

Open Gaze: Open Source eye tracker for smartphone devices using Deep Learning

Authors

Jyothi Swaroop Reddy Bommareddy
JNTUH College of Engineering Hyderabad
Hyderabad
Telangana
India

Sushmanth Reddy Mereddy
Amrita Vishwa Vidyapeetam
Amritapuri
Kerala
India

Corresponding author

Jyothi Swaroop Reddy Bommareddy
JNTUH College of Engineering Hyderabad
Hyderabad
Telangana
India
Email: bjsreddy742002@gmail.com

Abstract

Eye tracking has been a pivotal tool in diverse fields such as vision research, language analysis, and usability assessment. The majority of prior investigations, however, have concentrated on expansive desktop displays employing specialized, costly eye tracking hardware that lacks scalability. Remarkably little insight exists into ocular movement patterns on smartphones, despite their widespread adoption and significant usage. In this manuscript, we present an open-source implementation of a smartphone-based gaze tracker that emulates the methodology proposed by a [GooglePaper](#) (whose source code remains proprietary). Our focus is on attaining accuracy comparable to that attained through the [GooglePaper](#)'s methodology, without the necessity for supplementary hardware.

Through the integration of machine learning techniques, we unveil an accurate eye tracking solution that is native to smartphones. Our approach demonstrates precision akin to the state-of-the-art mobile eye trackers, which are characterized by a cost that is two orders of magnitude higher. Leveraging the vast MIT GazeCapture dataset, which is available through registration on the dataset's website, we successfully replicate crucial findings from previous studies concerning ocular motion behavior in oculomotor tasks and saliency analyses during natural image observation. Furthermore, we emphasize the applicability of smartphone-based gaze tracking in discerning reading comprehension challenges. Our findings exhibit the inherent potential to amplify eye movement research by significant proportions, accommodating participation from thousands of subjects with explicit consent. This scalability not only fosters advancements in vision research, but also extends its benefits to domains such as accessibility enhancement and healthcare applications.

Keywords

Gaze Capture, Deep Learning, Support Vector Regression, Computer Vision

1. Introduction

In recent decades, eye tracking has emerged as a fundamental tool across a spectrum of disciplines, encompassing vision research, linguistic analysis, and usability assessment. However, the predominant focus of prior investigations has been on conventional desktop displays, often necessitating the use of specialized and costly eye tracking hardware. Regrettably, such systems suffer from limitations in scalability and accessibility, impeding their widespread adoption. Concurrently, the ubiquity of smartphones has redefined human-computer interaction (HCI), prompting the need for a comprehensive understanding of eye movement behavior in this context. Despite the considerable time users spend interacting with smartphones, there remains a conspicuous dearth of research addressing ocular motion patterns on these devices. The untapped potential of smartphone-based eye tracking to unlock new avenues in vision research, accessibility improvements, and healthcare applications is a tantalizing prospect.

This paper introduces an open-source implementation of a gaze tracking solution for smartphones, inspired by the principles outlined in a proprietary [GooglePaper](#) [1]. The [GooglePaper](#) showcased a novel approach to accurate eye tracking, leveraging machine learning to achieve remarkable precision without the reliance on additional hardware. Specifically, the original implementation by Google reported that a deviation of 0.6–1 degree at a viewing range spanning 25 to 40 cm for smartphone usage indicates that the algorithm forecasts the position of a point on the device, observed from a distance of 25 to 40 cm, with a precision range of 0.46 ± 0.03 cm. Our efforts are aimed at reproducing and expanding upon their approach, specifically concentrating on adopting TensorFlow and PyTorch for implementing a Support Vector Regression (SVR) framework. The creators of the [GooglePaper](#) made the decision to keep their code private and have not made their trained models publicly available. Consequently, our main objective in this undertaking is to duplicate the documented levels of precision. Subsequently, we intend to extend the capabilities to include forecasting head position and other relevant factors.

To enable comprehensive evaluation, we utilize the extensive [MIT GazeCapture dataset](#), which provides a wealth of eye movement data [2]. This dataset, accessible through registration on its dedicated website, furnishes us with JSON files containing essential information such as bounding box coordinates for eyes and faces, alongside frame, face detection, and eye detection statistics. The inclusion of this dataset allows us to gauge the performance of our TensorFlow and PyTorch models in a standardized and rigorous manner. The overarching goal of this study is

to offer the open-source community a robust, accessible, and accurate gaze tracking solution that capitalizes on the widespread prevalence of smartphones. By demonstrating parity with the [GooglePaper](#)'s accuracy metrics, which typically require specialized and costly equipment, we aim to enable a new era of cost-effective, scalable eye movement research. This has the potential to catalyze advancements in fields ranging from vision research to improving digital accessibility and healthcare diagnostics.

In the subsequent sections of this paper, we delve into the technical details of our TensorFlow and PyTorch implementations, elucidating the core components of our models and their training processes. We then present a comprehensive analysis of the performance of our gaze tracking solution on the [MIT GazeCapture dataset](#), benchmarking it against the metrics established in the [GooglePaper](#). Additionally, we discuss the implications of our findings for vision research and other applications, emphasizing the scalability and accessibility that our smartphone-based approach brings to the table. This paper lays the foundation for a paradigm shift in eye tracking research by offering an open-source gaze tracking solution tailored for smartphones. With our approach, we aim to help democratize the access to high-precision eye movement analysis, fostering innovation and progress across a myriad of fields that benefit from insights into human visual behavior.

2. Dataset

The pivotal role of datasets in modern machine learning research cannot be understated, and the [MIT GazeCapture dataset](#) serves as a cornerstone for our smartphone-based gaze tracking project. This comprehensive dataset, originally released in 2016, is a treasure trove of eye movement data collected from a diverse set of participants engaged in various activities using smartphones. Accessible through registration on its dedicated website, the [MIT GazeCapture dataset](#) is meticulously curated to facilitate research on eye tracking methodologies. The dataset contains JSON files that are paired with corresponding images, creating a rich repository of information about eye and face movements. These files encompass a wealth of critical details, including bounding box coordinates for eyes and faces, the number of frames captured, as well as data on face and eye detection.

One of the dataset's strengths lies in its capacity to capture naturalistic scenarios. Participants were encouraged to interact with their smartphones in everyday activities, such as reading, browsing, and texting. This diversity of activities lends an authentic quality to the data, making it highly suitable for evaluating the performance of our smartphone-based gaze tracking models in real-world settings. Moreover, the dataset provides a unique opportunity to gauge the accuracy of our gaze tracking solution. By comparing our model's predictions with the ground truth eye movements recorded in the dataset, we are able to ascertain the extent to which our implementation aligns with actual gaze behavior. The [MIT GazeCapture dataset](#), with its extensive coverage of eye movement behaviors in smartphone interactions, constitutes an invaluable resource in our quest to create an open-source smartphone-based gaze tracker. Its inclusion in our study not only ensures robust evaluation, but also enables the replication of gaze-related findings in a controlled and standardized manner.

In the subsequent sections of this study, we delve into the technical intricacies of how we harnessed this dataset to train and validate our gaze tracking models in TensorFlow and PyTorch. Through this exploration, we offer insights into the efficacy of our approach in replicating the results reported in the [GooglePaper](#) and potentially expanding their scope to include additional variables such as head position prediction. The [MIT GazeCapture dataset](#) serves as the foundation upon which our research stands, enabling us to contribute to the democratization of accurate gaze tracking technology. This dataset is not just a repository of data points, but also a key enabler in our efforts to advance eye movement research, accessibility, and healthcare applications through innovative smartphone-based gaze tracking solutions.

Figure 1 illustrates the participant distribution across various devices, accompanied by the test/train/validation split as initially presented by the Team of GazeCapture.

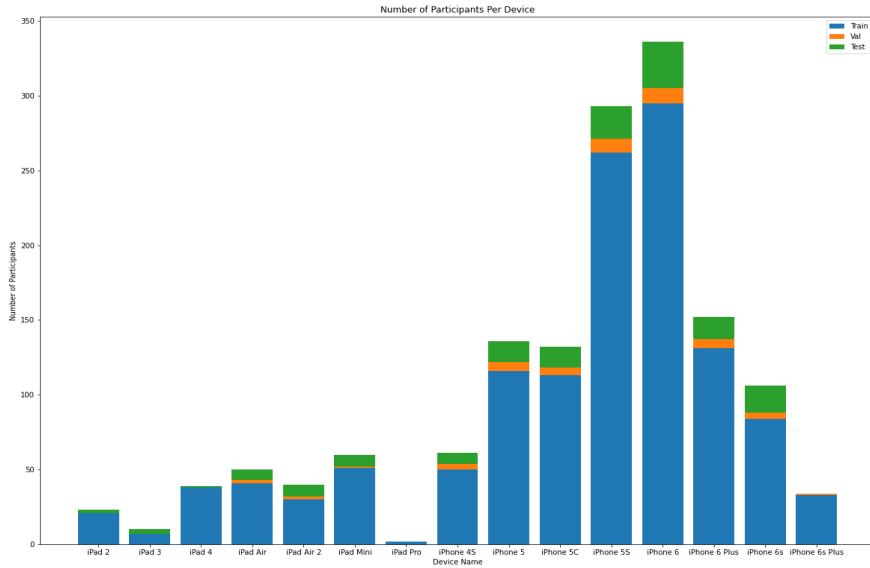


Fig. 1 Participant distribution across various devices

3. Dataset splits

3.1 MIT Split

Taking cues from the GazeCapture methodology, the MIT Split strategy upholds the practice of segregating training, validation, and testing datasets on a per-participant basis. This deliberate technique guarantees the exclusive allocation of a participant's data to the train, test, or validation sets. Such an approach effectively boosts the model's learning and generalization capabilities, as data from a single individual does not appear across all the splits.

The detailed information about this split is provided in the following sections.

3.2 Exclusively Mobile Devices in All Orientations

The subsequent dataset maintains the division of data following the guidance from GazeCapture, orientations encompassing all possibilities. The filters are subsequently employed:

- Face Detections that are legitimate
- Exclusively mobile data
- Eye detections that are legitimate.

The breakdown of frame quantities for each device is illustrated in Figure 2. The corresponding test/train/validation split is shown in Table 1.

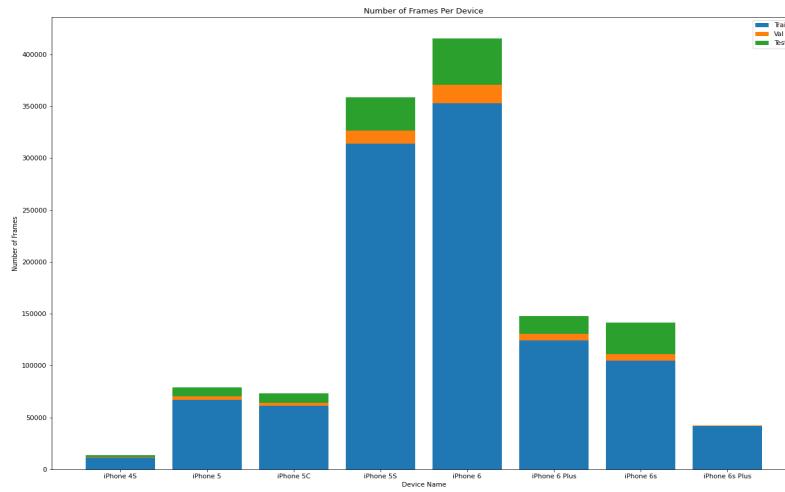


Fig. 2 Breakdown of frame quantities per device

Table 1 Train/validation/test split related to Fig. 2

Split	Number of Participants	Total Frames
Train	1,081	1,076,797
Validation	45	51,592
Test	121	143,796

3.3 Exclusively Mobile Devices in Portrait Mode

The initial dataset under consideration closely aligns with Google's training approach. The following filters are put into effect:

- Exclusively mobile device data
- Face detections that meet the criteria
- Orientation solely in portrait
- Eye detections that have been validated.

Figure 3 illustrates the dispersion of frame counts per individual device. The corresponding train/validation/test split is shown in Table 2.

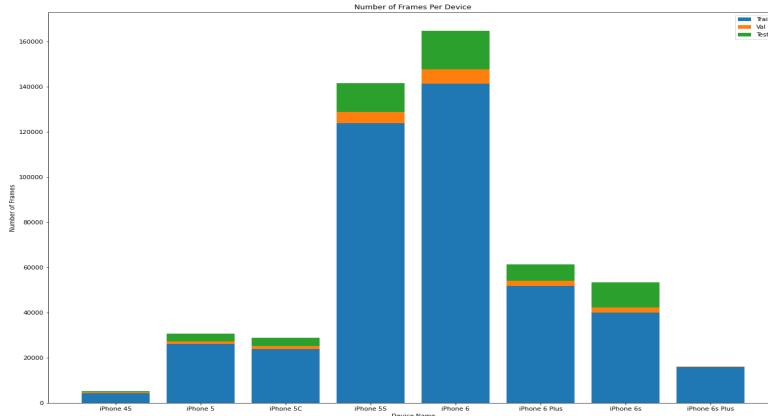


Fig. 3 Dispersion of frame counts per individual device

Table 2 Train/validation/test split related to Fig. 3

Split	Number of Participants	Total Frames
Train	1,075	427,092
Validation	45	19,102
Test	121	55,541

3.4 Google Split

In contrast, the Google Split is structured based on the unique ground truth points adopted by Google for their dataset division. Consequently, each participant's frames are included in the train, test, and validation sets.

Figure 4 illustrates the number of frames per individual device.

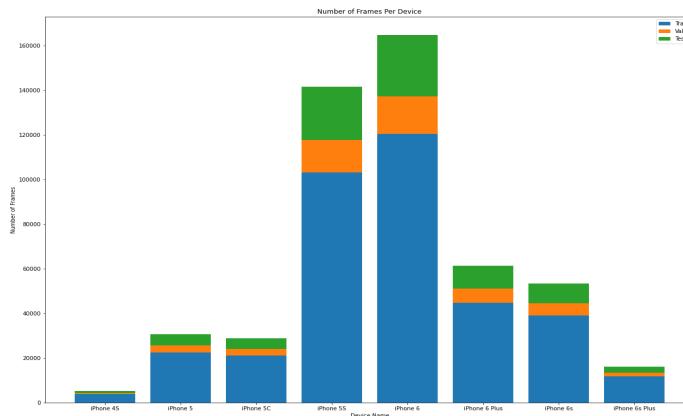


Fig. 4 Number of frames per individual device

The detailed information about the train/validation/test split is provided in Table 3.

Table 3 Train/validation/test split related to Fig. 4

Split	Number of Participants	Total Frames
Train	1,241	366,940
Validation	1,219	50,946
Test	1,233	83,849

The utilization of these distinct splits provides us with a robust framework to gauge the performance of our gaze tracking models across various conditions. It enables a rigorous evaluation process that showcases the efficacy of our models in emulating the accuracy and functionality of the [GooglePaper](#), while also potentially extending their capabilities to predict additional variables such as head position. Through this detailed exploration, we aim to provide a comprehensive understanding of our approach and its implications for advancing eye movement research and related applications.

4. Model Architecture

The gaze tracking model we've developed employs a multilayer feed-forward convolutional neural network (ConvNet) that encompasses several distinctive components for accurate gaze prediction. The process begins with extracting vital face features from input images, encompassing the face's bounding box and six distinct landmarks. These features are detected through a specialized face detector based on MobileNets, utilizing the Single Shot MultiBox Detector (SSD) technique to precisely identify faces. Subsequently, the base model is trained, prioritizing the face and its landmarks, using the [MIT GazeCapture dataset](#). This training equips the model to predict gaze locations based on facial cues effectively. Moving forward, we process the regions around each eye by cropping them with the guidance of corner landmarks. These cropped regions are then scaled to a consistent size of $128 \times 128 \times 3$ pixels. Each eye region is passed through individual identical ConvNet towers, each comprising three convolutional layers. The kernel sizes of these layers progressively decrease from 7×7 to 5×5 and then to 3×3 . The layers facilitate the model in learning intricate patterns within the eye regions. Deeper layers yield an increasing number of features, generating 32, 64, and 128 output channels respectively. Figure 5 depicts the gaze track network architecture.

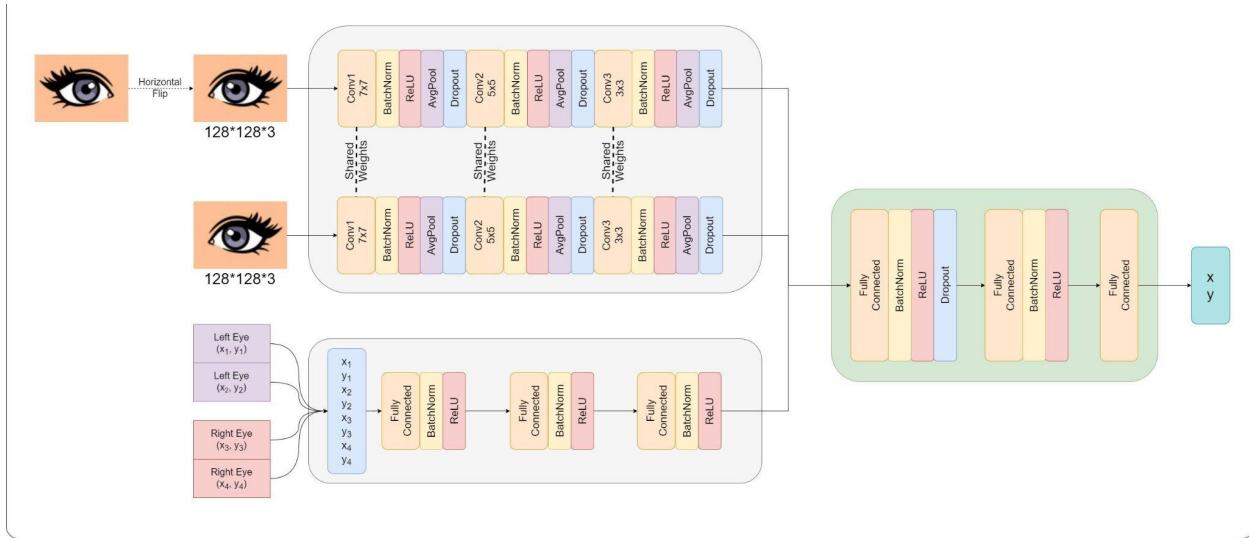


Fig. 5 Gaze track network architecture

Pooling layers are applied after each convolutional layer to condense information and alleviate complexity. Nonlinearities are introduced through Rectified Linear Units (ReLUs), enabling the model to capture complex

relationships within the data. To promote symmetry and simplify training, we horizontally flip the crop of the left eye and process it in a manner analogous to the right eye, allowing learned patterns to be shared between both eyes. The integration of eye and landmark data transpires as the outputs from the ConvNet towers are combined with the results from fully connected layers processing inner and outer eye corner landmarks. This fusion step enables effective incorporation of both eye and facial information. The final regression head of the model generates two coordinates – the x and y screen coordinates representing the gaze direction. Comprising multiple fully connected layers, each with distinct hidden units, this section processes the amalgamated eye and face information to predict gaze locations. To heighten model accuracy, we employ calibration data acquired as participants focus on a stimulus. This data refines the base model through fine-tuning, adapting its parameters to align more precisely with individual gaze behaviors. Subsequent to fine-tuning, a lightweight regression model, specifically support vector regression (SVR), is introduced. This model refines gaze predictions by minimizing the disparity between predicted and actual gaze locations, leveraging the outputs of the fine-tuned ConvNet. During real-world usage, the pre-trained base model and the regression model collaboratively generate personalized gaze estimates. This involves sequential processing of input images through both models, culminating in a precise gaze location determination. Our model's holistic architecture harnesses the power of ConvNet-based analysis of facial features and eye regions, synergizing with fine-tuning and a specialized regression technique. This intricate amalgamation of techniques results in robust and accurate gaze tracking performance, as substantiated by evaluations against participants' genuine gaze behavior.

5. Training

5.1 TensorFlow Implementation

In the process of adapting our model to the requirements of the .tflite format, we initiate with meticulous data preparation, structuring our data into the .tfrecs format, which seamlessly aligns with the model's .tflite input format. This preparatory phase ensures the streamlined and efficient feeding of data into the subsequent pipeline. An essential consideration during the conversion of our trained model into the .tflite version is the potential trade-off between model accuracy and the efficiency of the converted model. However, this challenge is proactively addressed through the implementation of post-training quantization. This technique, akin to Google's approach, optimizes the model to function effectively with reduced numerical precision, effectively sidestepping significant accuracy degradation.

To maintain the model's accuracy throughout this conversion process, we employ a technique called *post-training quantization*, which optimizes the model for reduced numerical precision while mitigating significant accuracy losses. Additionally, our training methodology encompasses a critical aspect of learning rate scheduling. Specifically, we leverage the "Reduce Learning Rate on Plateau" scheduler, which adapts the learning rate dynamically based on training progress. This differs from certain PyTorch versions, highlighting the context-specific nature of learning rate strategies. During the training phase, we utilize the Mean Squared Error (MSE) loss function to quantify the disparity between predicted and actual gaze locations, facilitating the model's accurate learning. The evaluation of our model's performance is executed using the Mean Euclidean Distance (MED) metric, which computes the average distance between predicted and actual gaze positions, offering a comprehensive measure of predictive precision.

Our exploration of various learning rate scheduling approaches led us to a pertinent insight. Among the Exponential Learning Rate, Reduce Learning Rate on Plateau, and no scheduler scenarios, the latter emerged as the most advantageous in terms of model convergence and performance improvement. The distinctiveness of our findings became apparent when they diverged from the outcomes of PyTorch implementations, where alternate learning rate scheduling methods were favored. This observation underscores the intricate interplay between models, datasets, and training paradigms, underscoring the necessity of tailored strategies that align with the specific context. In summation, our approach encompasses meticulous data preparation, proactive measures against accuracy reduction during model conversion, an adaptive learning rate strategy, distinctive loss functions and evaluation metrics, and insights gleaned from experimental disparities. These cumulative efforts are channeled towards the creation of a robust gaze tracking pipeline, analogous to Google's acclaimed approach.

5.2 PyTorch Implementation

In our pursuit of refining the gaze tracking model, we leverage the capabilities of PyTorch Lightning. This framework allows us to seamlessly train our model across multiple GPUs, fully utilizing the available resources. Given the relatively compact size of the network with approximately 140,000 parameters, our training speeds demonstrate commendable efficiency. Within our PyTorch Lightning environment, we experiment with various learning rate scheduling strategies. Of these, the Exponential Learning Rate (LR) scheduler emerges as the most effective. This strategy, aligned with Google's recommendations, dynamically adjusts the learning rate during training, contributing to improved convergence and model performance. Our approach involves a meticulous examination of Google's supplementary details, enriching our understanding of their implementation. With a focus on enhancing the model, we delve deeper into the specifics provided by Google. This comparative analysis allows us to pinpoint potential areas for improvement and experimentation with hyperparameters.

Through this in-depth study, we identify two key changes to augment our model, drawing inspiration from Google's approach:

Epsilon Value: We scrutinize the epsilon value in batch normalization layers. While PyTorch's default epsilon is 1e-5, Google employs 1e-3. Recognizing the impact of this hyperparameter, we adapt the model by setting `nn.BatchNorm2d('fill according to layer', momentum=0.9, eps=0.001)`.

Learning Rate Schedule Parameters: Google's implementation highlights the use of `tf.keras.optimizers.schedules.ExponentialDecay`. With an initial learning rate of 0.016, decay steps of 8000, decay rate of 0.64, and 'staircase' decay type, we reconfigure our optimizer and scheduler as follows:

```
optimizer = torch.optim.Adam(self.parameters(), lr=self.lr, betas=(0.9, 0.999), eps=1e-07)
scheduler = StepLR(optimizer, step_size=8000, gamma=0.64, verbose=True)
```

5.3 Refining Hyperparameters

These alterations, guided by Google's insights, serve as foundational improvements. By aligning our model more closely with their methodology, we position ourselves to observe potential enhancements in performance. Specifically, the changes in epsilon value and the learning rate schedule parameters promise to contribute to our model's precision and convergence. Our approach underscores the significance of meticulous study and iterative refinement. We amalgamate Google's recommendations with our prior implementation, ensuring a robust foundation for our gaze tracking model's advancement. This process showcases the iterative nature of research and underscores the dedication to adopting best practices to achieve optimal outcomes.

6. SVR Personalization after training (*TensorFlow and PyTorch*)

The subsequent objective involved a comparison of Support Vector Regression (SVR) outcomes with the prevailing implementations. Following the path laid out by Google, they extracted a specific output of dimensions (1,4) from the penultimate layer of a multi-layer feed-forward convolutional neural network (CNN) in their pipeline. This extracted output was then utilized, on an individualized basis, to construct a meticulously accurate personalized model. In pursuit of this methodology, the same approach was adopted.

To acquire the output from the penultimate layer, they introduced a mechanism termed a "hook" that intercepted the model's process. Once this mechanism furnished an output of dimensions (x,4) from the penultimate layer, a multi output regressor SVR was engaged. This SVR model, harnessed with the obtained

output, was trained and subsequently applied to the test set of an already trained model. The process of adjusting SVR parameters involved considering specific values: 'rbf' as the kernel, 'C' set to 20, and 'gamma' at 0.6. These choices closely mirrored the parameters described in Google's supplementary materials. To identify the most optimal parameter values, a sweep of the epsilon value in the Multi Output Regressor was conducted within the range of 0.01 to 1000. The process of fitting the SVR model included dividing the test set into two ratios: 70:30 and 2/3:1/3. Through the implementation of 3-fold and 5-fold grid searches, the best parameters for individual cases were identified. These optimal parameters were subsequently employed to train the SVR model.

In the realm of SVR personalization, two distinct versions were explored, both centered around individual cases:

1. The *Google Split* version incorporated individuals from the base model training set into the SVR's training and testing sets. This approach, while susceptible to data leakage, was found to result in minimized errors compared to other strategies.
2. The *MIT Split* version, in contrast, excluded individuals from the base model training set from the SVR's training and testing sets. This alternative was deemed more practical due to its avoidance of data leakage.

Within these versions, further variations emerged, each with its own characteristics:

- In the "Unique Ground Truth values" version, the dataset was divided into training and testing subsets based on distinct ground truth values. This ensured the sets possessed different ground truth values, achieved by selecting a single frame corresponding to each unique ground truth value. Consequently, this version contained 30 frames, matching the number of unique ground truth values.
- The "Random Data points/samples" version entailed random division of the dataset into training and testing subsets. This randomized approach encompassed data points from all screen positions in both sets. However, this approach was susceptible to poor generalization due to the potential for similarity between samples in both sets.

Additionally, an alternative division approach was introduced—the "No Shuffle split." This approach assigned the first 70% of test set points for SVR fitting and the remaining 30% for SVR testing. This approach replicated real-world usage, wherein the SVR is calibrated before user application. The selection of users for SVR evaluation was guided by the highest number of frames across the aforementioned splitting strategies. These users, who had not been part of the base model's training data, were selected for SVR fitting, signifying an assessment on previously unseen data.

6.1 Differential Impact based on Dataset Characteristics

Within the SVR 13 Point Calibration dataset, where training prominently involves ground truth points situated at the screen's periphery, the foundational model's inaccuracies are notably reduced compared to the dataset split into a 70/30 ratio. Nonetheless, when evaluating the SVR-trained model's effectiveness on the latter dataset, a more significant discrepancy between errors prior to and following SVR implementation becomes apparent. This differential underscores a heightened level of individualized optimization.

6.2 Division Based on 13-Point Calibration

The division based on the dataset Calibration of 13 Points for SVR is shown in Figure 6.

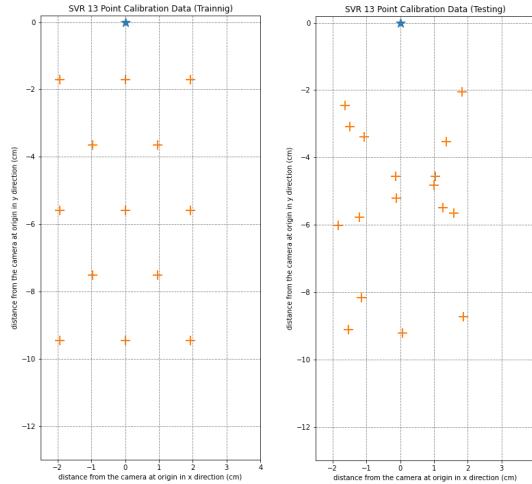


Fig. 6 Division based on SVR 13 Point Calibration dataset

6.3 Google's 70/30 Training and Testing Split Method

In this context, the application of the SVR approach notably boosts the model's performance, leading to a reduction in the overall average error to 1.87cm from 2.03cm. This demonstrates the potential of SVR for better personalization in scenarios with a broader array of gaze positions. Our SVR personalization strategy stands as a promising avenue for enhancing gaze tracking accuracy. The nuanced impact across different datasets underlines the importance of context-aware model adjustments, offering insight into the balance between generalization and personalized refinement.

The division based on the Google Calibration dataset is shown in *Figure 7*.

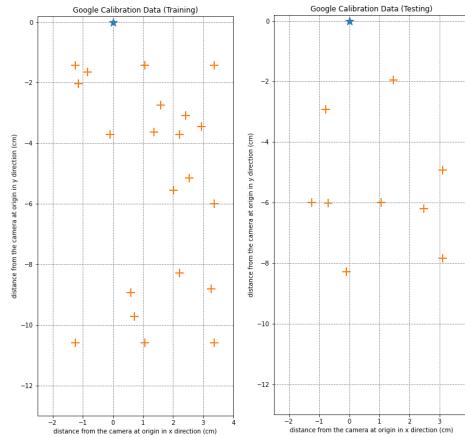


Fig. 7 Division based on Google Calibration dataset

7. Results

7.1 PyTorch Results

We present a comprehensive overview of the performance achieved through our PyTorch-trained models, encapsulating their predictive capacities. While our results exhibit a marginally higher error compared to the figures reported by Google, it is important to note that our data is subjected to fewer constraints. Table 4 furnishes a detailed account of the outcomes derived from our trained models across diverse datasets.

Table 4 Outcomes derived from trained models.

PyTorch Model	Images considered	Non Calibration Test Points count	Non Calibration Points Test Error (cm)	Test files count	Overall Test Error(cm)
Massachusetts Institute of TechnologySplit _GoogleCheckpoint.checkpoint	Portrait Mode, Iphone 6S, MIT Split	7,777	1.75cm	11,178	1.91cm
	Portrait Mode, On Iphone 6, MIT Split	9,683	1.98cm	16,993	2.05cm
	Portrait Mode, On Iphone 5C, MIT Split	2,072	1.70cm	3,643	1.87cm
	Portrait Mode, On Iphone 5S, MIT Split	7,158	1.85cm	12,648	1.95cm
	Portrait Mode, On Iphone 5, MIT Split	1,883	1.61cm	3,362	1.75cm
	Portrait Mode, On Iphone 4S, MIT Split	NA	NA	522	2.04cm
	Portrait Mode, All Phones, MIT Split	33,674	1.92cm	55,541	2.038cm
GoogleSplit_GoogleCheckpoint.checkpoint	Portrait Mode, All Phones, Google Split	NA	NA	83,849	1.86cm

It is evident that our results, while slightly deviating from Google's reported errors, maintain their validity under our experimental settings. These outcomes reflect our comprehensive approach to gaze tracking accuracy through the lens of PyTorch-based models.

7.1.1 PyTorch SVR Personalization Results

Our focus on enhancing gaze tracking precision led us to implement Support Vector Regression (SVR) personalization. By utilizing the result of the second-to-last rectified linear unit layer from the foundational pytorch model, we proceeded to train distinct SVR models to attain personalized outcomes. This approach has shown promising results, although there are instances where its impact isn't uniformly positive. A synopsis of our SVR-based individualization results is provided in *Table 5*.

Table 5 Synopsis of SVR-based individualization results.

Image Data	Test files count	Average Error after SVR(cm)	Average Error of Base Model (cm)	Enhancement (cm)
Splitting Points for SVR into 70/30 Ratio	16,642	1.868	2.029	0.156
Calibration of 13 Points for SVR	33,761	1.787	1.908	0.119

A distinction arises due to the nature of the datasets. The SVR 13 Point Calibration dataset primarily contains ground truth points at the screen's edges. As a result, the base model error is lower compared to the 70/30 split dataset, which tends to have ground truth points more centered. Nonetheless, the SVR benefits from a larger pool of data in the 70/30 split dataset, resulting in a more substantial improvement. This difference suggests enhanced individualization, as the contrast between errors before and after SVR implementation becomes more conspicuous.

7.1.2 Division Based on 13-Point Calibration

The application of SVR leads to a significant reduction in the overall average error to 1.81cm from 1.92cm (see Figure 8). This achievement is particularly remarkable, especially when acknowledging the infrequent instances where the SVR demonstrated performance limitations.

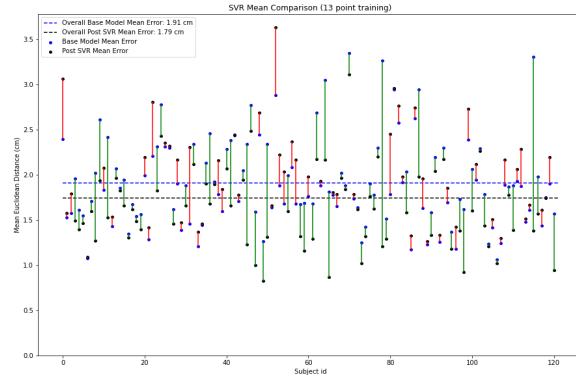


Fig. 8 Reduction in the average overall error in SVR 13-Point Calibration dataset

7.1.3 Google's 70/30 Division

Our integration of SVR leads to a decrease in the overall average error to 1.87cm from 2.03cm, thus reinforcing its capability to improve gaze tracking accuracy across diverse situations (see Figure 9).

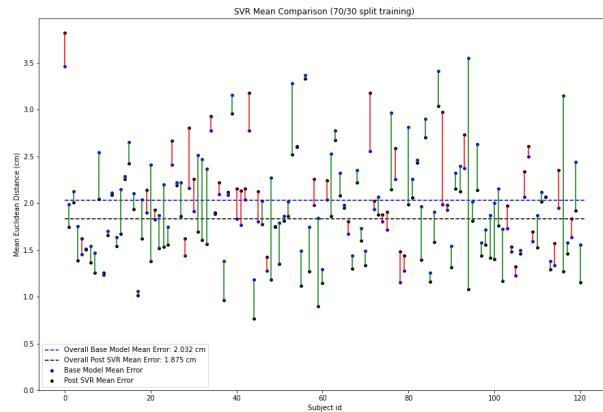


Fig. 9 Reduction in the average overall error in Google's 70/30 Split dataset

7.1.4 PyTorch SVR Output

The visual representations below provide a clear illustration of how our SVR implementation enhances the predictions made by the base model. These visuals effectively highlight the improvements brought about by SVR personalization.

In these graphics:

- The '+' symbols represent the actual ground truth gaze locations.
- The dots signify the predictions made by the network.
- Dots sharing the same color correspond to the '+' symbols of that color.
- Each gaze location corresponds to multiple frames and, consequently, multiple predictions.
- Chromatic classification is employed to associate predictions with their respective ground truth.
- Location of the cam lens is marked by a star.

Moreover, to enhance our grasp of SVR's influence, we plot connecting lines between predictions of network and their associated SVR-informed forecasts for a more holistic perspective. Of particular significance is the visual representation of training the SVR, demonstrated using the calibration of 13 points dataset. Insights obtained from this calibration procedure are then applied to the broader dataset. To facilitate interpretation, a centroid plot is provided, displaying the average of all predictions associated with a particular ground truth. This effectively demonstrates the central point of predictions aligned with that specific ground truth. These visual depictions collectively emphasize the tangible enhancements brought about by our SVR-driven approach, illuminating the intricacies of personalized gaze tracking enhancements.

7.1.5 Instances where SVR impacts is positive:

The instances where SVR leads to positive consequences are as follows (see Figure 10):

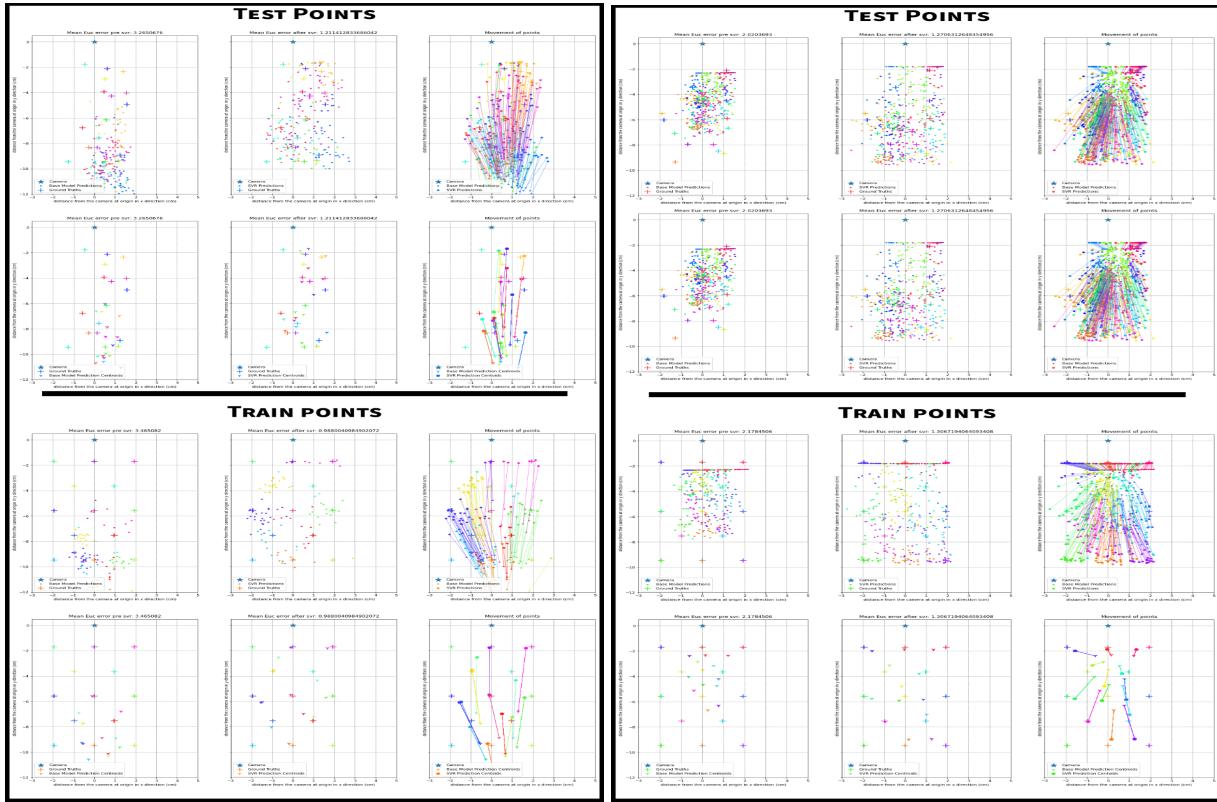


Fig. 10 Instances of positive SVR impact

7.1.6 Instances where SVR's Impact is Negative

The instances where SVR leads to negative consequences are as follows (see Figure 11):

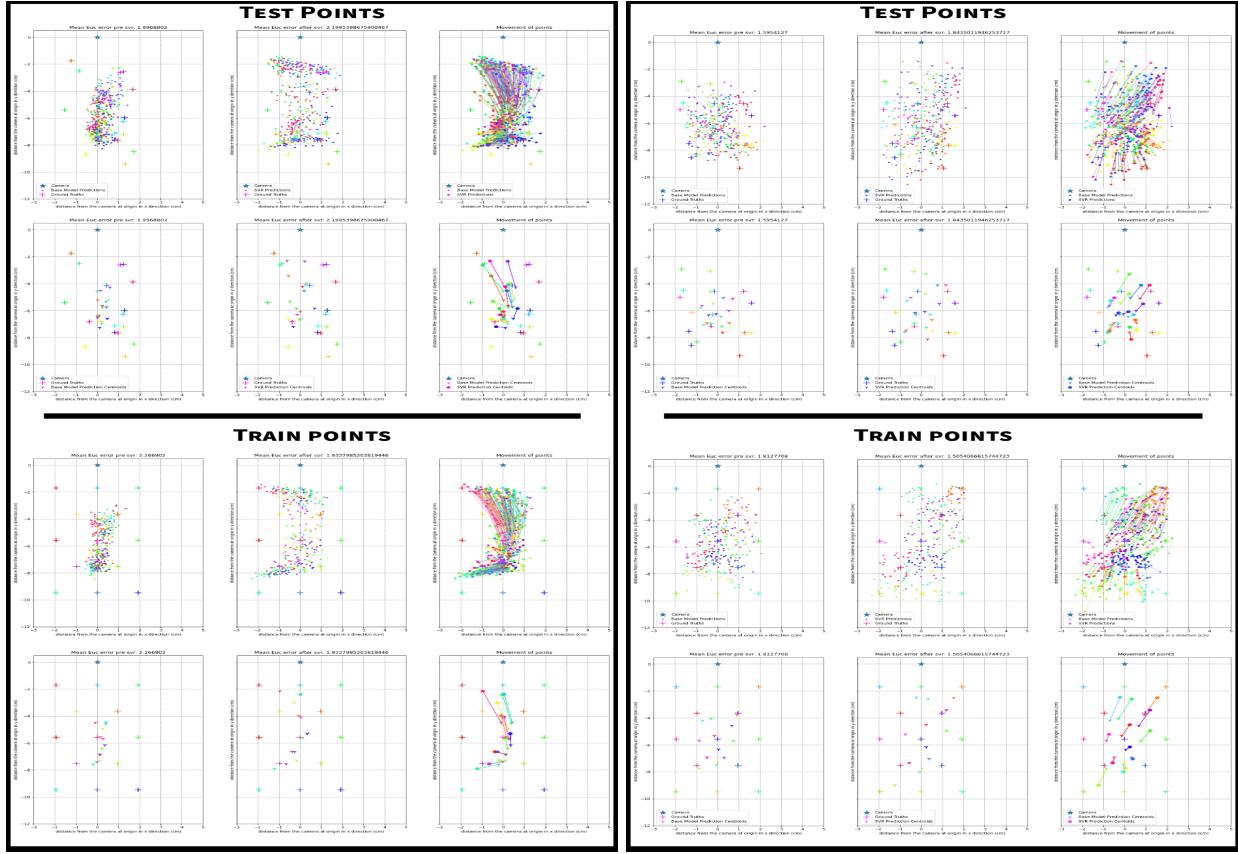


Fig. 11 Instances of negative SVR impact

8. Individualization Using Affine Transformations

Upon closely examining the forecasts generated by the network, it becomes evident that the introduction of a simple affine transform could potentially bolster accuracy. Drawing insights from these observations, we embarked on a journey to incorporate an affine transform-based approach, leveraging predictions generated by the network on frames aligned with the 13-point calibration data. Within the realm of an affine transform, permissible operations encompass shifting, scaling, and rotating.

These transformations present an avenue for refining prediction accuracy. Remarkably, our exploration into affine transforms yielded promising outcomes. The application of such transforms led to a noteworthy reduction in the foundational model's error, reducing it from 1.91cm to 1.859cm (see Figure 12). During this process, improvement may not be as substantial as what we achieved through SVR training, it remains a significant advancement.

It's important to emphasize that the influence of the affine transform tends to be more notable in scenarios where the initial error of the foundational model is relatively higher. This observation underscores the nuanced interplay between the nature of the base model's performance and the effectiveness of the affine transform. In essence, our pursuit of affine transform-based personalization underscores an intriguing avenue for refining gaze tracking accuracy, offering yet another tool in our quest to enhance the precision of our model.

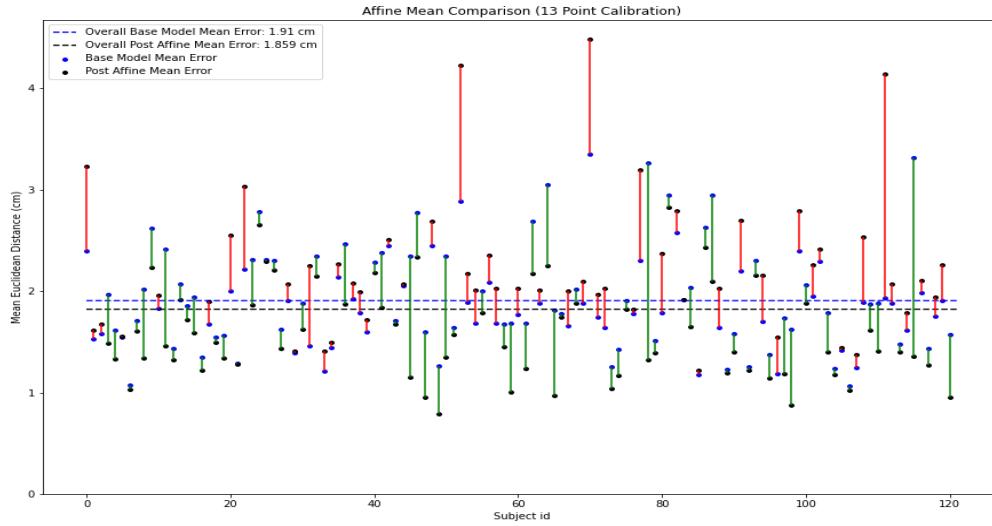
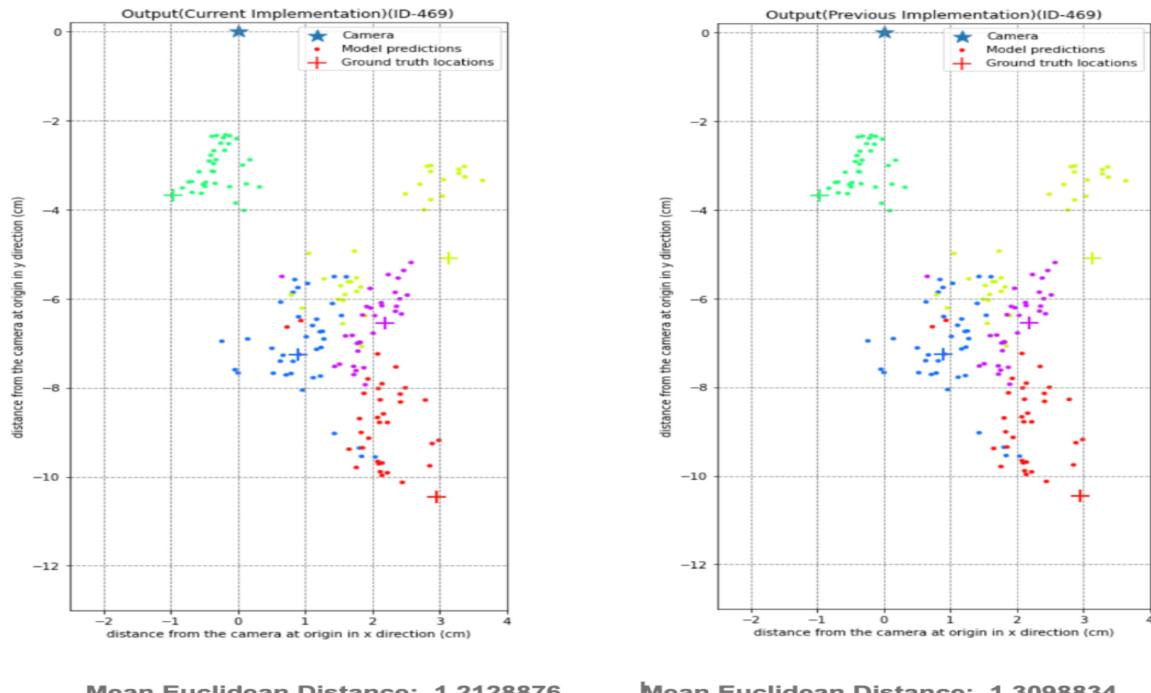


Fig. 12 Reduction in foundational model mean error using Affine transformation

8.1 Upgraded model of PyTorch

Notably, the epsilon (eps) value for Batch Normalization has been increased from $1\text{e-}5$ to $1\text{e-}3$, reflecting a shift towards greater numerical stability in the normalization process. Additionally, in the learning rate scheduling domain, the initial learning rate has been raised to 0.016 , and the decay rate has increased to 0.64 . These changes indicate a more aggressive reduction in learning rates over training epochs, potentially allowing for finer convergence control. It's evident that these increases have been strategically chosen to influence the model's training dynamics and stability, thereby potentially enhancing overall performance and convergence speed (see Figure 13).



(a) and (b)

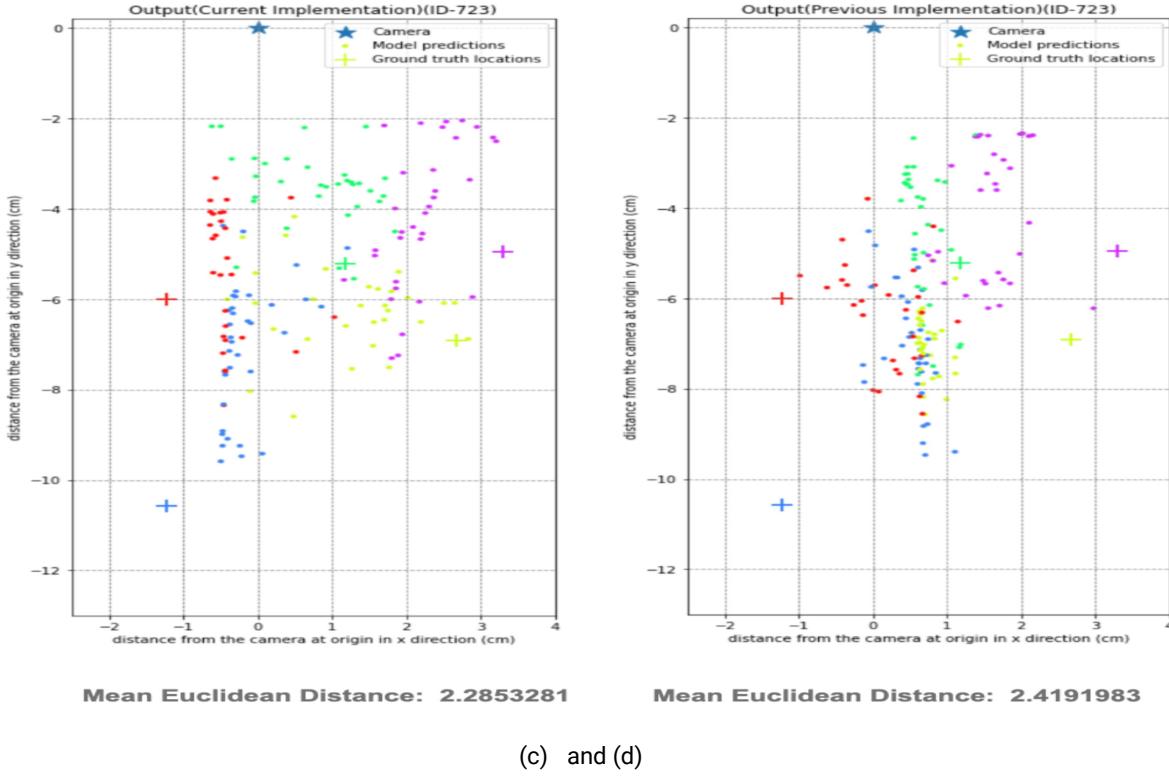


Fig. 13 Outputs of current and previous implementations from two sets: (a) and (b) ID-469; (c) and (d) ID-723

Indeed, examining the two sets of outputs for the current model with adjusted hyperparameters and the previous implementation highlights certain differences in the clustering of results. The reduced clustering observed in the current model's outputs might suggest a broader distribution of predictions, indicating that the model's responses are more spread out across the data range. However, this dispersion may not always be advantageous.

While decreased clustering can imply a better model fit in some cases, it's essential to consider the context and the ultimate goal of the model. A model that exhibits less clustering might have a tendency to produce more diverse predictions, potentially capturing a wider range of scenarios or capturing finer variations in the data. However, this might also lead to predictions that deviate more from the ground truth in some instances, which could be undesirable if accuracy and reliability are primary concerns. Balancing between tightly clustered outputs and more dispersed outputs is a nuanced challenge. It requires a trade-off between precision and generalization. It's crucial to evaluate the model's performance comprehensively, taking into account metrics such as RMSE, accuracy, and the specific application's requirements. If the broader distribution of outputs is leading to predictions that are farther from the ground truth in certain cases, it might be worth further analysis and perhaps tuning the hyperparameters to strike a better balance between precision and generalization.

8.2 SVR Results

The Euclidean Distance (MED) was employed as the metric to assess the model's performance. This distance metric was utilized to measure the dissimilarity between predicted values and actual ground truth values. The formula for calculating the Euclidean Distance involves computing the square root of the sum of squared differences between corresponding elements in the arrays of predicted and actual values. The resulting value indicates the degree of alignment between predictions and ground truth, with smaller distances indicating better agreement and larger distances reflecting greater discrepancies.

Among all the users, the lowest loss² is observed in User ID 2078, measuring at 0.982cm when considering all frames for training. However, when focusing on User ID 2301 and utilizing only 30 unique points for training, the

loss is slightly higher at 0.99cm. Notably, the most significant reduction in loss occurs for User ID 2301, displaying a remarkable decrease of approximately 35%. Calculating the average loss across all users before implementing Support Vector Regression (SVR) with the entire frame dataset yields a value of 1.82cm. After application of SVR, the loss decreases to 1.47cm when considering all frames, and 1.76cm for 30 unique frames. This showcases an overall enhancement of approximately 19% in performance. These results are obtained through the utilization of the upgraded model and its evaluation on the MIT split dataset (see Table 6), as outlined in the preceding section.

Table 6 SVR results as per current and previous PyTorch models and evaluation on MIT split dataset

User ID	Number of frames	Current PyTorch Model			Previous PyTorch Model		
		MED(MIT Split)	After SVR (3 fold) (including all frames)	After SVR (3 fold) (including only 30 unique frames)	MED (MIT split)	After SVR(3 fold) (including all frames)	After SVR (3 fold) (including only 30 unique frames)
3183	874	1.86cm	1.30cm	1.15cm	1.67cm	1.41cm	1.43cm
1877	860	2.09cm	1.23cm	1.19cm	2.08cm	1.35cm	1.40cm
1326	784	1.78cm	1.36cm	2.02cm	1.69cm	1.23cm	1.93cm
3140	783	1.71cm	1.68cm	2.47cm	1.72cm	1.38cm	1.83cm
2091	788	1.86cm	1.77cm	1.87cm	1.72cm	1.54cm	2.17cm
2301	864	1.69cm	1.08cm	0.99cm	1.72cm	1.2cm	1.16cm
2240	801	1.69cm	1.26cm	1.40cm	1.63cm	1.22cm	1.30cm
382	851	2.57cm	2.37cm	3.01cm	2.67cm	2.33cm	3.13cm
2833	796	1.68cm	1.61cm	2.42cm	1.71cm	1.56cm	1.88cm
2078	786	1.23cm	0.98cm	1.11cm	1.22cm	1.03cm	1.38cm

The subsequent step involved training the model on the Google Split dataset, implementing the same parameter adjustments as before. Specifically, the model's parameters remained consistent with the changes previously outlined. For experimentation, a set of 10 individuals was carefully selected based on having the highest number of frames within the train, test, and validation sets. It is likely that this selection process aimed to ensure a robust evaluation of the model's performance by focusing on individuals with substantial representation in the dataset.

The utilization of the Google Split dataset allows for an assessment of how well the model generalizes to new data, particularly emphasizing individuals with varying numbers of frames. On selecting individuals with the maximum number of frames for experimentation, the evaluation becomes more comprehensive, capturing a broader range of scenarios and potential challenges that the model might encounter when faced with real-world data (see Table 7).

Table 7 SVR results as per current and previous PyTorch models and evaluation on Google split dataset

User ID	Number of Frames	Current PyTorch Model(Google Split)	Previous PyTorch Model(Google Split)
503	965	1.32cm	1.50cm
1866	1018	0.99cm	1.16cm
2459	1006	0.97cm	1.15cm
1816	989	0.86cm	1.25cm
3004	983	1.43cm	1.4cm
3253	978	0.94cm	1.25cm
1231	968	1.38cm	1.40cm
2152	957	1.24cm	1.26cm
2051	947	1.15cm	1.38cm
1046	946	1.25cm	1.24cm

Among all the users evaluated, the lowest loss was observed for User ID 2078, registering at 1.03cm when considering all frames for training. Alternatively, for User ID 2301, utilizing only 30 unique points for training, the loss was 1.16cm. Notably, the most significant reduction in loss occurred for User ID 1877, demonstrating a remarkable decrease of approximately 35%. Calculating the average loss across all users prior to applying Support Vector Regression (SVR) using all frames yielded a value of 1.78cm. However, after the application of SVR, the loss decreased to 1.43cm when considering all frames and to 1.76cm when only 30 unique frames were taken into account. This demonstrates an overall enhancement of around 20% in model performance. Comparing the decrease in loss between the two models trained on the MIT Split dataset, both experienced an approximate 20% reduction in error (cm) after SVR application. This signifies the positive impact of the SVR technique on both models' predictive capabilities, resulting in significantly improved performance and indicating its effectiveness in refining the models' predictions.

8.3 TensorFlow Results

When evaluating on the MIT Split dataset, the TensorFlow implementation achieved a performance of 2.03cm, closely aligning with previous PyTorch result of 2.03cm and current PyTorch result of 2.06cm. Similarly, on the Google Split dataset, the TF implementation achieved a performance of 1.80cm, comparable to previous PyTorch performance of 1.86cm and notably improved from current PyTorch performance 1.68cm (see Table 8).

Table 8 TensorFlow implementations compared with current and previous PyTorch models when evaluating on MIT Split and Google Split datasets

Split	TF Implementation	Current PyTorch Model	Previous PyTorch Model
MIT Split	2.03	2.03	2.06
Google Split	1.80	1.86	1.86

These results demonstrate that by following the TensorFlow pipeline, the project achieved results on par with the previous PyTorch implementations. This comparison serves as a useful benchmark for future evaluations. It's especially pertinent when planning to compare the project's TensorFlow Lite (TFLite) version with the TFLite binary provided by Google, enabling a comprehensive assessment of the model's performance in a standardized context. The project's TF model checkpoints are available in the project repository for both the MIT Split and Google Split versions. Additionally, visualizations of gaze predictions were generated for this TensorFlow Implementation. In the visualizations, the ground truth gaze locations are represented by '+' signs, base model predictions by dots, and the mean of base model predictions for each ground truth gaze location by Tri-Down symbols. These visualizations incorporate color coding to correlate predictions with corresponding ground truth, enhancing the interpretability of the model's predictions. Notably, the star (*) symbol indicates the camera position, located at the origin in the visualization. These visualizations provide a clear insight into the model's performance and its relationship with actual gaze locations (see Figure 14)

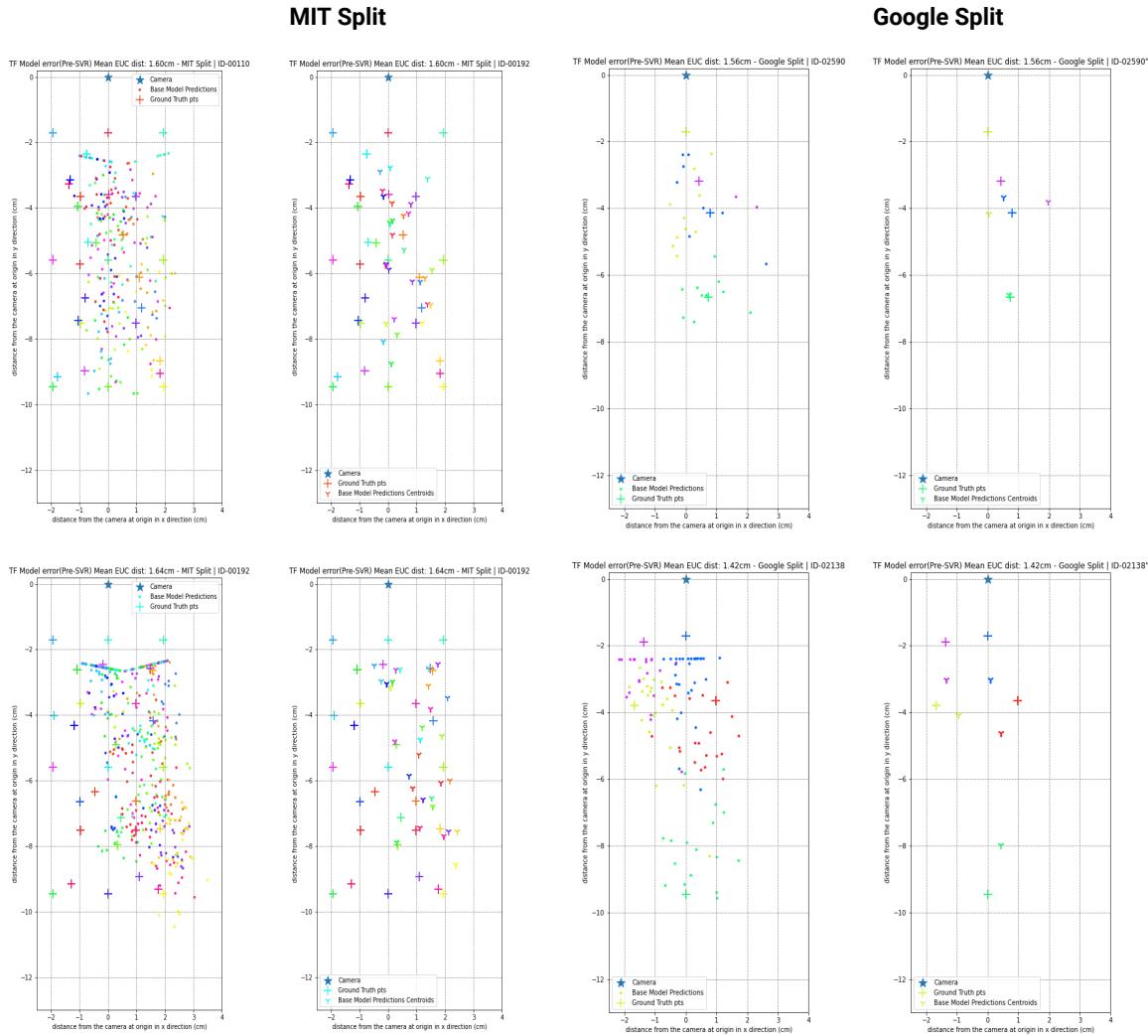


Fig. 14 TF Model errors: MIT Split versus Google Split

8.4 Google Split Version

In this version, individuals from the base model's train set are included in the SVR train/test sets. This approach, while not very practical, could potentially result in data leakage and consequently yield lower errors compared to other splits.

8.5 MIT Split Version

Here, individuals from the base model's train set are excluded from the SVR train/test sets, aligning with a more practical scenario. Within both of these versions, two sub-versions are examined:

8.6 Unique Ground Truth Values

The dataset is split into train and test based on unique ground truth values. This means that each set comprises different ground truth values. Specifically, one frame corresponding to each unique ground truth value is randomly chosen, resulting in a set of 30 frames (reflecting 30 unique ground truth values). This sub-version mirrors real-life scenarios more closely, as it reduces the risk of poor generalization due to similar samples in both train and test sets.

8.7 Random Data Points/Samples

The dataset is randomly split into train and test, potentially leading to data points from all screen positions in both sets. However, there's a risk that this approach could result in similar samples in both train and test sets, hampering generalization. Additionally, a "No Shuffle" split was attempted. In this setup, the first 70% of the test set points are utilized for SVR fitting, and the remaining 30% for SVR testing. This closely simulates the real use-case scenario where SVR calibration precedes subject use. From each of the described splits, 10 users were selected based on having the highest number of frames. These users, unseen by the base model, serve as the data on which SVR is fitted. This personalized SVR approach seeks to enhance the model's performance by tailoring predictions to individual users' characteristics.

8.8 MIT Split Results

The mean results comparison for MIT Split is shown in Table 9.

Table 9 Mean results comparison for MIT split

Implementation	MED
PyTorch Current Model	1.82cm
TensorFlow Model	1.68cm

8.8.1 Post-SVR Results

The random data points/samples (all frames) for MIT Split are shown in Table 10.

Table 10 Random data points/samples (all frames) for MIT Split

Implementation	70 & 30 Split		$\frac{2}{3}$ & $\frac{1}{3}$ Split	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
PyTorch Current Model	1.46cm	-	-	-
TensorFlow Model	1.48cm	1.69cm	1.49cm	1.64cm

The unique ground truth values (30 points) for MIT Split are shown in Table 11.

Table 11 Unique ground truth values (30 points) for MIT Split

Implementation	70 & 30 Split		$\frac{2}{3}$ & $\frac{1}{3}$ Split	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
PyTorch Current Model	1.76cm	-	-	-
TensorFlow Model	1.73cm	1.75cm	1.84cm	1.72cm

8.8.2 Base Model MED versus Post SVR MED

The random data points/samples (all frames) for Base Model MED versus Post-SVR MED are shown in Table 12.

Table 12 Random data points/samples (all frames) for Base Model MED versus Post-SVR MED

Version	70 & 30 Split		$\frac{2}{3}$ & $\frac{1}{3}$ Split	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
TF Base Model MED	1.79cm	1.88cm	1.79cm	1.86cm
Post SVR MED	1.48cm	1.69cm	1.49cm	1.64cm

The unique ground truth values (30 points) for Base Model MED versus Post-SVR MED are shown in Table 13.

Table 13 Unique ground truth values (30 points) for Base Model MED versus Post-SVR MED

Version	70 & 30 Split		$\frac{2}{3}$ & $\frac{1}{3}$ Split	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
TF Base Model MED	1.78cm	1.75cm	1.83cm	2cm
Post SVR MED	1.73cm	1.75cm	1.84cm	1.72cm

8.8.3 Per-Individual Comparison

The random data points/samples (all frames) for per-individual comparison are shown in Table 14.

Table 14 Random data points/samples (all frames) for per-individual comparison

User ID	Number of frames	MED(MIT Split)	SVR -3CV (70 & 30)		SVR -3CV (2/3 & 1/3)	
			Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
3183	874	1.38cm	1.34cm	1.42cm	1.35cm	1.32cm
1877	860	2.03cm	1.28cm	1.13cm	1.32cm	1.09cm
1326	784	1.53cm	1.31cm	1.47cm	1.29cm	1.44cm
3140	783	1.54cm	1.54cm	1.44cm	1.56cm	1.45cm
2091	788	1.70cm	1.80cm	1.98cm	1.81cm	1.92cm
2301	864	1.86cm	1.36cm	1.75cm	1.34cm	1.69cm
2240	801	1.46cm	1.24cm	1.52cm	1.23cm	1.46cm
382	851	2.38cm	2.44cm	2.89cm	2.44cm	2.75cm
2833	796	1.71cm	1.68cm	1.86cm	1.67cm	1.87cm
2078	786	1.24cm	0.82cm	1.42cm	0.83cm	1.37cm

The unique ground truth values (30 points) for per-individual comparison are shown in Table 15.

Table 15 Unique ground truth values (30 points) for per-individual comparison

User ID	SVR -3CV (70 & 30)		SVR -3CV (2/3 & 1/3)	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
3183	1.83cm	0.85cm	1.88cm	1.58cm
1877	1.82cm	1.46cm	1.64cm	1.47cm
1326	2.39cm	2.10cm	2.12cm	2.09cm
3140	1.20cm	1.30cm	1.73cm	1.58cm
2091	1.81cm	1.99cm	1.94cm	1.73cm
2301	1.43cm	1.5cm	1.77cm	1.61cm
2240	1.26cm	1.62cm	1.16cm	1.73cm
382	2.43cm	2.52cm	2.69cm	2.39cm
2833	1.82cm	1.82cm	1.89cm	1.79cm
2078	1.27cm	1.28cm	1.09cm	1.20cm

Firstly, when taking into account all frames, the mean losses are consistently lower compared to the scenario involving unique ground truth values. This disparity can be attributed to the discussed data leakage issue. Additionally, the post-SVR results demonstrate notably reduced mean errors compared to the base model errors. Notably, refraining from shuffling the set during the split leads to increased loss, as it more accurately mirrors real-life situations where users may encounter new ground truth points.

Furthermore, the consideration of frames with unique ground truth values highlights a significant variability in errors per individual, ultimately leading to nearly identical mean errors. This outcome is particularly evident due to the model being trained on a limited 30-frame dataset. Consequently, the SVR struggles to generalize effectively, potentially capturing unwanted features. It's noteworthy that addressing this limitation is part of the ongoing work to refine the model's performance in the future.

8.9 Google Split Results

The mean results comparison for Google Split results is shown in Table 16.

Table 16 Mean results comparison for Google Split results

Implementation	MED
PyTorch Current Model	1.15cm
TensorFlow Model	1.24cm

8.9.1 Base Model MED versus Post-SVR MED

The random data points/samples (all frames) for Base Model MED versus Post-SVR MED are shown in Table 17.

Table 17 Random data points/samples (all frames) for Base Model MED versus Post-SVR MED

Version	70 & 30 Split		$\frac{2}{3}$ & $\frac{1}{3}$ Split	
	Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
TF Base Model MED	1.31cm	1.31cm	1.32cm	1cm
Post SVR MED	1.04cm	1.14cm	1.12cm	1.04cm

8.9.2 Per-Individual Comparison

The random data points/samples (all frames) for per-individual comparison are shown in Table 18.

Table 18 Random data points/samples (all frames) for per-individual comparison

User ID	Number of frames	MED(MIT Split)	SVR -3CV (70 & 30)		SVR -3CV (2/3 & 1/3)	
			Shuffle=True	Shuffle=False	Shuffle=True	Shuffle=False
503	965	1.38cm	1.37cm	1.35cm	1.32cm	1.41cm
1866	1018	1.34cm	0.86cm	1.24cm	1.18cm	0.88cm
2459	1006	1.48cm	0.69cm	0.81cm	0.81cm	0.68cm
1816	989	1.04cm	0.92cm	0.93cm	0.92cm	0.94cm
3004	983	1.22cm	1.18cm	1.07cm	1.05cm	1.16cm
3253	978	1.26cm	0.84cm	1.07cm	0.98cm	0.84cm
1231	968	1.39cm	1.09cm	1.33cm	1.36cm	1.06cm
2152	957	1.28cm	1.36cm	1.28cm	1.27cm	1.38cm
2015	947	1.27cm	1.12cm	1.23cm	1.20cm	1.11cm
1046	946	1.24cm	0.97cm	1.07cm	1.07cm	0.97cm

The contrasting performance between the Google split and the MIT split can be attributed to the nature of their dataset compositions. The Google split includes frames of each individual in both the training and testing sets, contributing to notably lower errors on the dataset of 10 individuals, when compared to the MIT split. This design choice likely plays a significant role in driving Google's reported mean errors to an impressively low level of 0.46 ± 0.03 cm. The utilization of this particular split version potentially provides Google's model with more comprehensive and representative training data, contributing to its superior performance.

9. App

Data collection was facilitated through an Android application, where a systematic approach was adopted. During various instances when circles or dots were displayed on the screen, users' photos were captured at random intervals. The focal point of the circle or dot was meticulously recorded as the corresponding (x, y) coordinate of the user's gaze. These frames were then associated with the specific coordinate based on their corresponding time stamps, ensuring accurate alignment of the data. The process involved the utilization of an Android app that efficiently captured the data in a controlled manner, enabling the collection of valuable information about users' gaze positions.

10. Future Scope and Improvements

Our investigative efforts encompassed several essential aspects. Firstly, we thoroughly assessed the accuracy of our queries to the Google model binary, ensuring precise interaction with the model. Additionally, we delved into comprehending the distinctive patterns of the Support Vector Regression (SVR) across various model versions, aiming to gauge the extent of our model's generalization capability. To enhance our model's performance,

we embarked on training it using the normalization function adopted by Google. Subsequently, the model was subjected to rigorous testing using data obtained from our proprietary app, allowing for a comprehensive comparison with Google's binary model outputs. We also embarked on a comparative analysis with alternative implementations such as iTracker, exploring potential avenues to further enhance the model's efficacy through network expansion. It's important to note that during the process of fitting and evaluating the SVR on the Google split version, instances of data leakage were identified, a concern that will be addressed and rectified in our future work. Our comprehensive approach embraced these critical steps to ensure a robust evaluation of our model's performance and potential for improvement.

11. Code and Data Availability

The complete set of analytical code and de-identified, summarized data tables necessary to reproduce the figures and statistical analyses presented in this manuscript have been made publicly available. You can access these resources on GitHub: PyTorch Model-1 (https://github.com/lcandoit07/PyTorch_Model_1), PyTorch Model-2 (https://github.com/lcandoit07/PyTorch_Model_2), and TensorFlowModel (https://github.com/lcandoit07/TensorFlow_Model). Our study utilized the MIT Gaze Capture Dataset for experimentation (<https://gazecapture.csail.mit.edu/>). It's important to note that we do not possess any personalized data from Google, as such data was not made available in open sources. Instead, we leveraged the MIT Gaze Capture Dataset and introduced personalized data for a subset of random users in our work.

Declarations

Funding

No funds, grants, or other support were received.

Competing Interests

The authors have no financial or non-financial interests to declare.

References

1. Valliappan N, Dai N, Steinberg E, He J, Rogers K, Ramachandran V, Xu P, Shojaeizadeh M, Guo L, Kohlhoff K, Navalpakkam V (2020) Accelerating eye movement research via accurate and affordable smartphone eye tracking. *Nat Commun* 11: 1–12. <https://doi.org/10.1038/s41467-020-18360-5>
2. Kafka K, Khosla A, Kellnhofer P, Kannan H, Bhandarkar S, Matusik W, Torralba A (2016) Eye tracking for everyone. IEEE conference on computer vision and pattern recognition (CVPR).