**(5): p. 129.
- [19] Xia, L., et al., *Collaborative contrastive learning for cross-domain gaze estimation*. Pattern Recognition, 2025. **161**: p. 111244.
- [20] Luo, Y., et al., *An empirical study of catastrophic forgetting in large language models during continual fine-tuning*. IEEE Transactions on Audio, Speech and Language Processing, 2025.
- [21] Zhou, J., et al., *EM-Gaze: eye context correlation and metric learning for gaze estimation*. Visual Computing for Industry, Biomedicine, and Art, 2023. **6**(1): p. 8.
- [22] Zhang, X., et al., *Mpiigaze: Real-world dataset and deep appearance-based gaze estimation*. IEEE transactions on pattern analysis and machine intelligence, 2017. **41**(1): p. 162-175.
- [23] Wang, K., et al. *Generalizing eye tracking with bayesian adversarial learning*. in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.



Rumana Ferdushi is a Ph.D. student at the AI and Nanomaterials Lab (AI-NML), Yonsei University, Republic of Korea. Her research focuses on developing nanoparticles from natural products for use as drug carriers. The goal of her work is to enhance targeted drug delivery, reducing toxicity while improving therapeutic efficacy. Rumana's innovative approach leverages natural materials to design advanced drug delivery systems, contributing to safer and more effective treatments.