

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

<sup>3</sup> Zachary J. Cole<sup>1</sup>, Karl M. Kuntzelman<sup>1</sup>, Michael D. Dodd<sup>1</sup>, & Matthew R. Johnson<sup>1</sup>

<sup>4</sup> <sup>1</sup> University of Nebraska-Lincoln

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. Work  
done to develop the analysis approach was supported by NSF/EPSCoR grant #1632849  
(MRJ and MDD) and NIH grant GM130461 awarded to MRJ and colleagues. Additionally,  
this work was supported by the National Institute of General Medical Sciences of the  
National Institutes of Health [grant number P20 GM130461] and the Rural Drug Addiction  
Research Center at the University of Nebraska-Lincoln. The content is solely the  
responsibility of the authors and does not necessarily represent the official views of the  
National Institutes of Health or the University of Nebraska.

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

17

## Abstract

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

40        *Keywords:* deep learning, eye tracking, convolutional neural network, cognitive state,  
41 endogenous attention

42        Word count: 7960

43

## Introduction

44 The association between eye movements and mental activity is a fundamental topic of  
45 interest in attention research that has provided a foundation for developing a wide range of  
46 human assistive technologies. Early work by Yarbus (1967) showed that eye movement  
47 patterns appear to differ qualitatively depending on the task-at-hand (for a review of this  
48 work, see Tatler, Wade, Kwan, Findlay, & Velichkovsky, 2010). A replication of this work by  
49 DeAngelus and Pelz (2009) showed that the differences in eye movements between tasks can  
50 be quantified, and appear to be somewhat generalizable. Technological advances and  
51 improvements in computing power have allowed researchers to make inferences regarding the  
52 mental state underlying eye movement data, also known as the “inverse Yarbus process”  
53 (Haji-Abolhassani & Clark, 2014).

54 Current state-of-the-art machine learning and neural network algorithms are capable of  
55 identifying diagnostic patterns for the purpose of decoding a variety of data types, but the  
56 inner workings of the resulting model solutions are difficult or impossible to interpret.  
57 Algorithms that provide such solutions are referred to as *black box* models. Dissections of  
58 black box models have been largely uninformative (Zhou, Bau, Oliva, & Torralba, 2019),  
59 limiting the potential for researchers to apply the mechanisms underlying successful  
60 classification of the data. Still, black box models provide a powerful solution for  
61 technological applications such as human-computer interfaces (HCI; for a review, see  
62 Lukander, Toivanen, & Puolamäki, 2017). While the internal operations of the model  
63 solutions used for HCI applications do not necessarily need to be interpretable to serve their  
64 purpose, Lukander et al. (2017) pointed out that the inability to interpret the mechanisms  
65 underlying the function of black box solutions impedes the generalizability of these methods,  
66 and increases the difficulty of expanding these findings to real life applications. To ground  
67 these solutions, researchers guide decoding efforts by using eye movement data and/or  
68 models with built-in theoretical assumptions. For instance, eye movement data is processed

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

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

Currently, there is not a clearly established upper limit to how well cognitive task can be classified from eye movement data. 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 to relatively robust classification (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 related 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 decoded to the extent to which the data are capable of representing differences between

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%)

<sup>150</sup> categories. Given that the cognitive processes underlying distinct tasks are often overlapping<sup>151</sup> (Coco & Keller, 2014), decreasing the granularity of the data may actually limit the potential

152 of the algorithm to make fine-grained distinctions between diagnostic components underlying  
153 the tasks to be decoded.

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

168 The current study explored pragmatic solutions to the problem of classifying task from  
169 eye movement data by submitting unprocessed x-coordinate, y-coordinate, and pupil size  
170 data to a convolutional neural network (CNN) model. Instead of transforming the data into  
171 theoretically defined units, we allowed the network to learn meaningful patterns in the data  
172 on its own. CNNs have a natural propensity to develop low-level feature detectors similar to  
173 the primary visual cortex (e.g., Seeliger et al., 2018); for this reason, they are commonly  
174 implemented for image classification. To test the possibility that the image data are better  
175 suited to the CNN classifier, the data were also transformed from raw timelines into simple  
176 image representations. To our knowledge, no study has attempted to address the inverse  
177 Yarbus problem using any combination of the following methods: (1) Non-aggregated data,

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

187

## Methods

188 **Participants**

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

197 **Materials and Procedures**

198 Each participant viewed a series of indoor and outdoor scene images while carrying out  
199 a search, memorization, or rating task. For the memorization task, participants were  
200 instructed to memorize the image for a forced choice recognition test. At the end of each  
201 Memorize trial, the participants were prompted to indicate which of two images was just  
202 presented. For the rating task, participants were asked to think about how they would rate

203 the image on a scale from 1 (very unpleasant) to 7 (very pleasant). The participants were  
204 prompted to provide a rating immediately after viewing the image. For the search task,  
205 participants were instructed to find a small “Z” or “N” embedded in the image. Trials  
206 containing a target ( $n = 5$ ) were not analyzed but were included in the experiment design to  
207 encourage searching behavior on other Search trials. Trials containing the target were  
208 excluded because search behavior was likely to stop if the target was found, adding  
209 considerable noise to the eye movement data. For consistency between trial types,  
210 participants were prompted to indicate if they found a “Z” or “N” at the end of each Search  
211 trial.

212 The same materials were used in both experiments with a minor variation in the  
213 procedures. In the Confirmatory experiment, participants were directed as to where search  
214 targets might appear in the image (e.g., on flat surfaces). No such instructions were provided  
215 in the Exploratory experiment.

216 In both experiments, participants completed three mixed or uniform blocks of 40 trials  
217 ( $n = 120$  trials). Block type was assigned in counterbalanced order. When the blocks were  
218 mixed, the trial types were randomly intermixed within the block. For uniform blocks, each  
219 block consisted entirely of one of the three conditions (Search, Memorize, Rate), with block  
220 types presented in random order. Each stimulus image was presented for 8 seconds. The  
221 pictures were presented in color, with a size of 1024 x 768 pixels, subtending a visual angle of  
222  $23.8^\circ \times 18.0^\circ$ .

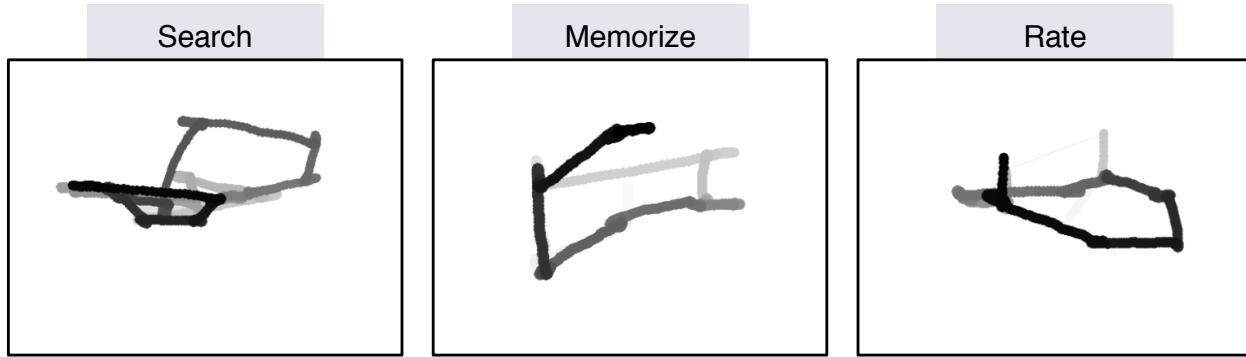
223 Eye movements were recorded using an SR Research EyeLink 1000 eye tracker with a  
224 sampling rate of 1000Hz. Only the right eye was recorded. The system was calibrated using  
225 a nine-point accuracy and validity test. Errors greater than  $1^\circ$  or averaging greater than  $0.5^\circ$   
226 in total were re-calibrated.

**227 Datasets**

228 On some trials, a probe was presented on the screen six seconds after the onset of the  
229 trial. To avoid confounds resulting from the probe, only the first six seconds of the data for  
230 each trial was analyzed. Trials that contained fewer than 6000 samples within the first six  
231 seconds of the trial were excluded before analysis. For both datasets, the trials were pooled  
232 across participants. After excluding trials, the Exploratory dataset consisted of 12,177 of the  
233 16,740 total trials, and the Confirmatory dataset consisted of 9,301 of the 10,395 total trials.

234 The raw x-coordinate, y-coordinate, and pupil size data collected at every sampling  
235 time point in the trial were used as inputs to the deep learning classifier. These data were  
236 also used to develop plot image datasets that were classified separately from the raw timeline  
237 datasets. For the plot image datasets, the timeline data for each trial were converted into  
238 scatterplot diagrams. The x- and y- coordinates and pupil size were used to plot each data  
239 point onto a scatterplot (e.g., see Figure 1). The coordinates were used to plot the location  
240 of the dot, pupil size was used to determine the relative size of the dot, and shading of the  
241 dot was used to indicate the time-course of the eye movements throughout the trial. The  
242 background of the plot images and first data point were white. Each subsequent data point  
243 was one shade darker than the previous data point until the final data point was reached.  
244 The final data point was black. For standardization, pupil size was divided by 10, and one  
245 unit was added. The plots were sized to match the dimensions of the data collection monitor  
246 (1024 x 768 pixels) and then shrunk to (240 x 180 pixels) in an effort to reduce the  
247 dimensionality of the data.

248 **Data Subsets.** The full timeline dataset was structured into three columns  
249 representing the x- and y- coordinates, and pupil size for each data point collected in the  
250 first six seconds of each trial. To systematically assess the predictive value of each XYP (i.e.,  
251 x-coordinates, y-coordinates, pupil size) component of the data, the timeline and image  
252 datasets were batched into subsets that excluded one of the components (i.e., XYØ, XØP,



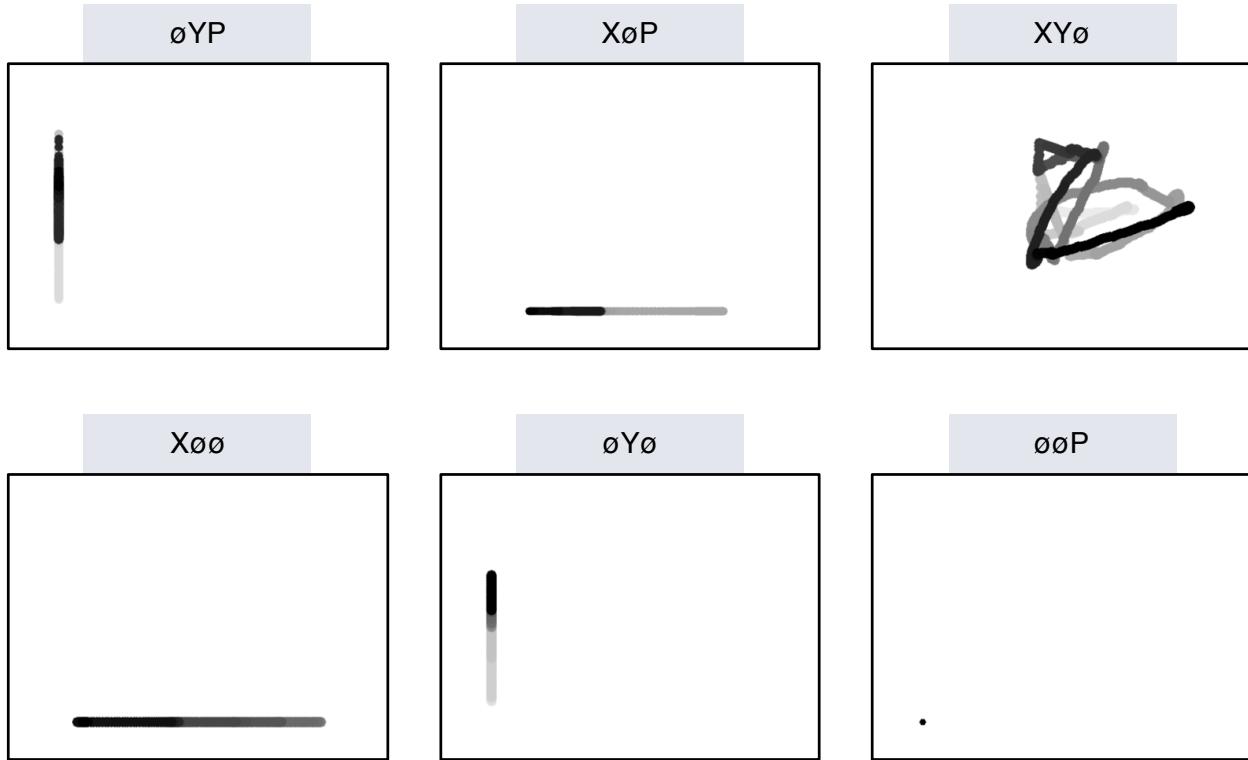
*Figure 1.* Each trial was represented as an image. Each sample collected within the trial was plotted as a dot in the image. Pupil size was represented by the size of the dot. The time course of the eye movements was represented by the gradual darkening of the dot over time.

253      $\emptyset YP$ ), or contained only one of the components (i.e.,  $X\emptyset\emptyset$ ,  $\emptyset Y\emptyset$ ,  $\emptyset\emptyset P$ ). For the timeline  
 254     datasets, this means that the columns to be excluded in each data subset were replaced with  
 255     zeros. The data were replaced with zeros because removing the columns would change the  
 256     structure of the data. The same systematic batching process was carried out for the image  
 257     dataset. See Figure 2 for an example of each of these image data subsets.

## 258     Classification

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

265         The development of these models was not guided by any formal theoretical  
 266     assumptions regarding the patterns or features likely to be extracted by the classifier. Like  
 267     many HCI models, the development of these models followed general intuitions concerned  
 268     with building a model architecture capable of transforming the data inputs into an  
 269     interpretable feature set that would not overfit the dataset. The models were developed  
 270     using version 0.3b of the DeLINEATE toolbox, which operates over a Keras backend



*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.

271 (<http://delineate.it>; Kuntzman et al., under review). Each training/test iteration randomly  
 272 split the data so that 70% of the trials were allocated to training, 15% to validation, and  
 273 15% to testing. Training of the model was stopped when validation accuracy did not improve  
 274 over the span of 100 epochs. Once the early stopping threshold was reached, the resulting  
 275 model was tested on the held-out test data. This process was repeated 10 times for each  
 276 model, resulting in 10 classification accuracy scores for each model. The resulting accuracy  
 277 scores were used for the comparisons against chance and other datasets or data subsets.

278 The models were developed and tested on the Exploratory dataset. Model  
 279 hyperparameters were adjusted until the classification accuracies appeared to peak. The  
 280 model architecture with the highest classification accuracy on the Exploratory dataset was  
 281 trained, validated, and tested independently on the Confirmatory dataset. This means that  
 282 the model that was used to analyze the Confirmatory dataset was not trained on the

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

285 **Analysis**

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

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

304 **Results**

305 **Timeline Data Classification**

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

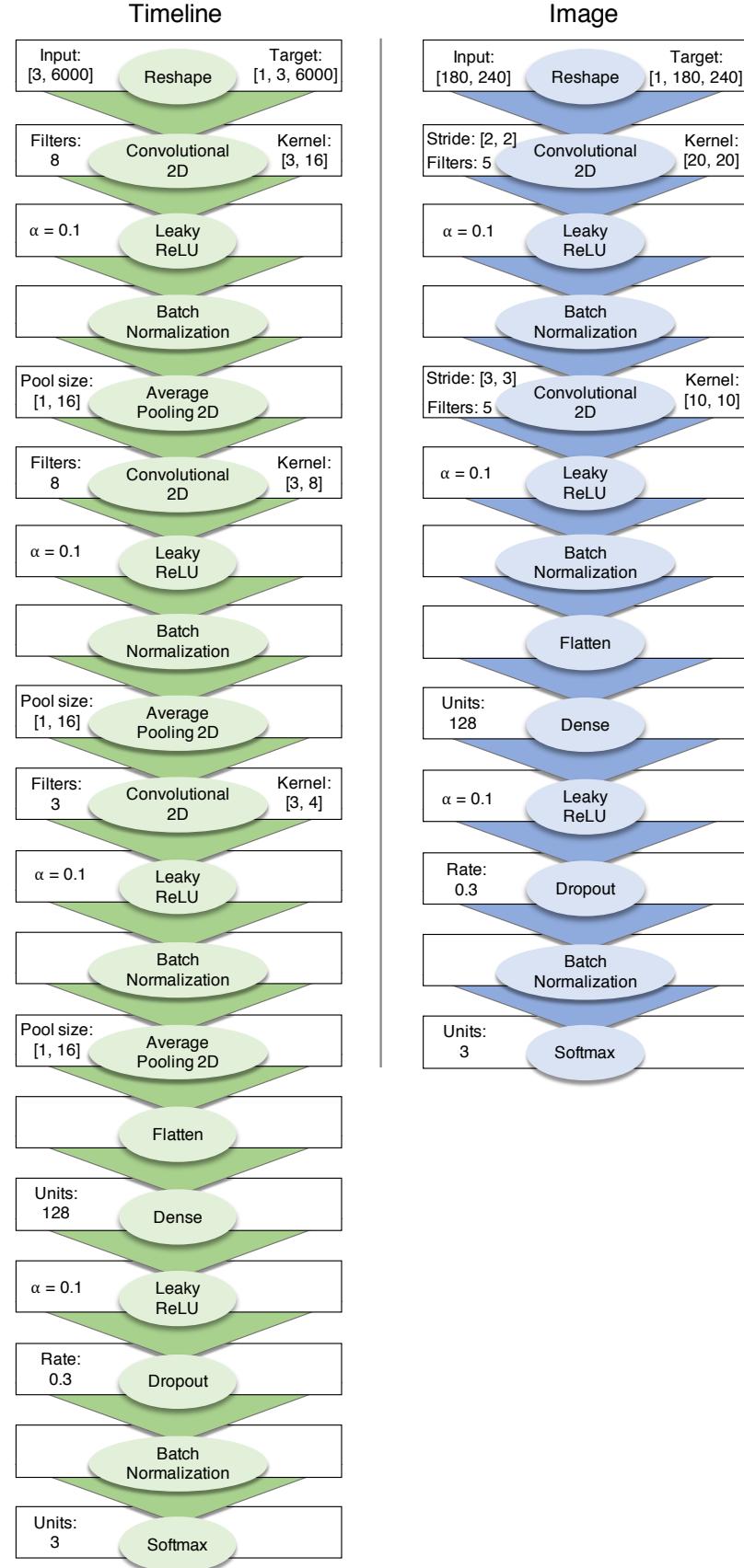


Figure 3. Two different model architectures were used to classify the timeline and image data. Both models were compiled using a categorical crossentropy loss function, and optimized with the Adam algorithm.

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

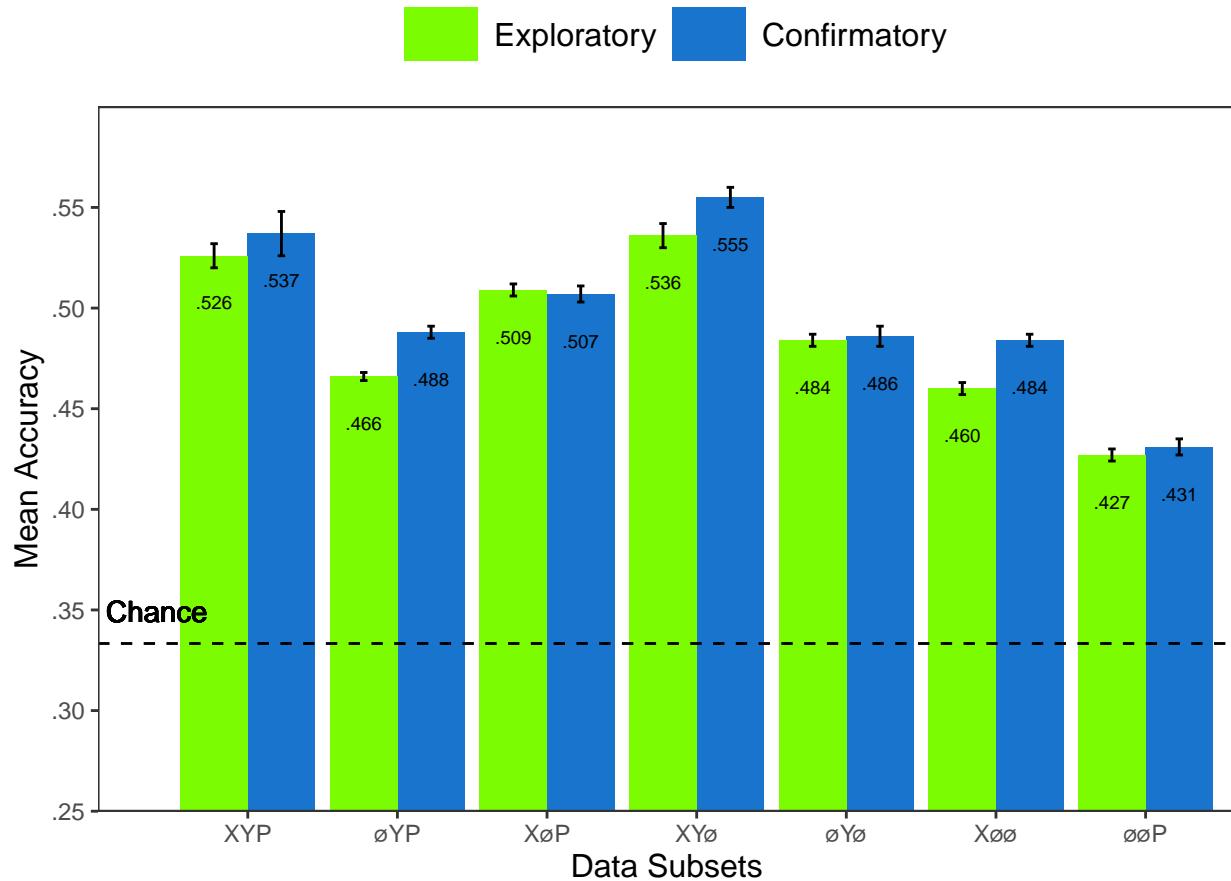


Figure 4. All of the data subsets were decoded at levels better than chance (.33). Each subset is labeled with the mean accuracy. The error bars represent standard errors.

313 There was a difference in classification accuracy for the XYP dataset and the subsets  
 314 that had the pupil size, x-coordinate, and y-coordinate data systematically removed ( $F_{(3,36)}$   
 315 = 47.471,  $p < .001$ ,  $\eta^2 = 0.798$ ). Post-hoc comparisons against the XYP dataset showed that  
 316 classification accuracies were not affected by the removal of pupil size or y-coordinate data  
 317 (see Table 2). The null effect present when pupil size was removed suggests that the pupil  
 318 size data were not contributing unique information that was not otherwise provided by the x-

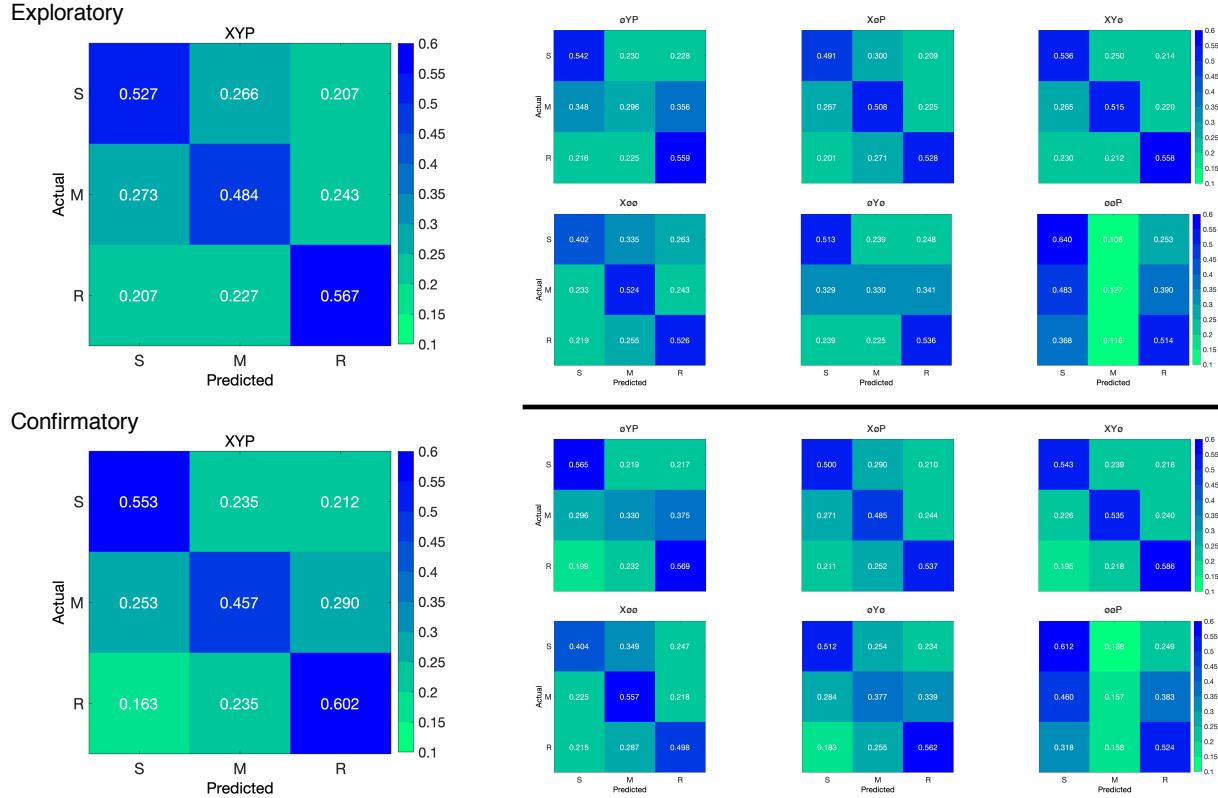


Figure 5. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents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

319 and y-coordinates. A strict significance threshold of  $\alpha = .05$  implies the same conclusion for  
 320 the y-coordinate data, but the relatively low degrees of freedom ( $df = 18$ ) and the borderline  
 321 observed  $p$ -value ( $p = .056$ ) afford the possibility that there exists a small effect. However,  
 322 classification for the  $\emptyset YP$  subset was significantly lower than the  $XYP$  dataset, showing that  
 323 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. $\emptyset YP$	9.420	< .001	5.210	< .001
XYP vs. $X\emptyset P$	2.645	.056	3.165	.016
XYP vs. $XY\emptyset$	1.635	.372	1.805	.288
$X\emptyset\emptyset$ vs. $\emptyset Y\emptyset$	5.187	< .001	0.495	.874
$X\emptyset\emptyset$ vs. $\emptyset\emptyset P$	12.213	< .001	10.178	< .001
$\emptyset Y\emptyset$ vs. $\emptyset\emptyset P$	7.026	< .001	9.683	< .001

324 There was also a difference in classification accuracies for the X $\emptyset\emptyset$ ,  $\emptyset Y\emptyset$ , and  $\emptyset\emptyset P$

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

326 classification accuracy for the  $\emptyset\emptyset P$  subset was lower than the X $\emptyset\emptyset$  and  $\emptyset Y\emptyset$  subsets.

327 Classification accuracy for the X $\emptyset\emptyset$  subset was higher than the  $\emptyset Y\emptyset$  subset. Altogether,

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

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

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

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

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

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

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

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

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

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

338 and Rate tasks (see Figure 5).

339 To test the generalizability of the model architecture, classification accuracies for the

340 XYP Exploratory and Confirmatory timeline datasets were compared. The Shapiro-Wilk

341 test for normality indicated that the Exploratory ( $W = 0.937, p = .524$ ) and Confirmatory

342 ( $W = 0.884, p = .145$ ) datasets were normally distributed, but Levene's test indicated that

343 the variances were not equal,  $F_{(1,18)} = 8.783, p = .008$ . Welch's unequal variances  $t$ -test did

344 not show a difference between the two datasets,  $t_{(13.045)} = 0.907, p = .381$ , Cohen's  $d =$

345 0.406. These findings indicate that the deep learning model decoded the Exploratory and

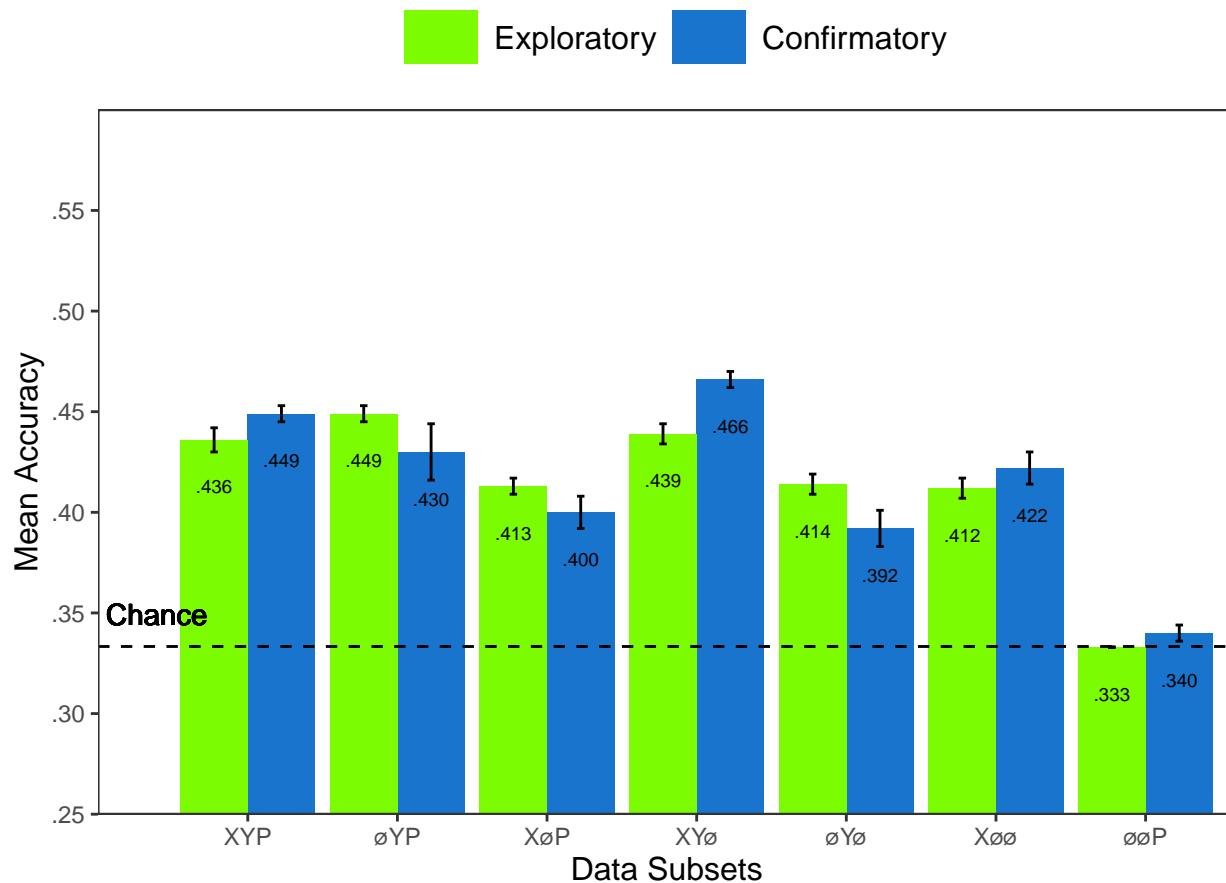
346 Confirmatory timeline datasets equally well, but the Confirmatory dataset classifications

347 were less consistent across training/test iterations (as indicated by the increase in standard

348 deviation).

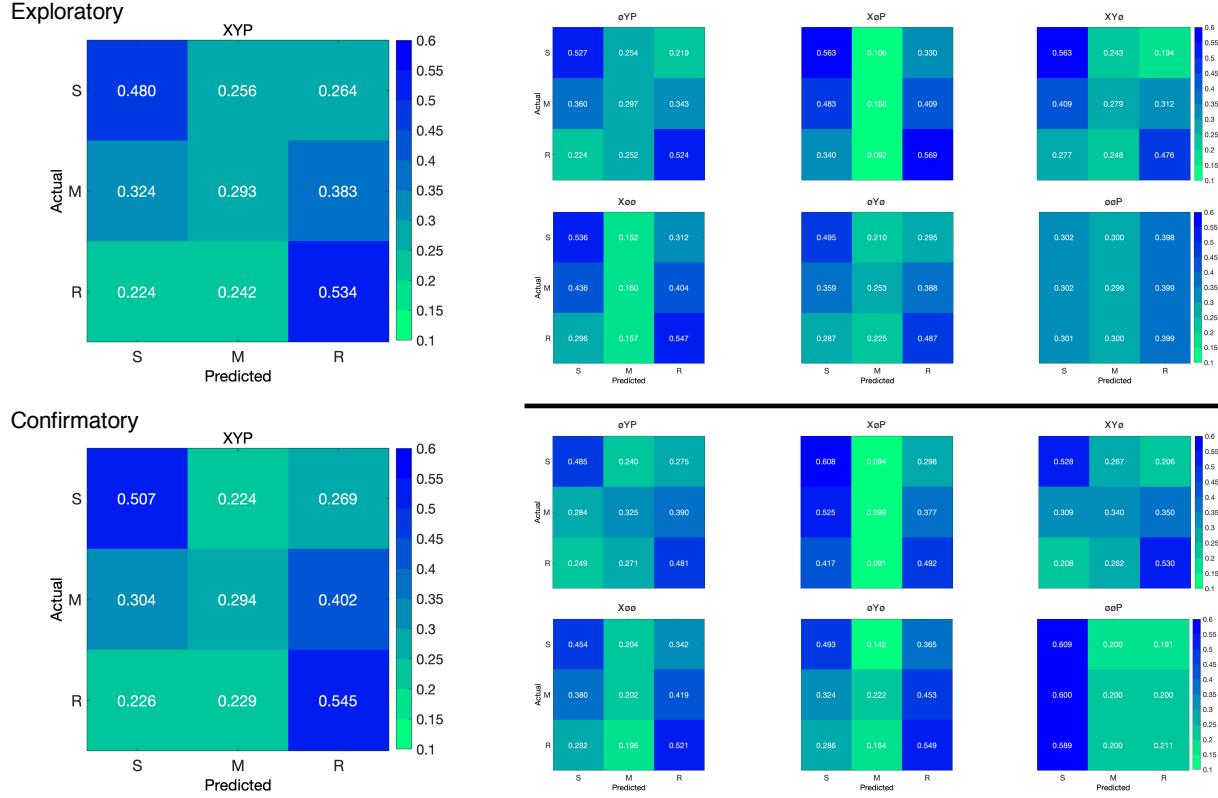
<sup>349</sup> **Plot Image Classification**

<sup>350</sup> **Exploratory.** Classification accuracies for the XYP plot image data were better  
<sup>351</sup> than chance ( $M = .436$ ,  $SD = .020$ ,  $p < .001$ ), but were less accurate than the classifications  
<sup>352</sup> for the XYP Exploratory timeline data ( $t_{(18)} = 10.813$ ,  $p < .001$ ). Accuracies for the  
<sup>353</sup> classifications for all subsets of the plot image data except the  $\emptyset\emptyset P$  subset were better than  
<sup>354</sup> chance (see Figure 6). Following the pattern expressed by the timeline dataset, the confusion  
<sup>355</sup> matrices showed that the Memorize condition was misclassified more often than the other  
<sup>356</sup> conditions, and appeared to be equally mis-identified as a Search or Rate condition (see  
<sup>357</sup> Figure 7).



<sup>358</sup> *Figure 6.* 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.

<sup>358</sup> There was a difference in classification accuracy between the XYP dataset and the data



*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.

359 subsets ( $F_{(4,45)} = 7.093, p < .001, \eta^2 = .387$ ). Post-hoc comparisons showed that compared  
 360 to the XYP dataset, there was no effect of removing pupil size or the x-coordinates, but  
 361 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. $\emptyset$ Y $\emptyset$	0.423	.906	2.807	.204
XoO vs. $\emptyset$ OoP	13.569	< .001	5.070	< .001
$\emptyset$ Y $\emptyset$ vs. $\emptyset$ OoP	13.235	< .001	7.877	< .001

362 There was also a difference in classification accuracies between the XoO,  $\emptyset$ Y $\emptyset$ , and  
 363  $\emptyset$ OoP 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 showed that there was no difference in classification accuracies for the XØØ and ØYØ subsets, but classification for the ØØP subset were less accurate than the XØØ and ØYØ subsets.

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

To test the generalizability of the model architecture, the classification accuracies for the XYP Exploratory and Confirmatory plot image datasets were compared. The independent samples *t*-test comparing the classification accuracies for the Exploratory and Confirmatory plot image datasets did not show a significant difference,  $t_{(18)} = 1.777, p = .092$ , Cohen's *d* = 0.795.

## Discussion

The present study aimed to produce a practical and reliable example of a black box solution to the inverse Yarbus problem. To implement this solution, we classified raw timeline and minimally processed plot image data using a CNN model architecture. To our knowledge, this study was the first to provide a solution to determining mental state from eye movement data using each of the following: (1) Non-aggregated eye tracking data (i.e., raw x-coordinates, y-coordinates, pupil size), (2) timeline and image data formats (see Figure 2), and (3) a black box CNN architecture. This study probed the relative predictive value of the x-coordinate, y-coordinate, and pupil size components of the eye movement data

390 using a CNN. The CNN was able to decode the timeline and plot image data better than  
391 chance, although only the timeline datasets were decoded with accuracies comparable to  
392 other state-of-the-art approaches. Datasets with lower classification accuracies were not able  
393 to differentiate the cognitive processes underlying the Memorize task from the cognitive  
394 processes underlying the Search and Rate tasks. Decoding subsets of the data revealed that  
395 pupil size was the least uniquely informative component of the eye movement data. This  
396 pattern of findings was consistent between the Exploratory and Confirmatory datasets.

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

411 When considering the superior performance of the timeline data (vs., plot image data),  
412 we must also consider the differences in the model architectures. Because the structures of  
413 the timeline and plot image data formats were different, the models decoding those data  
414 structures also needed to be different. Both model architectures were optimized individually  
415 on the Exploratory dataset before being tested on the Confirmatory dataset. For both

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

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

438 To more clearly understand how the different tasks influenced the decodability of the  
439 eye movement data, additional analyses were conducted on the Exploratory and  
440 Confirmatory timeline datasets (see Appendix). For the main supplementary analysis, the  
441 data subsets were re-submitted to the CNN and re-classified as 2-category task sets. In

addition to the main supplementary analysis, the results from the primary analysis were re-calculated from 3-category task sets to 2-category task sets. In the primary analyses, the Memorize condition was predicted with the lowest accuracy, but mis-classifications of the Search and Rate trials were most often categorized as Memorize. As a whole, this pattern of results indicated a general bias for uncertain trials to be categorized as Memorize. The re-calculation analysis generally replicated the pattern of results seen in the main supplementary analysis but with larger variance, suggesting that including lower-accuracy trial types can decrease the consistency of classifier performance. Overall, the findings from this supplemental analysis show that conclusions drawn from comparisons between approaches that do not use the same task sets, or the same number of tasks, could be potentially uninterpretable because the features underlying the task categories are interpreted differently by the neural network algorithm.

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

469 randomly assigned, there is the possibility that uncontrolled stimulus properties known to  
470 affect pupil size impeded the CNN's capacity to detect patterns in the pupil size data.

471 The findings from the current study support the notion that black box CNNs are a  
472 viable approach to determining task from eye movement data. In a recent review, Lukander  
473 et al. (2017) expressed concern regarding the lack of generalizability of black box approaches  
474 when decoding eye movement data. Overall, the current study showed a consistent pattern  
475 of results for the XYP timeline and image datasets, but some minor inconsistencies in the  
476 pattern of results for the x- and y- coordinate subset comparisons. These inconsistencies may  
477 be a product of overlap in the cognitive processes underlying the three tasks. When the data  
478 are batched into subsets, at least one dimension (i.e., x-coordinates, y-coordinates, or pupil  
479 size) is removed, leading to a potential loss of information. When the data provide fewer  
480 meaningful distinctions, finer-grained inferences are necessary for the tasks to be  
481 distinguishable. As shown by Coco and Keller (2014), eye movement data can be more  
482 effectively decoded when the cognitive processes underlying the tasks are explicitly  
483 differentiable. While the cognitive processes distinguishing memorizing, searching, or rating  
484 an image are intuitively different, the eye movements elicited from these cognitive processes  
485 are not easily differentiated. To correct for potential mismatches between the distinctive  
486 task-diagnostic features in the data and the level of distinctiveness required to classify the  
487 tasks, future research could more definitively conceptualize the cognitive processes  
488 underlying the task-at-hand.

489 Classifying mental state from eye movement data is often carried out in an effort to  
490 advance technology to improve educational outcomes, strengthen the independence of  
491 physically and mentally handicapped individuals, or improve HCI's (Koochaki &  
492 Najafizadeh, 2018). Given the previous questions raised regarding the reliability and  
493 generalizability of black-box CNN classification, the current study first tested models on an  
494 exploratory dataset, then confirmed the outcome using a second independent dataset.

495 Overall, the findings of this study indicate that this black-box approach is capable of  
496 producing a stable and generalizable outcome. Additionally, the supplementary analyses  
497 showed that different task sets, or a different number of tasks, could lead the algorithm to  
498 interpret features differently, which should be taken into account when comparing task  
499 classification approaches. Future studies that incorporate features from the stimulus might  
500 have the potential to surpass current state-of-the-art classification. According to Bulling,  
501 Weichel, and Gellersen (2013), incorporating stimulus feature information into the dataset  
502 may improve accuracy relative to decoding gaze location data and pupil size. Alternatively,  
503 Borji and Itti (2014) suggested that accounting for salient features in the the stimulus might  
504 leave little to no room for theoretically defined classifiers to consider mental state. Future  
505 research should examine the potential for the inclusion of stimulus feature information in  
506 addition to the eye movement data to boost black-box CNN classification accuracy of image  
507 data beyond that of timeline data.

508

## References

- 509 Boisvert, J. F. G., & Bruce, N. D. B. (2016). Predicting task from eye movements: On the  
510 importance of spatial distribution, dynamics, and image features. *Neurocomputing*,  
511 207, 653–668. <https://doi.org/10.1016/j.neucom.2016.05.047>
- 512 Borji, A., & Itti, L. (2014). Defending Yarbus: Eye movements reveal observers' task.  
513 *Journal of Vision*, 14(3), 1–21. <https://doi.org/10.1167/14.3.29>
- 514 Bulling, A., Weichel, C., & Gellersen, H. (2013). EyeContext: Recognition of high-level  
515 contextual cues from human visual behaviour. In *Proceedings of the SIGCHI  
516 Conference on Human Factors in Computing Systems - CHI '13* (p. 305). Paris,  
517 France: ACM Press. <https://doi.org/10.1145/2470654.2470697>
- 518 Castelhano, M. S., Mack, M. L., & Henderson, J. M. (2009). Viewing task influences eye  
519 movement control during active scene perception. *Journal of Vision*, 9(3), 1–15.  
520 <https://doi.org/10.1167/9.3.6>
- 521 Coco, M. I., & Keller, F. (2014). Classification of visual and linguistic tasks using  
522 eye-movement features. *Journal of Vision*, 14(3), 1–18.  
523 <https://doi.org/10.1167/14.3.11>
- 524 DeAngelus, M., & Pelz, J. B. (2009). Top-down control of eye movements: Yarbus revisited.  
525 *Visual Cognition*, 17(6-7), 790–811. <https://doi.org/10.1080/13506280902793843>
- 526 Greene, M. R., Liu, T., & Wolfe, J. M. (2012). Reconsidering Yarbus: A failure to predict  
527 observers' task from eye movement patterns. *Vision Research*, 62, 1–8.  
528 <https://doi.org/10.1016/j.visres.2012.03.019>
- 529 Haji-Abolhassani, A., & Clark, J. J. (2014). An inverse Yarbus process: Predicting observers'  
530 task from eye movement patterns. *Vision Research*, 103, 127–142.

- 531           <https://doi.org/10.1016/j.visres.2014.08.014>
- 532   Henderson, J. M., Shinkareva, S. V., Wang, J., Luke, S. G., & Olejarczyk, J. (2013).  
533           Predicting Cognitive State from Eye Movements. *PLoS ONE*, 8(5), e64937.  
534           <https://doi.org/10.1371/journal.pone.0064937>
- 535   Kahneman, D., & Beatty, J. (1966). Pupil Diameter and Load on Memory. *Science*,  
536           154(3756), 1583–1585. Retrieved from <http://www.jstor.org/stable/1720478>
- 537   Kanan, C., Ray, N. A., Bseiso, D. N. F., Hsiao, J. H., & Cottrell, G. W. (2014). Predicting  
538           an observer's task using multi-fixation pattern analysis. In *Proceedings of the*  
539           *Symposium on Eye Tracking Research and Applications - ETRA '14* (pp. 287–290).  
540           Safety Harbor, Florida: ACM Press. <https://doi.org/10.1145/2578153.2578208>
- 541   Karatekin, C., Couperus, J. W., & Marcus, D. J. (2004). Attention allocation in the  
542           dual-task paradigm as measured through behavioral and psychophysiological  
543           responses. *Psychophysiology*, 41(2), 175–185.  
544           <https://doi.org/10.1111/j.1469-8986.2004.00147.x>
- 545   Koochaki, F., & Najafizadeh, L. (2018). Predicting Intention Through Eye Gaze Patterns.  
546           In *2018 IEEE Biomedical Circuits and Systems Conference (BioCAS)* (pp. 1–4).  
547           <https://doi.org/10.1109/BIOCAS.2018.8584665>
- 548   Król, M. E., & Król, M. (2018). The right look for the job: Decoding cognitive processes  
549           involved in the task from spatial eye-movement patterns. *Psychological Research*, 84,  
550           245–258. <https://doi.org/10.1007/s00426-018-0996-5>
- 551   Lukander, K., Toivanen, M., & Puolamäki, K. (2017). Inferring Intent and Action from Gaze  
552           in Naturalistic Behavior: A Review. *International Journal of Mobile Human*  
553           *Computer Interaction*, 9(4), 41–57. <https://doi.org/10.4018/IJMHCI.2017100104>

- 554 MacInnes, W., Joseph, Hunt, A. R., Clarke, A. D. F., & Dodd, M. D. (2018). A Generative  
555 Model of Cognitive State from Task and Eye Movements. *Cognitive Computation*,  
556 10(5), 703–717. <https://doi.org/10.1007/s12559-018-9558-9>
- 557 Mills, M., Hollingworth, A., Van der Stigchel, S., Hoffman, L., & Dodd, M. D. (2011).  
558 Examining the influence of task set on eye movements and fixations. *Journal of*  
559 *Vision*, 11(8), 1–15. <https://doi.org/10.1167/11.8.17>
- 560 Porter, G., Troscianko, T., & Gilchrist, I. D. (2007). Effort during visual search and  
561 counting: Insights from pupillometry. *Quarterly Journal of Experimental Psychology*  
562 (2006), 60(2), 211–229. <https://doi.org/10.1080/17470210600673818>
- 563 Seeliger, K., Fritzsche, M., Güçlü, U., Schoenmakers, S., Schoffelen, J.-M., Bosch, S. E., &  
564 van Gerven, M. A. J. (2018). Convolutional neural network-based encoding and  
565 decoding of visual object recognition in space and time. *NeuroImage*, 180, 253–266.  
566 <https://doi.org/10.1016/j.neuroimage.2017.07.018>
- 567 Tatler, B. W., Wade, N. J., Kwan, H., Findlay, J. M., & Velichkovsky, B. M. (2010). Yarbus,  
568 Eye Movements, and Vision. *I-Perception*, 1(1), 7–27. <https://doi.org/10.1068/i0382>
- 569 Wang, C.-A., Baird, T., Huang, J., Coutinho, J. D., Brien, D. C., & Munoz, D. P. (2018).  
570 Arousal Effects on Pupil Size, Heart Rate, and Skin Conductance in an Emotional  
571 Face Task. *Frontiers in Neurology*, 9, 1029. <https://doi.org/10.3389/fneur.2018.01029>
- 572 Yarbus, A. (1967). *Eye Movements and Vision*. New York, NY: Plenum Press.
- 573 Zhou, B., Bau, D., Oliva, A., & Torralba, A. (2019). Comparing the Interpretability of Deep  
574 Networks via Network Dissection. In W. Samek, G. Montavon, A. Vedaldi, L. K.  
575 Hansen, & K.-R. Müller (Eds.), *Explainable AI: Interpreting, Explaining and*  
576 *Visualizing Deep Learning* (pp. 243–252). Cham: Springer International Publishing.  
577 [https://doi.org/10.1007/978-3-030-28954-6\\_12](https://doi.org/10.1007/978-3-030-28954-6_12)

578

## Appendix

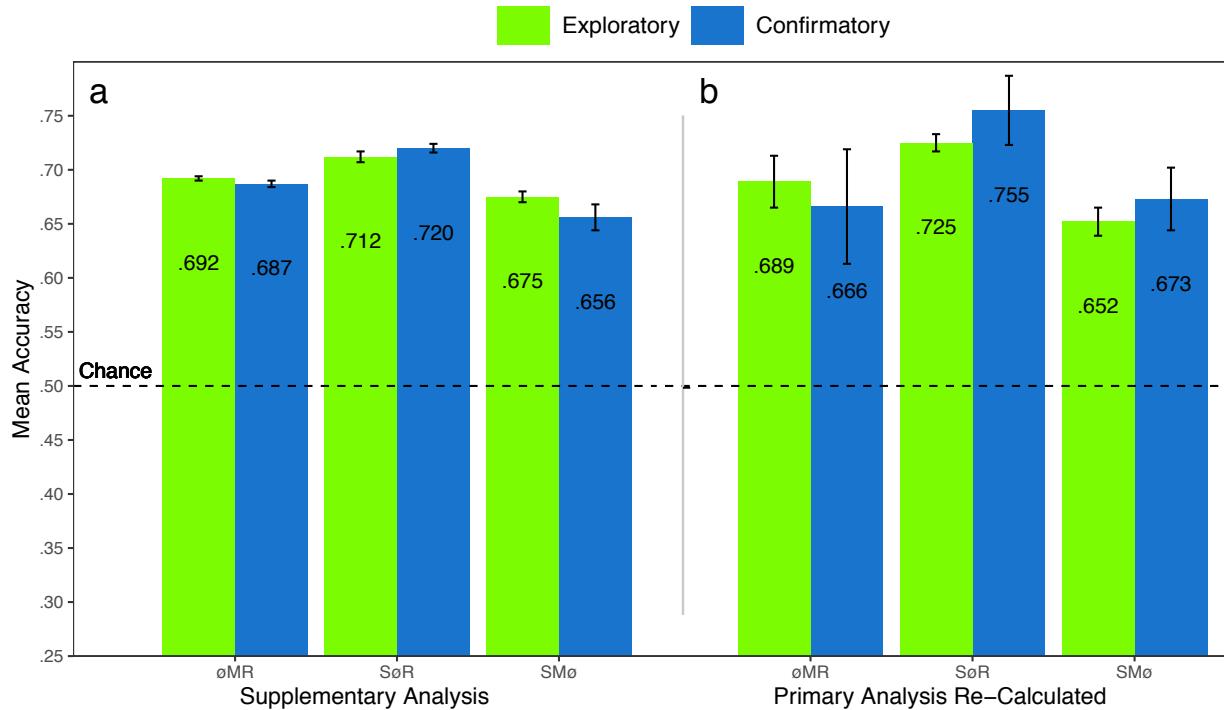
579        Additional analyses were conducted in an attempt to clarify the effect of task on  
 580 classification accuracy. These supplementary analyses were not seen as central to the current  
 581 study, but could prove to be informative to researchers attempting to replicate or extend  
 582 these findings in the future. The results from the primary analysis showed that classification  
 583 accuracies were the lowest for the Memorize condition. To further understand why  
 584 classification accuracy was lower for the Memorize condition than it was for the Search or  
 585 Rate condition, the Exploratory and Confirmatory timeline datasets were systematically  
 586 batched into subsets with the Search (S), Memorize (M), or Rate (R) condition removed (i.e.,  
 587  $\emptyset$ MR, S $\emptyset$ R, SM $\emptyset$ ), and then run through the CNN classifier using the same methods as the  
 588 primary analysis, but with only two classes.

589        All of the data subsets analyzed in this supplementary analysis were decoded with  
 590 better than chance accuracy (see Figure 8a). The same pattern of results was observed in  
 591 both the Exploratory and Confirmatory datasets. When the Memorize condition was  
 592 removed, classification accuracy improved (see Table 4, Figure 8a). When the Rate condition  
 593 was removed, classification was the worst. When the Memorize condition was included (i.e.,  
 594 SM $\emptyset$  and  $\emptyset$ MR), mis-classifications were biased toward Memorize, and the Memorize  
 595 condition was more accurately predicted than the Search and Rate conditions (see Figure 9).

Table 4  
*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

596        The accuracies for all of the data subsets observed in the supplementary analysis were  
 597 higher than the accuracies observed in the main analysis. Although there is a clear difference  
 598 in accuracy, the primary analysis was classifying three categories (chance = .33) and the



*Figure 8.* The graph represents the average accuracy reported for each subset of the Exploratory and Confirmatory timeline data for (a) the supplementary analysis, and the (b) re-calculated accuracies from the primary analysis. All of the data subsets were decoded at levels better than chance (.50). The error bars represent standard errors.

supplementary analysis was classifying two categories (chance = .50). Because the baseline chance performance was different for the primary and supplemental analyses, any conclusions drawn from a comparison of the results of analyses could be misleading. For this reason, we revisited the results from the primary analysis and re-calculated the predictions to be equivalent to a 50% chance threshold. Because the cross-validation scheme implemented by the DeLINEATE toolbox (<http://delineate.it>; Kuntzelman et al., under review) guaranteed an equal number of trials in the test set were assigned to each condition for each dataset, we were able to re-calculate 2-category predictions from the 3-category predictions presented in the confusion matrices from the primary analysis (see Figure 5). The predictions were re-calculated using the following formula:  $\text{Prediction}_{(A,A,A \otimes C)} = \text{Prediction}_{(A,A,ABC)} / (\text{Prediction}_{(A,A,ABC)} + \text{Prediction}_{(A,C,ABC)})$ . For example, accuracy for the Search classification for S $\otimes$ R would be calculated with the following:  $\text{Prediction}_{(S,S,S \otimes R)} = \text{Prediction}_{(S,S,SMR)} / (\text{Prediction}_{(S,S,SMR)} + \text{Prediction}_{(S,R,SMR)})$ , where  $\text{Prediction}_{(S,R,S \otimes R)}$  is

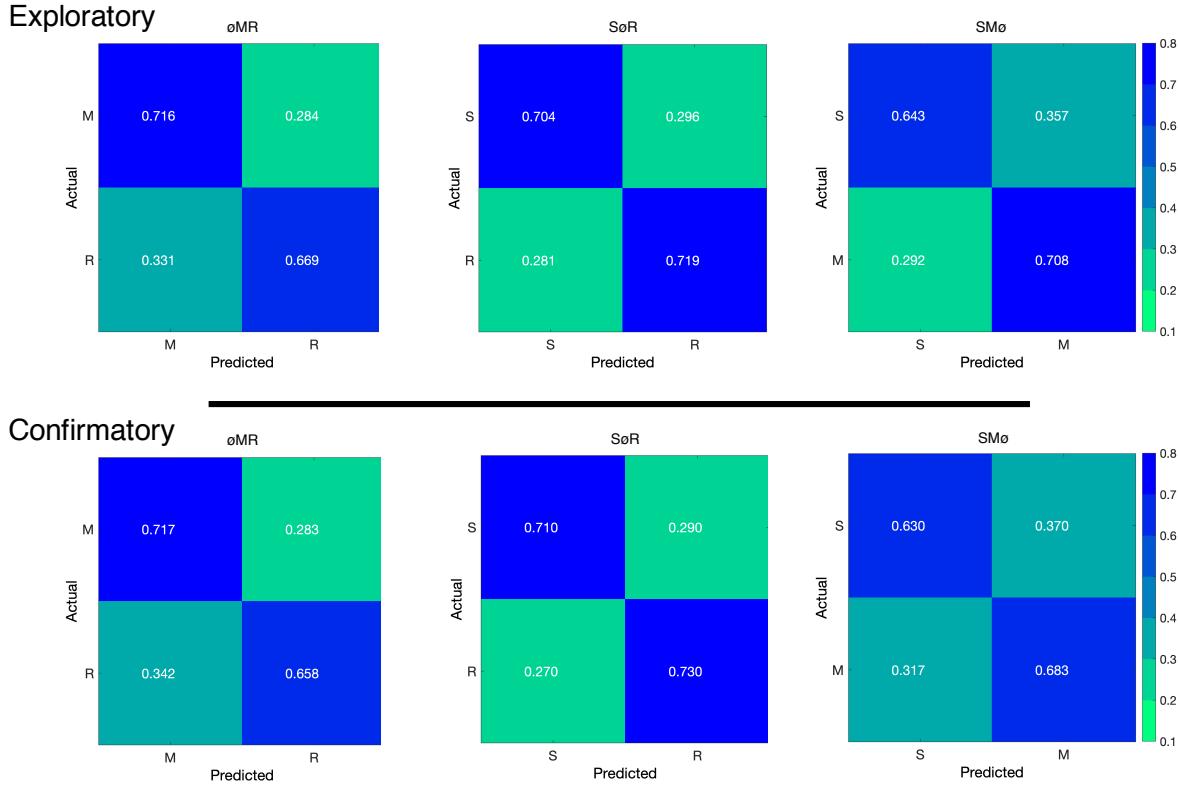


Figure 9. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents the condition that was predicted by the model.

612 the ratio of Search trials that were misclassified as Rate.

613 The results for the re-calculated predictions followed a pattern similar to the main  
 614 supplementary analysis (see Figure 8b). Looking back at the primary analysis, the 3-category  
 615 classifications predicted the Memorize conditions with the lowest accuracy (c.f., Search and  
 616 Rate conditions), and mis-classifications of the Search and Rate conditions were most often  
 617 categorized as Memorize (see Figure 5). Because the Memorize condition was mis-classified  
 618 more often than the other conditions in the primary analysis, the removal of the third class  
 619 in the re-calculated SMø and øM subsets resulted in a disproportionate amount of  
 620 mis-classified Memorize trials being removed from the data subset, somewhat eliminating the  
 621 tendency to mis-classify Search and Rate trials as Memorize (see Figure 10). Nevertheless,  
 622 the re-calculated SMø and øMR subsets were classified less accurately than SøR.

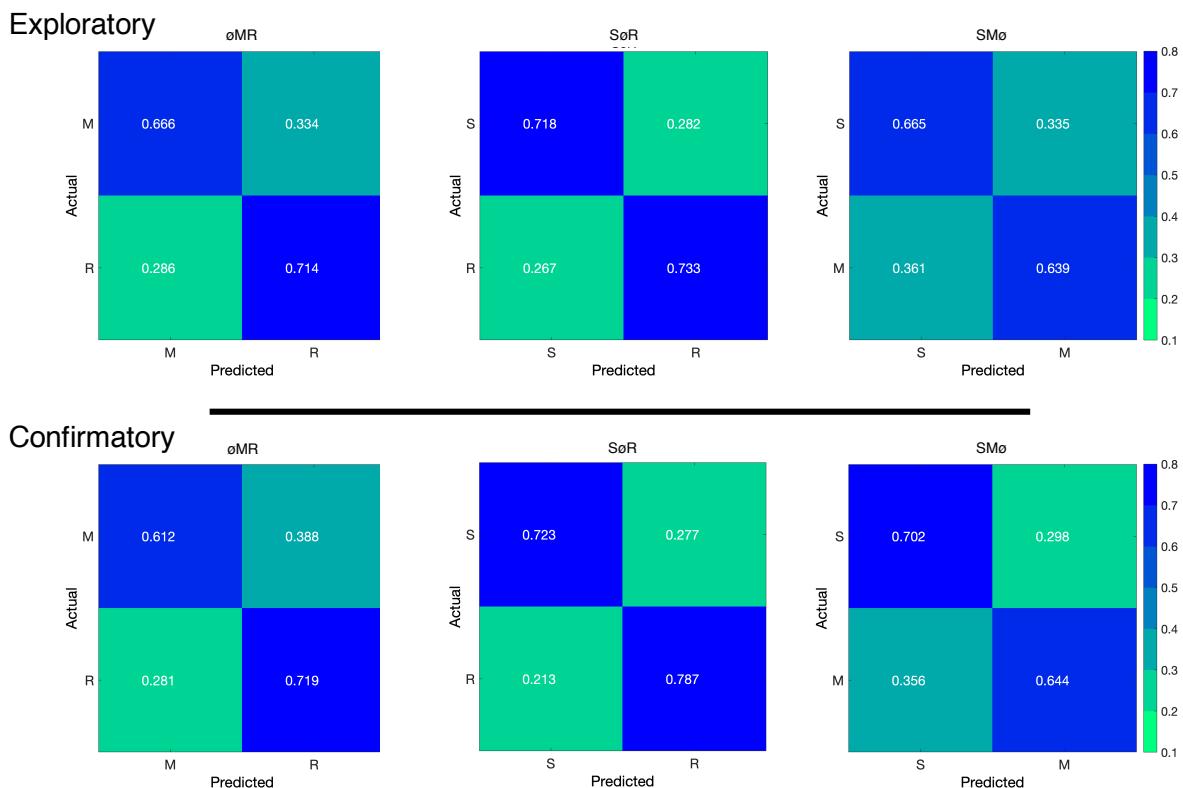


Figure 10. The confusion matrices represent a re-calculation of the classification accuracies for each category from the primary analysis. This re-calculation is meant to make the accuracies presented in the primary analysis (chance = .33) equivalent to the classification accuracies presented in the supplementary analysis (chance = .50).