

INNOVATIVE METHODOLOGY | *Sensory Processing*

Human-level saccade detection performance using deep neural networks

Marie E. Bellet,^{1*} Joachim Bellet,^{2,3,4*} Hendrikje Nienborg,²  Ziad M. Hafed,^{2,4*} and  Philipp Berens^{1,2,5*}

¹Institute for Ophthalmic Research, University of Tübingen, Tübingen, Germany; ²Werner Reichardt Centre for Integrative Neuroscience, University of Tübingen, Tübingen, Germany; ³International Max Planck Research School for Cognitive and Systems Neuroscience, Tübingen, Germany; ⁴Hertie Institute for Clinical Brain Research, University of Tübingen, Tübingen, Germany; and ⁵Bernstein Center for Computational Neuroscience, Tübingen, Germany

Submitted 5 September 2018; accepted in final form 17 December 2018

Bellet ME, Bellet J, Nienborg H, Hafed ZM, Berens P. Human-level saccade detection performance using deep neural networks. *J Neurophysiol* 121: 646–661, 2019. First published December 19, 2018; doi:10.1152/jn.00601.2018.—Saccades are ballistic eye movements that rapidly shift gaze from one location of visual space to another. Detecting saccades in eye movement recordings is important not only for studying the neural mechanisms underlying sensory, motor, and cognitive processes, but also as a clinical and diagnostic tool. However, automatically detecting saccades can be difficult, particularly when such saccades are generated in coordination with other tracking eye movements, like smooth pursuits, or when the saccade amplitude is close to eye tracker noise levels, like with microsaccades. In such cases, labeling by human experts is required, but this is a tedious task prone to variability and error. We developed a convolutional neural network to automatically detect saccades at human-level accuracy and with minimal training examples. Our algorithm surpasses state of the art according to common performance metrics and could facilitate studies of neurophysiological processes underlying saccade generation and visual processing.

NEW & NOTEWORTHY Detecting saccades in eye movement recordings can be a difficult task, but it is a necessary first step in many applications. We present a convolutional neural network that can automatically identify saccades with human-level accuracy and with minimal training examples. We show that our algorithm performs better than other available algorithms, by comparing performance on a wide range of data sets. We offer an open-source implementation of the algorithm as well as a web service.

algorithm; deep neural network; eye movements; microsaccade; saccade

INTRODUCTION

Eye tracking is widely used in both animals and humans to study the mechanisms underlying perception, cognition, and action, and it is useful for investigating neurological and neurodegenerative diseases in human patients (Carpenter 1988; Kowler 2011; Leigh and Kennard 2004; Leigh and Zee 2015; MacAskill and Anderson 2016). This is in part due to practical reasons: recording eye movements is relatively easy (Duch-

owski 2007), while, at the same time, eye movements can be highly informative about brain state (Borji and Itti 2014; Haji-Abolhassani and Clark 2014).

The most prominent type of eye movement, in terms of eyeball rotation speed, is a ballistic shift in gaze position, called saccade. This type of eye movement occurs 3–5 times per second, and it can realign the fovea with interesting scene locations within only ~50 ms. Naturally, saccades cause dramatic changes in visual input when they occur, and they therefore impact neural processing in different visual areas and also in a variety of ways (Burr et al. 1994; Crevecoeur and Kording 2017; Duhamel et al. 1992; Golan et al. 2017; Ross et al. 1997; Reppas et al. 2002; Sommer and Wurtz 2008; Yao et al. 2018; Zirnsak et al. 2014). This even happens for the tiniest of saccades, called microsaccades, that occur when gaze is fixed (Bellet et al. 2017; Bosman et al. 2009; Chen and Hafed 2017; Hafed 2011; Hafed et al. 2015; Hass and Horwitz 2011; Herrington et al. 2009; Gur et al. 1997; Leopold and Logothetis 1998; Yu et al. 2017). Therefore, studies not quantitatively analyzing microsaccades can miss important behavioral and neural modulations in experiments (Hafed 2013). Saccades and microsaccades are, additionally, key discrete events in eye tracking traces that can be useful for parsing other eye movement epochs (e.g., smooth pursuits, ocular drifts, ocular tremors) for further analysis. Therefore, detecting saccades is typically the first step in any quantitative analysis of behavior or neural activity that might be impacted by these eye movements.

Several algorithms have been proposed for automating the task of saccade detection (reviewed in Andersson et al. 2017). For example, Engbert and Mergenthaler (2006) developed a method for classifying saccades and microsaccades based on an adaptive threshold. This algorithm (which we refer to here as EM) is particularly popular because of its simple implementation and ease of use, as well as its ability to detect even microsaccades. However, this algorithm, like others, may still mislabel some microsaccades due to high eye tracker noise (as is typical with video-based eye trackers) as well as small catch-up saccades occurring during smooth pursuit. Other existing algorithms (Larsson et al. 2013; Pekkanen and Lappi 2017) have the added advantage of providing additional labels for fixations and postsaccadic oscillations (PSO) in eye position.

* M.E.B. and J.B. contributed equally to this work; Z. M. Hafed and P. Berens are co-senior authors.

Address for reprint requests and other correspondence: Z. M. Hafed, Werner Reichardt Centre for Integrative Neuroscience, Tübingen, Germany (e-mail: ziad.m.hafed@cin.uni-tuebingen.de).

Despite their success, several shortcomings still render the use of existing algorithms either less reliable than desired or, at the very least, cumbersome. While the performance of many published algorithms is promising (Andersson et al. 2017; Pekkanen and Lappi 2017), it does not reach the level of trained human experts. Also, none of the existing algorithms show convincing performance for all eye movement-related events that may need to be analyzed (e.g., fixations, saccades, PSO, blinks, smooth pursuits). In addition, equipment-dependent hyperparameters, such as thresholds, need to be chosen for most algorithms, a fact that renders broad usability difficult. For example, even simple changes in eye tracking hardware, involving changes in sampling frequency or measurement noise, require retuning of such parameters. Retuning is also needed when the ranges of eye movement amplitudes being studied are modified (e.g., microsaccades vs. larger saccades). Perhaps most importantly, objective parameter estimation in existing algorithms is currently a challenging task because of a limited amount of available reliably labeled data. Finally, in many cases, applying available online resources is not straightforward. As a result of all of the above shortcomings, current laboratory practice often still involves experimenters spending substantial amounts of time to carefully relabel at least parts of their data after automatic saccade detection.

Here we propose a convolutional neural network (CNN) for classifying eye movements. The architecture of the network is inspired by U-Net, which has successfully been used for image segmentation (Ronneberger et al. 2015). We evaluated our network (U'n'Eye) on four challenging data sets containing small saccades occurring during fixations or smooth pursuits. On these data sets, U'n'Eye reached the performance level of human experts in labeling saccades and microsaccades, while being much faster. The network also beat state-of-the-art algorithms on a benchmark data set not just for saccade detection, but also for PSO. As we show here, our network can be trained quickly, even on a standard laptop, and with minimal amounts of training data. More importantly, our network's adaptability to different data sets makes U'n'Eye the novel state-of-the-art eye movement detection algorithm. We provide an easily accessible web service for running U'n'Eye (<http://uneye.berenslab.org>), as well as an open source implementation (<https://github.com/berenslab/uneye>). Our labeled data sets will also be freely available upon publication.

METHODS

Data sets. All experiments used for collecting the data sets were approved by ethics committees at Tübingen University. Human subjects provided informed, written consent in accordance with the Declaration of Helsinki. Monkey experiments were approved by the regional governmental offices of the city of Tübingen.

Data set 1 was collected from human subjects using the Eyelink 1000 video-based eye tracker (SR Research) sampling eye position at 1 kHz. The data set contains mostly microsaccades and small-amplitude memory-guided saccades. It contains 2,000 trials of 1 s. Out of these 2,000 trials, 1,000 were selected to compare U'n'Eye to other algorithms via cross-validation (Fig. 4). We named these trials "set1A." When testing for the impact of missing labels on performance (Fig. 7B), we used the other 1,000 trials, "set1B," to train networks and tested them on set1A.

Data set 2 was collected from three male rhesus macaque monkeys implanted with scleral search coils (in one eye for each of the monkeys). The data set contains catch-up saccades generated during

smooth pursuit. Eye position was again sampled at 1 kHz. For the trials containing smooth pursuit of sinusoidal target motion trajectories in this data set, the data were obtained from the experiments described in Hafed et al. (2008) and Hafed and Krauzlis (2008). For the trials containing pursuit of constant speed, the experimental conditions are described in Buonocore et al. (2018). Eye movement calibration for search coil data was done according to the procedures in Tian et al. (2016). The overall data set consists of 2,000 segments of 1 s of eye traces. Like in the case of data set 1, we split the set into two sets of 1,000 segments each, "set2A" and "set2B." set2A was used to compare U'n'Eye to Daye and Optican's (2014) algorithm (Fig. 4).

Data set 3 was collected from a single male macaque monkey using the Eyelink 1000 video-based system sampling eye position at 500 Hz. The data set contains microsaccades generated during fixation. The data were obtained from experiments described in Kawaguchi et al. (2018). It consists of 403 segments of 1.438 s. Similarly to data sets 1 and 2, we split the data in two subsets, "set3A" and "set3B." Set3A contains 350 segments and set3B 53 segments. Set3A was used for comparing U'n'Eye's performance to other that of algorithms (Fig. 4).

For the results shown in Fig. 7D, we used setA of all data sets for training and the respective setB for testing.

Data set 4 was collected from the same eye tracker as data set 1 but with different sets of subjects. It comes from a recently published study (Bellet et al. 2017) in which subjects had to keep fixation at the center of the screen before a peripheral target appearance. We selected 630 segments of 750 ms from each of 10 subjects (4,725 s in total). The data set includes not only successful trials, in which subjects maintained fixation, but also trials containing blinks or saccades outside of the fixation window. Again, we split the data set into two subsets. Set4A contained 330 segments per subject and was used to train networks. Set4B contained 300 segments per subject and was used to test the performance of the networks.

In all data sets, we manually detected saccades using a custom-made graphical user interface (GUI) in MATLAB. The GUI displayed horizontal and vertical eye position traces, as well as filtered radial eye velocity. The GUI internally estimated saccade onset and end times using a combination of velocity and acceleration thresholds (Chen and Hafed 2013). The user then manually interacted with the GUI to delete false alarms, correct false negatives, and adjust estimation of onset and offset timing.

Simulated saccades. To test the performance of our network on noisy labeled data, we designed artificial eye traces for which we knew the ground truth. Saccades ranging from 0.5 to 60° were simulated using an adaptation of a model for saccade waveforms (Dai et al. 2016). The model is a sum of soft ramp functions, which follows the relationship between amplitude and peak velocity observed in real saccades (Dai et al. 2016). Since the model is originally one dimensional, we adapted it so that it generates two-dimensional trajectories. Saccade generation in time was made to follow a Poisson process with λ equal to 3 saccades/s. Simulated blinks were also added by inducing sharp transients in the eye traces. Finally, a Gaussian white noise with a standard deviation of 0.02° was added to the trace. Then, as described in RESULTS, we trained U'n'Eye under a variety of conditions in which we intentionally removed a subset of saccade labels during training, to explore robustness to missing labels (Fig. 7).

U'n'Eye: our convolutional neural network. The architecture of CNN was inspired by U-Net, a CNN first used for image segmentation (Ronneberger et al. 2015). Here we modified U-Net to meet the requirements of an eye movement classifier. The network was built of seven convolutional layers with kernel size 5, each followed by a linear-rectifying unit (ReLU) and a BatchNorm layer, both described in detail in RESULTS. Batches consisted of samples of the same duration. The input to the network was eye velocity which was computed as the first-order difference of the eye position signal. The input was of dimension $N \times T \times 2$, where N is the batch size, T the number of time points, and 2 the number of coordinates (horizontal

and vertical eye velocity). The number of input time points could be variable but had to be a multiple of 25 bins due to the max pooling operations. The output of the network was a matrix of dimension $N \times K \times T$, where K was the user-defined number of classes. For example, we could have a “saccade” and “fixation” class in the networks of Fig. 3 and we could also add other classes like “PSO” in the network of Fig. 6.

We applied a softmax (Bishop 2016) activation function to the output of the last convolutional layer x :

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}} \quad (1)$$

where x_i is the layer corresponding to *class i*. Thus, the network’s output y represented the sample-by-sample conditional probability of each class (e.g., “fixation” or “saccade”) given the eye-velocity x and the network weights w :

$$Y_k = p(k = 1|x, w) \quad (2)$$

The final prediction of the algorithm represented the class that maximized this conditional probability:

$$\hat{k} = \text{argmax}_k p(k = 1|x, w) \quad (3)$$

We chose the kernel sizes of the convolutional and max pooling operations in a way to capture a relevant signal range around each time point. Based on the given kernel sizes of the network, it can be shown that the prediction of one time bin is influenced by the preceding and following 89 time bins of the velocity signal (Fig. 2B, red color).

Network training. We trained the network with minibatches whose size depended on the total number of training samples. We performed 10 training iterations in each epoch. Overfitting on the training set was prevented by computing the loss on a validation set and stopping training when the validation loss increased for three successive epochs. We used a multiclass error function, which, for two classes, equals the cross entropy loss. Weight-regularization was done with L2-penalty (Bishop 2016), which corresponds to a Gaussian prior with zero mean over the network weights. The optimal parameter λ was determined to be 0.01. The loss function was thus defined as:

$$L = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y(x_n, w) + \lambda \|w\|_2^2 \quad (4)$$

where N is the number of time points and K the number of classes. The ground truth label t_{nk} equals 1 if the *time point n* belongs to *class k*. Gradient computation was done with PyTorch autograd method.

We used the Adam optimizer (Kingma and Ba 2014) with an initial learning rate of 0.001. Adam is a stochastic gradient-based optimizer that uses adaptive learning rates for different weights of the network. An additional step decay by a factor of 2 was applied to the current learning rates when the loss on the validation set increased during one epoch.

Postprocessing. In the case of binary prediction into the classes fixation and saccade, we provided the possibility to define thresholds for minimum saccade duration and minimum saccade distance. If thresholds were given, saccades closer than the minimum distance were merged and saccades shorter than the minimum duration were removed. We obtained the results reported here with a minimum saccade distance threshold of 10 ms for data set 1 and of 3, 4, and 5 ms for data set 2, because we previously observed that some saccades occurred very close in time in this data set. For data sets 1–4, we used a minimum saccade duration threshold of 6 ms. The same thresholds were used for the algorithm Engbert and Mergenthaler (2006).

Data augmentation. U’n’Eye performs better with a bigger training set. However, we aimed to reduce the amount of saccades that a user should provide to train U’n’Eye. In this study, to increase the number

of training samples, the input eye positions were rotated and added to the original training samples:

$$x_2 = x \cos(\theta) + y \sin(\theta) \quad (5)$$

$$y_2 = -x \sin(\theta) + y \cos(\theta) \quad (6)$$

where x and y are the horizontal and vertical eye positions. We used $\theta = (1/4\pi, 3/4\pi, 5/4\pi, 7/4\pi)$ radians. Thus, we could increase by fivefold the size of our training set without causing overfitting.

Performance measures. To evaluate the eye movement detection performance of our network, we used the following metrics: Cohen’s kappa, F1 score, and onset and offset time differences.

Cohen’s kappa is a sample-based statistic. It reflects how much two coders agree on the class that each time bin belongs to, while controlling for chance agreement of the two coders. It is given by

$$\kappa = \frac{p_0 - p_e}{1 - p_e} \quad (7)$$

where p_0 is the proportion of time bins for which two coders agree, and p_e is the proportion of time bins for which agreement can be expected by chance.

For a binary classification of fixation vs. saccades, the Cohen’s kappa value p_e is given by

$$P_e = \frac{1}{N^2} \times \sum_{k=1}^K nk_{\text{coder1}} \times nk_{\text{coder2}} \quad (8)$$

where $nk_{\text{coder}X}$ is the number of time bins coder X assigned to *class k*.

The F1 score is a measure of classification accuracy that combines precision and recall of a predictor. Precision is defined as the proportion of correctly classified saccades over all predicted saccades. Recall is defined as the proportion of correctly classified saccades over all saccades in the ground truth. The F1 score is the harmonic mean of these two measures. It is given by

$$F1 = 2 \times \frac{TP}{2 * TP + FN + FP} \quad (9)$$

where TP is the number of true positives, FN the number of false negatives, and FP the number of false positives. For all true positive saccades, we compared saccade timing between the ground truth and prediction by calculating the absolute time differences between true and predicted saccade onsets and offsets.

Evaluation on a benchmark data set. We evaluated U’n’Eye performance on a benchmark data set by Andersson et al. (2017). This data set comprises 500 Hz eye-tracking data from humans looking at images, movies, or moving dots. It contains human labels for the events fixations, smooth pursuits, saccades, PSO, and blinks. Events that the human experts did not assign to any of these classes were labeled as “others.” For some trials, the data set contained labels from two different human coders. For other trials, only one label was available. We trained 20 independent networks with different random initializations on the data with labels from one human coder (coder RA). Performance was then tested on the trials with labels from two coders, which makes our result comparable with previously reported results (Pekkanen and Lappi 2017). Note that we were not able to reproduce the interrater measures reported by Andersson et al. (2017) in line with the results of Pekkanen and Lappi (2017). For comparability with the NSLR-HMM algorithm (Pekkanen and Lappi 2017), we excluded the event labels “other” for the calculation of Cohen’s kappa scores. The performance on the class “blinks” was not compared with other algorithms since it was not reported.

Evaluation of other algorithms. We compared U’n’Eye performance on our data sets to several already published algorithms. For data sets 1 and 3, which contain microsaccades occurring during fixation of a static target, we evaluated the performance of three algorithms designed for microsaccade detection (Engbert and Mer-

genthaler 2006; Otero-Millan et al. 2014; Sheynikhovich et al. 2018) and one algorithm designed for saccade detection in a high-noise regime (Pekkanen and Lappi 2017).

The algorithm by Engbert and Mergenthaler (2006) is commonly used as an unsupervised method to detect microsaccades. It selects saccades based on a threshold that depends on the level of the noise in the velocity. One parameter, called λ , can be also be fit to the data to obtain better results. This λ is multiplied with the velocity noise to determine a threshold for saccade selection. Here we chose λ values that maximized the metric of interest on the training data from our cross-validations. This was done to give this algorithm the benefit of the doubt in our comparisons. Importantly, before saccade detection, we smoothed the eye traces using a five-point average independently of the sampling frequency, as described by Engbert and Mergenthaler (2006).

The approach from Otero-Millan et al. (2014) is an unsupervised method. It gives an estimate of saccade onset and offset timing and thus can be compared in terms of both the Cohen's kappa and F1 metrics.

Another unsupervised method has been developed by Sheynikhovich et al. (2018). This algorithm only gives an estimate of microsaccade occurrence at one point in time without determining onset and offset. We thus compared only the performance in terms of F1 score for this algorithm. We considered as true positive any saccade detected ± 10 ms away from a ground truth saccade.

For data set 2, we evaluated the performance of the method by Daye and Optican (2014), which uses particle filters to detect saccades embedded in high-velocity eye movements. The algorithm was kindly provided by the authors. To increase performance, we detected sac-

ades independently in the horizontal and vertical channel and then merged the predictions. This is because the Daye and Optican algorithm only considers as a saccade an event crossing a threshold both in horizontal and vertical components at the same time, which increases the number of false negatives. To increase performance, the parameters were tuned differently for trials containing sinusoidal pursuit than for those containing linear pursuit, again to give the algorithm the benefit of the doubt when comparing to U'n'Eye. To detect saccades in sinusoidal pursuit, the parameters were set to $\Omega = 10^{-4}$, $\xi = 3 \cdot 10^{-4}$, $N = 100$, $m = 20$, $\lambda = 5 \cdot 10^{-3}$, $\psi = 2 \cdot 10^{-3}$. To detect saccades in linear pursuit, the parameters were set to $\Omega = 10^{-3}$, $\xi = 3 \cdot 10^{-3}$, $N = 100$, $m = 20$, $\lambda = 5 \cdot 10^{-3}$, $\psi = 10^{-4}$.

For all unsupervised algorithms, the 10 testing subsets from the cross-validation data were evaluated at once to yield better clustering estimates.

The results of the algorithm by Pekkanen and Lappi (2017) on data sets 1 and 3 were obtained by estimating the model's parameters via cross-validation using the same training folds as for U'n'Eye. The estimated parameters were kindly provided by the authors.

Compute time. The computation times of our algorithm reported here were achieved on a personal computer with a 2 GHz Intel Core i5 processor at 8 GB RAM running on Mac OS X 10.11.6.

Code and data availability. A web service for running the algorithm is available at <http://uneye.berenslab.org>. All code is available from <https://github.com/berenslab/uneye>. Data will be available upon publication.

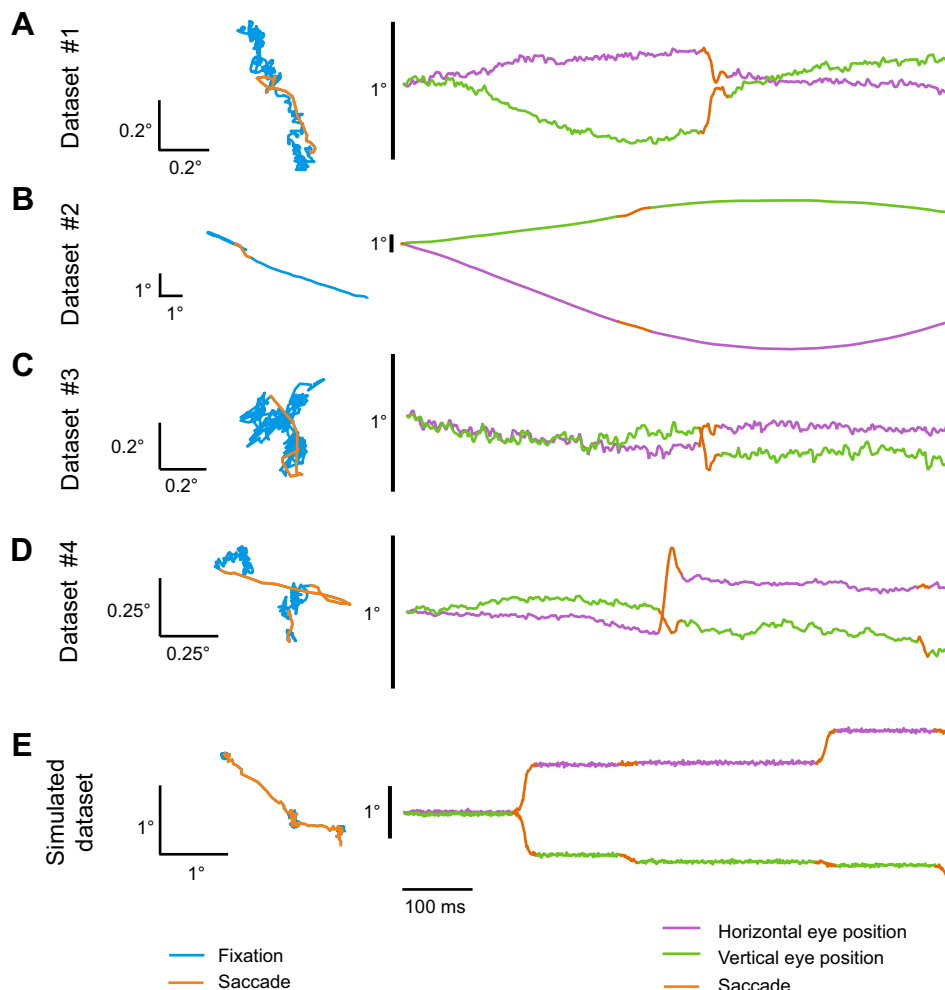


Fig. 1. Examples of eye traces containing saccades for detection. *A*: microsaccades during fixation recorded with a video-based eye tracker. *B*: catch-up saccades during smooth tracking recorded with scleral search coils. *C*: microsaccades during fixation recorded with a video-based eye tracker. *D*: microsaccades during fixation recorded with the same video-based eye tracker as in *A* but for different sets of subjects. *E*: simulated saccades. In all panels, 2D plots on the left are the 2D representation of the eye trajectory over 1 s of recording, and to the right of them are the horizontal and vertical components of the corresponding traces presented as a function of time; in this case, an upward deflection in the shown traces corresponds to a rightward or upward eye movement for the horizontal and vertical components, respectively. Note that, in *B*, we refer to the nonsaccadic smooth change in eye position as “fixation” for simplicity, since the primary goal of our algorithm was to detect saccades, irrespective of whether they happened during fixation or embedded in smooth pursuit eye movements.

RESULTS

Design of a convolutional neural network for eye movement classification. We developed a CNN that predicts the state of the eye for each time point of an eye trace. The aim of the network was to segment eye movement recordings (Fig. 1) into epochs containing saccades/microsaccades (orange highlights) vs. epochs not containing these eye movements (but see also *U'n'Eye: new state-of-the-art eye movement classifier* below for additionally classifying PSO using our network). Our primary goal was to have a network that can seamlessly handle the challenging scenarios of tiny microsaccades during fixation (Fig. 1A), small catch-up saccades embedded in relatively high smooth pursuit eye velocities (Fig. 1B), and microsaccades and saccades occurring in recordings with higher noise levels associated with video-based eye trackers when compared with, say, scleral search coil techniques (Fuchs and Robinson 1966; Judge et al. 1980) (Fig. 1C). We therefore trained and tested the network on three different challenging data sets (see METHODS and Table 1), which contain labels for fixations and saccades manually determined by human experts. To test the ability of our network to generalize across eye movement traces recorded from different individuals, we also included a fourth data set (Fig. 1D), which was obtained from 10 different subjects using the same eye tracker as in data set 1 (see METHODS). Testing generalizability was also achieved using with a fifth and final data set containing artificially generated noisy eye movement traces, in which the ground truth for saccade times was known (see METHODS) (Fig. 1E). Finally, we compared our network's performance to different existing algorithms, both on our data sets and also on a publicly available benchmark data set (Larsson et al. 2013) (http://dev.humlab.lu.se/www-transfer/people/marcus-nystrom/annotated_data.zip).

The network operates on the eye velocity signal and requires no other preprocessing. Eye velocity is computed as the differential of eye position (see METHODS), and chunks of eye

velocity signals are then input to the network. Briefly, the network's architecture is based on the U-Net, a CNN for pixel-by-pixel image segmentation (Ronneberger et al. 2015), which we modified to process one-dimensional signals and output a predictive probability for each eye movement class at every time point (Fig. 2A). A major change compared with the original U-Net architecture is that we introduced batch normalization (BatchNorm) layers (Klambauer et al. 2017). BatchNorm layers subtract a mean from their input and divide it by a standard deviation. Both of these parameters are estimated for each layer over minibatches of training samples during learning. This method normalizes the distribution of activations across the network layers, allowing for higher learning rates and reducing overfitting (Ioffe and Szegedy 2015). We also applied a rectified linear unit (ReLU) function between each convolutional and batch normalization layer. The ReLU function, or heaviside step function, introduces nonlinearities in the network, allowing it to apply arbitrary-shaped functions to the input data. Finally, the U-shaped architecture of the network leads to temporal downsampling and upsampling in the hidden layer representations (Fig. 2). Downsampling is achieved by max pooling (MaxPool) operations that reduce the dimensionality of the network content, extracting relevant features. Upsampling is realized by transposed convolution. Convolutional kernels and max pooling operations together lead to the integration of information over time. Due to the network design, the probability assigned to each time bin can be influenced by ± 89 preceding and following time bins (Fig. 2B). Thus, U'n'Eye takes into account a large enough signal to make point predictions of the correct eye movement class.

U'n'Eye achieves human-level performance. Our network achieved human-level performance after training on our data sets. We first illustrate this with three example scenarios for detecting saccades (Fig. 3). For illustrative purposes, we also show how the commonly used EM algorithm might perform for the examples; we later provide an exhaustive quantitative

Table 1. Data set characteristics

| Data Set | 1 | 2 | 3 | 4 |
|-------------------------------------|-----------------------------------|--------------------------------|--------------------|----------------------------|
| Subjects | Humans | Monkeys | Monkeys | Humans |
| Eye tracker | Eyelink 1000 | Search coil | Eyelink 1000 | Eyelink 1000 |
| Sampling frequency, Hz | 1,000 | 1,000 | 500 | 1,000 |
| Saccade type | Microsaccades and memory saccades | Saccades during smooth pursuit | Microsaccades | Microsaccades and saccades |
| Duration | | | | |
| Mean \pm SD, ms | 44.58 \pm 15.42 | 37.51 \pm 8.81 | 23.12 \pm 6.52 | 31.66 \pm 8.93 |
| Median, ms | 42 | 36 | 22 | 31 |
| Minimum, ms | 11 | 18 | 8 | 8 |
| Maximum, ms | 169 | 97 | 54 | 110 |
| Amplitude | | | | |
| Mean \pm SD, $^{\circ}$ | 0.69 \pm 0.93 | 1.07 \pm 0.70 | 0.22 \pm 0.13 | 0.33 \pm 0.28 |
| Median, $^{\circ}$ | 0.43 | 0.96 | 0.20 | 0.24 |
| Minimum, $^{\circ}$ | 0.02 | 0.04 | 0.010 | 0.008 |
| Maximum, $^{\circ}$ | 11.34 | 7.03 | 1.27 | 2.66 |
| Velocity | | | | |
| Mean peak \pm SD, $^{\circ}$ /s | 102.46 \pm 68.82 | 68.23 \pm 42.98 | 208.41 \pm 65.95 | 61.93 \pm 35.18 |
| Median peak, $^{\circ}$ /s | 81.91 | 56.59 | 198.86 | 52.96 |
| Minimum peak, $^{\circ}$ /s | 17.81 | 11.49 | 85.28 | 15.32 |
| Maximum peak, $^{\circ}$ /s | 547.72 | 450.44 | 560.28 | 423.18 |
| Median instantaneous, $^{\circ}$ /s | 5.63 | 15.70 | 10.20 | 5.63 |

All statistics refer to saccades. Note that minimum saccade amplitude may appear very low due to the existence of some saccades that had very strong dynamic overshoot (a substantial saccadic movement followed by one lobe of a postsaccadic oscillation almost to the original eye position before saccade onset). The statistics of the simulated data set are described in METHODS.

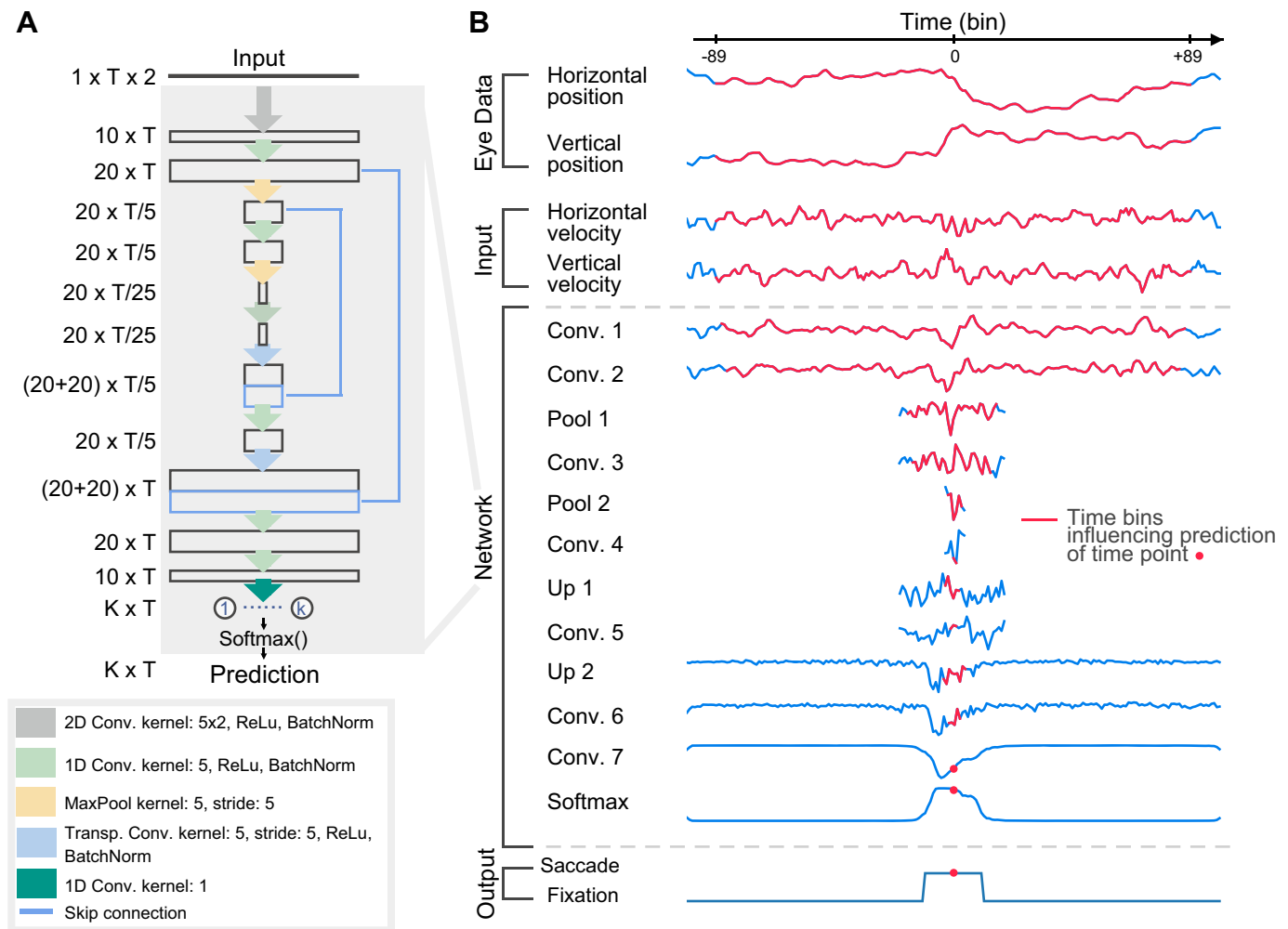


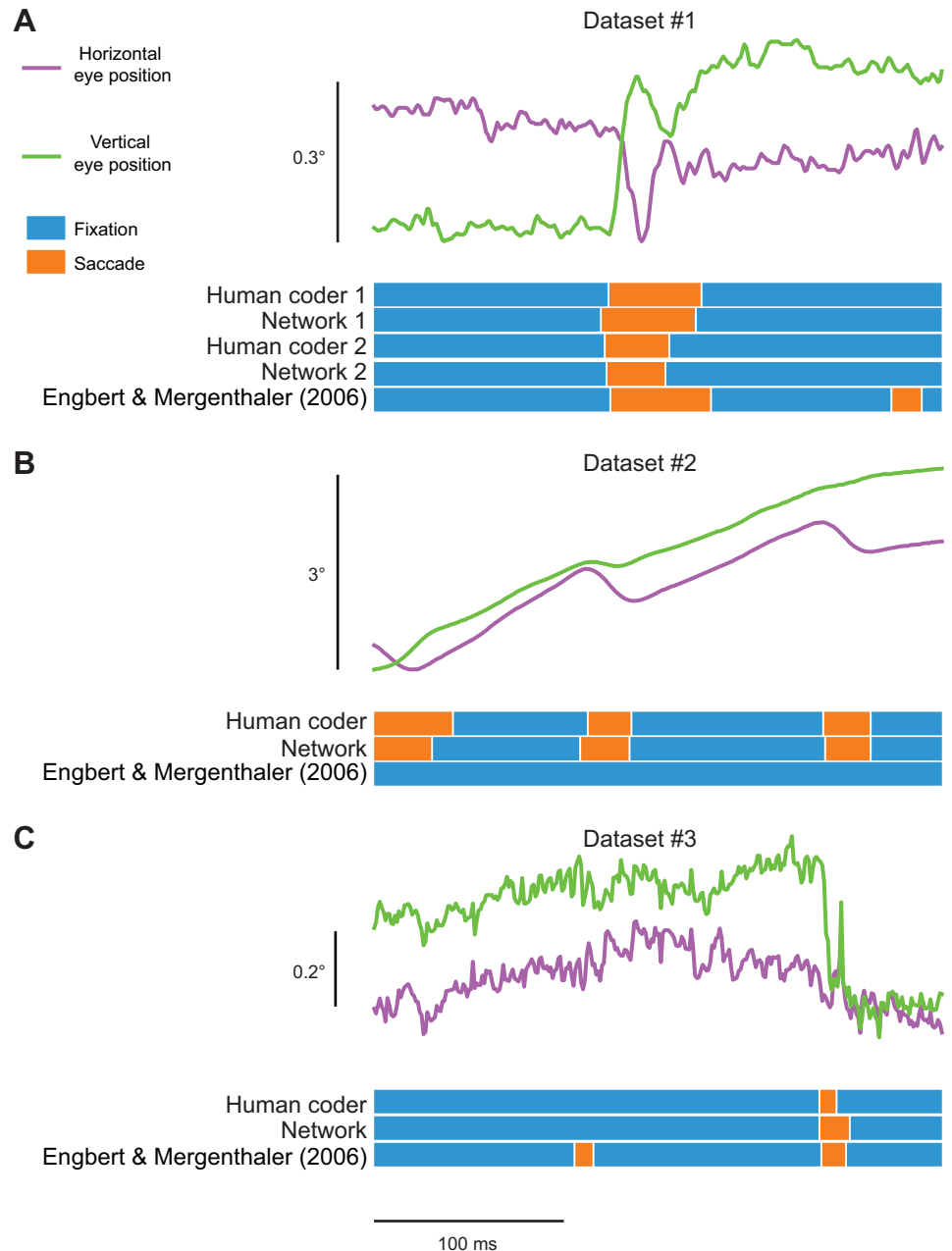
Fig. 2. U'n'Eye. **A**: network architecture. The input matrix contains horizontal and vertical eye velocity. T is the number of input time points (see **B**), and K is the user-defined number of eye movement classes (e.g., “fixation” vs. “saccade” in a binary classifier). The different network layers are described in *Design of a convolutional neural network for eye movement classification* under RESULTS. **B**: the output probability of one time bin is influenced by 89 time samples before and after this time bin. For each layer of the network, the red color indicates the range of influence of the time bin indicated by the red dot in the output. Traces show the projection of the layer’s output onto its first principal component. The outputs of convolutional (Conv) layers 6 and 7 resemble the final classifier’s output probability (Softmax), whereas early convolutional layers 1 and 2 seem to perform noise reduction.

comparison with several more algorithms (Fig. 4). In the first example, a small microsaccade occurred with substantial oscillation in eye position toward movement end, and with the amplitude of the movement being near the eye tracker noise level (Fig. 3A). Human coder 1 considered the postsaccadic oscillation as part of the saccade, and so did our network trained on his training set (compare the binary classification output of the coder and network 1 below the eye movement traces in Fig. 3A). On the other hand, coder 2 determined that the saccade ended earlier, and our network trained on his training set did the same (again, compare the classification output for human coder 2 and network 2). Thus, our network could match the criterion used by an individual human coder very well. Moreover, our network successfully avoided a false detection by the EM algorithm on these traces. In the second example, the EM algorithm missed all three saccades, which is perfectly reasonable since this algorithm was never designed to work in association with smooth pursuit eye movements, but our network successfully flagged them (Fig. 3B). Finally, the eye movement in the third example was collected with a

video-based eye tracker having substantially more noise (Fig. 3C). In this case, one false detection made by the EM algorithm was successfully excluded by our network.

To present more quantitative performance measures, we first tested our network on our in-house data sets (Fig. 1) and compared its performance to that of commonly used or recently published algorithms. For our network, we performed 10-fold cross-validation separately for data sets 1–3. In each cross-validation round, 90% of the data were used for training the network, and the remaining 10% were used to test performance. A separate validation set from each data set was used to detect overfitting of the network. To prevent such overfitting, we regularized the weights of the network using the L2 penalty (Bishop 2016) (see METHODS), preventing the parameters of the network from deviating excessively from zero. Furthermore, we made use of early stopping. For this, a separate validation set was used, and the validation set error was computed in each epoch. Training was stopped at the point of smallest validation set error. For data sets 1 and 2, 950 s of eye traces were used for cross-validation and 50 s for valida-

Fig. 3. Examples of eye traces from our first three data sets. Saccades are labeled by either human coders, different instances of U'n'Eye, or a popular algorithm from the literature included here for illustrative purposes (also see Fig. 4 for detailed performance comparisons to several algorithms). *A*: an example microsaccade exhibiting substantial postsaccadic oscillation (PSO). The top two traces show eye position as a function of time in an identical format to Fig. 1. Below the eye position traces, we show labels for “fixation” or “saccade” made by two human experts (human coder 1 and human coder 2) as well as predictions of two separate networks. Network 1 was trained on labels from human coder 1, and network 2 was trained on labels from human coder 2. Note how each network matched the performance of its corresponding human coder. The very bottom row shows the performance of the Engbert and Mergenthaler (2006) algorithm (EM), which suffered from a false alarm later in the trace due to eye tracker noise. *B*: saccades embedded in smooth pursuit eye movements. Here, our network successfully detected three catch-up saccades, all of which were missed by the EM algorithm, which was not designed to work with eye movement records containing smooth pursuit. The reason that these saccades were missed is that the saccades were directed opposite to the ongoing pursuit, resulting in momentary reductions in eye speed, as opposed to increases. *C*: an example microsaccade embedded in high eye tracker noise. Once again, the EM algorithm suffered from a false alarm due to eye tracker noise.



tion. Thus, each training set contained 855 s of data. For data set 3, 330 s were used for cross-validation and 23 s for validation, resulting in 297 s of data in each training set. For the other algorithms that we tested, we used the same cross-validation approach in the case of supervised algorithms [EM (Engbert and Mergenthaler 2006), Pekkanen and Lappi (2017)]. Note that we used the EM algorithm as a supervised method since we fitted its single parameter on our training data (see METHODS). For unsupervised methods (Daye and Optican 2014; Otero-Millan et al. 2014; Sheynikhovich et al. 2018), the identical 10 test sets were evaluated without using the training set (see METHODS).

Finally, similarity of the algorithms' predictions to human labels was evaluated using three metrics. First, we calculated Cohen's kappa, which is a sample-by-sample similarity measure that takes chance agreement of two predictors into account (Cohen 1960). Second, we calculated the F1 score, which is an

accuracy measure that considers precision and recall of a classifier. Recall corresponds to the number of correctly detected saccades divided by the number of saccades that were labeled by the human expert. Precision, on the other hand, is the number of correctly classified saccades divided by the total number of saccades detected by the classifier (see METHODS). The F1 score is defined as the harmonic mean of both, and it thus only measures how accurate saccades were detected without taking into account their timing (i.e., exact saccade onset and offset times). Correctly labeling saccade onset and offset can be crucial for further analyses. Therefore, for our third and final metric, we additionally computed the absolute time difference in onset and offset of correctly predicted saccades and of saccades labeled by the human experts. This measure reflects how well an algorithm agrees with the human coder in terms of saccade start and end.

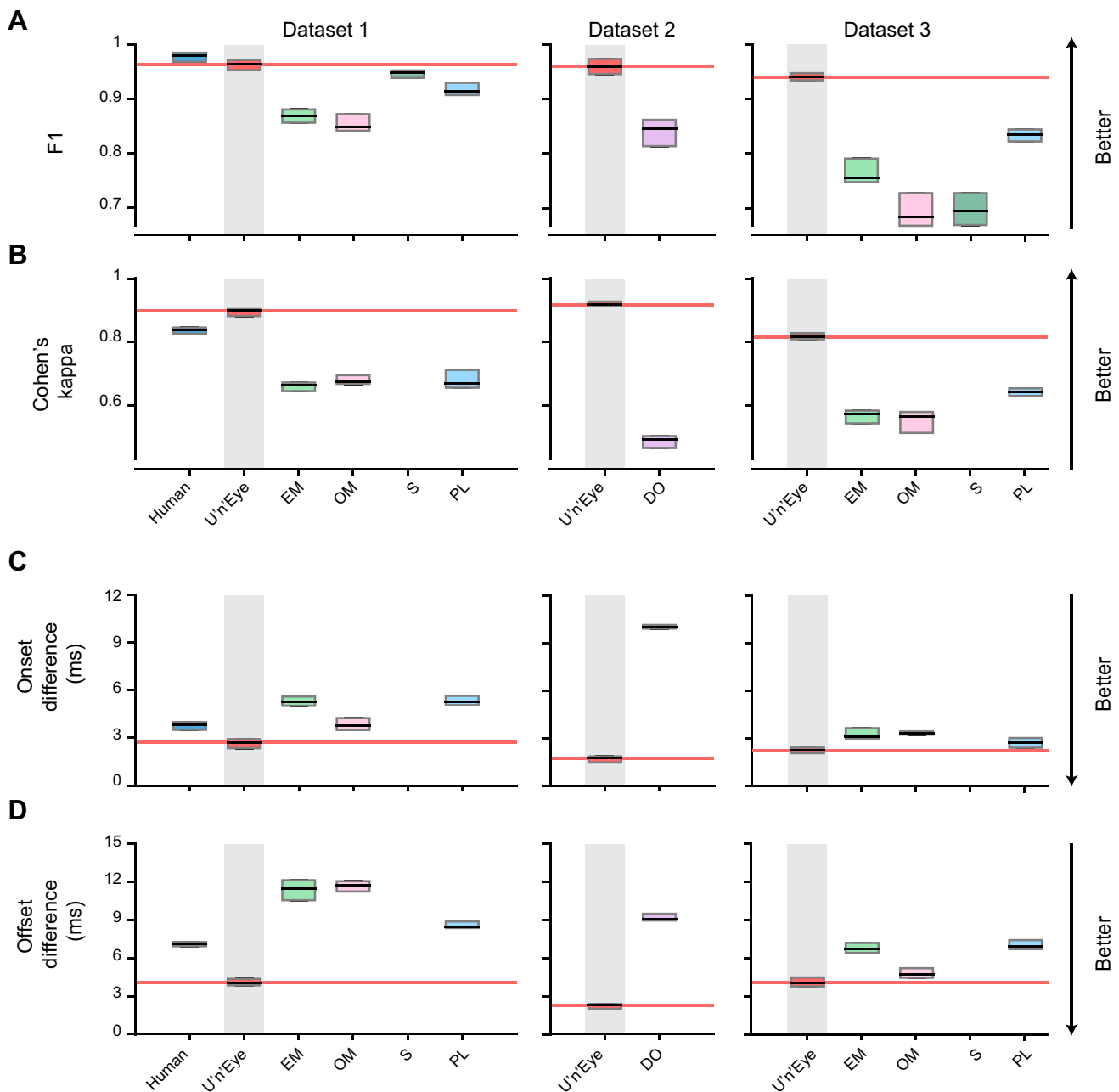


Fig. 4. High performance of U'n'Eye. Each panel shows results from one performance metric described in the text, and on each of our first three data sets. For each metric, we show the median across 10 different cross-validation runs. The boxes show two quartiles of the distributions. *A*: F1 score (See Eq. 9 and the surrounding text) summarizing precision and recall performance between two predictors. The first predictor was always a human coder (considered as ground truth). Therefore, the first column indicates agreement between a second coder (labeled Human in the figure) to the original coder used to train our network. *B*: Cohen's kappa measuring sample to sample agreement. *C*: average absolute difference in the timing of saccade onset times. *D*: same as *C* but for saccade offset times. In all cases, our network (highlighted by gray rectangles) demonstrated superior performance (the arrows on the far right side indicate the direction of superior performance for each metric). EM, Engbert and Mergenthaler (2006); OM, Otero-Millan et al. (2014); S, Sheynikhovich et al. (2018); PL, Pekkanen and Lappi (2017); DO, Daye and Optican (2014). Note that we only tested data set 2 on DO because only this algorithm was explicitly designed to deal with smooth pursuit eye movements.

U'n'Eye reached high similarity to the human coder (Fig. 4, *A* and *B*) and outperformed all the other compared algorithms (Fig. 4, *A* and *B*). U'n'Eye also detected saccade onset and offset in high agreement with the human labels. On average, saccade onset differences to human labels were smaller than 3 ms, and saccade offset differences were smaller than 4 ms. Saccade onset and offset labels by the other algorithms deviated more strongly from the human-labeled saccades (Fig. 4, *C*

and *D*; Table 2). This indicates that U'n'Eye's saccade predictions were more humanlike.

In the more challenging data set 2, in which saccades occurred during smooth pursuit eye movements, U'n'Eye outperformed the algorithm by Daye and Optican (2014), which was designed to overcome this difficulty. Here, saccade peak velocity was close to the instantaneous velocity of the ongoing smooth pursuit movements. In fact, the minimum saccade peak

Table 2. Comparison of U'n'Eye performance to other algorithms on data sets 1, 2, and 3

| Data Set | Algorithm | F1 | Cohen's Kappa | Δ Onset, ms | Δ Offset, ms |
|----------|-----------|------------------------|------------------------|------------------------|------------------------|
| 1 | U'n'Eye | 0.96 \pm 0.01 | 0.89 \pm 0.02 | 2.66 \pm 0.34 | 4.11 \pm 0.41 |
| | EM | 0.87 \pm 0.03 | 0.66 \pm 0.02 | 5.39 \pm 0.49 | 11.28 \pm 1.00 |
| | OM | 0.85 \pm 0.03 | 0.68 \pm 0.03 | 3.80 \pm 0.48 | 11.50 \pm 0.77 |
| | S | 0.95 \pm 0.02 | | | |
| | PL | 0.92 \pm 0.02 | 0.68 \pm 0.03 | 5.51 \pm 0.88 | 8.53 \pm 0.60 |
| 2 | U'n'Eye | 0.96 \pm 0.01 | 0.92 \pm 0.01 | 1.70 \pm 0.29 | 2.19 \pm 0.37 |
| | DO | 0.84 \pm 0.03 | 0.49 \pm 0.03 | 9.99 \pm 0.19 | 9.22 \pm 0.35 |
| 3 | U'n'Eye | 0.94 \pm 0.01 | 0.82 \pm 0.02 | 2.23 \pm 0.22 | 3.99 \pm 0.60 |
| | EM | 0.77 \pm 0.04 | 0.58 \pm 0.04 | 3.22 \pm 0.53 | 6.87 \pm 0.66 |
| | OM | 0.68 \pm 0.07 | 0.55 \pm 0.06 | 3.48 \pm 0.55 | 4.90 \pm 0.64 |
| | S | 0.70 \pm 0.04 | | | |
| | PL | 0.83 \pm 0.02 | 0.64 \pm 0.03 | 2.74 \pm 0.39 | 6.94 \pm 0.50 |

In bold are the best performances for each data set. In all cases, U'n'Eye achieved highest performance. Values report mean and standard deviation across validation sets. EM, Engbert and Mergenthaler (2006); OM, Otero-Millan et al. (2014); PL, Pekkanen and Lappi (2017); S, Sheynikhovich et al. (2018).

velocity in this data set was smaller than the median instantaneous velocity during pursuit (Table 1). Yet U'n'Eye succeeded in detecting such saccades. This was because the network architecture utilized a substantial time window (Fig. 2), allowing it to infer changes in the state of the eye even if the instantaneous velocity is low compared with the surrounding eye trace.

We next addressed the question of whether U'n'Eye can achieve a similar level of interhuman agreement when multiple human experts analyze the same data. For this, we used data set 1 because, among the four data sets, it contained saccades with the widest range of amplitudes (from as small as 0.02° up to a size of 11° ; see Table 1 for a reason why saccades as small as 0.02° were possible). We could thus assess interrater agreement for a broad range of saccades. Data set 1 was labeled by a second independent human coder (Fig. 3, top panel; Fig. 4, data set 1). Coder 1 estimated saccade timing based on a combination of the raw eye traces and the smoothed radial velocity, whereas coder 2 used the raw eye traces only. We trained independent networks either with labels from coder 1 or coder 2 (network 1 and network 2, respectively), and we tested the networks' performance on the 10 test sets from the 10-fold cross-validation routine described above, both against ground truth labels from coder 1 or coder 2. U'n'Eye's saccade labels were as similar to both human coders as the human labels were to each other (Table 3). In terms of the F1 score, the interhuman agreement was not significantly different from the network-human agreement (Table 4). Interestingly, network 1 showed higher similarity scores than coder 2 when both were compared with labels of coder 1 in the test sets, and vice versa for network 2 and coder 2. This is reflected by larger Cohen's kappa scores and smaller onset and offset differences (Table 4, all $P < 5 \times 10^{-5}$ after Bonferroni correction for multiple comparisons, Student's paired samples *t*-test for Cohen's kappa and F1 scores, and independent samples *t*-test for on- and offset differences). This indicates that U'n'Eye's saccade estimation surpasses interrater consistency.

U'n'Eye misses only a small fraction of microsaccades. We then analyzed the patterns of agreement and disagreement between U'n'Eye and human labeling. For true positive sac-

Table 3. Interrater comparison

| | Cohen's kappa | F1 | Δ Onset, ms | Δ Offset, ms |
|-----------------------|------------------------|------------------------|------------------------|------------------------|
| Coder 1 vs. coder 2 | 0.83 \pm 0.02 | 0.98 \pm 0.01 | 3.72 \pm 0.39 | 7.10 \pm 0.34 |
| Network 1 vs. coder 1 | 0.89 \pm 0.02 | 0.96 \pm 0.01 | 2.65 \pm 0.34 | 4.11 \pm 0.41 |
| Network 2 vs. coder 2 | 0.89 \pm 0.01 | 0.96 \pm 0.01 | 2.00 \pm 0.11 | 4.81 \pm 0.33 |
| Network 2 vs. coder 1 | 0.85 \pm 0.01 | 0.96 \pm 0.01 | 3.34 \pm 0.34 | 5.58 \pm 0.33 |
| Network 1 vs. coder 2 | 0.86 \pm 0.01 | 0.96 \pm 0.01 | 2.82 \pm 0.32 | 6.57 \pm 0.53 |

The first row shows the similarity measures between labels from two human experts (coder 1 and coder 2). Network 1 was trained on labels from coder 1, and network 2 was trained on labels from coder 2. In bold are comparisons leading to best performances. Values report mean and standard deviation across cross-validations. Interrater agreement was evaluated on the 10 test samples from cross-validation.

cadets, the two dimensional histogram of detected movements reflected the typical main sequence relationship between peak velocity and amplitude of saccades (Fig. 5, A, D, and G) (Zuber et al. 1965). A few false positives were present within the range of the main sequence, suggesting that the human coder forgot to label some saccades (for example, see the movement in the inset in Fig. 5B). Concerning the rare false negatives that occurred, some of them had fairly large amplitudes (beyond eye tracker noise). Closer inspection revealed that there were pairs of successive saccades that had very short intersaccadic intervals. The network lumped them into one movement, whereas the human coders separated them. Most remaining disagreements between the human and the network were associated with the smallest microsaccades, closest to eye tracker noise levels.

U'n'Eye: new state-of-the-art eye movement classifier. To compare our algorithm to state-of-the-art methods for eye movement classification, we next evaluated its performance on a benchmark data set (Larsson et al. 2013), which has previously been used for the comparison of 12 eye movement classifiers (Andersson et al. 2017; Pekkanen and Lappi 2017). The data set comprises 500-Hz eye tracking recordings from humans watching videos, images, or moving dots, and it contains human labels for fixations, smooth pursuits, saccades,

Table 4. Statistical tests in interrater comparison

| Metric | Comparison | Test | <i>t</i> -Value | <i>P</i> Value |
|--------------------------------|------------|----------------------------|-----------------|-----------------------|
| Kappa to C1 | N1 vs. C2 | paired <i>t</i> -test | 18.38 | $2.98 \cdot 10^{-7}$ |
| Kappa relative to C2 | N2 vs. C1 | paired <i>t</i> -test | 10.88 | $3.69 \cdot 10^{-10}$ |
| F1 relative to C1 | N1 vs. C2 | paired <i>t</i> -test | -3.52 | $5.08 \cdot 10^{-2}$ |
| F1 relative to C2 | N2 vs. C1 | paired <i>t</i> -test | -3.7 | $5.19 \cdot 10^{-2}$ |
| Onset distance relative to C1 | N1 vs. C2 | independent <i>t</i> -test | -6.6 | $2.98 \cdot 10^{-5}$ |
| Onset distance relative to C2 | N2 vs. C1 | independent <i>t</i> -test | -13.6 | $5.26 \cdot 10^{-10}$ |
| Offset distance relative to C1 | N1 vs. C2 | independent <i>t</i> -test | -17.9 | $5.28 \cdot 10^{-12}$ |
| Offset distance relative to C2 | N2 vs. C1 | independent <i>t</i> -test | -15.3 | $7.33 \cdot 10^{-11}$ |

Network 1 (N1) was trained on labels from coder 1 (C1), and network 2 (N2) was trained on labels from coder 2 (C2). All *P* values were Bonferroni corrected for multiple comparisons.

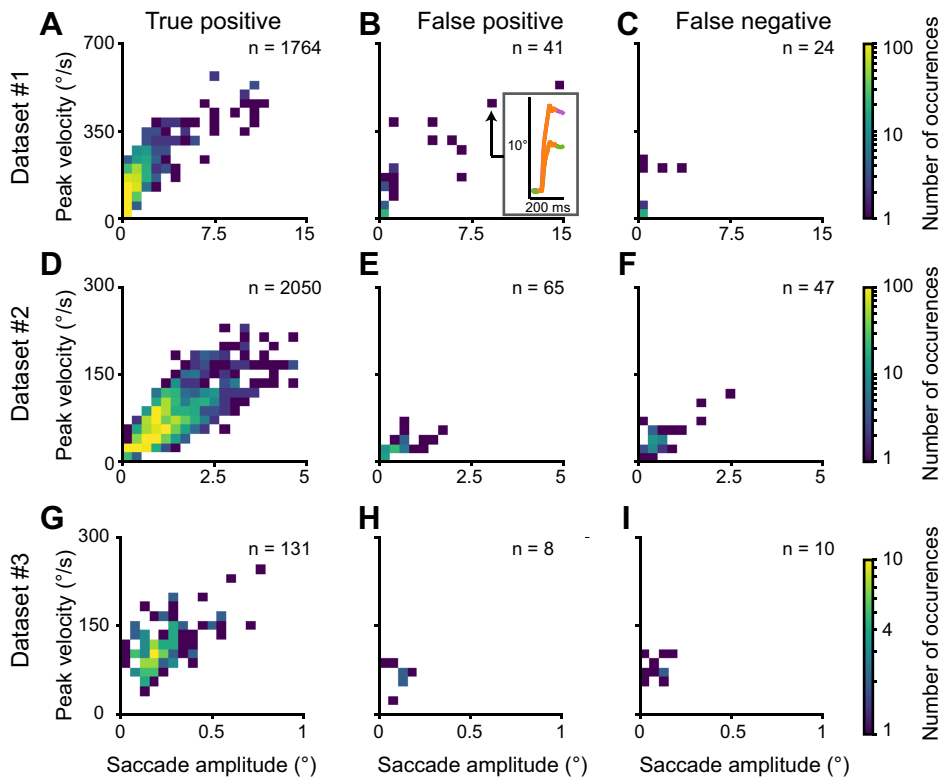


Fig. 5. Location on the main sequence of detected and undetected saccades. *A, D, and G*: saccades that were detected both by a human expert and U'n'Eye. The detected saccades expectedly followed the main sequence relationship between peak velocity and movement amplitude. *B, E, and H*: saccades that were detected only by U'n'Eye. Most saccades were small and close to the eye tracker noise, likely being cautiously unlabeled by human coders. In the inset, a large saccade was detected by U'n'Eye but not by the human coder, suggesting a possible lapse by the latter. *C, F, and I*: saccades missed by U'n'Eye. Most of these were very small.

PSO (Fig. 6A), and blinks. We therefore used U'n'Eye as a multiclass classifier to predict saccades, PSOs, and blinks (Fig. 6B). Fixations and smooth-pursuit eye movements were both assigned to the fixation class. U'n'Eye output a predictive probability for each class (Fig. 6D), with the prediction value

corresponding to the class that maximized this predictive probability (Fig. 6C). We trained U'n'Eye on one part of the data and evaluated its performance on the test trials listed in Andersson et al. (2017; their Table 11). When considering the whole benchmark data set, U'n'Eye outperformed the state-of-

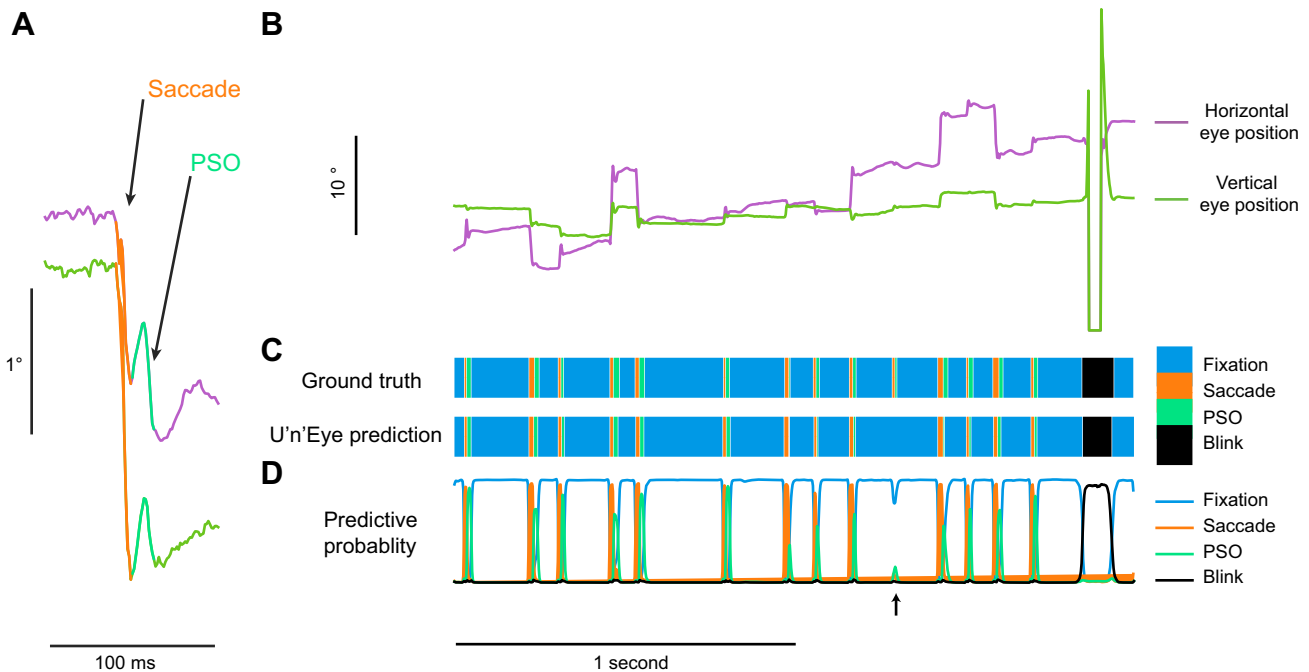


Fig. 6. Multiclass labeling by U'n'Eye. *A*: an example saccade showing substantial postsaccadic oscillation (PSO) from the data in Larsson et al. (2013). *B*: an example full trace from the same data set showing sequences of saccades, PSOs, and blinks. *C*: for the trace in *B*, ground truth labels are shown, in addition to labels by U'n'Eye. The latter successfully classified all ground truth labels, except for one instance marked by a black vertical arrow. *D*: nonetheless, the predictive probability of the network still showed a transient for the missed microsaccade (upward black arrow), suggesting that additional postprocessing may be used to improve the performance of U'n'Eye even more. For example, the user could manually inspect significant transients in predictive probability.

the-art classifiers for saccades and PSOs (Table 5). Moreover, U'n'Eye's performance lay within the range of the intercoder agreement of the two human experts who labeled the data set (Table 5). This result indicates that U'n'Eye is very well suited for multiclass eye movement classification.

Practical considerations for U'n'Eye usage. To better understand the practical aspects of using our approach, we additionally assessed how U'n'Eye performs under different training scenarios. The results of this section can be used as good practice guidelines by the users in their own applications that employ our algorithm.

First, we studied how the amount of training data impacts saccade detection performance. In practice, the available number of annotated training samples might be limited. To achieve good performance of U'n'Eye, small training sets were sufficient (Fig. 7A). Even with only 50 s of labeled data, our network outperformed other algorithms. Using more training samples led to a further increase of performance. Training time was also no limiting factor, since training a new network even on a CPU took only ~2 min for every minute of training data (Fig. 7A).

In machine learning, the quality of the training data is also crucial for the performance of a classifier, since the latter directly learns from the human ground truth labels. Human labeling, however, is prone to mistakes and lapses: saccades might be missed by the human coder, leading to noisy labels. We therefore assessed how U'n'Eye's performance was influenced by noise-corrupted labels. We evaluated the network's performance when trained on real data (data set 1) from which we artificially removed a fixed fraction of saccade labels. U'n'Eye was robust to the presence of noisy labels in the training data: even with 20% of missing labels, our network outperformed other algorithms (Fig. 7B). We also trained the network on simulated data (Fig. 1E) for which we knew the ground truth. While noise-corrupted labels in the training data impaired saccade detection performance as expected, this effect could be compensated for by using a larger amount of training data (Fig. 7C). This indicates that U'n'Eye can achieve good performance even if the human coder misses some saccades in the training set.

Table 5. Performance of U'n'Eye compared with state-of-the-art algorithms

| Event | | Coder MN | U'n'Eye | NSLR-HMM | LNS |
|----------|-------|----------|-------------|----------|------|
| Saccades | Image | 0.91 | 0.89 | | 0.81 |
| | Dot | 0.80 | 0.79 | | 0.75 |
| | Video | 0.88 | 0.89 | | 0.81 |
| | All | 0.89 | 0.88 | 0.82 | 0.81 |
| PSOs | Image | 0.76 | 0.72 | | 0.64 |
| | Dot | 0.59 | 0.59 | | 0.53 |
| | Video | 0.73 | 0.68 | | 0.63 |
| | All | 0.73 | 0.70 | 0.53 | 0.64 |
| Blinks | Image | 0.92 | 0.84 | | |
| | Dot | 0.77 | 0.71 | | |
| | Video | 0.82 | 0.84 | | |
| | All | 0.91 | 0.83 | | |

Naive Segmented Linear Regression-Hidden Markov models (NSLR-HMM) and Larsson, Nyström, and Stridh (LNS) values were taken from Pekkanen and Lappi (2017) and Andersson et al. (2017), respectively. For U'n'Eye, values are the median across 20 independent networks. In bold are the values reached by the best performing algorithm. MN, initial of the expert labeling the dataset (anonymous); PSO, postsaccadic oscillation.

Next, we studied how well an already-trained network can be applied to label saccades in new data, for which no training labels are available. Our results show that this is possible if the new data has broadly the same signal characteristics as the data used for training the network (i.e., if it was sampled with the same eye tracker during a sufficiently similar task). To illustrate this, we trained networks on our first three data sets and evaluated their performance on each data set plus an additional data set 4 on which none of the networks was trained. Data set 4 was similar to data set 1 in that it was recorded with the same eye tracker in human subjects performing fixations (Table 1). Therefore, a network trained on data set 1 performed very well not only in detecting saccades in the same data set, but also in data set 4 (Fig. 7D). Overall, good performance was guaranteed when the test set exhibited similar statistics as the training set or was exposed in training to a sufficiently wide variety of training samples (Fig. 7D).

Likewise, our network extrapolated well over subjects, for example in large cohort studies with many different observers (as is often the case in clinical investigations of neurological diseases). We studied whether a network trained on data from one subject was able to detect saccades well in data from another subject. To this end, we trained separate networks on data from 10 individual human subjects in data set 4 and applied them to all other subjects. Overall, performance on data that came from the same subject as the training data were only marginally higher than performance on data that came from a different subject (F1 mean and SD: 0.96 and 0.01 vs. 0.92 and 0.08, Fig. 7E). The higher standard deviation of intersubject performance was due to the apparent difference between data from certain subjects (Fig. 7E). We therefore advise users to combine training data from a few subjects to obtain a network that is able to deal with different signal statistics (Fig. 7E, network trained on all). Note that for the network trained on a combination of subjects, we made sure to keep the number of training samples the same as for networks trained on individual subjects. Thus, the better performance was a result of having more variable samples in the training set and not of more training examples being available.

Eye movement representation becomes disentangled along network layers. We finally had a closer look at how the network achieves the separation of two eye states (e.g., fixations and saccades; Fig. 8A). In the velocity domain, saccades and fixations can show highly overlapping distributions (Fig. 8B). This explains why velocity threshold-based algorithms can fail to distinguish fixations from saccades (Fig. 4). Here, we showed that U'n'Eye can differentiate between fixations and saccades with high accuracy (Fig. 4). The classification was based on the output layer of the network. To illustrate how this decision arises throughout the hidden layers, we performed principal component analysis (PCA) on the features of each convolutional layer. The fraction of explained variance by the first two principal components (PCs) reflects the U-shaped architecture of the network (Fig. 8C): in the middle layers, information is distributed across more components than in early and late layers. We projected the hidden layer activations onto the PC space and labeled time bins according to their ground truth labels (fixation or saccade, Fig. 8D). We observed in higher layers that the two classes were better separated (Fig. 8D). Finally, in the output layer, fixations and saccades became linearly separable (Fig. 8E). Thus, through training, the net-

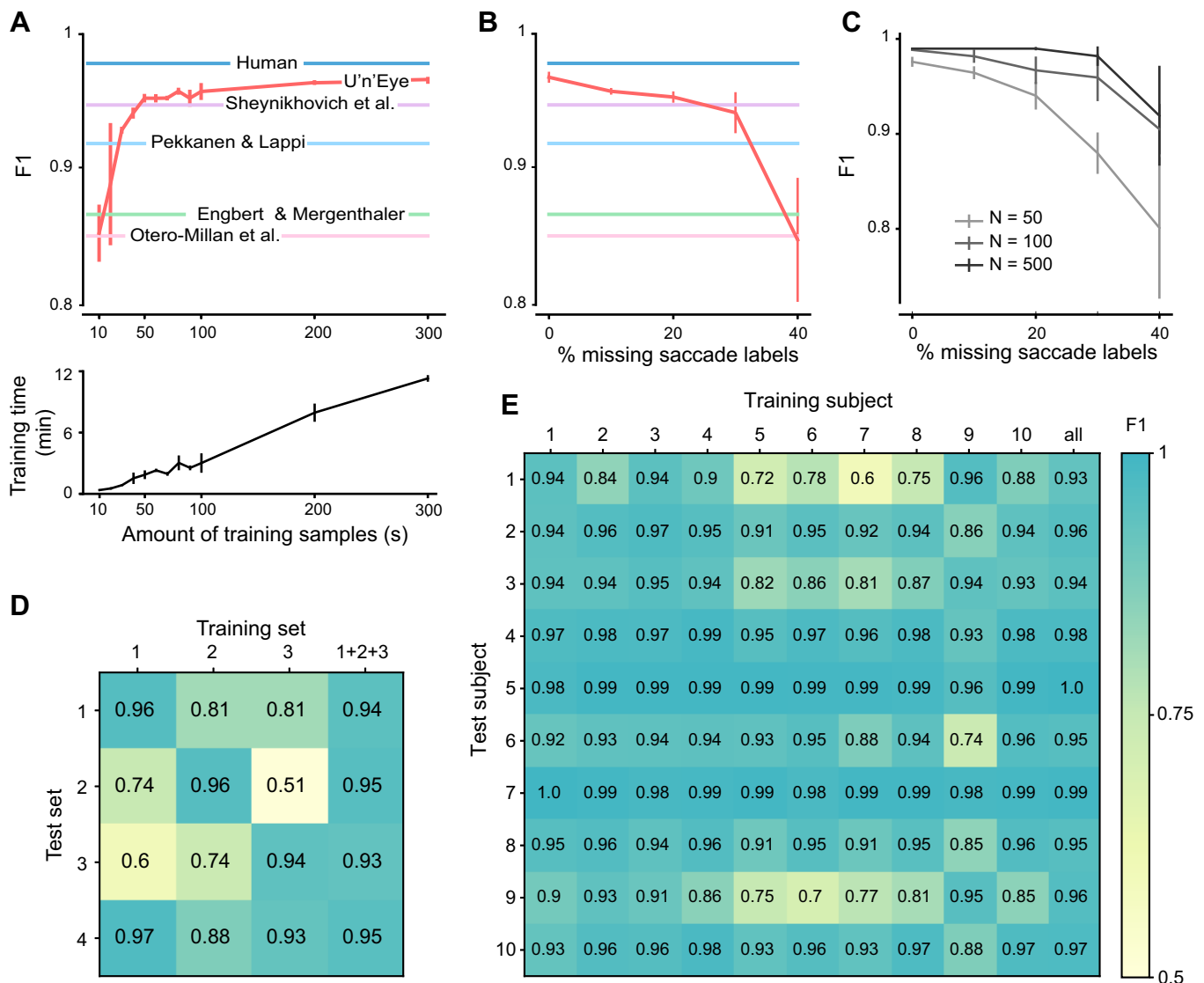


Fig. 7. Robustness of U'n'Eye performance under a variety of training regimes. *A, top*: U'n'Eye saccade detection performance (red) as a function of amount of training samples. *Bottom*: linear increase of training time with the number of training samples on a CPU. Data shows mean \pm standard deviation across the same training epochs as in the *top* panel. The colored horizontal lines with labels refer to the performance of other algorithms from the literature that we tested. *B*: U'n'Eye saccade detection performance as a function of the fraction of saccade labels missing in 300 s of training data. The colored horizontal lines refer to the same algorithms as in *A*. Also, the red line in *A* and *B* shows mean \pm standard deviation across networks trained on three different subsets of data set 1. Performance of U'n'Eye and other algorithms was evaluated on 1,000 s of test data from data set 1. *C*: U'n'Eye saccade detection performance in simulated data with missing labels for different amounts of training samples N in seconds. *D*: U'n'Eye saccade detection performance for different combinations of training and test sets. Each number in a square (and its associated color code) indicates the F1 score for training on one data set and testing on another. The column labeled 1+2+3 on the far right shows results when the network was trained on all three data sets simultaneously (but ensuring the same amount of training data as in the other columns of the figure). *E*: U'n'Eye saccade detection performance for combinations of different human subjects in training and test data from data set 4. Each column shows the results of training on a single subject from the data set, in a similar format to *D*. The final column on the right indicates that training the network on a (small) population of subjects yields best performance, and the rest of the figure indicates that additional subjects can then be tested with the pretrained network without much loss in performance. The same color scale was used for *D* and *E*.

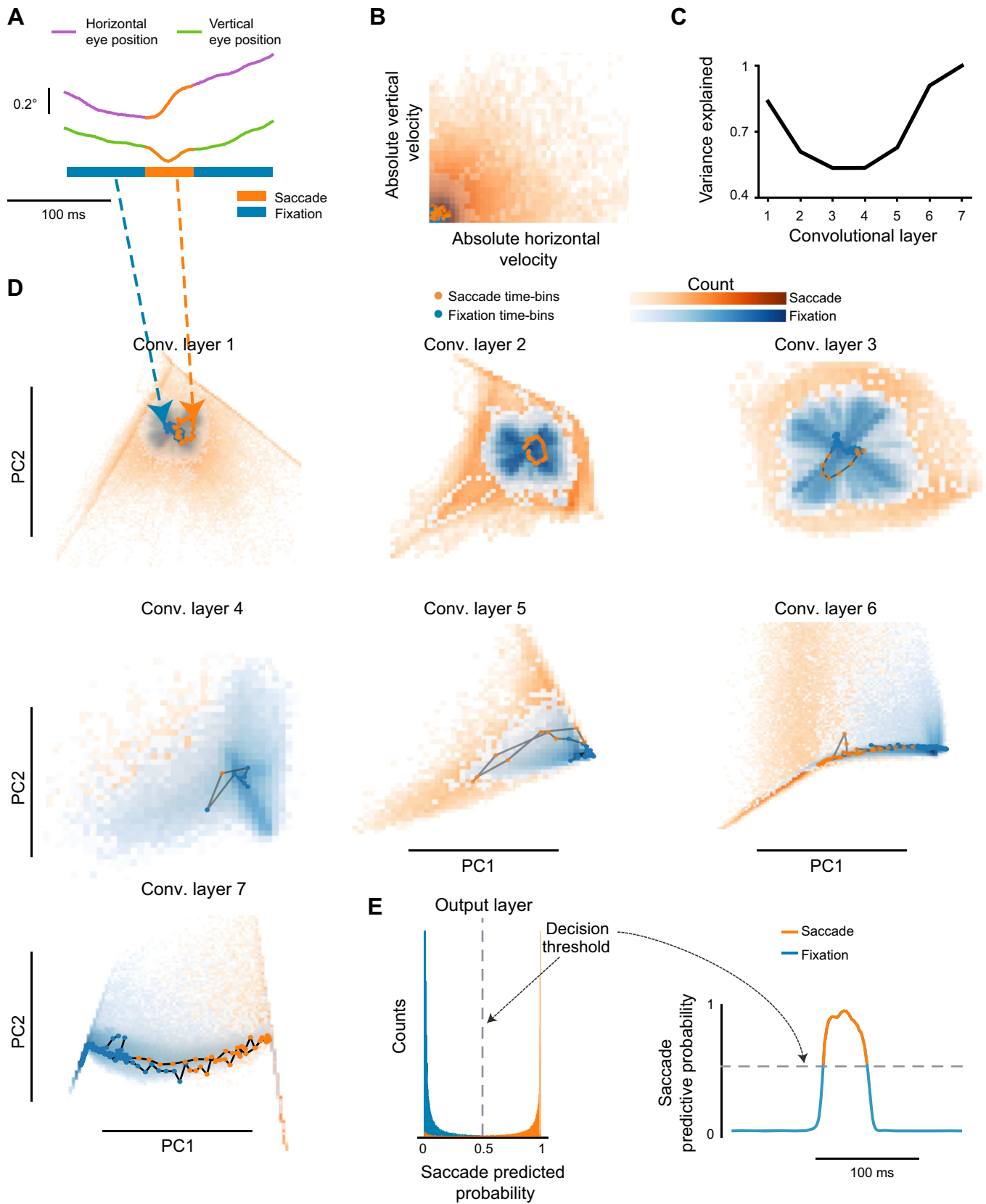
work effectively learns to extract relevant features and to project those onto a plane where the two eye movement classes are linearly separable.

DISCUSSION

In this study, we presented U'n'Eye, a convolutional neural network for eye movement classification. We demonstrated that U'n'Eye achieved human-level performance in the detection of saccades and microsaccades. In addition, the network was able to predict other classes of eye movements, which we

exemplified with the detection of blinks and PSOs in a benchmark data set.

Furthermore, we showed that U'n'Eye achieved excellent performance both when trained on a single type of data with labels from one coder and when trained on different data sets with labels from two coders. While data sets 1 and 3 used in this study contained data with only one type of visual task and labels from one coder each, data set 2 was composed of two different pursuit tasks and contained labels from two different human coders. Moreover, data set 4 allowed us to conclude that



our algorithm has generalization properties and can therefore be used when training on a subset of individuals and then testing with large cohorts of subsequent subjects measured with similar eye tracking technology. The data set by Andersson et al. (2017) also contained greatly varying types of saccades and other eye movements. Still, U'n'Eye achieved good performance when trained and tested on this data set. Note that the network might fail to detect eye movements when tested on data that show a very different distribution than the data it was trained on. We therefore recommend to either train a network with a variety of data or to train separate specialized networks for each task.

In this regard, our approach falls in the class of supervised learning algorithms, as opposed to methods not requiring parameter estimation based on annotated data (Engbert and Mergenthaler 2006; Otero-Millan et al. 2014; Sheynikhovich et al. 2018). However, we typically see in different scenarios that casting an algorithmic issue as a supervised problem helps in terms of performance. For example, we recently showed that supervised techniques perform as well as, or better than, unsupervised ones for spike inference from calcium imaging data (Berens et al. 2018; Theis et al. 2016). Similarly, Mathis et al. (2018) recently showed that supervised learning provides superior animal tracking with few annotated samples. We showed here that the situation is similar for eye movement detection. Importantly, we showed that performance generalizes to new unseen data sets and subjects, yielding better performance than any of the unsupervised algorithms. Of course, there is some manual work involved in preparing the training samples for our network, but we posit here that this amount of manual work is significantly less intensive than the manual postprocessing that we typically perform with other saccade-detection algorithms.

U'n'Eye is publicly available and provides a user friendly interface as well as a web service in which users can upload their data and receive classification outputs (see METHODS). No parameter tuning is needed even for training (e.g., learning rate, and so on) since the standard settings were found to work well across data sets. Instead, an experimenter just needs to provide a few hundred seconds of labeled data to train the network once. Even if some labels are missing in the training data, U'n'Eye can still reach high performance. We recommend, however, to use only carefully annotated data for training, as this should improve results.

Of the few algorithms that are capable of detecting saccades as well as PSO (Larsson et al. 2013; Pekkanen and Lappi 2017; Zemblys et al. 2018b), U'n'Eye achieves highest performance. Note that Zemblys et al. (2018a) also recently proposed a deep learning method for eye movement detection. Their approach consists of generating a large training set out of a small human-labeled data set using a generative neural network. A second network is then trained on this data to classify eye movements. This method reports performance similar to that of

U'n'Eye in a subset of the benchmark data set by Andersson et al. (2017), but it remains to be seen how this algorithm performs on more exhaustive tasks like the ones that we reported here. For example, the applicability to data containing smooth pursuit has not been demonstrated. Conversely Startsev et al. (2018) recently published a deep learning approach showing reasonable performance, but again they tested only on a subset of the benchmark data set containing smooth pursuit.

Recently, a Bayesian approach for the detection of microsaccades based on a generative model has been proposed (Mihali et al. 2017). Inherently, Bayesian methods provide estimates of uncertainty, in addition to estimates of the quantity of interest. Indeed, it is an interesting future perspective to combine U'n'Eye with Bayesian Deep Learning techniques to provide uncertainty estimates for the detected eye movements (Gal and Ghahramani 2015).

Future work should include combining data sets with different characteristics, such as different sampling frequencies, to obtain a network that can generalize on a large range of data. Such a network could be used by a large part of the scientific community, which would allow for reproducibility of scientific results. We recommend that anyone who uses our algorithm to publish the weights of the trained network so that eye movement detection can be reproduced. For our own trained networks, all weights have been published online (<https://github.com/berenslab/uneye>) along with the code of the network. This has the advantage that users with similar data characteristics to one of our three data sets (e.g., microsaccades during fixation with a video-based eye tracker as in data set 3) can directly use our weights from the proper data set without having to retrain their own network. We also intend to make all three data sets publicly available, facilitating the further development for eye movement detection algorithms.

Of course, it should be noted that some prediction errors may still occur with U'n'Eye. However, such errors fall within the range of interrater variability across humans anyway. Also, even when U'n'Eye does make mistakes, the predictive probability that it outputs can be used to retrieve missed events (e.g., see the upward black arrow in the bottom of Fig. 6D). For example, detecting peaks in the predictive probability output that did not cross the threshold can accelerate eventual manual postprocessing.

Finally, U'n'Eye's capacity to learn nonlinear relationships between an eye trace and some annotated labels opens new horizons in neuroscience: the network could be used to understand the properties of neural activities related in a complex manner to eye movements. For example, the disentanglement in later layers (Fig. 8) could be used to quantitatively analyze the activity patterns of premotor neurons in the brain stem, which themselves ultimately transform brain processing into individual ocular muscle innervations. Furthermore, U'n'Eye could be turned into a generative model for eye movements, as was shown for neural networks that are used for image clas-

Fig. 8. Disentanglement of fixations and saccades throughout the network. *A*: example eye trace with a microsaccade. *B*: distribution of data set 2 in the velocity domain. Fixations and saccades (shown in bluish and orangish colors, respectively) showed overlapping distributions. *C*: fraction of explained variance by the two first principal components (PCs) of the network's convolutional (Conv.) layers. There was a reduction in the middle layers followed by a peak at the final seventh layer. *D*: projection of hidden layer activations by eye traces of data set 2 onto the first two principal components. Fixation and saccade classes became better separated throughout the hidden layers. *B–D*: Dots indicate the time points of the example eye trace in *A*, and the rest of the background data show the entire data set time samples. *E*: the probability output allowed for a linear separation of the two classes. Time points with a saccade predictive probability above 0.5 were classified as a saccade.

sification (Gatys et al. 2015). The information about eye movements that is contained in the network architecture might in the future be used to identify variations in eye movement characteristics that could hint at underlying pathologies.

ACKNOWLEDGMENTS

We thank Konstantin-Friedrich Willeke and Antimo Buonocore for providing help with labeling saccade data, Murat Ayhan for input on deep neural networks, and Adam von Daranyi for setting up the web service.

GRANTS

This work was funded by the German Ministry of Education and Research (FKZ 01GQ1601) and the German Research Foundation (EXC307, EXC2064/1–Project ID 390727645, SFB 1233–Project ID 276693517, BE5601/4-1).

DISCLOSURES

No conflicts of interest, financial or otherwise, are declared by the authors.

AUTHOR CONTRIBUTIONS

M.E.B., J.B., Z.M.H., and P.B. conceived and designed research; M.E.B. and J.B. analyzed data; M.E.B., J.B., Z.M.H., and P.B. interpreted results of experiments; M.E.B. prepared figures; M.E.B. and J.B. drafted manuscript; M.E.B., J.B., Z.M.H., and P.B. edited and revised manuscript; M.E.B., J.B., H.N., Z.M.H., and P.B. approved final version of manuscript; J.B. and H.N. performed experiments.

ENDNOTE

At the request of the authors, readers are herein alerted to the fact that additional materials related to this manuscript may be found at the institutional website of one of the authors, which at the time of publication they indicate is: at <http://uneye.berenslab.org>. All code is available from <https://github.com/berenslab/uneye>. These materials are not a part of this manuscript and have not undergone peer review by the American Physiological Society (APS). APS and the journal editors take no responsibility for these materials, for the website address, or for any links to or from it.

REFERENCES

- Andersson R, Larsson L, Holmqvist K, Stridh M, Nyström M. One algorithm to rule them all? An evaluation and discussion of ten eye movement event-detection algorithms. *Behav Res Methods* 49: 616–637, 2017. doi:10.3758/s13428-016-0738-9.
- Bellet J, Chen CY, Hafed ZM. Sequential hemifield gating of α - and β -behavioral performance oscillations after microsaccades. *J Neurophysiol* 118: 2789–2805, 2017. doi:10.1152/jn.00253.2017.
- Berens P, Freeman J, Deneux T, Chenkov N, McColgan T, Speiser A, Macke JH, Turaga SC, Mineault P, Rupprecht P, Gerhard S, Friedrich RW, Friedrich J, Paninski L, Pachitariu M, Harris KD, Bolte B, Machado TA, Ringach D, Stone J, Rogerson LE, Sofroniew NJ, Reimer J, Froudarakis E, Euler T, Román Rosón M, Theis L, Tolias AS, Bethge M. Community-based benchmarking improves spike rate inference from two-photon calcium imaging data. *PLOS Comput Biol* 14: e1006157, 2018. doi:10.1371/journal.pcbi.1006157.
- Bishop CM. *Pattern Recognition and Machine Learning* New York: Springer, 2016.
- Borji A, Itti L. Defending Yarbus: eye movements reveal observers' task. *J Vis* 14: 29, 2014. doi:10.1167/14.3.29.
- Bosman CA, Womelsdorf T, Desimone R, Fries P. A microsaccadic rhythm modulates gamma-band synchronization and behavior. *J Neurosci* 29: 9471–9480, 2009. doi:10.1523/JNEUROSCI.1193-09.2009.
- Buonocore A, Skinner J, Hafed ZM. Eye-position error influence over “open-loop” smooth pursuit initiation. *bioRxiv*: 404491, 2018. doi:10.1101/404491.
- Burr DC, Morrone MC, Ross J. Selective suppression of the magnocellular visual pathway during saccadic eye movements. *Nature* 371: 511–513, 1994. doi:10.1038/371511a0.
- Carpenter RH. *Movements of the Eyes* (2nd rev. ed.) London: Pion, 1988.
- Chen CY, Hafed ZM. Postmicrosaccadic enhancement of slow eye movements. *J Neurosci* 33: 5375–5386, 2013. doi:10.1523/JNEUROSCI.3703-12.2013.
- Chen CY, Hafed ZM. A neural locus for spatial-frequency specific saccadic suppression in visual-motor neurons of the primate superior colliculus. *J Neurophysiol* 117: 1657–1673, 2017. doi:10.1152/jn.00911.2016.
- Cohen J. A coefficient of agreement for nominal scales. *Educ Psychol Meas* 20: 37–46, 1960. doi:10.1177/001316446002.000104.
- Crevecoeur F, Kording KP. Saccadic suppression as a perceptual consequence of efficient sensorimotor estimation. *eLife* 6: e25073, 2017. doi:10.7554/eLife.25073.
- Dai W, Selesnick I, Rizzo JR, Rucker J, Hudson T. A parametric model for saccadic eye movement. Signal Processing in Medicine and Biology Symposium (SPMB). Philadelphia, PA, December 3, 2016. doi:10.1109/SPMB.2016.7846860.
- Daye PM, Optican LM. Saccade detection using a particle filter. *J Neurosci Methods* 235: 157–168, 2014. doi:10.1016/j.jneumeth.2014.06.020.
- Duchowski AT. *Eye Tracking Methodology: Theory and Practice* (2nd ed.). London: Springer, 2007.
- Duhamel JR, Colby CL, Goldberg ME. The updating of the representation of visual space in parietal cortex by intended eye movements. *Science* 255: 90–92, 1992. doi:10.1126/science.1553535.
- Engbert R, Mergenthaler K. Microsaccades are triggered by low retinal image slip. *Proc Natl Acad Sci USA* 103: 7192–7197, 2006. doi:10.1073/pnas.0509557103.
- Fuchs AF, Robinson DA. A method for measuring horizontal and vertical eye movement chronically in the monkey. *J Appl Physiol* 21: 1068–1070, 1966. doi:10.1152/jappl.1966.21.3.1068.
- Gal Y, Ghahramani Z. Bayesian convolutional neural networks with Bernoulli approximate variational inference (Preprint). *arXiv* 1506.02158, 2015.
- Gatys LA, Ecker AS, Bethge M. A neural algorithm of artistic style (Preprint). *arXiv* 1508.06576, 2015.
- Golan T, Davidesco I, Meshulam M, Groppe DM, Mégevand P, Yeagle EM, Goldfinger MS, Harel M, Melloni L, Schroeder CE, Deouell LY, Mehta AD, Malach R. Increasing suppression of saccade-related transients along the human visual hierarchy. *eLife* 6: e27819, 2017. doi:10.7554/eLife.27819.
- Gur M, Beylin A, Snodderly DM. Response variability of neurons in primary visual cortex (V1) of alert monkeys. *J Neurosci* 17: 2914–2920, 1997. doi:10.1523/JNEUROSCI.17-08-02914.1997.
- Hafed ZM. Mechanisms for generating and compensating for the smallest possible saccades. *Eur J Neurosci* 33: 2101–2113, 2011. doi:10.1111/j.1460-9568.2011.07694.x.
- Hafed ZM. Alteration of visual perception prior to microsaccades. *Neuron* 77: 775–786, 2013. doi:10.1016/j.neuron.2012.12.014.
- Hafed ZM, Chen CY, Tian X. Vision, perception, and attention through the lens of microsaccades: mechanisms and implications. *Front Syst Neurosci* 9: 167, 2015. doi:10.3389/fnsys.2015.00167.
- Hafed ZM, Goffart L, Krauzlis RJ. Superior colliculus inactivation causes stable offsets in eye position during tracking. *J Neurosci* 28: 8124–8137, 2008. doi:10.1523/JNEUROSCI.1317-08.2008.
- Hafed ZM, Krauzlis RJ. Goal representations dominate superior colliculus activity during extrafoveal tracking. *J Neurosci* 28: 9426–9439, 2008. doi:10.1523/JNEUROSCI.1313-08.2008.
- Haji-Abolhassani A, Clark JJ. An inverse Yarbus process: predicting observers' task from eye movement patterns. *Vision Res* 103: 127–142, 2014. doi:10.1016/j.visres.2014.08.014.
- Hass CA, Horowitz GD. Effects of microsaccades on contrast detection and V1 responses in macaques. *J Vis* 11: 3, 2011. doi:10.1167/11.3.3.
- Herrington TM, Masse NY, Hachmeh KJ, Smith JE, Assad JA, Cook EP. The effect of microsaccades on the correlation between neural activity and behavior in middle temporal, ventral intraparietal, and lateral intraparietal areas. *J Neurosci* 29: 5793–5805, 2009. doi:10.1523/JNEUROSCI.4412-08.2009.
- Ioffe S, Szegedy C. Batch normalization: accelerating deep network training by reducing internal covariate shift (Preprint). *arXiv* 1502.03167, 2015.
- Judge SJ, Richmond BJ, Chu FC. Implantation of magnetic search coils for measurement of eye position: an improved method. *Vision Res* 20: 535–538, 1980. doi:10.1016/0042-6989(80)90128-5.
- Kawaguchi K, Clery S, Pourriahi P, Seillier L, Haefner RM, Nienborg H. Differentiating between models of perceptual decision making using pupil size inferred confidence. *J Neurosci* 38: 8874–8888, 2018. doi:10.1523/JNEUROSCI.0735-18.2018.

- Kingma DP, Ba J.** Adam: a method for stochastic optimization (Preprint). *arXiv* 1412.6980, 2014.
- Klibisz A, Rose D, Eicholtz M, Blundon J, Zakharenko S.** Fast, simple calcium imaging segmentation with fully convolutional networks. In: *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, edited by Cardoso MJ, Arbel T. Cham, Switzerland: Springer, 2017, p. 285–293. doi:10.1007/978-3-319-67558-9_33.
- Kowler E.** Eye movements: the past 25 years. *Vision Res* 51: 1457–1483, 2011. doi:10.1016/j.visres.2010.12.014.
- Larsson L, Nyström M, Stridh M.** Detection of saccades and postsaccadic oscillations in the presence of smooth pursuit. *IEEE Trans Biomed Eng* 60: 2484–2493, 2013. doi:10.1109/TBME.2013.2258918.
- Leigh RJ, Kennard C.** Using saccades as a research tool in the clinical neurosciences. *Brain* 127: 460–477, 2004. doi:10.1093/brain/awh035.
- Leigh RJ, Zee DS.** *The Neurology of Eye Movements*. New York: Oxford University Press, 2015. Contemporary Neurology Series 90.
- Leopold DA, Logothetis NK.** Microsaccades differentially modulate neural activity in the striate and extrastriate visual cortex. *Exp Brain Res* 123: 341–345, 1998. doi:10.1007/s002210050577.
- MacAskill MR, Anderson TJ.** Eye movements in neurodegenerative diseases. *Curr Opin Neurol* 29: 61–68, 2016. doi:10.1097/WCO.0000000000000274.
- Mathis A, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, Bethge M.** Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nat Neurosci* 21: 1281–1289, 2018. doi:10.1038/s41593-018-0209-y.
- Mihali A, van Opheusden B, Ma WJ.** Bayesian microsaccade detection. *J Vis* 17: 13, 2017. doi:10.1167/17.1.13.
- Otero-Millan J, Castro JLA, Macknik SL, Martinez-Conde S.** Unsupervised clustering method to detect microsaccades. *J Vis* 14: 18, 2014. doi:10.1167/14.2.18.
- Pekkanen J, Lappi O.** A new and general approach to signal denoising and eye movement classification based on segmented linear regression. *Sci Rep* 7: 17726, 2017. doi:10.1038/s41598-017-17983-x.
- Reppas JB, Usrey WM, Reid RC.** Saccadic eye movements modulate visual responses in the lateral geniculate nucleus. *Neuron* 35: 961–974, 2002. doi:10.1016/S0896-6273(02)00823-1.
- Ronneberger O, Fischer P, Brox T.** U-net: convolutional networks for biomedical image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, edited by Navab N, Hornegger J, Wells W, Frangi A. Cham, Switzerland: Springer, 2015, p. 234–241. Lecture Notes in Computer Science 9351.
- Ross J, Morrone MC, Burr DC.** Compression of visual space before saccades. *Nature* 386: 598–601, 1997. doi:10.1038/386598a0.
- Sheynikhovich D, Bécu M, Wu C, Arleo A.** Unsupervised detection of microsaccades in a high-noise regime. *J Vis* 18: 19, 2018. doi:10.1167/18.6.19.
- Sommer MA, Wurtz RH.** Brain circuits for the internal monitoring of movements. *Annu Rev Neurosci* 31: 317–338, 2008. doi:10.1146/annurev.neuro.31.060407.125627.
- Startsev M, Agtzidis I, Dorr M.** 1D CNN with BLSTM for automated classification of fixations, saccades, and smooth pursuits. *Behav Res Methods*, 2018. doi:10.3758/s13428-018-1144-2.
- Theis L, Berens P, Froudarakis E, Reimer J, Román Rosón M, Baden T, Euler T, Tolias AS, Bethge M.** Benchmarking spike rate inference in population calcium imaging. *Neuron* 90: 471–482, 2016. doi:10.1016/j.neuron.2016.04.014.
- Tian X, Yoshida M, Hafed ZM.** A microsaccadic account of attentional capture and inhibition of return in posner cueing. *Front Syst Neurosci* 10: 23, 2016. doi:10.3389/fnsys.2016.00023.
- Yao T, Treue S, Krishna BS.** Saccade-synchronized rapid attention shifts in macaque visual cortical area MT. *Nat Commun* 9: 958, 2018. doi:10.1038/s41467-018-03398-3.
- Yu G, Yang M, Yu P, Dorris MC.** Time compression of visual perception around microsaccades. *J Neurophysiol* 118: 416–424, 2017. doi:10.1152/jn.00029.2017.
- Zemblys R, Niehorster DC, Holmqvist K.** gazeNet: End-to-end eye-movement event detection with deep neural networks. *Behav Res Methods*, 2018a. doi:10.3758/s13428-018-1133-5.
- Zemblys R, Niehorster DC, Komogortsev O, Holmqvist K.** Using machine learning to detect events in eye-tracking data. *Behav Res Methods* 50: 160–181, 2018b. doi:10.3758/s13428-017-0860-3.
- Zirnsak M, Steinmetz NA, Noudoost B, Xu KZ, Moore T.** Visual space is compressed in prefrontal cortex before eye movements. *Nature* 507: 504–507, 2014. doi:10.1038/nature13149.
- Zuber BL, Stark L, Cook G.** Microsaccades and the velocity-amplitude relationship for saccadic eye movements. *Science* 150: 1459–1460, 1965. doi:10.1126/science.150.3702.1459.