

NVGaze: An Anatomically-Informed Dataset for Low-Latency, Near-Eye Gaze Estimation

Supplementary Material

Joohwan Kim*
NVIDIA

Michael Stengel*
NVIDIA

Alexander Majercik
NVIDIA

Shalini De Mello
NVIDIA

David Dunn
UNC

Samuli Laine
NVIDIA

Morgan McGuire
NVIDIA

David Luebke
NVIDIA

ACM Reference Format:

Joohwan Kim*, Michael Stengel*, Alexander Majercik, Shalini De Mello, David Dunn, Samuli Laine, Morgan McGuire, and David Luebke. 2019. NVGaze: An Anatomically-Informed Dataset, for Low-Latency, Near-Eye Gaze Estimation, Supplementary Material. In *CHI Conference on Human Factors in Computing Systems Proceedings (CHI 2019)*, May 4–9, 2019, Glasgow, Scotland UK. ACM, New York, NY, USA, 20 pages. <https://doi.org/10.1145/3290605.3300780>

Dataset Publication Page

The NVGaze eye ball model and datasets are available on the project page <https://sites.google.com/nvidia.com/nvgaze>.

Resolution and Network Complexity

Though neural network architectures for gaze estimation vary widely in form and complexity [3, 11, 15], the minimum required network sophistication for accurately estimating a viewer’s gaze has not previously been analyzed. For a camera that is stably fixed with respect to the eye, gaze direction might be a simple function of the position of iris or pupil, easily computed by a simple network. However, blinks, different users, and slippage are common factors that often challenge such simple algorithms and make a more complex mechanism necessary. This section examines to what extent these variations have an effect on the size and complexity of the neural network required for accurate gaze estimation.

Network Architecture and Training. Our template neural network consists of a variable number of convolutional layers, followed by a single fully-connected layer that outputs the

two result values for horizontal and vertical gaze. All convolutional layers use 3×3 kernels and have a stride of 2×2 pixels, i.e., they approximately halve the image resolution at each layer. No pooling or padding were used. Motivated by Laine et al. [9], the output channel count is increased by a factor of 1.5 at each convolutional layer. The convolutional layers use ReLU activation, whereas the final fully-connected layer uses linear activation. Dropout layers [13] with $p = 0.1$ were added after each convolutional layer to prevent overfitting. Finally, when calibrating for multiple subjects, a learned, per-subject affine transformation is performed as a post-processing step.

To evaluate how well a trained network generalizes on a novel subject, we define *generalization error* as the absolute error between the test labels and inferred values transformed according to an optimal affine calibration transform, computed between the set of inferred values and the set of test labels.

When training, the network weights were randomly initialized following He et al. [5], and Adam optimizer [8] was run for 1500 epochs with parameters $\beta_1 = 0.9$, $\beta_2 = 0.99$ and $\epsilon = 10^{-8}$. Every minibatch contained 10 images from each training subject, allowing the minibatch size to vary between different experimental setups, but keeping the number of updates to the network constant. Learning rate was ramped up to $\lambda = 10^{-3}$ during the first 10 training epochs, and ramped down to zero during the last 150 epochs to ensure convergence to a local optimum.

Experimental Conditions. There are three conditions, each of which was characterized by a dataset representing different degrees of input data complexity. The first and simplest dataset consists of data generated using a single synthetic head model. The distance and rotation between the camera and the face are constant for all samples and the eye has no animated blinks. The second and modestly complex dataset was again generated by using a single synthetic model but includes randomized blinks and randomized slippage between the camera and the face for each image sample. The last and most complex dataset was generated using all 10 synthetic

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CHI 2019, May 4–9, 2019, Glasgow, Scotland UK

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-5970-2/19/05...\$15.00

<https://doi.org/10.1145/3290605.3300780>

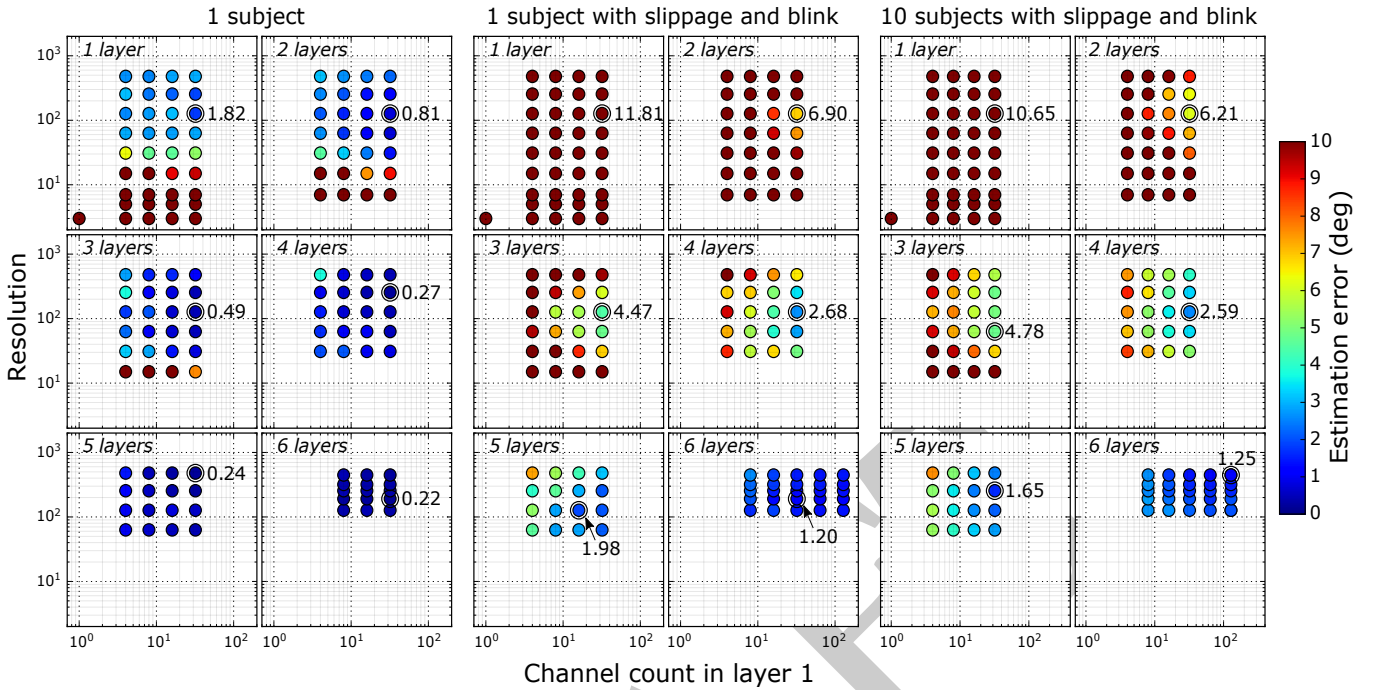


Figure 1: Effect of data and network complexity on gaze estimation accuracy. Using our synthetic data, we explore the effect of complexity in input data on the required amount of information and size of network for accurate gaze estimation. Left, middle, and right columns are, respectively, when input data consist of 1 subject’s images without simulating slippage and blink, 1 subject’s images with slippage and blink simulated, and 10 subjects’ images with slippage and blink simulated. Each plot represents a set of gaze estimation accuracy with a network with certain number of layers. The network yielding the best accuracy in each plot is denoted by a black circle with its accuracy labelled next to it. The unit of estimation error was degree in visual angle.

head models and included blinks and slippage randomization. The slippage value is randomly selected within $[-0.5, 0.5]$ cm around the original camera position in horizontal and vertical directions with uniform distribution. The pupil size is uniformly randomized within the normal range (2 to 8 mm) for all the datasets.

In each experimental condition, we trained approximately 150 networks with different degrees of network complexity and input resolution. The number of convolutional layers varies from 1 to 6 and, independently, the number of output channels of the first convolutional layer is varied. As we grow the number of channels by a fixed factor of 1.5 for each layer, the total number of activations in the network scales linearly with respect to the output channel count of the first layer. Input resolution R is selected among the values that meet the criterion in Eq. 1, which describes the condition where every pixel in the input data contributes to the output for a given stride size s , number of layers l , kernel size k , and an arbitrary positive integer N .

$$R = s^l \times N + (k - s) \frac{s^l - 1}{s - 1} \quad (1)$$

Experimental Results. Fig. 1 summarizes the results of our investigation. Note that the supported minimum resolution of our network architecture grows with the number of used convolutional layers due to the criterion of Eq. 1. As expected, rather shallow network architectures are sufficient for the simplest case (Fig. 1, left) to perform gaze estimation with a surprisingly high degree of accuracy. Even a single-layer network can estimate gaze with less than 2 degrees angular error. A two-layer network brings the error down to below 1 deg. Increasing the number of convolutional layers improves accuracy, but the improvement is marginal after 4 layers. Interestingly, higher input resolution was not helpful for the simple single subject case.

The middle plots in Fig. 1 show the results for modestly complex dataset (1 subject with slippage and blink). Note the dramatic increase in the estimation error in shallow networks. In this case, we need at least 5 convolutional layers in order to get down below 2 deg of estimation error. Such complex network architecture is probably required for dealing with the complexity in input induced by the simulated slippage and blink. As before, higher resolution input does not guarantee improvement in estimation accuracy.

Finally, the right plots in Fig. 1 show the result from the most complex dataset (10 subjects with slippage and blink). Note that we report generalization error in this condition. As for a single subject with slippage and blink, we clearly need a more complex architecture for an accurate gaze estimation. The best network configurations are quite consistent as in the single subject condition when the number of convolutional layers are small. However, deeper networks benefit from high input resolutions more than in the previous test.

We conclude from our investigation that a multi-layer convolutional network with high feature count is necessary to become robust to realistic challenges such as slippage and blinks. A shallow network architecture can accurately perform gaze estimation only when there is not much variation in input (e.g. no slippage or blinks). More channels in each convolutional layer generally improve results, suggesting a trade-off between performance and computational resources. However, increasing the input resolution does not guarantee improvement in accuracy; we advise choosing the optimal resolution through experiments. We hope our synthetic dataset is helpful in executing such experiments.

A good reference for computational complexity analysis in terms of FLOPS and milliwatts can be found in Zemblys et al. [18].

Synthetic Dataset Evaluation

In addition to the evaluation in the main manuscript, we additionally compared our network trained on the UnityEyes model of Wood et al.’s [16] with rasterized images and path traced images generated from our anatomically-informed extensions to the SynthesEyes model of Wood et al. [17], holding both resolution and number of images (1M) constant. In this experiment, path tracing the images and improved anatomical accuracy decreased the generalization error from 3.7° to 3.1° when validated on real images. This result provides additional evidence of the superior performance of our proposed synthetic model for the task of near-eye gaze estimation under IR lighting.

Pupil Localization: Training

The network architecture is equivalent to the gaze estimation experiment, except that we use 7 convolutional layers as shown below.

Layer index	1	2	3	4	5	6	7
Kernel size	9×9	7×7	5×5	5×5	3×3	3×3	3×3
Output channels	24	36	52	80	124	256	512

We perform various augmentation steps during training as we did for the gaze estimation network, making a subset of our synthetic data sufficient for convergence of our network. We use a resolution of 293x293 as input for pupil location

estimation and always rescale the input to this resolution using bicubic filtering.

We train on 80% of all dataset samples for 150 epochs using Adam [8] with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$ and use the remaining samples for testing. The learning rate is kept at $\lambda = 10^{-4}$ until a fairly long ramp-down to zero during the last 50 epochs. There is no learning rate ramp-up.

Our augmentations are as follows:

First, we transform the image and label using a random affine transformation (translation up to 50 pixels, rotation about origin up to 10°, rescale with a randomized factor within the range [0.9,1.5]) while making sure that pupil location is within a reasonable range of [0.1,0.9] of horizontal and vertical image size. Second, we apply pixel-wise intensity noise with a maximum offset of ± 10 (probability $p=0.5$). Third, we apply a global intensity offset of ± 40 ($p=0.5$). Fourth, we apply Gaussian filtering with $\sigma \in [0.6,1.6]$ and kernel size of 7 ($p=0.5$). Fifth, using bicubic filtering we apply random shrinking with a scale factor $s \in [0.25,1.0]$ followed by upscaling again to the full input resolution ($p=0.5$). Sixth, to simulate environment reflections in the eye we randomly overlay the image with images out of 326 natural photographs from the dataset published in [7] ($p=0.25$). Before blending, the typically larger overlay image B is converted to grayscale and randomly cropped the eye image resolution. For eye image E and a randomly chosen opacity value $o \in [0.1, 0.2]$ we use a soft blending function $(1 - ((1 - E) * (1 - B))) * o + E * (1 - o)$ and clip the resulting pixel intensities to a maximum of 255. Seventh, we apply histogram equalization implemented in OpenCV [1]. As a last step, we rescale the image intensities to a range of [-1,1] and randomly shift the mean about ± 0.15 and the min and max intensities about ± 0.1 . During inference of test images only histogram equalization and normalization to [-1,1] are applied.

Pupil Localization: Comparison to previous methods

We compare performance of our network for pupil center detection against Park et al. (Fig.6 on the PupilNet datasets from Fuhl et al. (Fig.2,3,4,5; Table 1,2).

Our comparison shows performance of each approach as a percentage of samples within a give pixel range of values from 0.1–15, increasing by 0.1 pixels at each step. To measure Park et al’s approach, we used their pretrained network with an input resolution of 180x108. To make the most charitable comparison, we cropped images from the PupilNet dataset in a 180x108 crop centered on the pupil position label and clipped to the edges of the image. This approach provides the highest resolution possible and the clearest pupil image, thus ensuring the best performance of Park et al.’s method on the PupilNet dataset.

REFERENCES

- [1] Ivan Culjak, David Abram, Tomislav Pribanic, Hrvoje Dzapo, and Mario Cifrek. 2012. A brief introduction to OpenCV. In *MIPRO, 2012 proceedings of the 35th international convention*. IEEE, 1725–1730.
- [2] Wolfgang Fuhl, Thomas Kübler, Katrin Sippel, Wolfgang Rosenstiel, and Enkelejda Kasneci. 2015. Excuse: Robust pupil detection in real-world scenarios. In *International Conference on Computer Analysis of Images and Patterns*. Springer, 39–51.
- [3] Wolfgang Fuhl, Thiago Santini, Gjergji Kasneci, Wolfgang Rosenstiel, and Enkelejda Kasneci. 2017. PupilNet v2.0: Convolutional Neural Networks for CPU based real time Robust Pupil Detection. *CoRR* abs/1711.00112 (2017). arXiv:1711.00112 <http://arxiv.org/abs/1711.00112>
- [4] Wolfgang Fuhl, Thiago C Santini, Thomas Kübler, and Enkelejda Kasneci. 2016. Else: Ellipse selection for robust pupil detection in real-world environments. In *Proceedings of the Ninth Biennial ACM Symposium on Eye Tracking Research & Applications*. ACM, 123–130.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. *CoRR* abs/1502.01852 (2015).
- [6] Amir-Homayoun Javadi, Zahra Hakimi, Morteza Barati, Vincent Walsh, and Lili Tcheang. 2015. SET: a pupil detection method using sinusoidal approximation. *Frontiers in neuroengineering* 8 (2015), 4.
- [7] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2008. Hamming embedding and weak geometric consistency for large scale image search. In *European conference on computer vision*. Springer, 304–317.
- [8] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. *CoRR* abs/1412.6980 (2014).
- [9] Samuli Laine, Tero Karras, Timo Aila, Antti Herva, Shunsuke Saito, Ronald Yu, Hao Li, and Jaakko Lehtinen. 2017. Production-Level Facial Performance Capture Using Deep Convolutional Neural Networks. *Proc. Symposium on Computer Animation (SCA)*.
- [10] Dongheng Li, David Winfield, and Derrick J Parkhurst. 2005. Starburst: A hybrid algorithm for video-based eye tracking combining feature-based and model-based approaches. In *null*. IEEE, 79.
- [11] Seonwook Park, Adrian Spurr, and Otmar Hilliges. 2018. Deep Pictorial Gaze Estimation. *arXiv preprint arXiv:1807.10002* (2018).
- [12] Seonwook Park, Xucong Zhang, Andreas Bulling, and Otmar Hilliges. 2018. Learning to find eye region landmarks for remote gaze estimation in unconstrained settings. *arXiv preprint arXiv:1805.04771* (2018).
- [13] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* 15 (2014), 1929–1958.
- [14] Lech Świrski and Neil A. Dodgson. 2014. Rendering synthetic ground truth images for eye tracker evaluation. In *Proceedings of ETRA 2014*. 219–222. <http://www.cl.cam.ac.uk/research/rainbow/projects/eyerender/>
- [15] Marc Tonsen, Julian Steil, Yusuke Sugano, and Andreas Bulling. 2017. InvisibleEye: Mobile Eye Tracking Using Multiple Low-Resolution Cameras and Learning-Based Gaze Estimation. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 3, Article 106 (Sept. 2017), 21 pages. <https://doi.org/10.1145/3130971>
- [16] Erroll Wood, Tadas Baltrusaitis, Louis-Philippe Morency, Peter Robinson, and Andreas Bulling. 2016. Learning an appearance-based gaze estimator from one million synthesised images. In *ETRA*.
- [17] Erroll Wood, Tadas Baltrusaitis, Xucong Zhang, Yusuke Sugano, Peter Robinson, and Andreas Bulling. 2015. Rendering of Eyes for Eye-Shape Registration and Gaze Estimation. In *Proc. of the IEEE International Conference on Computer Vision (ICCV 2015)* (2015-12-12).
- [18] Raimondas Zemblys and Oleg Komogortsev. 2018. Making stand-alone PS-OG technology tolerant to the equipment shifts. In *Proceedings of the 7th Workshop on Pervasive Eye Tracking and Mobile Eye-Based Interaction*. ACM, 2.

Dataset	ElSe[4]	ExCuSe[2]	Fuhl et al.[3]	Fuhl et al.[3]	Fuhl et al.[3]	Park et al.[12]	Ours
			SK_8P_8	$F_{CK_XP_Y}$	$F_{SK_XP_Y}$		
I	0.86	0.72	0.77	0.78	0.82	0.35	0.76
II	0.65	0.40	0.80	0.79	0.79	0.65	0.78
III	0.64	0.38	0.62	0.60	0.66	0.40	0.87
IV	0.83	0.80	0.90	0.90	0.92	0.72	0.92
V	0.85	0.76	0.91	0.89	0.92	0.42	0.95
VI	0.78	0.60	0.73	0.78	0.79	0.60	0.94
VII	0.60	0.49	0.73	0.80	0.73	0.39	0.78
VIII	0.68	0.55	0.84	0.83	0.81	0.47	0.91
IX	0.87	0.76	0.86	0.86	0.86	0.41	0.76
X	0.79	0.79	0.80	0.78	0.81	0.53	0.88
XI	0.75	0.58	0.85	0.74	0.91	0.77	0.73
XII	0.79	0.80	0.87	0.85	0.85	0.61	0.85
XIII	0.74	0.69	0.79	0.81	0.83	0.66	0.74
XIV	0.84	0.68	0.91	0.94	0.95	0.53	0.93
XV	0.57	0.56	0.81	0.71	0.81	0.35	0.88
XVI	0.60	0.35	0.80	0.72	0.80	0.65	0.86
XVII	0.90	0.79	0.99	0.87	0.97	0.66	0.96
XVIII	0.57	0.24	0.55	0.44	0.62	0.08	0.85
XIX	0.33	0.23	0.34	0.20	0.37	0.19	0.68
XX	0.78	0.58	0.79	0.73	0.79	0.31	0.85
XXI	0.47	0.52	0.81	0.67	0.83	0.25	0.91
XXII	0.53	0.26	0.50	0.52	0.58	0.30	0.82
XXIII	0.94	0.93	0.86	0.87	0.90	0.87	0.88
XXIV	0.53	0.46	0.46	0.55	0.55	0.30	0.71
New I	0.62	0.22	0.69	0.56	0.69	0.31	0.69
New II	0.26	0.16	0.44	0.35	0.45	0.08	0.69
New III	0.39	0.34	0.45	0.44	0.49	0.21	0.74
New IV	0.54	0.48	0.83	0.77	0.82	0.34	0.92
New V	0.75	0.59	0.78	0.76	0.81	0.25	0.88
Average	0.67	0.54	0.74	0.71	0.76	0.44	0.83

Table 1: Five pixel error on PupilNet dataset. We estimate the probability of finding the pupil within a distance of max. 5 pixels with respect to the ground truth label of the 384x288 input images. Values are estimated using our 293x293 network images rescaled to 293x293. Images for Park et al.[12] are cropped to 180x108 as described in the text. Other values are directly copied from original PupilNet paper by Fuhl et al. [3].

Pixel Error	Starburst [10] 384x288	SET [6] 384x288	Swirski [14] 384x288	ExCuSe [2] 384x288	ElSe [4] 384x288	Fuhl et al. [3] $F_{CK \times P_Y}$ 384x288	Fuhl et al. [3] $F_{SK \times P_Y}$ 384x288	Park et al. [3] cropped 180x108	Ours 384x288	Ours 293x293
0	0	0	0	0	0.01	0	0.01	0	0	0
1	0.007	0.0278	0.0528	0.108	0.1306	0.105	0.1287	0.01284	0.1040	0.1296
2	0.03	0.0853	0.1548	0.2959	0.3693	0.2945	0.3477	0.09298	0.3295	0.3963
3	0.06	0.148	0.2343	0.4376	0.541	0.499	0.5564	0.19988	0.5565	0.6304
4	0.09	0.1849	0.2757	0.5045	0.6236	0.6279	0.6856	0.35068	0.7259	0.7876
5	0.115	0.2068	0.3012	0.5418	0.6709	0.7059	0.7671	0.43719	0.8308	0.8736
6	0.135	0.2159	0.318	0.5656	0.6973	0.7481	0.7991	0.48127	0.8910	0.9218
7	0.151	0.2229	0.3348	0.5824	0.717	0.7709	0.8225	0.52939	0.9282	0.9532
8	0.162	0.2279	0.3463	0.5933	0.7307	0.7857	0.8367	0.59343	0.9528	0.9719
9	0.173	0.231	0.3599	0.6034	0.742	0.7961	0.8465	0.67657	0.9708	0.9848
10	0.1819	0.2324	0.3694	0.6113	0.7502	0.8028	0.854	0.71770	0.9829	0.9931
11	0.1891	0.2332	0.3794	0.6187	0.756	0.8052	0.8589	0.73676	0.9906	0.9952
12	0.1965	0.2349	0.3888	0.6248	0.7606	0.8067	0.8624	0.75653	0.9950	0.9963
13	0.205	0.2368	0.4005	0.632	0.7637	0.8101	0.8652	0.78478	0.9972	0.9971
14	0.21	0.237	0.4093	0.6384	0.7689	0.8143	0.8679	0.81807	0.9975	0.9977
15	0.215	0.24	0.4171	0.6406	0.7738	0.8167	0.8711	0.83605	0.9978	0.9979

Table 2: Average Pupil Tracking Error on PupilNet dataset. The average pixel error distribution is given for our network and other pupil trackers. Error is computed with respect to the given input resolution. PupilNet images have native size of 384x288. For fair comparison to Park et al. we crop the image to the network input size of 180x108. Crop is performed as centric as possible to the ground truth pupil location. No rescaling of image is performed. For our network we rescale the image to 293x293 and compute the error with respect to 384x288 as well as 293x293. The results show that our network is performing better even with respect to 384x288 although inference is only using 293x293 input.

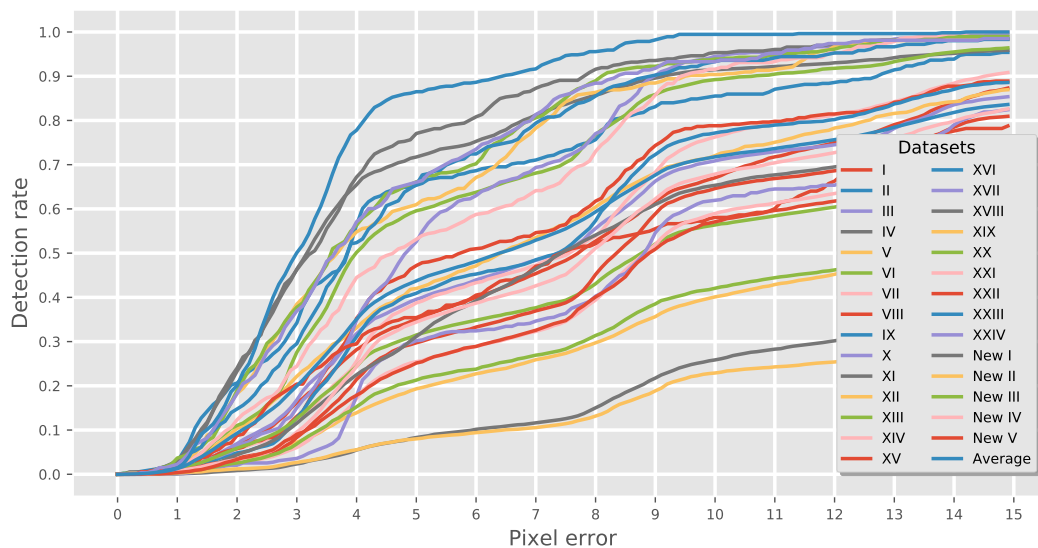


Figure 6: Park et al. Average Pupil/Iris Tracking Error on PupilNet dataset with respect to dataset resolution. Results shown for 180x108 network of Park et al. for Pupil/Iris Center Tracking. Please consider explanation in Table 2 for details on the shown values.

DRAFT

NVGaze: An Anatomically-Informed Dataset for Low-Latency, Near-Eye Gaze Estimation

CHI 2019, May 4-9, 2019, Glasgow, Scotland UK

Pixel error	I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII	XIII	XIV	XV	XVI	XVII	XVIII	XIX	XX	XXI	XXII	XXIII	XXIV	XXV	New I	New II	New III	New IV	New V	Average																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																									
7.1	0.918966	0.903509	0.906251	0.906813	0.908336	0.918710	0.924353	0.935376	0.947872	0.962357	0.979245	0.997692	1.027024	1.064881	1.112034	1.167579	1.230471	1.300824	1.378485	1.463638	1.556084	1.655926	1.772283	1.905266	2.054747	2.220615	2.403996	2.605188	2.834799	3.093179	3.379179	3.693478	4.037111	4.410858	4.815608	5.251447	5.719447	6.220858	6.756884	7.328447	7.936884	8.583147	9.268447	10.003884	10.789447	11.626147	12.514147	13.454447	14.447884	15.495447	16.608447	17.788447	19.036884	20.355447	21.745447	23.208447	24.747884	26.364884	28.061447	29.839447	31.699884	33.645447	35.668447	37.771447	39.956884	42.214884	44.556447	46.971447	49.463884	52.034884	54.686447	57.421447	60.232884	63.114884	66.070447	69.103884	72.218447	75.417884	78.704884	82.081447	85.549884	89.111447	92.768447	96.523884	100.378447	104.343884	108.461447	112.732884	117.159884	121.743884	126.386447	131.188447	136.150447	141.273884	146.558447	152.005884	157.616447	163.391447	169.330447	175.434884	181.704884	188.240447	194.943884	201.815884	208.856447	216.065884	223.443884	230.990447	238.704884	246.586447	254.635884	262.853884	271.240447	279.796447	288.521447	297.414884	306.476447	315.705884	325.103884	334.570447	343.605884	352.809447	362.181447	371.721447	381.329447	391.005884	400.850447	410.863884	421.045884	431.396447	441.916447	452.505884	463.263884	474.191447	485.288447	496.554884	507.990447	519.595884	531.370447	543.314884	555.428447	567.711447	580.263884	593.085884	606.177447	619.540447	633.174884	647.079447	661.254884	675.696447	690.403884	705.376447	720.615884	736.120447	751.892447	767.931447	784.237447	800.810447	817.550447	834.554884	851.822447	869.353884	887.147447	905.203884	923.425884	941.813884	960.367447	979.086447	998.070447	1017.319447	1036.833884	1056.612447	1076.656447	1096.965447	1117.539447	1138.378447	1159.482447	1180.850447	1202.482447	1224.378447	1246.538447	1268.961447	1291.647447	1314.591447	1337.792447	1361.250447	1384.965447	1408.937447	1433.166447	1457.653447	1482.397447	1507.397447	1532.653447	1558.165447	1583.933447	1609.957447	1636.237447	1662.771447	1689.560447	1716.604447	1743.903447	1771.457447	1799.266447	1827.330447	1855.648447	1884.221447	1913.049447	1942.132447	1971.470447	2001.068447	2030.926447	2061.044447	2091.412447	2122.030447	2152.908447	2184.046447	2215.444447	2247.092447	2278.990447	2311.138447	2343.536447	2376.184447	2409.082447	2442.228447	2475.624447	2509.269447	2543.164447	2577.308447	2611.701447	2646.343447	2681.234447	2716.374447	2751.762447	2787.397447	2823.280447	2859.411447	2895.789447	2932.414447	2969.285447	3006.402447	3043.765447	3081.374447	3119.228447	3157.326447	3195.668447	3234.254447	3273.084447	3312.156447	3351.470447	3391.026447	3430.824447	3470.864447	3511.146447	3551.670447	3592.436447	3633.444447	3674.692447	3716.180447	3757.908447	3799.876447	3842.084447	3884.532447	3927.220447	3970.148447	4013.316447	4056.724447	4100.372447	4144.259447	4188.385447	4232.751447	4277.356447	4322.199447	4367.280447	4412.598447	4458.153447	4503.944447	4549.971447	4596.243447	4642.760447	4689.522447	4736.529447	4783.775447	4831.260447	4878.983447	4926.944447	4975.143447	5023.580447	5072.253447	5121.162447	5170.306447	5219.684447	5269.306447	5319.171447	5369.279447	5419.629447	5470.221447	5521.054447	5572.128447	5623.451447	5675.023447	5726.844447	5778.914447	5831.233447	5883.800447	5936.614447	5989.675447	6042.983447	6096.536447	6150.343447	6204.403447	6258.715447	6313.278447	6368.092447	6423.156447	6478.468447	6534.024447	6589.833447	6645.894447	6702.206447	6758.768447	6815.580447	6872.641447	6929.951447	6987.510447	7045.317447	7103.372447	7161.684447	7220.251447	7279.072447	7338.147447	7397.475447	7457.056447	7516.889447	7576.973447	7637.308447	7697.894447	7758.731447	7819.818447	7881.155447	7942.742447	8004.579447	8066.666447	8128.903447	8191.290447	8253.927447	8316.814447	8379.951447	8443.338447	8506.974447	8570.860447	8634.996447	8699.382447	8764.018447	8828.904447	8894.040447	8959.426447	9025.062447	9090.948447	9157.084447	9223.470447	9290.106447	9356.992447	9424.128447	9491.514447	9559.150447	9627.036447	9695.172447	9763.558447	9832.194447	9901.080447	9970.216447	10059.602447	10149.238447	10239.124447	10329.260447	10419.646447	10510.282447	10601.168447	10692.304447	10783.690447	10875.326447	10967.212447	11059.348447	11151.734447	11244.370447	11337.256447	11430.392447	11523.778447	11617.414447	11711.300447	11805.436447	11899.822447	11994.458447	12089.344447	12184.480447	12279.866447	12375.502447	12471.388447	12567.524447	12663.910447	12760.546447	12857.432447	12954.568447	13051.954447	13149.590447	13247.476447	13345.612447	13443.998447	13542.634447	13641.520447	13740.656447	13839.942447	13939.378447	14039.064447	14138.900447	14238.886447	14339.022447	14439.308447	14539.744447	14640.330447	14741.066447	14841.952447	14942.988447	15044.174447	15145.510447	15247.096447	15348.932447	15450.918447	15553.054447	15655.340447	15757.776447	15860.362447	15963.098447	16065.984447	16169.020447	16272.206447	16375.546447	16479.039447	16582.684447	16686.489447	16790.445447	16894.552447	16998.810447	17103.219447	17207.779447	17312.490447	17417.351447	17522.362447	17627.523447	17732.834447	17838.295447	17943.906447	18049.667447	18155.578447	18261.639447	18367.850447	18474.211447	18580.722447	18687.383447	18794.194447	18901.155447	19008.266447	19115.527447	19222.938447	19330.499447	19438.210447	19546.071447	19654.082447	19762.243447	19870.554447	19979.015447	20087.626447	20196.387447	20305.298447	20414.359447	20523.570447	20632.931447	20742.442447	20852.103447	20961.914447	21071.874447	21181.985447	21292.246447	21402.657447	21513.218447	21623.929447	21734.790447	21845.801447	21956.962447	22068.273447	22179.734447	22291.345447	22403.106447	22515.017447	22627.078447	22739.290447	22851.651447	22964.162447	23076.823447	23189.634447	23302.595447	23415.706447	23528.967447	23642.378447	23755.939447	23869.650447	23983.511447	24097.522447	24211.683447	24325.994447	24440.455447	24555.066447	24669.827447	24784.738447	24899.799447	25014.910447	25130.171447	25245.582447	25361.143447	25476.854447	25592.715447	25708.726447	25824.887447	25941.198447	26057.659447	26174.270447	26291.031447	26407.942447	26525.003447	26642.214447	26759.575447	26877.086447	26994.747447	27112.558447	27230.519447	27348.630447	27466.891447	27585.302447	27703.863447	27822.574447	27941.435447	28060.446447	28179.607447	28298.918447	28418.379447	28537.990447	28657.751447	28777.662447	28897.723447	29017.934447	29138.295447	29258.806447	29379.467447	29499.278447	29619.239447	29739.350447	29859.611447	29979.922447	30099.383447	30219.994447	30340.755447	30461.666447	30582.727447	30703.938447	30825.299447	30946.810447	31068.471447	31190.282447	31312.243447	31434.354447	31556.615447	31679.026447	31801.587447	31924.298447	32047.159447	32170.170447	32293.331447	32416.642447	32540.103447	32663.714447	32787.475447	32911.386447	33035.447447	33159.658447	33284.020447	33408.531447	33533.192447	33657.913447	33782.794447	33907.835447	34033.036447	34158.397447	34283.918447	34409.599447	34535.440447	34661.441447	34787.592447	34913.893447	35040.344447	35166.945447	35293.696447	35420.597447	35547.648447	35674.849447	35802.100447	35929.541447	36057.172447	36184.993447	36312.914447	36441.035447	36569.356447	36697.877447	36826.598447	36955.519447	37084.640447	37213.861447	37343.282447	37472.903447	37602.724447	37732.745447	37862.966447	37993.387447	38123.908447	38254.529447	38385.250447	38516.071447	38647.082447	38778.243447	38909.554447	39040.915447	39172.426447	39304.087447	39435.998447	39568.159447	39699.570447	39831.131447	39962.842447	40094.703447	40226.814447	40359.075447	40491.486447	40624.047447	40756.758447	40889.619447	41022.630447	41155.791447	41289.102447	41422.663447	41556.374447	41690.235447	41824.246447	41958.407447	42092.718447	42227.179447	42361.790447	42496.551447	42631.462447	42766.523447	42901.734447	43037.095447	43172.606447	43308.267447	43444.078447	43580.039447	43716.150447	43852.411447	43988.822447	44125.383447	44262.094447	44398.955447	44535.966447	44673.127447	44810.438447	44947.899447	45085.510447	45223.271447	45361.182447	45499.243447	45637.454447	45775.815447	45914.326447	46052.987447	46191.800447	46330.761447	46470.872447	46611.133447	46751.544447	46892.105447	47032.816447	47174.677447	47316.688447	47458.849447	47601.160447	47743.621447	47886.232447	48028.993447	48171.904447	48314.965447	48458.176447	48601.537447	48745.058447	48888.729447	49032.550447	49176.571447	49320.792447	49465.213447	49609.834447	49754.555447	49899.376447	50044.297447	50189.418447	50334.739447	50480.260447	50625.981447	50771.902447	50918.023447	51064.344447	51210.865447	51357.586447	51504.507447	51651.628447	51798.949447	51946.470447	52094.191447	52242.112447	52390.233447	52538.554447	52687.075447	52835.796447	52984.717447	53133.838447	53283.159447	53432.680447	53582.401447	53732.322447	53882.443447	54032.764447	54183.285447	54333.906447	54484.627447	54635.548447	54786.669447	54937.990447	55089.511

Pixel error	I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII	XIII	XIV	XV	XVI	XVII	XVIII	XIX	XX	XXI	XXII	XXIII	XXIV	New I	New II	New III	New IV	New V	Average												
0.1	0.0008	0.0022	0.0012	0.0034	0.0024	0.0032	0.0004	0.0033	0.0011	0.0014	0.0000	0.0041	0.0043	0.0043	0.0057	0.0028	0.0021	0.0011	0.0017	0.0010	0.0017	0.0014	0.0016	0.0006	0.0017	0.0002	0.0016	0.0006	0.0021	0.0002	0.0023											
0.2	0.0041	0.0022	0.0051	0.0121	0.0099	0.0108	0.0023	0.0136	0.0065	0.0116	0.0007	0.0082	0.0042	0.0042	0.0056	0.0062	0.0056	0.0062	0.0059	0.0079	0.0057	0.0057	0.0056	0.0036	0.0017	0.0002	0.0053	0.004	0.0068	0.0004	0.0064											
0.3	0.0082	0.0066	0.0131	0.0261	0.0218	0.0235	0.009	0.0232	0.0115	0.0159	0.0054	0.0121	0.0115	0.0122	0.0198	0.0139	0.0112	0.0118	0.0122	0.0138	0.0108	0.0108	0.0078	0.0063	0.0036	0.0017	0.0063	0.0082	0.016	0.0166	0.0128											
0.4	0.0132	0.011	0.0245	0.0396	0.0399	0.0447	0.0147	0.0447	0.0184	0.0224	0.0108	0.0224	0.0238	0.0238	0.0368	0.0279	0.0217	0.0217	0.0222	0.0268	0.0178	0.0178	0.0148	0.0098	0.0028	0.0123	0.0148	0.0308	0.0288	0.023												
0.5	0.0194	0.0132	0.0366	0.0594	0.0539	0.0602	0.0247	0.0613	0.0274	0.0448	0.0161	0.0346	0.0224	0.0302	0.0424	0.0325	0.0335	0.0335	0.0343	0.0434	0.0308	0.0308	0.0245	0.0162	0.0042	0.0342	0.0366	0.0729	0.0608	0.0493												
0.6	0.0282	0.0263	0.052	0.0894	0.07	0.088	0.0377	0.0859	0.065	0.085	0.0287	0.0432	0.0395	0.0585	0.0597	0.0592	0.0333	0.0333	0.0343	0.0434	0.0355	0.0355	0.0289	0.0162	0.0042	0.0366	0.0425	0.0929	0.0608	0.0493												
0.7	0.0409	0.0417	0.0729	0.1328	0.1189	0.0571	0.0877	0.0533	0.091	0.043	0.075	0.0449	0.0648	0.0765	0.0769	0.0769	0.0683	0.0683	0.0683	0.0765	0.0683	0.0683	0.0545	0.0325	0.0054	0.0366	0.0425	0.1018	0.0813	0.0669												
0.8	0.0536	0.0526	0.1139	0.1964	0.1604	0.1093	0.1533	0.074	0.1093	0.0663	0.1199	0.0865	0.0612	0.0864	0.0963	0.0864	0.0864	0.0864	0.0963	0.0864	0.0864	0.0726	0.0425	0.0054	0.0366	0.0425	0.1218	0.1059	0.0864	0.0669												
0.9	0.0672	0.0746	0.1439	0.2494	0.1864	0.1299	0.1864	0.0946	0.1439	0.0813	0.1439	0.1145	0.0788	0.1145	0.1439	0.1145	0.1145	0.1145	0.1145	0.1439	0.1145	0.1145	0.0929	0.0545	0.0054	0.0366	0.0425	0.1439	0.1145	0.0864	0.0669											
1	0.0878	0.0897	0.1665	0.2897	0.2285	0.1581	0.2285	0.1126	0.1665	0.1126	0.1665	0.1337	0.0864	0.1337	0.1665	0.1337	0.1337	0.1337	0.1665	0.1337	0.1337	0.1010	0.0638	0.0054	0.0366	0.0425	0.1665	0.1337	0.1010	0.0864	0.0669											
1.1	0.1078	0.1097	0.1865	0.3129	0.2517	0.1813	0.2517	0.1358	0.1865	0.1358	0.1865	0.1579	0.1100	0.1579	0.1865	0.1579	0.1579	0.1579	0.1865	0.1579	0.1579	0.1262	0.0886	0.0054	0.0366	0.0425	0.1865	0.1579	0.1262	0.1010	0.0864	0.0669										
1.2	0.1278	0.1297	0.2085	0.3320	0.2708	0.2004	0.2708	0.1711	0.2085	0.1711	0.2085	0.1826	0.1358	0.1826	0.2085	0.1826	0.1826	0.1826	0.2085	0.1826	0.1826	0.1511	0.1039	0.0054	0.0366	0.0425	0.2085	0.1826	0.1511	0.1262	0.1010	0.0864	0.0669									
1.3	0.1478	0.1497	0.2285	0.3555	0.2943	0.2239	0.2943	0.2004	0.2285	0.2004	0.2285	0.2142	0.1673	0.2142	0.2285	0.2142	0.2142	0.2142	0.2285	0.2142	0.2142	0.1826	0.1358	0.0054	0.0366	0.0425	0.2285	0.2142	0.1826	0.1511	0.1262	0.1010	0.0864	0.0669								
1.4	0.1678	0.1697	0.2485	0.3845	0.3233	0.2529	0.3233	0.2285	0.2485	0.2285	0.2485	0.2341	0.1872	0.2341	0.2485	0.2341	0.2341	0.2341	0.2485	0.2341	0.2341	0.2026	0.1557	0.0054	0.0366	0.0425	0.2485	0.2341	0.2026	0.1711	0.1463	0.1212	0.0963	0.0712								
1.5	0.1878	0.1897	0.2685	0.4135	0.3523	0.2819	0.3523	0.2577	0.2685	0.2577	0.2685	0.2542	0.2073	0.2542	0.2685	0.2542	0.2542	0.2542	0.2685	0.2542	0.2542	0.2227	0.1758	0.0054	0.0366	0.0425	0.2685	0.2542	0.2227	0.1913	0.1664	0.1413	0.1164	0.0913	0.0664							
1.6	0.2078	0.2097	0.2875	0.4425	0.3813	0.3109	0.3813	0.2865	0.2875	0.2865	0.2875	0.2730	0.2261	0.2730	0.2875	0.2730	0.2730	0.2730	0.2875	0.2730	0.2730	0.2415	0.1946	0.0054	0.0366	0.0425	0.2875	0.2730	0.2415	0.2068	0.1819	0.1570	0.1321	0.1072	0.0823							
1.7	0.2278	0.2297	0.3065	0.4715	0.4103	0.3399	0.4103	0.3155	0.3065	0.3155	0.3065	0.2920	0.2451	0.2920	0.3065	0.2920	0.2920	0.2920	0.3065	0.2920	0.2920	0.2595	0.2126	0.0054	0.0366	0.0425	0.3065	0.2920	0.2595	0.2247	0.1998	0.1749	0.1500	0.1251	0.0992	0.0743						
1.8	0.2478	0.2497	0.3265	0.5005	0.4393	0.3689	0.4393	0.3445	0.3265	0.3445	0.3265	0.3120	0.2651	0.3120	0.3265	0.3120	0.3120	0.3120	0.3265	0.3120	0.3120	0.2765	0.2296	0.0054	0.0366	0.0425	0.3265	0.3120	0.2765	0.2417	0.2168	0.1919	0.1670	0.1421	0.1172	0.0923	0.0674					
1.9	0.2678	0.2697	0.3465	0.5295	0.4683	0.3979	0.4683	0.3741	0.3465	0.3741	0.3465	0.3320	0.2851	0.3320	0.3465	0.3320	0.3320	0.3320	0.3465	0.3320	0.3320	0.2965	0.2496	0.0054	0.0366	0.0425	0.3465	0.3320	0.2965	0.2617	0.2368	0.2119	0.1870	0.1621	0.1372	0.1123	0.0874	0.0625				
2	0.2878	0.2897	0.3665	0.5585	0.4973	0.4269	0.4973	0.4041	0.3665	0.4041	0.3665	0.3520	0.3051	0.3520	0.3665	0.3520	0.3520	0.3520	0.3665	0.3520	0.3520	0.3165	0.2696	0.0054	0.0366	0.0425	0.3665	0.3520	0.3165	0.2817	0.2568	0.2319	0.2070	0.1821	0.1572	0.1323	0.1074	0.0825				
2.1	0.3078	0.3097	0.3865	0.5875	0.5263	0.4559	0.5263	0.4341	0.3865	0.4341	0.3865	0.3720	0.3251	0.3720	0.3865	0.3720	0.3720	0.3720	0.3865	0.3720	0.3720	0.3365	0.2896	0.0054	0.0366	0.0425	0.3865	0.3720	0.3365	0.3017	0.2768	0.2519	0.2270	0.2021	0.1772	0.1523	0.1274	0.1025				
2.2	0.3278	0.3297	0.4065	0.6165	0.5553	0.4849	0.5553	0.4641	0.4065	0.4641	0.4065	0.3940	0.3471	0.3940	0.4065	0.3940	0.3940	0.3940	0.4065	0.3940	0.3940	0.3585	0.3116	0.0054	0.0366	0.0425	0.4065	0.3940	0.3585	0.3237	0.2988	0.2739	0.2490	0.2241	0.1992	0.1743	0.1494	0.1245				
2.3	0.3478	0.3497	0.4265	0.6455	0.5843	0.5139	0.5843	0.4931	0.4265	0.4931	0.4265	0.4146	0.3677	0.4146	0.4265	0.4146	0.4146	0.4146	0.4265	0.4146	0.4146	0.3791	0.3322	0.0054	0.0366	0.0425	0.4265	0.4146	0.3791	0.3443	0.3194	0.2945	0.2696	0.2447	0.2198	0.1949	0.1700	0.1451				
2.4	0.3678	0.3697	0.4465	0.6745	0.6133	0.5429	0.6133	0.5221	0.4465	0.5221	0.4465	0.4346	0.3877	0.4346	0.4465	0.4346	0.4346	0.4346	0.4465	0.4346	0.4346	0.3991	0.3522	0.0054	0.0366	0.0425	0.4465	0.4346	0.3991	0.3643	0.3394	0.3145	0.2896	0.2647	0.2398	0.2149	0.1900	0.1651				
2.5	0.3878	0.3897	0.4665	0.7035	0.6423	0.5719	0.6423	0.5511	0.4665	0.5511	0.4665	0.4546	0.4077	0.4546	0.4665	0.4546	0.4546	0.4546	0.4665	0.4546	0.4546	0.4191	0.3722	0.0054	0.0366	0.0425	0.4665	0.4546	0.4191	0.3843	0.3594	0.3345	0.3096	0.2847	0.2598	0.2349	0.2100	0.1851				
2.6	0.4078	0.4097	0.4865	0.7325	0.6713	0.6009	0.6713	0.5801	0.4865	0.5801	0.4865	0.4746	0.4277	0.4746	0.4865	0.4746	0.4746	0.4746	0.4865	0.4746	0.4746	0.4391	0.3922	0.0054	0.0366	0.0425	0.4865	0.4746	0.4391	0.4043	0.3794	0.3545	0.3296	0.3047	0.2798	0.2549	0.2300	0.2051	0.1802			
2.7	0.4278	0.4297	0.5065	0.7615	0.7003	0.6299	0.7003	0.6101	0.4278	0.6101	0.4278	0.4158	0.3689	0.4158	0.4278	0.4158	0.4158	0.4158	0.4278	0.4158	0.4158	0.4391	0.3922	0.0054	0.0366	0.0425	0.5065	0.4946	0.4598	0.4250	0.3902	0.3653	0.3404	0.3155	0.2906	0.2657	0.2408	0.2159	0.1910			
2.8	0.4478	0.4497	0.5265	0.7905	0.7293	0.6589	0.7293	0.6301	0.4478	0.6301	0.4478	0.4358	0.3889	0.4358	0.4478	0.4358	0.4358	0.4358	0.4478	0.4358	0.4358	0.4591	0.4122	0.0054	0.0366	0.0425	0.5265	0.5146	0.4798	0.4450	0.4102	0.3853	0.3604	0.3355	0.3106	0.2857	0.2608	0.2359	0.2110	0.1861		
2.9	0.4678	0.4697	0.5465	0.8195	0.7583	0.6879	0.7583	0.6501	0.4678	0.6501	0.4678	0.4558	0.4089	0.4558	0.4678	0.4558	0.4558	0.4558	0.4678	0.4558	0.4558	0.4791	0.4322	0.0054	0.0366	0.0425	0.5465	0.5346	0.4998	0.4650	0.4302	0.4053	0.3804	0.3555	0.3306	0.3057	0.2808	0.2559	0.2310	0.2061	0.1812	
3	0.4878	0.4897	0.5665	0.8485	0.7873	0.7169	0.7873	0.6901	0.4878	0.6901	0.4878	0.4758	0.4289	0.4758	0.4878	0.4758	0.4758	0.4758	0.4878	0.4758	0.4758	0.4991	0.4522	0.0054	0.0366	0.0425	0.5665	0.5546	0.5198	0.4850	0.4502	0.4253	0.4004	0.3755	0.3506	0.3257	0.3008	0.2759	0.2510	0.2261	0.2012	0.1763
3.1	0.5078	0.5097	0.5865	0.8775	0.8163	0.7459	0.8163	0.7101	0.5078	0.7101	0.50																															

NVGaze: An Anatomically-Informed Dataset for Low-Latency, Near-Eye Gaze Estimation

CHI 2019, May 4–9, 2019, Glasgow, Scotland UK

Pixel error	I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII	XIII	XIV	XV	XVI	XVII	XVIII	XIX	XX	XXI	XXII	XXIII	XXIV	New I	New II	New III	New IV	New V	Average	
7	0.9476	0.9276	0.9733	0.9832	0.9939	0.9857	0.9412	0.9652	0.9138	0.961	0.9211	0.975	0.9653	0.9892	0.9745	0.9889	1	0.955	0.8796	0.9585	0.9746	0.9485	0.9734	0.9505	0.9892	0.9925	0.9091	0.978	0.9703	0.9551	
7.1	0.9516	0.9342	0.9751	0.9861	0.9948	0.9862	0.9439	0.9669	0.9192	0.9624	0.9211	0.9788	0.9673	0.9892	0.9802	0.9889	1	0.957	0.8857	0.9608	0.9756	0.9514	0.9734	0.9522	0.9075	0.9868	0.9108	0.9785	0.9719	0.9572	
7.2	0.9536	0.9386	0.9764	0.9871	0.9948	0.9869	0.9463	0.9699	0.9212	0.9668	0.9319	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9584	0.8921	0.963	0.9769	0.9544	0.9734	0.9522	0.9079	0.9865	0.9147	0.9799	0.9735	0.9593	
7.3	0.9588	0.9408	0.9771	0.9898	0.9948	0.9878	0.9497	0.9685	0.9227	0.9685	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9611	0.8969	0.9646	0.9779	0.9565	0.9769	0.9529	0.9041	0.9713	0.9804	0.9751	0.9613	0.9593	
7.4	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
7.5	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
7.6	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
7.7	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
7.8	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
7.9	0.9607	0.9408	0.9782	0.9902	0.9948	0.9887	0.9523	0.9685	0.9227	0.9668	0.9359	0.9788	0.9735	0.9935	0.9802	0.9889	1	0.9625	0.9016	0.9668	0.9787	0.9587	0.9769	0.9529	0.9041	0.9716	0.9808	0.9751	0.9626	0.9593	
8	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.1	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.2	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.3	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.4	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.5	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.6	0.9715	0.9559	0.9853	0.9928	0.9962	0.9928	0.9606	0.9818	0.9352	0.9709	0.9462	0.9865	0.9837	0.9978	0.9802	0.9944	1	0.945	0.9264	0.9762	0.9833	0.9683	0.984	0.9608	0.9531	0.9824	0.9854	0.9854	0.982	0.9726	
8.7	0.9862	0.9627	0.9924	0.9958	0.9981	0.9959	0.9827	0.9934	0.9516	0.9841	0.9677	0.9904	0.9878	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
8.8	0.9868	0.9649	0.9924	0.9958	0.9981	0.9959	0.9827	0.9934	0.9516	0.9841	0.9677	0.9904	0.9878	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
8.9	0.9886	0.9649	0.9924	0.9958	0.9981	0.9959	0.9827	0.9934	0.9516	0.9841	0.9677	0.9904	0.9878	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9	0.9901	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.1	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.2	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.3	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.4	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.5	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.6	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.7	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.8	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
9.9	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10.1	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10.2	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10.3	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10.4	0.9911	0.9671	0.9948	0.9977	0.9986	0.9968	0.9879	0.9934	0.9536	0.9855	0.9749	0.9923	0.9959	0.9978	0.9972	0.9944	1	0.9845	0.9539	0.9857	0.9871	0.9901	0.9785	0.9911	0.9659	0.9621	0.954	0.9544	0.9912	0.9899	0.9822
10.5	0.9911	0.9671	0.9948</																												



Figure 7: Successful pupil estimation using 7-layer network with 293x239 input resolution. The dataset includes challenging cases such as bad lighting conditions, occluded and partly occluded pupils, head motion, mascara darkening eye lashes, refraction and reflections off glasses and contacts. The estimated pupil location is shown in red.



Figure 8: Augmentation samples for pupil estimation network. Image augmentations include rescaling, affine transformation, brightness and contrast variation, gaussian noise, randomly overlaid images. The estimated pupil location is shown in red.

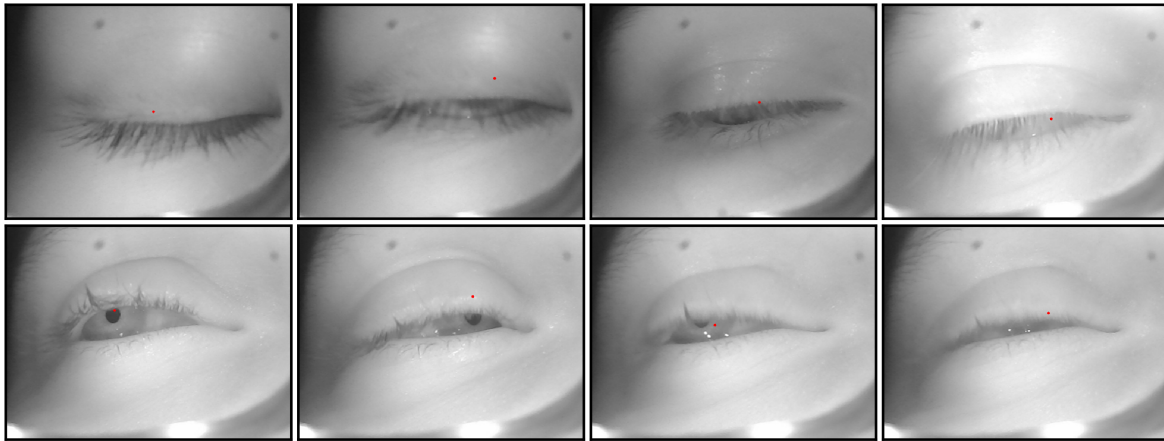


Figure 9: Limiting samples for pupil estimation network. Pupil estimation failed due to insufficient pupil visibility. The estimated pupil location is shown in red.

DRAFT

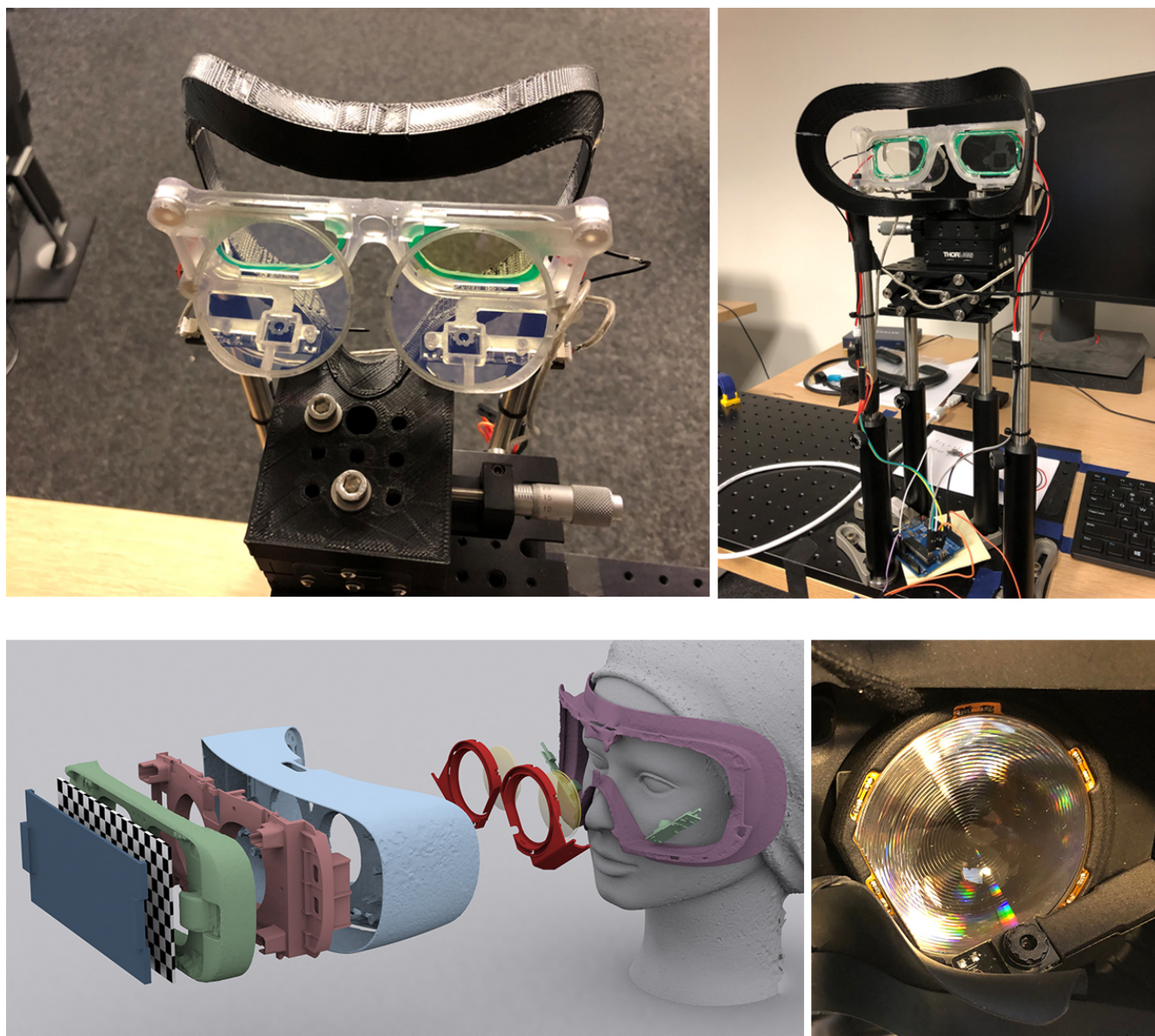


Figure 10: Eye Capture Hardware setup. Top row: The mounted AR glasses setup used for the first real dataset includes beam splitters to enable the on-axis view for the binocular eye tracking cameras. The infrared lighting units is driven by an Arduino microcontroller to enable temporally varying lighting conditions. The face mask stabilizes the user’s head during the capture experiment. The overall geometric setup is also used for generating the rendered synthetic dataset. Bottom row: Custom VR headset based on Samsung GearVR and Pupil Labs cameras. The headset is used for creating the off-axis eye image dataset. The headset includes clips attached to the lenses to hold the binocular eye cameras and infrared LEDs in place.

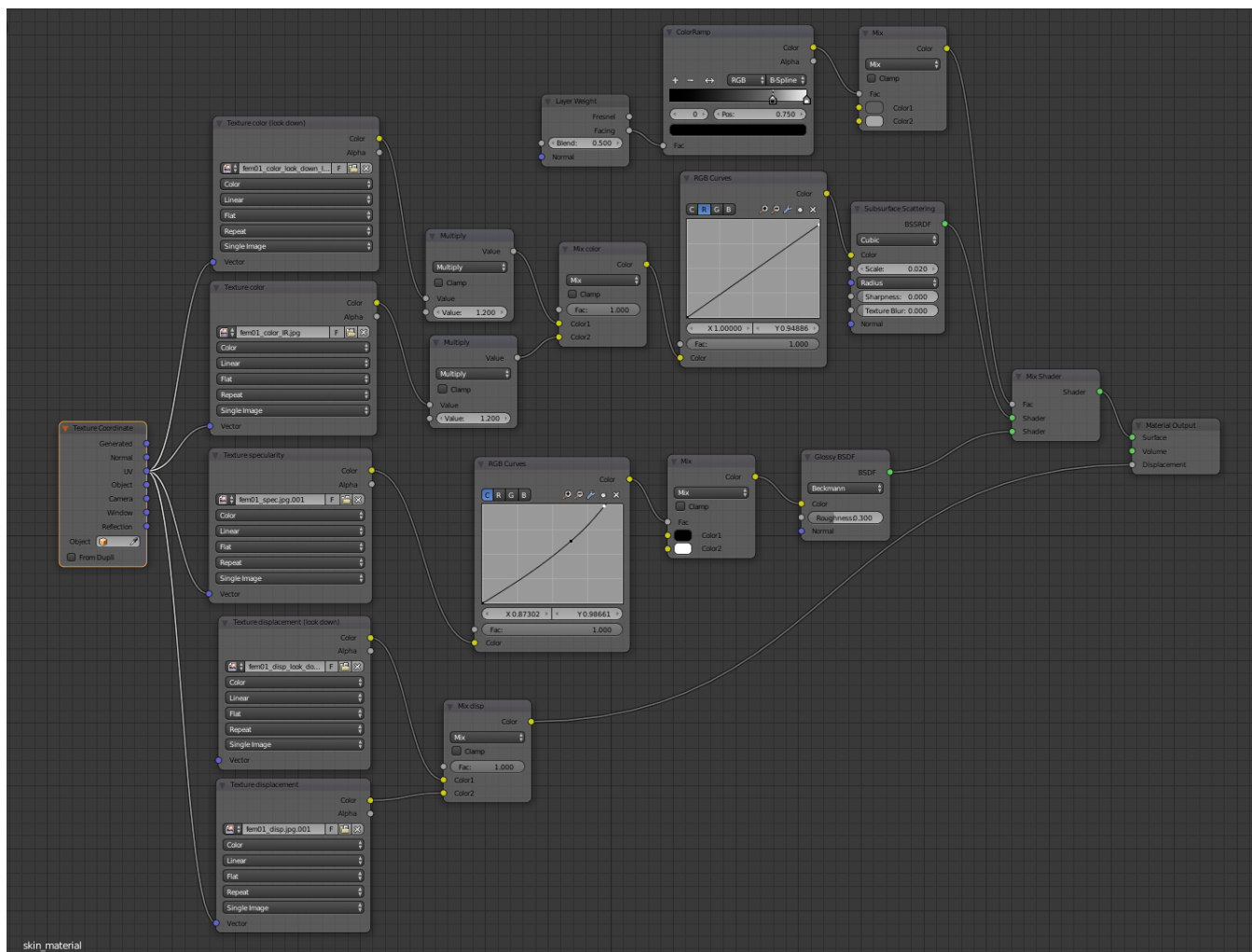


Figure 11: Shader graph of skin material for Cycles renderer in Blender.

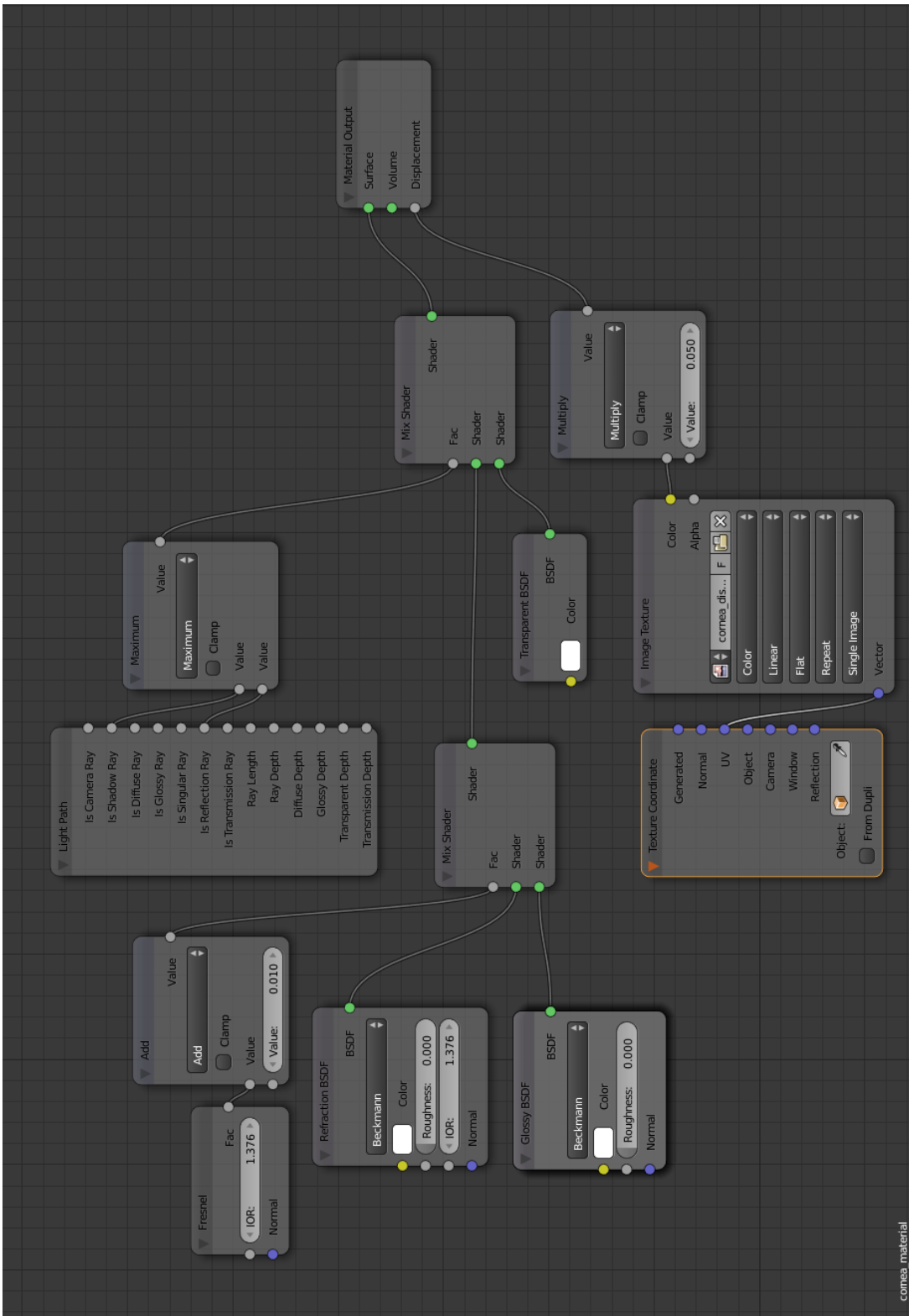


Figure 13: Shader graph of cornea material for Cycles renderer in Blender.