

1 Convolutional neural networks can decode eye movement data: A black box approach to
2 predicting task from eye movements

³ Zachary J. Cole¹, Karl M. Kuntzman¹, Michael D. Dodd¹, & Matthew R. Johnson¹

⁴ ¹ University of Nebraska-Lincoln

5 Author Note

The data used for the exploratory and confirmatory analyses in the present manuscript are derived from experiments funded by NIH/NEI Grant 1R01EY022974 to MDD. Additionally, work done to develop the analysis approach was supported by NSF/EPSCoR grant #1632849 and NIH grant GM130461 awarded to MRJ and colleagues.

¹⁰ Correspondence concerning this article should be addressed to Zachary J. Cole, 238
¹¹ Burnett Hall, Lincoln, NE 68588-0308. E-mail: z@neurophysicole.com

12

Abstract

13 Previous attempts to classify task from eye movement data have relied on model
14 architectures designed to emulate theoretically defined cognitive processes, and/or data that
15 has been processed into aggregate (e.g., fixations, saccades) or statistical (e.g., fixation
16 density) features. *Black box* convolutional neural networks (CNNs) are capable of identifying
17 relevant features in raw and minimally processed data and images, but difficulty interpreting
18 the mechanisms underlying these model architectures have contributed to challenges in
19 generalizing lab-trained CNNs to applied contexts. In the current study, a CNN classifier
20 was used to classify task from two eye movement datasets (Exploratory and Confirmatory)
21 in which participants searched, memorized, or rated indoor and outdoor scene images. The
22 Exploratory dataset was used to tune the hyperparameters of the model, and the resulting
23 model architecture was re-trained, validated, and tested on the Confirmatory dataset. The
24 data were formatted into raw timeline data (i.e., x-coordinate, y-coordinate, pupil size) and
25 minimally processed images. To further understand the relative informational value of the
26 raw components of the eye movement data, the timeline and image datasets were broken
27 down into subsets with one or more of the components of the data systematically removed.
28 Average classification accuracies were compared between datasets and subsets. Classification
29 of the timeline data consistently outperformed the image data. The Memorize condition was
30 most often confused with the Search and Rate conditions. Pupil size was the least uniquely
31 informative eye movement component when compared with the x- and y-coordinates. The
32 general pattern of results for the Exploratory dataset was replicated in the Confirmatory
33 dataset. Overall, the present study provides a practical and reliable black box solution to the
34 inverse Yarbus problem.

35 *Keywords:* deep learning, eye tracking, convolutional neural network, cognitive state,
36 endogenous attention

37 Word count: 7260

38 Convolutional neural networks can decode eye movement data: A black box approach to
39 predicting task from eye movements

40 **Background**

41 The association between eye movements and mental activity is a fundamental topic of
42 interest in attention research that has provided a foundation for developing a wide range of
43 human assistive technologies. Early work by Yarbus (1967) showed that eye movement
44 patterns appear to differ qualitatively depending on the task-at-hand (for a review of this
45 work, see Tatler, Wade, Kwan, Findlay, & Velichkovsky, 2010). A replication of this work by
46 DeAngelus and Pelz (2009) shows that the differences in eye movements between tasks can
47 be quantified, and appear to be somewhat generalizable. Technological advances and
48 improvements in computing power have allowed researchers to make inferences regarding the
49 mental state underlying eye movement data, also known as the “inverse Yarbus process”
50 (Haji-Abolhassani & Clark, 2014). Current state-of-the-art machine learning and neural
51 network algorithms are capable of identifying diagnostic patterns for the purpose of decoding
52 a variety of data types, but the inner workings of the resulting model solutions are difficult
53 or impossible to interpret. Algorithms that provide such solutions are referred to as *black box*
54 models. Dissections of black box models have been largely uninformative (Zhou, Bau, Oliva,
55 & Torralba, 2019), limiting the potential for researchers to apply the mechanisms underlying
56 successful classification of the data. Still, black box models provide a powerful solution for
57 technological applications such as human-computer interfaces (HCI; for a review, see
58 Lukander, Toivanen, & Puolamäki, 2017). While the internal operations of the model
59 solutions used for HCI applications do not necessarily need to be interpretable to serve their
60 purpose, Lukander et al. (2017) pointed out that the inability to interpret the mechanisms
61 underlying the function of black box solutions impedes the generalizability of these methods,
62 and increases the difficulty of expanding these findings to real life applications. To ground
63 these solutions, researchers guide decoding efforts by using eye movement data and/or
64 models with built-in theoretical assumptions. For instance, eye movement data is processed

65 into meaningful aggregate properties such as fixations or saccades, or statistical features such
66 as fixation density, and the models used to decode these data are structured based on the
67 current understanding of relevant cognitive or neurobiological processes (e.g., MacInnes,
68 Hunt, Clarke, & Dodd, 2018). Despite the proposed disadvantages of black box approaches
69 to classifying eye movement data, there is no clear evidence to support the notion that the
70 grounded solutions described above are actually more valid or definitive than a black box
71 solution.

72 The scope of theoretically informed solutions to decoding eye movement data are
73 limited to the extent of the current theoretical knowledge linking eye movements to cognitive
74 and neurobiological processes. As our theoretical understanding of these processes develops,
75 older theoretically informed models become outdated. Furthermore, these solutions are
76 susceptible to any inaccurate preconceptions that are built into the theory. Consider the case
77 of Greene, Liu, and Wolfe (2012), who were not able to classify task from commonly used
78 aggregate eye movement features (i.e., number of fixations, mean fixation duration, mean
79 saccade amplitude, percent of image covered by fixations) using correlations, a linear
80 discriminant model, and a support vector machine (see Table 1). This led Greene and
81 colleagues to question the robustness of Yarbus's (1967) findings, inspiring a slew of
82 responses that successfully decoded the same dataset by aggregating the eye movements into
83 different feature sets or implementing different model architectures (see Table 1; i.e.,
84 Haji-Abolhassani & Clark, 2014; Borji & Itti, 2014; Kanan, Ray, Bseiso, Hsiao, & Cottrell,
85 2014). The subsequent re-analyses of these data support Yarbus (1967) and the notion that
86 mental state can be decoded from eye movement data using a variety of combinations of
87 data features and model architectures. Collectively, these re-analyses did not point to an
88 obvious global solution capable of clarifying future approaches to the inverse Yarbus problem
89 beyond what could be inferred from black box model solutions, but did provide a rigorous
90 test of a variety of methodological features which can be applied to theoretical or black box
91 approaches to the inverse Yarbus problem.

Eye movements can only delineate tasks to the extent that the cognitive processes underlying the tasks can be differentiated (Król & Król, 2018). Every task is associated with a unique set of cognitive processes (Coco & Keller, 2014; Król & Król, 2018), but in some cases, the cognitive processes for different tasks may produce indistinguishable eye movement patterns. To differentiate the cognitive processes underlying task-evoked eye movements, some studies have chosen to classify tasks that rely on stimuli that prompt easily distinguishable eye movements, such as reading text (e.g., Henderson, Shinkareva, Wang, Luke, & Olejarczyk, 2013). The eye movements elicited by salient stimulus features facilitate task classifications, however, because these eye movements are the consequence of a feature (or features) inherent to the stimulus rather than the task, it is unclear if these classifications are attributable to the stimulus or a complex mental state (e.g., Henderson et al., 2013; Boisvert & Bruce, 2016). Additionally, the distinct nature of exogenously elicited eye movements prompts decoding algorithms to prioritize these bottom-up patterns in the data over higher-level top-down effects (Borji & Itti, 2014). This means that these models are identifying the type of information that is being processed, but are not necessarily reflecting the mental state of the individual observing the stimulus. Eye movements that are the product of bottom-up processes have been reliably decoded, which is relevant for some HCI applications, but does not fit the nature of the inverse Yarbus problem, which is concerned with decoding high-level abstract mental operations that are not dependent on particular stimuli.

Currently, an upper limit to how well cognitive task can be classified from eye movement data has not been clearly established. Prior evidence has shown that the task-at-hand is capable of producing distinguishable eye movement features such as the total scan path length, total number of fixations, and the amount of time to the first saccade (Castelhano, Mack, & Henderson, 2009; DeAngelus & Pelz, 2009). Decoding accuracies within the context of determining task from eye movements typically range from chance performance (between 14.29% and 33%) to 59.64% (when chance was 25%; see Table 1). In

one case, Coco and Keller (2014) categorized the same eye movement features used by Greene et al. (2012) with respect to the relative contribution of latent visual or linguistic components of three tasks (visual search, name the picture, name objects in the picture) with 84% accuracy (chance = 33%). While this manipulation is reminiscent of other experiments relying on the bottom-up influence of words and pictures (e.g., Henderson et al., 2013; Boisvert & Bruce, 2016) the eye movements in the Coco and Keller (2014) tasks can be attributed to the occurrence of top-down attentional processes. A conceptually similar follow-up to this study classified tasks along two spatial and semantic dimensions, resulting in 51% classification accuracy (chance = 25%; Król & Król, 2018). A closer look at these results showed that the categories within the semantic dimension were consistently misclassified, suggesting that this level of distinction may require a richer dataset, or a more powerful decoding algorithm. Altogether, there is no measurable index of relative top-down or bottom-up influence, but this body of literature suggests that the relative influence of top-down and bottom-up attentional processes may have a role in determining the decodability of the eye movement data.

As shown in Table 1, when eye movement data are prepared for classification, fixation and saccade statistics are typically aggregated along spatial or temporal dimensions, resulting in variables such as fixation density or saccade amplitude (Castelhano et al., 2009; MacInnes et al., 2018; Mills, Hollingworth, Van der Stigchel, Hoffman, & Dodd, 2011). The implementation of these statistical methods is meant to explicitly provide the decoding algorithm with characteristics of the eye movement data that are representative of theoretically relevant cognitive processes. For example, MacInnes et al. (2018) attempted to provide an algorithm with data designed to be representative of inputs to the frontal eye fields. In some instances, such as the case of Król and Król (2018), grounding the data using theoretically driven aggregation methods may require sacrificing granularity in the dataset. This means that aggregating the data has the potential to wash out certain fine-grained distinctions that could otherwise be detected. Data structures of any kind can only be

Table 1

Previous Attempts to Classify Cognitive Task Using Eye Movement Data

Study	Tasks	Features	Model Architecture	Accuracy (Chance)
Greene et al. (2012)	memorize, decade, people, wealth	number of fixations, mean fixation duration, mean saccade amplitude, percent of image covered by fixations, dwell times	linear discriminant, correlation, SVM	25.9% (25%)
Haji-Abolhassani & James (2014)	Greene et al. tasks	fixation clusters	Hidden Markov Models	59.64% (25%)
Kanan et al. (2014)	Greene et al. tasks	mean fixation durations, number of fixations	multi-fixation pattern analysis	37.9% (25%)
Borji & Itti (2014)	Greene et al. tasks	number of fixations, mean fixation duration, mean saccade amplitude, percent of image covered by fixations, first five fixations, fixation density	kNN, RUSBoost	34.34% (25%)
Borji & Itti (2014)	Yarbus tasks (i.e., view, wealth, age, prior activity, clothes, location, time away)	number of fixations, mean fixation duration, mean saccade amplitude, percent of image covered by fixations, first five fixations, fixation density	kNN, RUSBoost	24.21% (14.29%)
Coco & Keller (2014)	search, name picture, name object	Greene et al. features, latency of first fixation, first fixation duration, mean fixation duration, total gaze duration, initiation time, mean saliency at fixation, entropy of attentional landscape	MM, LASSO, SVM	84% (33%)
MacInnes et al. (2018)	view, memorize, search, rate	saccade latency, saccade duration, saccade amplitude, peak saccade velocity, absolute saccade angle, pupil size	augmented Naive Bayes Network	53.9% (25%)
Król & Król (2018)	people, indoors/outdoors, white/black, search	eccentricity, screen coverage	feed forward neural network	51.4% (25%)

¹⁴⁶ decoded to the extent to which the data are capable of representing differences between¹⁴⁷ categories. Given that the cognitive processes underlying distinct tasks are often overlapping

¹⁴⁸ (Coco & Keller, 2014), decreasing the granularity of the data may actually limit the potential
¹⁴⁹ of the algorithm to make fine-grained distinctions between diagnostic components underlying
¹⁵⁰ the target task and the other tasks.

¹⁵¹ The current state of the literature does not provide any firm guidelines for determining
¹⁵² what eye movement features are most meaningful, or what model architectures are most
¹⁵³ suited for determining mental state from eye movements. The examples provided in Table 1
¹⁵⁴ used a variety of eye movement features and model architectures, most of which were
¹⁵⁵ effective to some extent. A proper comparison of these outcomes is difficult because these
¹⁵⁶ datasets vary in levels of chance and data quality. Datasets with more tasks to be classified
¹⁵⁷ have lower levels of chance, lowering the threshold for successful classification. Additionally,
¹⁵⁸ datasets with a lower signal-to-noise ratio will have a lower achievable classification accuracy.
¹⁵⁹ For these reasons, outside of re-analyzing the same datasets, there is no consensus on how to
¹⁶⁰ establish direct comparisons of these model architectures. Given the inability to directly
¹⁶¹ compare the relative effectiveness of the various theoretical approaches present in the
¹⁶² literature, the current study addressed the inverse Yarbus problem by allowing a black box
¹⁶³ model to self-determine the most informative features from minimally processed eye
¹⁶⁴ movement data.

¹⁶⁵ The current study explored pragmatic solutions to the problem of classifying task from
¹⁶⁶ eye movement data by submitting unprocessed x-coordinate, y-coordinate, and pupil size
¹⁶⁷ data to a convolutional neural network (CNN) model. Instead of transforming the data into
¹⁶⁸ theoretically defined units, we allowed the network to learn meaningful patterns in the data
¹⁶⁹ on its own. CNNs have a natural propensity to develop low-level feature detectors similar to
¹⁷⁰ primary visual cortex (e.g., Seeliger et al., 2018); for this reason, they are commonly
¹⁷¹ implemented for image classification. To test the possibility that the image data are better
¹⁷² suited to the CNN classifier, the data were also transformed into raw timeline and simple
¹⁷³ image representations. To our knowledge, no study has attempted to address the inverse

174 Yarbus problem using any combination of the following methods: (1) Non-aggregated data,
175 (2) image data format, and (3) a black-box CNN architecture. Given that CNN architectures
176 are capable of learning features represented in raw data formats, and are well-suited to
177 decoding multidimensional data that have a distinct spatial or temporal structure, we
178 expected that a non-theoretically-constrained CNN architecture could be capable of decoding
179 data at levels consistent with the current state of the art. Furthermore, despite evidence that
180 black box approaches to the inverse Yarbus problem can impede generalizability (Lukander
181 et al., 2017), we expected that when testing the approach on an entirely separate dataset,
182 providing the model with minimally processed data and the flexibility to identify the unique
183 features within each dataset would result in the replication of our initial findings.

184

Methods

185 **Participants**

186 Two separate datasets were used to develop and test the deep CNN architecture. The
187 two datasets were collected from two separate experiments, which we refer to as Exploratory
188 and Confirmatory. The participants for both datasets consisted of college students
189 (Exploratory $N = 124$; Confirmatory $N = 77$) from the University of Nebraska-Lincoln who
190 participated in exchange for class credit. Participants who took part in the Exploratory
191 experiment did not participate in the Confirmatory experiment. All procedures and
192 materials were approved by the University of Nebraska-Lincoln Institutional Review Board
193 prior to data collection.

194 **Materials and Procedures**

195 Each participant viewed a series of indoor and outdoor scene images while carrying out
196 a search, memorization, or rating task. For the search task, participants were instructed to
197 find a small “Z” or “N” embedded in the image. Trials containing a target ($n = 5$) were not
198 analyzed but were included in the experiment design to encourage searching behavior on

199 other Search trials. If the letter was found, the participants were instructed to press a
200 button, which terminated the trial. For the memorization task, participants were instructed
201 to memorize the image for a test that would take place when the task was completed. For
202 the rating task, participants were asked to think about how they would rate the image on a
203 scale from 1 (very unpleasant) to 7 (very pleasant). The participants were prompted for their
204 rating immediately after viewing the image. The same materials were used in both
205 experiments with a minor variation in the procedures. In the Confirmatory experiment,
206 participants were directed as to where search targets might appear in the image (e.g., on flat
207 surfaces). No such instructions were provided in the Exploratory experiment. In actuality,
208 none of the images in either experiment actually contained any search targets.

209 In both experiments, trials were presented in one mixed block, and three separate task
210 blocks. For the mixed block, the trial types were randomly intermixed within the block. For
211 the three separate task blocks, each block was 35 trials consisting entirely of one of the three
212 conditions (Search, Memorize, Rate). Each stimulus image was presented for 8 seconds. The
213 pictures were presented in color, with a size of 1024 x 768 pixels, subtending a visual angle of
214 23.84° x 17.99°.

215 Eye movements were recorded using an SR Research EyeLink 1000 eye tracker with a
216 sampling rate of 1000Hz. Only the right eye was recorded. The system was calibrated using
217 a nine-point accuracy adn validity test. Errors greater than 1° or averaging greater than 0.5°
218 in total were re-calibrated. When eye movement velocities remained below 30°/s for 10
219 consecutive samples, movement offset was detected.

220 Datasets

221 On some of the search trials, a probe was presented on the screen six seconds following
222 the onset of the trial. To avoid confounds resulting from the probe, only the first six seconds
223 of the data in all three conditions were analyzed. Trials that contained fewer than 6000

²²⁴ samples within the first six seconds of the trial were excluded before analysis. For both
²²⁵ datasets, the trials were pooled across participants. After excluding trials, the Exploratory
²²⁶ dataset consisted of 12,177 of the 16,740 total trials, and the Confirmatory dataset consisted
²²⁷ of 9,301 of the 10,395 total trials.

²²⁸ The raw x-coordinate, y-coordinate, and pupil size data collected at every sampling
²²⁹ time point in the trial were used as inputs to the deep learning classifier. These data were
²³⁰ also used to develop plot image datasets that were classified separately from the raw timeline
²³¹ datasets. For the plot image datasets, the timeline data for each trial were converted into
²³² scatterplot diagrams. The x- and y- coordinates and pupil size were used to plot each data
²³³ point onto a scatterplot (e.g., see Figure 1). The coordinates were used to plot the location
²³⁴ of the dot, pupil size was used to determine the relative size of the dot, and shading of the
²³⁵ dot was used to indicate the time-course of the eye movements throughout the trial. The
²³⁶ background of the plot images and first data point were white. Each subsequent data point
²³⁷ was one shade darker than the previous data point until the final data point was reached.
²³⁸ The final data point was black. For standardization, pupil size was divided by 10, and one
²³⁹ unit was added. The plots were sized to match the dimensions of the data collection monitor
²⁴⁰ (1024 x 768 pixels) and then shrunk to (240 x 180 pixels) in an effort to reduce the
²⁴¹ dimensionality of the data.

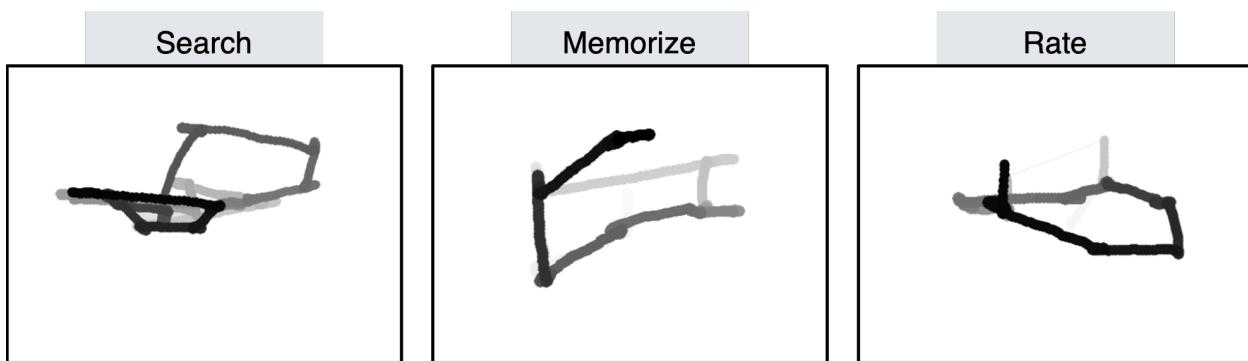


Figure 1. Each trial was represented as an image. Each sample collected within the trial was plotted as a dot in the image. Pupil size was represented by the size of the dot. The time course of the eye movements was represented by the gradual darkening of the dot over time.

242 **Data Subsets.** The full timeline dataset was structured into three columns

243 representing the x- and y- coordinates, and pupil size for each data point collected in the
244 first six seconds of each trial. To systematically assess the predictive value of each XYP (i.e.,
245 x-coordinates, y-coordinates, pupil size) component of the data, the timeline and image
246 datasets were batched into subsets that excluded one of the components (i.e., XYØ, XØP,
247 ØYP), or contained only one of the components (i.e., XØØ, ØYØ, ØØP). For the timeline
248 datasets, this means that the columns to be excluded in each data subset were replaced with
249 zeros. The data were replaced with zeros because removing the columns would change the
250 structure of the data. The same systematic batching process was carried out for the image
251 dataset. See Figure 2 for an example of each of these image data subsets.

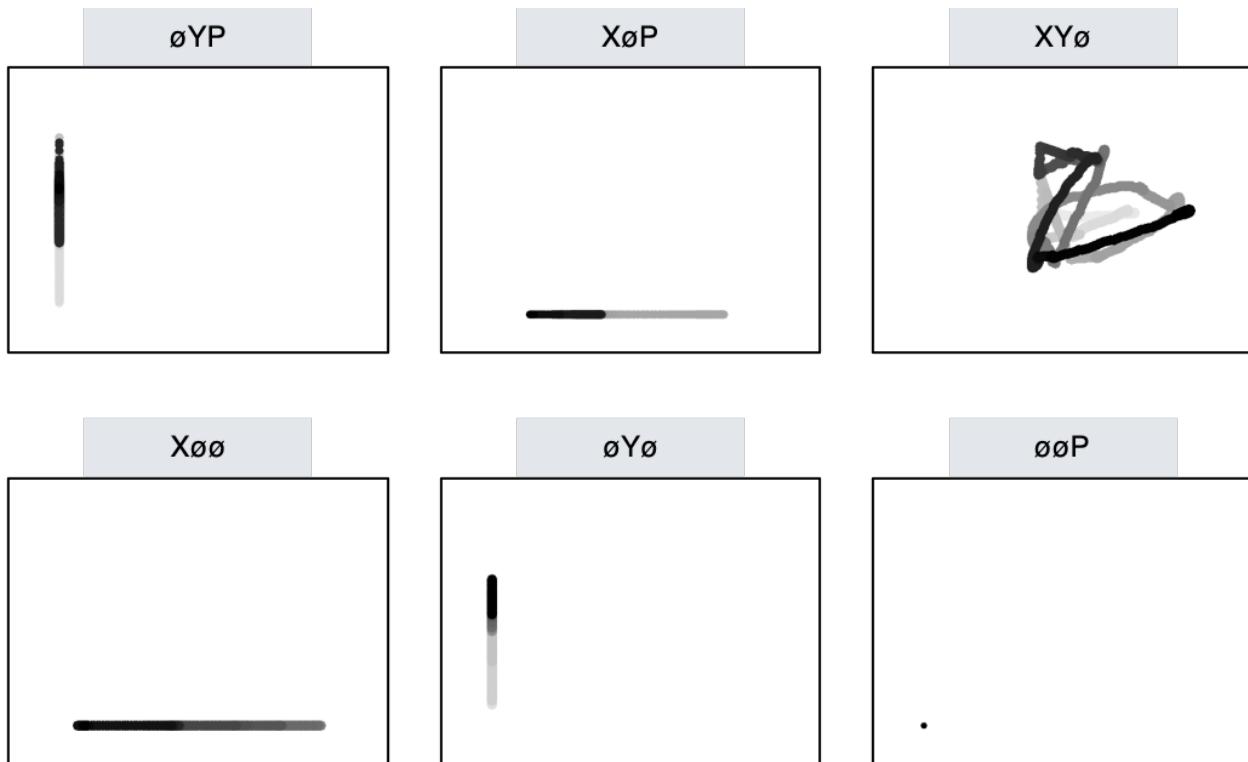


Figure 2. Plot images were used to represent each type of data subset. As with the trials in the full XYP dataset, the time course of the eye movements was represented by the shading of the dot. The first sample of each trial was white, and the last sample was black.

252 **Classification**

253 Deep CNN model architectures were implemented to classify the trials into Search,
254 Memorize, or Rate categories. Because CNNs act as a digital filter sensitive to the number of
255 features in the data, the differences in the structure of the timeline and image data formats
256 necessitated separate CNN model architectures. The model architectures were developed
257 with the intent of establishing a generalizable approach to classifying cognitive processes
258 from eye movement data.

259 The development of these models was not guided by any formal theoretical
260 assumptions regarding the patterns or features likely to be extracted by the classifier. Like
261 many HCI models, the development of these models followed general intuitions concerned
262 with building a model architecture capable of transforming the data inputs into an
263 interpretable feature set that would not overfit the dataset. The models were developed
264 using version 0.3b of the DeLINEATE toolbox, which operates over a Keras backend
265 (<http://delineate.it>). Each training/test iteration randomly split the data so that 70% of the
266 trial data were allocated to training, 15% of the trial data were allocated to validation, and
267 15% of the trial data were allocated to testing. Training of the model was stopped when
268 validation accuracy did not improve over the span of 100 epochs. Once the early stopping
269 threshold was reached, the resulting model was tested on the held-out test data. This
270 process was repeated 10 times for each model, resulting in 10 classification accuracy scores
271 for each model. The average of the resulting accuracy scores were the subject of comparisons
272 against chance and other datasets or data subsets.

273 The models were developed and tested on the Exploratory dataset. Model
274 hyperparameters were adjusted until the classification accuracies appeared to peak. The
275 model architecture with the highest classification accuracy on the Exploratory dataset was
276 trained, validated, and tested independently on the Confirmatory dataset. This means that
277 the model that was used to analyze the Confirmatory dataset was not trained on the

278 Exploratory dataset. The model architectures used for the timeline and plot image datasets
279 are shown in Figure 3.

280 **Analysis**

281 Results for the CNN architecture that resulted in the highest accuracy on the
282 Exploratory dataset are reported below. For every dataset tested, a one-sample two-tailed
283 *t*-test was used to compare the CNN accuracies against chance (33%). The Shapiro-Wilk test
284 was used to assess the normality for each dataset. When normality was assumed, the mean
285 accuracy for that dataset was compared against chance using Student's one-sample
286 two-tailed *t*-test. When normality could not be assumed, the median accuracy for that
287 dataset was compared against chance using Wilcoxon's Signed Rank test.

288 To determine the relative value of the three components of the eye movement data, the
289 data subsets were compared within the timeline and plot image data types. If classification
290 accuracies were lower when the data were batched into subsets, the component that was
291 removed was assumed to have some unique contribution that the model was using to inform
292 classification decisions. To determine the relative value of the contribution from each
293 component, the accuracies from each subset with one component of the data removed were
294 compared to the accuracies for the full dataset (XYP) using a one-way between-subjects
295 Analysis of Variance (ANOVA). To further evaluate the decodability of each component
296 independently, the accuracies from each subset containing only one component of the eye
297 movement data were compared within a separate one-way between-subjects ANOVA. All
298 post-hoc comparisons were corrected using Tukey's HSD.

299 **Results**

300 **Timeline Data Classification**

301 **Exploratory.** Classification accuracies for the XYP timeline dataset were well above
302 chance (chance = 33%; $M = .526$, $SD = .018$; $t_{(9)} = 34.565$, $p < .001$). Accuracies for

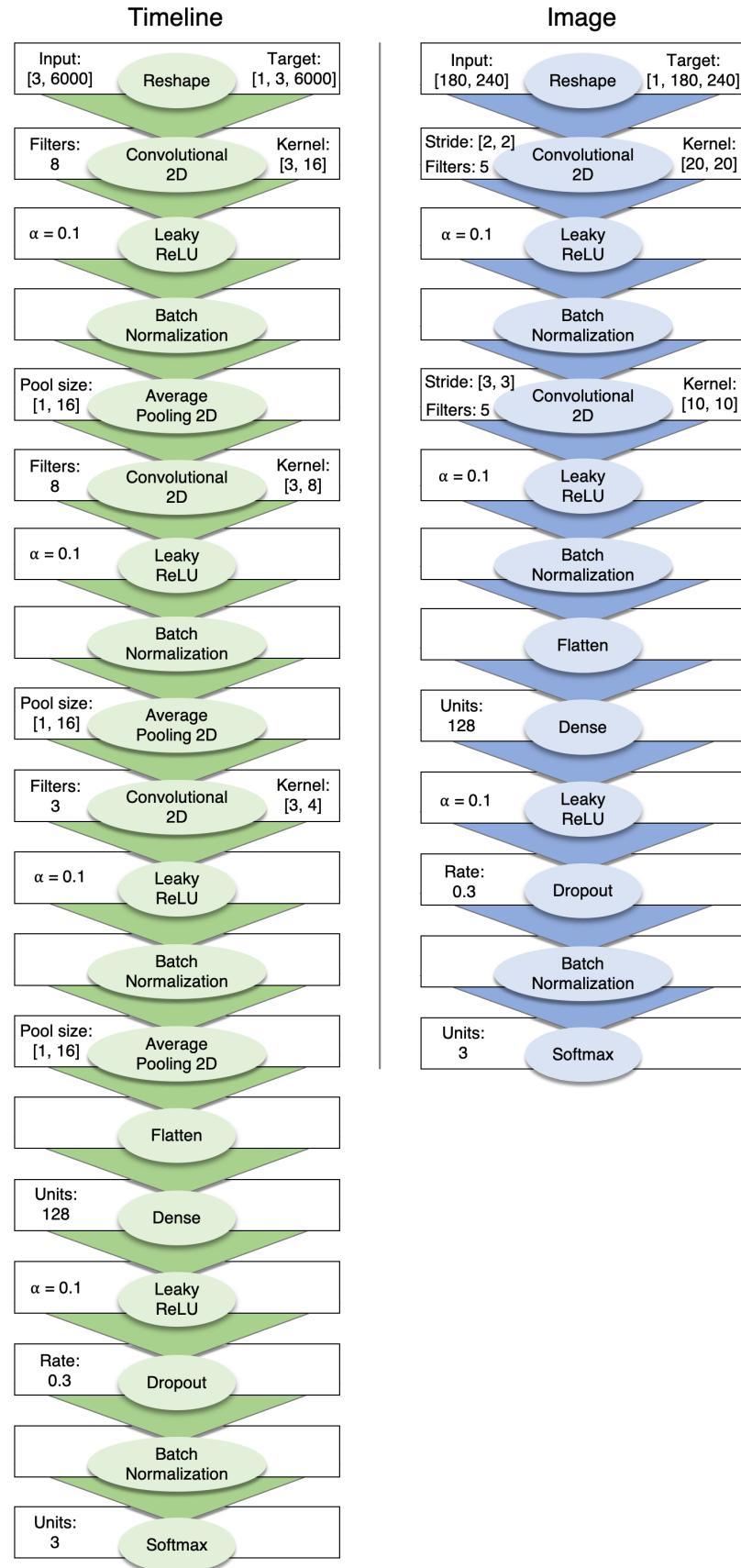


Figure 3. Two different model architectures were used to classify the timeline and image data. Both models were compiled using a categorical crossentropy loss function, and optimized with the Adam algorithm.

303 classifications of the batched data subsets were all better than chance (see Figure 4). As
 304 shown in the confusion matrices displayed in Figure 5, the data subsets with lower overall
 305 classification accuracies almost always classified the Memorize condition at or below chance
 306 levels of accuracy. Misclassifications of the Memorize condition were split relatively evenly
 307 between the Search and Rate conditions.

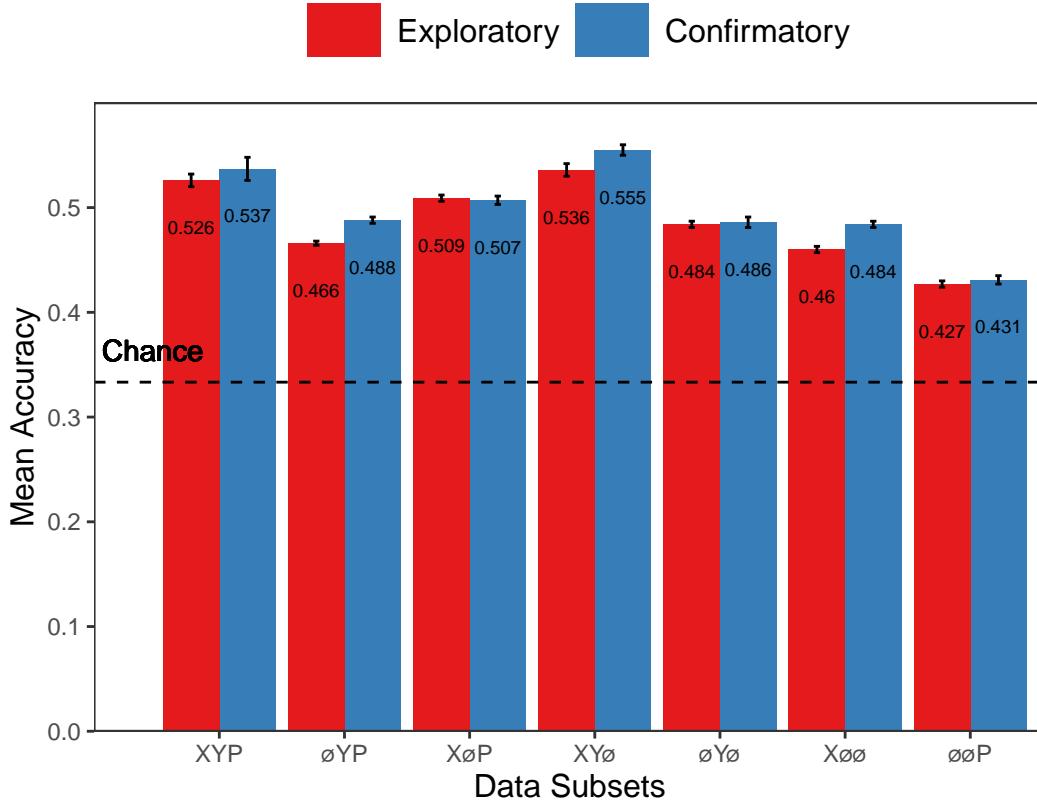


Figure 4. The graph represents the average accuracy reported for each subset of the timeline data. All of the data subsets were decoded at levels better than chance (33%). Each subset is labeled with the mean accuracy. The error bars represent standard errors.

308 There was a difference in classification accuracy for the XYP dataset and the subsets
 309 that had the pupil size, x-coordinate, and y-coordinate data systematically removed ($F_{(3,36)}$
 310 $= 47.471, p < .001, \eta^2 = 0.798$). Post-hoc comparisons against the XYP dataset showed that
 311 classification accuracies were not affected by the removal of pupil size or y-coordinate data
 312 (see Table 2). The null effect present when pupil size was removed suggests that the pupil
 313 size data were not contributing unique information that was not otherwise provided by the x-

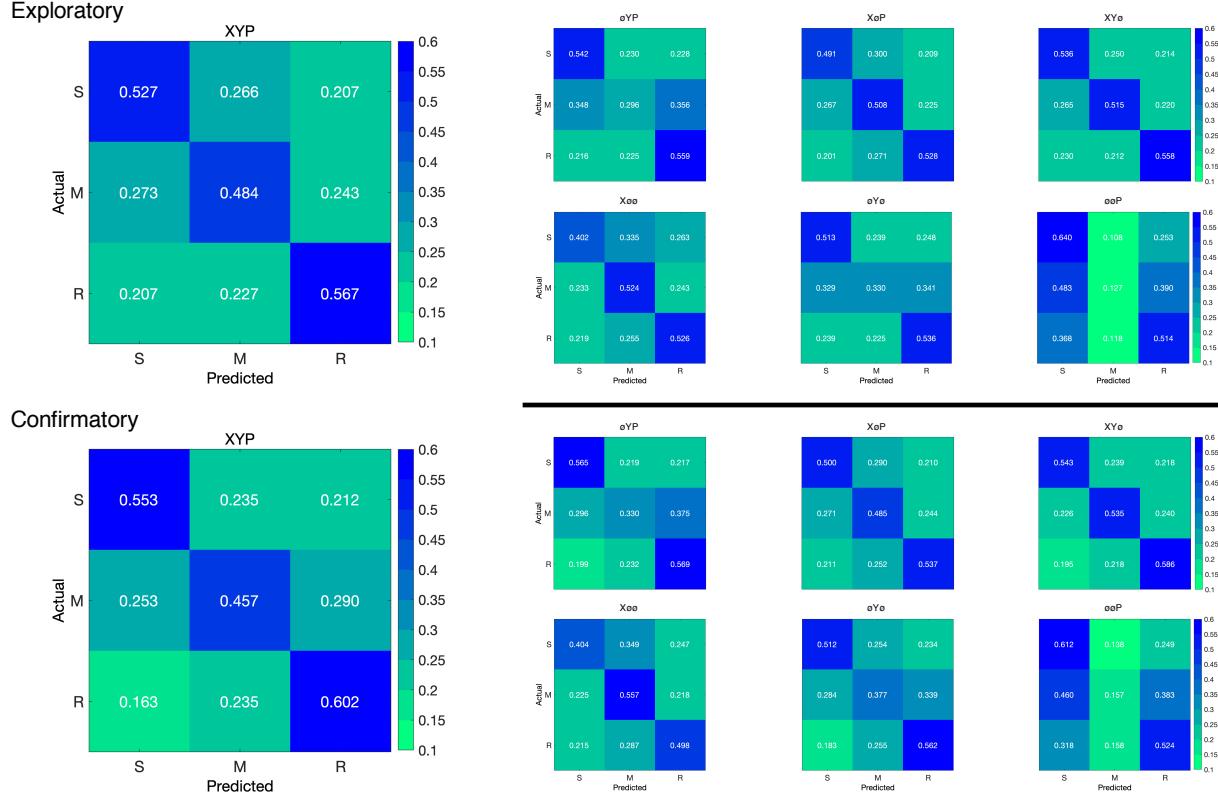


Figure 5. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for trial. The horizontal axis of the confusion matrices represents the condition that was predicted by the model.

and y-coordinates. A strict significance threshold of $\alpha = .05$ implies the same conclusion for the y-coordinate data, but the relatively low degrees of freedom ($df = 18$) and the borderline observed p -value ($p = .056$) afford the possibility that there exists a small effect. However, classification for the \emptyset Y \emptyset subset was significantly lower than the XYP dataset, showing that the x-coordinate data were uniquely informative to the classification.

Table 2
Timeline Subset Comparisons

Comparison	Exploratory		Confirmatory	
	<i>t</i>	<i>p</i>	<i>t</i>	<i>p</i>
XYP vs. \emptyset Y \emptyset	9.420	< .001	5.210	< .001
XYP vs. X \emptyset P	2.645	.056	3.165	.016
XYP vs. XYø	1.635	.372	1.805	.288
X \emptyset ø vs. \emptyset Yø	5.187	< .001	0.495	.874
X \emptyset ø vs. \emptyset øP	12.213	< .001	10.178	< .001
\emptyset Yø vs. \emptyset øP	7.026	< .001	9.683	< .001

319 There was also a difference in classification accuracies for the XØØ, ØYØ, and ØØP

320 subsets ($F_{(2,27)} = 75.145, p < .001, \eta^2 = 0.848$). Post-hoc comparisons showed that

321 classification accuracy for the ØØP subset was lower than the XØØ and ØYØ subsets.

322 Classification accuracy for the XØØ subset was higher than the ØYØ subset. Altogether,

323 these findings suggest that pupil size data was the least uniquely informative to classification

324 decisions, while the x-coordinate data was the most uniquely informative.

325 **Confirmatory.** Classification accuracies for the Confirmatory XYP timeline dataset

326 were well above chance ($M = .537, SD = 0.036, t_{(9)} = 17.849, p < .001$). Classification

327 accuracies for the data subsets were also better than chance (see Figure 4). Overall, there

328 was high similarity in the pattern of results for the Exploratory and Confirmatory datasets

329 (see Figure 4). Furthermore, the general trend showing that pupil size was the least

330 informative eye tracking data component was replicated in the Confirmatory dataset (see

331 Table 2). Also in concordance with the Exploratory timeline dataset, the confusion matrices

332 for these data revealed that the Memorize task was mis-classified more often than the Search

333 and Rate tasks (see Figure 5).

334 To test the generalizability of the model to other eye tracking data, classification

335 accuracies for the XYP Exploratory and Confirmatory timeline datasets were compared. The

336 Shapiro-Wilk test for normality indicated that the Exploratory ($W = 0.937, p = .524$) and

337 Confirmatory ($W = 0.884, p = .145$) datasets were normally distributed, but Levene's test

338 indicated that the variances were not equal, $F_{(1,18)} = 8.783, p = .008$. Welch's unequal

339 variances *t*-test did not show a difference between the two datasets, $t_{(13.045)} = 0.907, p =$

340 .381, Cohen's *d* = 0.406. These findings indicate that the deep learning model decoded the

341 Exploratory and Confirmatory timeline datasets equally well, but the Confirmatory dataset

342 classifications were less consistent across training/test iterations (as indicated by the increase

343 in standard deviation).

³⁴⁴ **Plot Image Classification**

³⁴⁵ **Exploratory.** Classification accuracies for the XYP plot image data were better
³⁴⁶ than chance ($M = .436$, $SD = .020$, $p < .001$), but were less accurate than the classifications
³⁴⁷ for the XYP Exploratory timeline data ($t_{(18)} = 10.813$, $p < .001$). Accuracies for the
³⁴⁸ classifications for all subsets of the plot image data except the $\emptyset\emptyset P$ subset were better than
³⁴⁹ chance (see Figure 6). Following the pattern expressed by the timeline dataset, the confusion
³⁵⁰ matrices showed that the Memorize condition was misclassified more often than the other
³⁵¹ conditions, and appeared to be evenly mis-identified as a Search or Rate condition (see
³⁵² Figure 7).

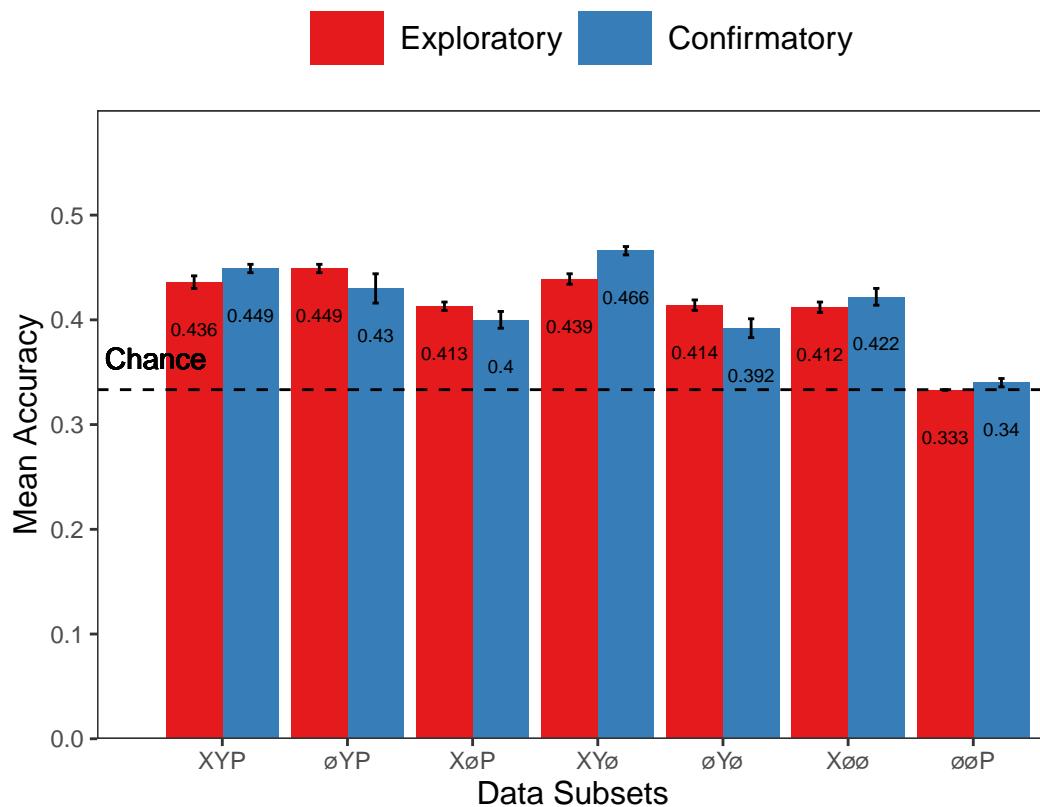


Figure 6. The graph represents the average accuracy reported for each subset of the image data. All of the data subsets except for the Exploratory $\emptyset\emptyset P$ dataset were decoded at levels better than chance (33%). Each subset is labeled with the mean accuracy. The error bars represent standard errors.

³⁵³ There was a difference in classification accuracy between the XYP dataset and the data
³⁵⁴ subsets ($F_{(4,45)} = 7.093$, $p < .001$, $\eta^2 = .387$). Post-hoc comparisons showed that compared

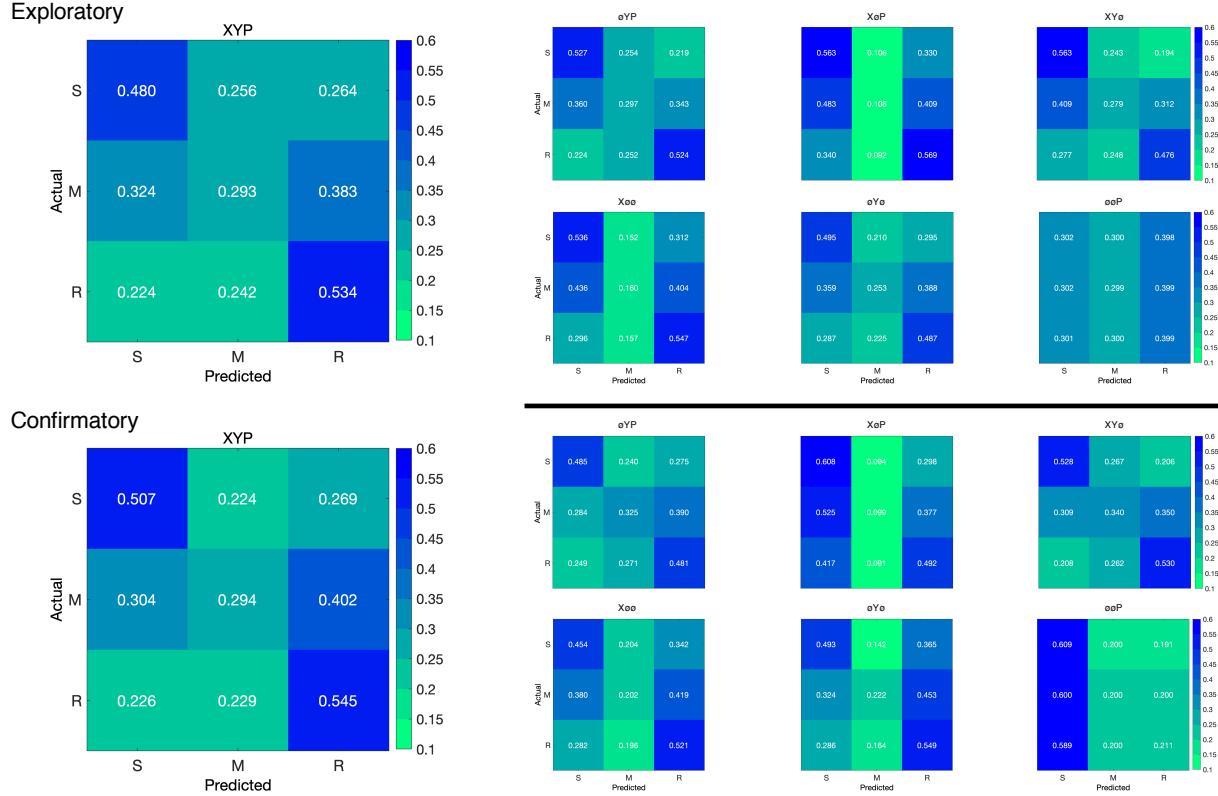


Figure 7. The confusion matrices represent the average classification accuracies for each condition of the image data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents the condition that was predicted by the model.

to the XYP dataset, there was no effect of removing pupil size or the x-coordinates, but classification accuracy was worse when the y-coordinates were removed (see Table 3).

Table 3
Image Subset Comparisons

Comparison	Exploratory		Confirmatory	
	<i>t</i>	<i>p</i>	<i>t</i>	<i>p</i>
XYP vs. \emptyset YP	1.792	.391	1.623	.491
XYP vs. XoP	2.939	.039	4.375	< .001
XYP vs. XYo	0.474	.989	1.557	.532
XoO vs. \emptyset Yo	0.423	.906	2.807	.204
XoO vs. $\emptyset\emptyset$ P	13.569	< .001	5.070	< .001
\emptyset Yo vs. $\emptyset\emptyset$ P	13.235	< .001	7.877	< .001

There was also a difference in classification accuracies between the XoO, \emptyset Yo, and $\emptyset\emptyset$ P subsets (Levene's test: $F_{(2,27)} = 3.815$, $p = .035$; Welch correction for lack of homogeneity of variances: $F_{(2,17.993)} = 228.137$, $p < .001$, $\eta^2 = .899$). Post-hoc comparisons

360 showed that there was no difference in classification accuracies for the XØØ and ØYØ
361 subsets, but classification for the ØØP subset were less accurate than the XØØ and ØYØ
362 subsets.

363 **Confirmatory.** Classification accuracies for the XYP confirmatory image dataset
364 were well above chance ($M = .449$, $SD = 0.012$, $t_{(9)} = 31.061$, $p < .001$), but were less
365 accurate than the classifications of the confirmatory timeline dataset ($t_{(18)} = 11.167$ $p <$
366 $.001$). Accuracies for classifications of the data subsets were also all better than chance (see
367 Figure 6). The confusion matrices followed the pattern showing that the Memorize condition
368 was confused most often, and was relatively evenly mis-identified as a Search or Rate trial
369 (see Figure 7). As with the timeline data, the general trend showing that pupil size data was
370 the least informative to the model was replicated in the Confirmatory dataset (see Table 3).

371 To test the generalizability of the model, the classification accuracies for the XYP
372 Exploratory and Confirmatory plot image datasets were compared. The independent samples
373 t -test showed that the deep learning model did equally well at classifying the Exploratory
374 and Confirmatory plot image datasets, $t_{(18)} = 1.777$, $p = .092$, Cohen's $d = 0.795$.

375 Discussion

376 The present study aimed to produce a practical and reliable example of a black box
377 solution to the inverse Yarbus problem. To implement this solution, we classified raw
378 timeline and minimally processed plot image data using a CNN model architecture. To our
379 knowledge, this study was the first to provide a solution to determining mental state from
380 eye movement data using each of the following: (1) Non-aggregated eye tracking data (i.e.,
381 raw x-coordinates, y-coordinates, pupil size), (2) timeline and image data formats (see
382 Figure 2), and (3) a black box CNN architecture. This study probed the relative predictive
383 value of the x-coordinate, y-coordinate, and pupil size components of the eye movement data
384 using a CNN. The CNN was able to decode the timeline and plot image data better than
385 chance, although only the timeline datasets were decoded with accuracies comparable to

386 other state-of-the-art approaches. Datasets with lower classification accuracies were not able
387 to differentiate the cognitive processes underlying the Memorize task from the cognitive
388 processes underlying the Search and Rate tasks. Decoding subsets of the data revealed that
389 pupil size was the least uniquely informative component of the eye movement data. This
390 pattern of findings was consistent between the Exploratory and Confirmatory datasets.

391 Although several aggregate eye movement features have been tested as task predictors,
392 to our knowledge, no other study has assessed the predictive value of the data format (viz.,
393 data in the format of a plot image). Our results suggest that although CNNs are robust
394 image classifiers, eye movement data is decoded in the standard timeline format more
395 effectively than in image format. This may be because the image data format contains less
396 decodable information than the timeline format. Over the span of the trial (six seconds), the
397 eye movements occasionally overlapped. When there was an overlap in the image data
398 format, the more recent data points overwrote the older data points. This resulted in some
399 information loss that did not occur when the data were represented in the raw timeline
400 format. Despite this loss of information, the plot image format was still decoded with better
401 than chance accuracy. To further examine the viability of classifying task from eye
402 movement image datasets, future research might consider representing the data in different
403 forms such as 3-dimensional data formats, or more complex color combinations capable of
404 representing overlapping data points.

405 When considering the superior performance of the timeline data (vs., plot image data),
406 we must also consider the differences in the model architectures. Because the structures of
407 the timeline and plot image data formats were different, the models decoding those data
408 structures also needed to be different. Both models were optimized individually on the
409 Exploratory dataset before being tested on the Confirmatory dataset. For both timeline and
410 plot image formats, there was good replicability between the Exploratory and Confirmatory
411 datasets, demonstrating that these architectures performed similarly from experiment to

412 experiment. An appropriately tuned CNN should be capable of learning any arbitrary
413 function, but given that the upper bound for decodability of these datasets is unknown,
414 there is the possibility that a model architecture exists that is capable of classifying the plot
415 image data format more accurately than the model used to classify the timeline data.
416 Despite this possibility, the convergence of these findings with other studies (see Table 1)
417 suggests that the results of this study are approaching a ceiling for the potential to solve the
418 inverse Yarbus problem with eye movement data. Although the true capacity to predict
419 mental state from eye movement data is unknown, standardizing datasets in the future could
420 provide a point for comparison that can more effectively indicate which methods are most
421 effective at solving the inverse Yarbus problem.

422 In the current study, the Memorize condition was most classified less accurately than
423 the Search and Rate conditions, especially for the datasets with lower overall accuracy. This
424 suggests that the eye movements associated with the Memorize task were potentially lacking
425 unique or informative features to decode. This means that eye movements associated with
426 the Memorize condition were interpreted as noise, or were sharing features of underlying
427 cognitive processes that were represented in the eye movements associated with the Search
428 and Rate tasks. Previous research (e.g., Król & Król, 2018) has attributed the inability to
429 differentiate one condition from the others to the overlapping of sub-features in the eye
430 movements between two tasks that are too subtle to be represented in the eye movement
431 data.

432 To more clearly understand how the different tasks influenced the decodability of the
433 eye movement data, additional analyses were conducted on the Exploratory and
434 Confirmatory timeline datasets (see Appendix). These analyses showed that classification
435 accuracy improved when the Memorize condition was removed. A closer look at these results
436 showed that when the Memorize condition was included in the subset, classification
437 accuracies of the Search and Rate conditions was lower. The analyses were complemented

438 with a re-calculation of the accuracies from the primary analysis to account for a 50%
439 threshold of chance performance. Altogether, these results could indicate that the eye
440 movement features underlying the Memorize condition are shared with the Search and Rate
441 conditions, or that the Memorize condition is contributing a substantial amount of noise.
442 Given the findings of these supplementary analyses, we believe that the Memorize trials were
443 more often mis-classified than the other trials due to an increased prevalence of noise in the
444 eye movement data for the Memorize trials.

445 When determining the relative contributions of the eye movement features used in
446 this study (x-coordinates, y-coordinates, pupil size), the pupil size data was consistently the
447 least uniquely informative. When pupil size was removed from the Exploratory and
448 Confirmatory timeline and plot image datasets, classification accuracy remained stable (vs.,
449 XYP dataset). Furthermore, classification of the $\emptyset\emptyset P$ subset was the lowest of all of the data
450 subsets, and in one instance, was no better than chance. Although these findings indicate
451 that, in this case, pupil size was a relatively uninformative component of the eye movement
452 data, previous research has associated changes in pupil size as indicators of working memory
453 load (Kahneman & Beatty, 1966; Karatekin, Couperus, & Marcus, 2004), arousal (Wang et
454 al., 2018), and cognitive effort (Porter, Troscianko, & Gilchrist, 2007). The results of the
455 current study indicate that the changes in pupil size associated with these underlying
456 processes are not useful in delineating the tasks being classified (i.e., Search, Memorize,
457 Rate), potentially because these tasks do not evoke a reliable pattern of changes in pupil size.

458 The findings from the current study support the notion that black box CNNs are a
459 viable approach to determining task from eye movement data. In a recent review, Lukander
460 et al. (2017) expressed concern regarding the lack of generalizability of black box approaches
461 when decoding eye movement data. Overall, the current study showed a consistent pattern
462 of results for the XYP timeline and image datasets, but some minor inconsistencies in the
463 pattern of results for the x- and y- coordinate subset comparisons. These inconsistencies may

464 be a product of overlap in the cognitive processes underlying the three tasks. When the data
465 are batched into subsets, at least one dimension (i.e., x-coordinates, y-coordinates, or pupil
466 size) is removed, leading to a potential loss of information. When the data provide fewer
467 meaningful distinctions, finer-grained inferences are necessary for the tasks to be
468 distinguishable. As shown by Coco and Keller (2014), eye movement data can be more
469 effectively decoded when the cognitive processes underlying the tasks are explicitly
470 differentiable. While the cognitive processes distinguishing memorizing, searching, or rating
471 an image are intuitively different, the eye movements elicited from these cognitive processes
472 are not easily differentiated. To correct for potential mismatches between the distinctive
473 task-diagnostic features in the data and the level of distinctiveness required to classify the
474 tasks, future research could more definitively conceptualize the cognitive processes
475 underlying the task-at-hand.

476 Classifying mental state from eye movement data is often carried out in an effort to
477 advance technology to improve educational outcomes, strengthen the independence of
478 physically and mentally handicapped individuals, or improve HCI's (Koochaki &
479 Najafizadeh, 2018). Given the previous questions raised regarding the reliability and
480 generalizability of black-box CNN classification, the current study first tested models on an
481 exploratory dataset, then confirmed the outcome using a second independent dataset.
482 Overall, the findings of this study indicate that this black-box approach is capable of
483 producing a stable and generalizable outcome. Future studies that incorporate stimulus
484 features might have the potential to surpass current state-of-the-art classification. According
485 to Bulling, Weichel, and Gellersen (2013), incorporating stimulus feature information into
486 the dataset may provide improve accuracy relative to decoding gaze location data and pupil
487 size. Alternatively, Borji and Itti (2014) suggested that accounting for salient features in the
488 the stimulus might leave little to no room for theoretically defined classifiers to consider
489 mental state. Future research should examine the potential for the inclusion of stimulus
490 feature information in addition to the eye movement data to boost black-box CNN

491 classification accuracy of image data beyond that of timeline data.

References

- 492
- 493 Boisvert, J. F. G., & Bruce, N. D. B. (2016). Predicting task from eye movements: On the
494 importance of spatial distribution, dynamics, and image features. *Neurocomputing*,
495 207, 653–668. <https://doi.org/10.1016/j.neucom.2016.05.047>
- 496 Borji, A., & Itti, L. (2014). Defending Yarbus: Eye movements reveal observers' task.
497 *Journal of Vision*, 14(3), 29–29. <https://doi.org/10.1167/14.3.29>
- 498 Bulling, A., Weichel, C., & Gellersen, H. (2013). EyeContext: Recognition of high-level
499 contextual cues from human visual behaviour. In *Proceedings of the SIGCHI
500 Conference on Human Factors in Computing Systems - CHI '13* (p. 305). Paris,
501 France: ACM Press. <https://doi.org/10.1145/2470654.2470697>
- 502 Castelhano, M. S., Mack, M. L., & Henderson, J. M. (2009). Viewing task influences eye
503 movement control during active scene perception. *Journal of Vision*, 9(3), 6–6.
504 <https://doi.org/10.1167/9.3.6>
- 505 Coco, M. I., & Keller, F. (2014). Classification of visual and linguistic tasks using
506 eye-movement features. *Journal of Vision*, 14(3), 11–11.
507 <https://doi.org/10.1167/14.3.11>
- 508 DeAngelus, M., & Pelz, J. B. (2009). Top-down control of eye movements: Yarbus revisited.
509 *Visual Cognition*, 17(6-7), 790–811. <https://doi.org/10.1080/13506280902793843>
- 510 Greene, M. R., Liu, T., & Wolfe, J. M. (2012). Reconsidering Yarbus: A failure to predict
511 observers' task from eye movement patterns. *Vision Res*, 62, 1–8.
512 <https://doi.org/10.1016/j.visres.2012.03.019>
- 513 Haji-Abolhassani, A., & Clark, J. J. (2014). An inverse Yarbus process: Predicting observers'
514 task from eye movement patterns. *Vision Research*, 103, 127–142.

515 <https://doi.org/10.1016/j.visres.2014.08.014>

516 Henderson, J. M., Shinkareva, S. V., Wang, J., Luke, S. G., & Olejarczyk, J. (2013).

517 Predicting Cognitive State from Eye Movements. *PLoS ONE*, 8(5), e64937.

518 <https://doi.org/10.1371/journal.pone.0064937>

519 Kahneman, D., & Beatty, J. (1966). Pupil Diameter and Load on Memory. *Science*,

520 154(3756), 1583–1585. Retrieved from <https://www.jstor.org/stable/1720478>

521 Kanan, C., Ray, N. A., Bseiso, D. N. F., Hsiao, J. H., & Cottrell, G. W. (2014). Predicting

522 an observer's task using multi-fixation pattern analysis. In *Proceedings of the*

523 *Symposium on Eye Tracking Research and Applications - ETRA '14* (pp. 287–290).

524 Safety Harbor, Florida: ACM Press. <https://doi.org/10.1145/2578153.2578208>

525 Karatekin, C., Couperus, J. W., & Marcus, D. J. (2004). Attention allocation in the

526 dual-task paradigm as measured through behavioral and psychophysiological

527 responses. *Psychophysiology*, 41(2), 175–185.

528 <https://doi.org/10.1111/j.1469-8986.2004.00147.x>

529 Koochaki, F., & Najafizadeh, L. (2018). Predicting Intention Through Eye Gaze Patterns.

530 In *2018 IEEE Biomedical Circuits and Systems Conference (BioCAS)* (pp. 1–4).

531 <https://doi.org/10.1109/BIOCAS.2018.8584665>

532 Król, M. E., & Król, M. (2018). The right look for the job: Decoding cognitive processes

533 involved in the task from spatial eye-movement patterns. *Psychological Research*.

534 <https://doi.org/10.1007/s00426-018-0996-5>

535 Lukander, K., Toivanen, M., & Puolamäki, K. (2017). Inferring Intent and Action from Gaze

536 in Naturalistic Behavior: A Review. *International Journal of Mobile Human*

537 *Computer Interaction*, 9(4), 41–57. <https://doi.org/10.4018/IJMHCI.2017100104>

- 538 MacInnes, W., Joseph, Hunt, A. R., Clarke, A. D. F., & Dodd, M. D. (2018). A Generative
539 Model of Cognitive State from Task and Eye Movements. *Cognitive Computation*,
540 10(5), 703–717. <https://doi.org/10.1007/s12559-018-9558-9>
- 541 Mills, M., Hollingworth, A., Van der Stigchel, S., Hoffman, L., & Dodd, M. D. (2011).
542 Examining the influence of task set on eye movements and fixations. *Journal of*
543 *Vision*, 11(8), 17–17. <https://doi.org/10.1167/11.8.17>
- 544 Porter, G., Troscianko, T., & Gilchrist, I. D. (2007). Effort during visual search and
545 counting: Insights from pupillometry. *Quarterly Journal of Experimental Psychology*
546 (2006), 60(2), 211–229. <https://doi.org/10.1080/17470210600673818>
- 547 Seeliger, K., Fritzsche, M., Güçlü, U., Schoenmakers, S., Schoffelen, J.-M., Bosch, S. E., &
548 van Gerven, M. A. J. (2018). Convolutional neural network-based encoding and
549 decoding of visual object recognition in space and time. *NeuroImage*, 180, 253–266.
550 <https://doi.org/10.1016/j.neuroimage.2017.07.018>
- 551 Tatler, B. W., Wade, N. J., Kwan, H., Findlay, J. M., & Velichkovsky, B. M. (2010). Yarbus,
552 Eye Movements, and Vision. *I-Perception*, 1(1), 7–27. <https://doi.org/10.1068/i0382>
- 553 Wang, C.-A., Baird, T., Huang, J., Coutinho, J. D., Brien, D. C., & Munoz, D. P. (2018).
554 Arousal Effects on Pupil Size, Heart Rate, and Skin Conductance in an Emotional
555 Face Task. *Frontiers in Neurology*, 9. <https://doi.org/10.3389/fneur.2018.01029>
- 556 Yarbus, A. (1967). Eye Movements and Vision. Retrieved January 24, 2019, from
557 [http://wexler.free.fr/library/files/yarbus%20\(1967\)%20eye%20movements%20and%20vision.pdf](http://wexler.free.fr/library/files/yarbus%20(1967)%20eye%20movements%20and%20vision.pdf)
- 559 Zhou, B., Bau, D., Oliva, A., & Torralba, A. (2019). Comparing the Interpretability of Deep
560 Networks via Network Dissection. In W. Samek, G. Montavon, A. Vedaldi, L. K.
561 Hansen, & K.-R. Müller (Eds.), *Explainable AI: Interpreting, Explaining and*

562 *Visualizing Deep Learning* (pp. 243–252). Cham: Springer International Publishing.

563 https://doi.org/10.1007/978-3-030-28954-6_12

Appendix

Supplementary Analysis

564 Additional analyses were conducted to clarify the effect of task on classification accuracy.
 565 These supplementary analyses were not seen as central to the current study, but could prove
 566 to be informative to researchers attempting to replicate or extend these findings in the
 567 future. The results from the primary analyses showed that classification accuracies were the
 568 lowest for the Memorize condition, but these findings did not indicate if the Memorize
 569 condition was adding noise to the data, or was providing redundant information to the
 570 model. To further understand why classification accuracy was lower for the Memorize
 571 condition than it was for the Search or Rate condition, the Exploratory and Confirmatory
 572 timeline datasets were systematically batched into subsets with the Search (S), Memorize
 573 (M), or Rate (R) condition removed (i.e., \emptyset MR, S \emptyset R, SM \emptyset).

574 All of the data subsets analyzed in this supplementary analysis were decoded with
 575 better than chance accuracy (see Figure A1). The same pattern of results was observed in
 576 both the Exploratory and Confirmatory datasets. When the Memorize condition was
 577 removed, classification accuracy improved (see Table A1). When the Rate condition was
 578 removed, classification was the worst. When the Memorize condition was included, the
 579 Memorize condition was more accurately predicted than the Search and Rate conditions (see
 580 Figure A2).

Table A1
Supplementary Subset Comparisons

Comparison	Exploratory		Confirmatory	
	<i>t</i>	<i>p</i>	<i>t</i>	<i>p</i>
\emptyset MR vs. S \emptyset R	3.248	.008	3.094	.012
\emptyset MR vs. SM \emptyset	2.875	.021	2.923	.018
S \emptyset R vs. SM \emptyset	6.123	< .001	6.017	< .001

581 Overall, the accuracies for all of the data subsets observed in the supplementary
 582 analysis were higher than the accuracies observed in the main analysis. Chance accuracy

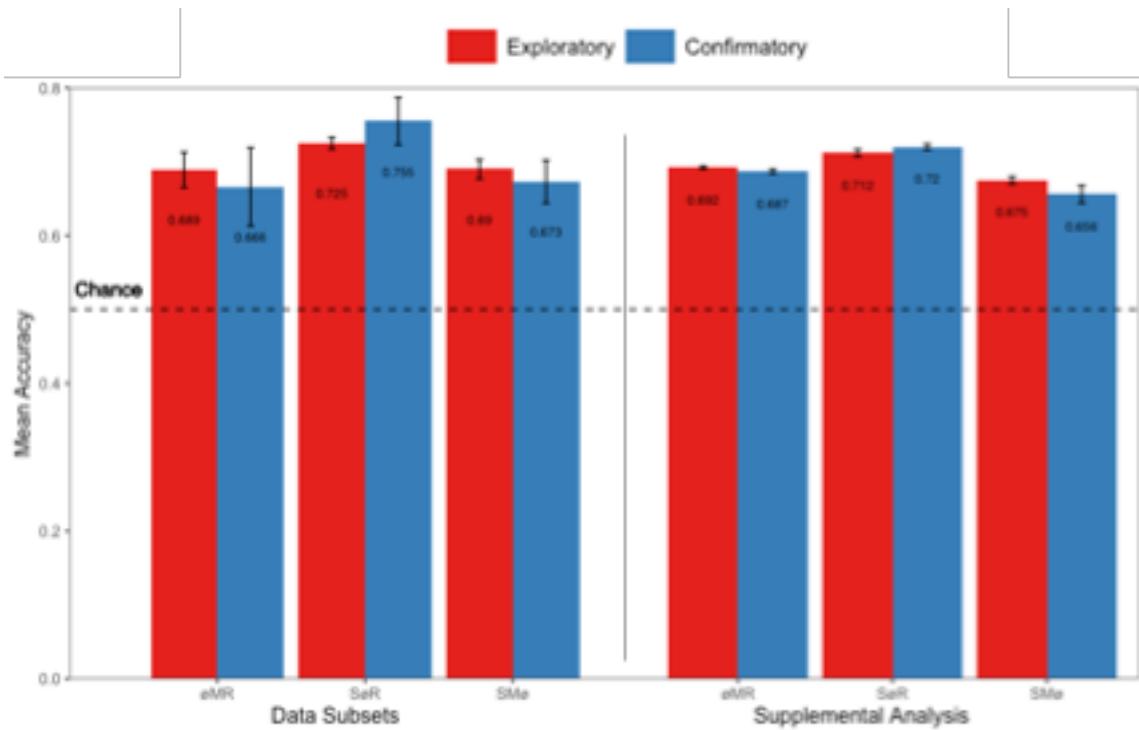


Figure A1. The graph represents the average accuracy reported for each subset of the Exploratory and Confirmatory timeline data for the supplementary analysis and the re-calculated accuracies from the primary analyses. All of the data subsets were decoded at levels better than chance (50%). The error bars represent standard errors.

583 levels for the primary analysis was 33%, but because one of the tasks was removed from each
 584 element observed in the supplementary analyses, chance accuracy for these analyses was 50%.
 585 Given the data analyzed for these supplementary purposes have different thresholds of
 586 chance performance, any conclusions drawn from a comparison between the primary and
 587 supplementary datasets could be misleading. For this reason, we revisited the results from
 588 the primary analyses and re-calculated the accuracies to be equivalent to a 50% chance
 589 threshold. The accuracies were recalculated using the following formula: $\text{Accuracy}_{\text{New}} =$
 590 $\text{Accuracy}_{\text{Original}} / (\text{Hits}_{\text{Original}} + \text{Misses}_{\text{Original}})$. For example, accuracy for the Search
 591 classification for SØR would be calculated with the following: $\text{Search}_{\text{Hits}} = \text{Hits}_{\text{Search}} /$
 592 $(\text{Hits}_{\text{Search}} + \text{Misses}_{\text{Search}})$, where $\text{Misses}_{\text{Search}}$ is the ratio of Search trials that were
 593 misclassified as Rate.

594 The re-calculated confusion matrices for each category are presented in Figure A3. The
 595 general pattern of findings presented in the re-calculated confusion matrices was the same as

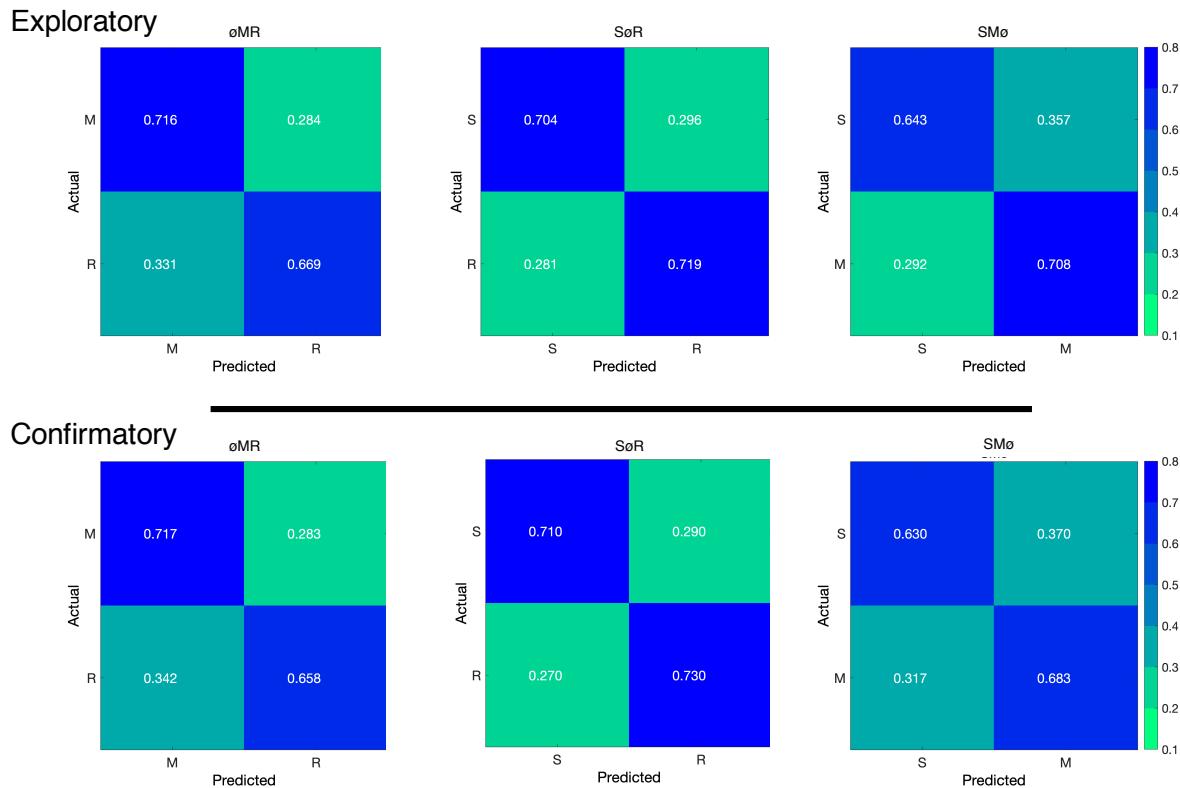


Figure A2. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for trial. The horizontal axis of the confusion matrices represents the condition that was predicted by the model.

596 the pattern that can be seen in the supplementary analysis. This is also true for the pattern
 597 of results for the overall accuracies (see Figure A1). In both analyses, the subsets that
 598 include the Memorize trials were classified with lower accuracy than the other data subsets.
 599 Furthermore, a closer look at the confusion matrices shows that the Memorize conditions
 600 were most often confused with the other two conditions. Because the Memorize condition
 601 was mis-classified to a similar extent when paired with Search or Rate conditions, and these
 602 mis-classifications affected the overall accuracies to a similar extent, we believe it is likely
 603 that there was an increased presence of noise in the eye movement data for the Memorize
 604 trials.

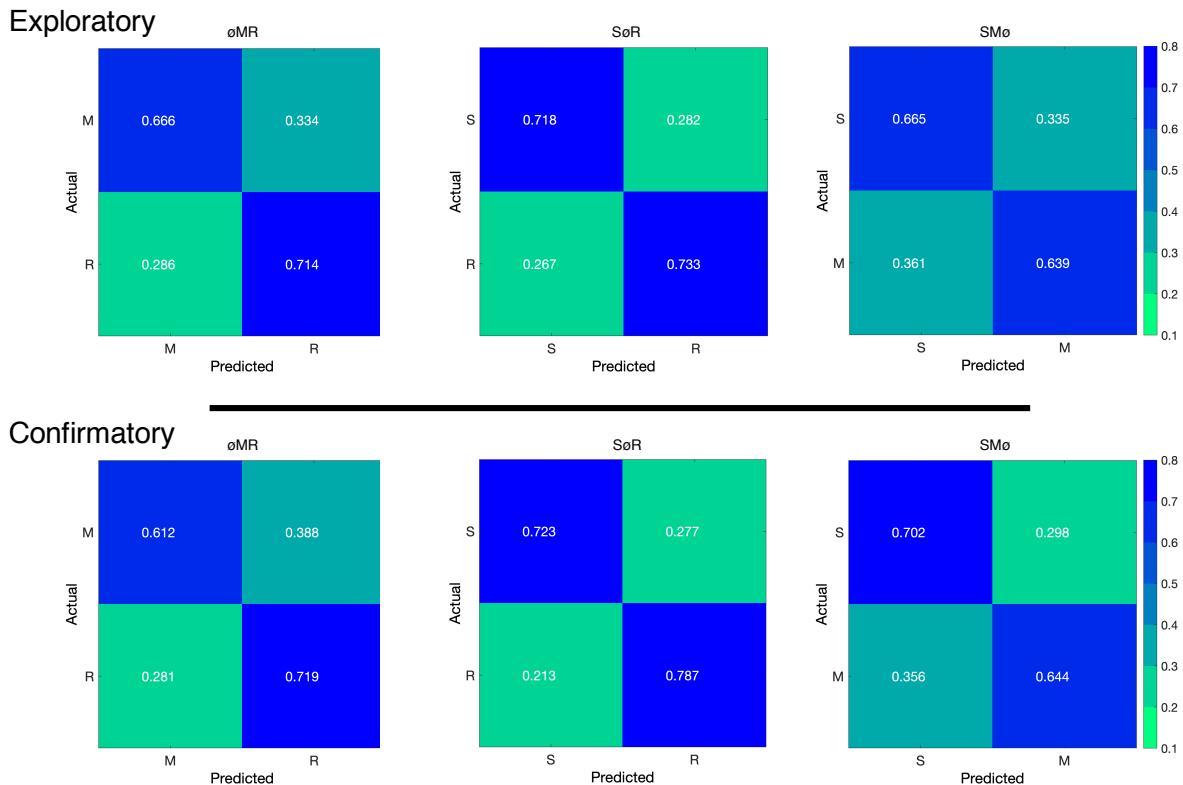


Figure A3. The confusion matrices represent a re-calculation of the classification accuracies for each category from the primary analysis. This re-calculation is meant to make the accuracies presented in the primary analysis (chance = 33%) equivalent to the classification accuracies presented in the supplementary analysis (chance = 50%).