

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

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

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

#### 5 Author Note

6 The data used for the exploratory and confirmatory analyses in the present manuscript  
7 are derived from experiments funded by NIH/NEI Grant 1R01EY022974 to MDD. Work  
8 done to develop the analysis approach was supported by NSF/EPSCoR grant #1632849  
9 (MRJ and MDD) and NIH grant GM130461 to MRJ and colleagues. Additionally,  
10 this work was supported by the National Institute of General Medical Sciences of the  
11 National Institutes of Health. The data used for the exploratory and confirmatory analyses in the present manuscript  
12 are derived from experiments funded by NIH/NEI Grant 1R01EY022974 to MDD.  
13 Additionally, work done to develop the analysis approach was supported by NSF/EPSCoR  
14 grant #1632849 and NIH grant GM130461 awarded to MRJ and colleagues.

15 Correspondence concerning this article should be addressed to Zachary J. Cole, 238  
16 Burnett Hall, Lincoln, NE 68588-0308. E-mail: z@neuophysicole.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  
34 classifying task from eye movement data.

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

37        Word count: 7260

## Introduction

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

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

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

71 The scope of theoretically informed solutions to decoding eye movement data is limited  
72 to the extent of the current theoretical knowledge linking eye movements to cognitive and  
73 neurobiological processes. As our theoretical understanding of these processes develops, older  
74 theoretically informed models become outdated. Furthermore, these solutions are susceptible  
75 to any inaccurate preconceptions that are built into the theory. Consider the case of Greene,  
76 Liu, and Wolfe (2012), who were not able to classify task from commonly used aggregate eye  
77 movement features (i.e., number of fixations, mean fixation duration, mean saccade  
78 amplitude, percent of image covered by fixations) using correlations, a linear discriminant  
79 model, and a support vector machine (see Table 1). This led Greene and colleagues to  
80 question the robustness of Yarbus's (1967) findings, inspiring a slew of responses that  
81 successfully decoded the same dataset by aggregating the eye movements into different  
82 feature sets or implementing different model architectures (see Table 1; Haji-Abolhassani &  
83 Clark, 2014; Borji & Itti, 2014; Kanan, Ray, Bseiso, Hsiao, & Cottrell, 2014). The  
84 subsequent re-analyses of these data support Yarbus (1967) and the notion that mental state  
85 can be decoded from eye movement data using a variety of combinations of data features and  
86 model architectures. Collectively, these re-analyses did not point to an obvious global  
87 solution capable of clarifying future approaches to the inverse Yarbus problem beyond what  
88 could be inferred from black box model solutions, but did provide a wide-ranging survey of a  
89 variety of methodological features that can be applied to theoretical or black box approaches  
90 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 (Castellano, Mack, & Henderson, 2009; DeAngelis & Peiz, 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 Table 1, Coco and Keller (2014) categorized

(2014) categorized the same eye movement features used by Green et al. (2012) with respect to the relative contribution of latent visual and linguistic components of three tasks (visual search, name the picture, name objects in the picture) with 84% accuracy (chance is 33%). While this manipulation is reminiscent of other experiments relying on the bottom-up influence of words and pictures (e.g., Henderson; Balay 2013; Boisvert & Bruce, 2016) the eye movements in the Kiboko (and Kellers, 2014) tasks can be attributed to the occurrence of top-down attentional processes. A conceptually similar follows up to this study classified tasks along two spatial and semantic dimensions, resulting in a 51% classification accuracy (Cháñ & Król, 2018). At closer looks 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, or by this body of literature suggests that the relative influence of top-down and bottom-up attentional may processes may have a role in determining the decoding ability 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>145</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<sup>146</sup> (Coco & Keller, 2014), decreasing the granularity of the data may actually limit the potential

- <sup>147</sup> of the algorithm to make fine-grained distinctions between diagnostic components underlying  
<sup>148</sup> the tasks to be decoded.

<sup>154</sup> The current state of the literature does not provide any firm guidelines for determining  
<sup>149</sup> The current state of the literature does not provide any firm guidelines for determining  
<sup>155</sup> what eye movement features are most meaningful, or what model architectures are best  
<sup>150</sup> what eye movement features are most meaningful, or what model architectures are best  
<sup>156</sup> suited for determining mental state from eye movements. The examples provided in Table 1  
<sup>151</sup> suited for determining mental state from eye movements. The examples provided in Table 1  
<sup>157</sup> used a variety of eye movement features and model architectures, most of which were  
<sup>152</sup> used a variety of eye movement features and model architectures, most of which were  
<sup>158</sup> effective to some extent. A proper comparison of these outcomes is difficult because these  
<sup>153</sup> effective to some extent. A proper comparison of these outcomes is difficult because these  
<sup>159</sup> datasets vary in levels of chance and data quality. Datasets with more tasks to be classified  
<sup>154</sup> datasets vary in levels of chance and data quality. Datasets with more tasks to be classified  
<sup>160</sup> have lower levels of chance, lowering the threshold for successful classification. Additionally,  
<sup>155</sup> have lower levels of chance, lowering the threshold for successful classification. Additionally,  
<sup>161</sup> datasets with a lower signal-to-noise ratio will have a lower achievable classification accuracy.  
<sup>156</sup> datasets with a lower signal-to-noise ratio will have a lower achievable classification accuracy.  
<sup>162</sup> For these reasons, outside of re-analyzing the same datasets, there is no consensus on how to  
<sup>157</sup> For these reasons, outside of re-analyzing the same datasets, there is no consensus on how to  
<sup>163</sup> establish direct comparisons of these model architectures. Given the inability to directly  
<sup>158</sup> establish direct comparisons of these model architectures. Given the inability to directly  
<sup>164</sup> compare the relative effectiveness of the various theoretical approaches present in the  
<sup>159</sup> compare the relative effectiveness of the various theoretical approaches present in the  
<sup>165</sup> literature, the current study addressed the inverse Yarbus problem by allowing a black box  
<sup>160</sup> literature, the current study addressed the inverse Yarbus problem by allowing a black box  
<sup>166</sup> model to self-determine the most informative features from minimally processed eye  
<sup>161</sup> model to self-determine the most informative features from minimally processed eye  
<sup>167</sup> movement data.  
<sup>162</sup> movement data.

<sup>168</sup> The current study explored pragmatic solutions to the problem of classifying task from  
<sup>163</sup> The current study explored pragmatic solutions to the problem of classifying task from  
<sup>169</sup> eye movement data by submitting unprocessed x-coordinate, y-coordinate, and pupil size  
<sup>164</sup> eye movement data by submitting unprocessed x-coordinate, y-coordinate, and pupil size  
<sup>170</sup> data to a convolutional neural network (CNN) model. Instead of transforming the data into  
<sup>165</sup> data to a convolutional neural network (CNN) model. Instead of transforming the data into  
<sup>171</sup> theoretically defined units, we allowed the network to learn meaningful patterns in the data  
<sup>166</sup> theoretically defined units, we allowed the network to learn meaningful patterns in the data  
<sup>172</sup> on its own. CNNs have a natural propensity to develop low-level feature detectors similar to  
<sup>167</sup> on its own. CNNs have a natural propensity to develop low-level feature detectors similar to  
<sup>173</sup> the primary visual cortex (e.g., Seeliger et al., 2018); for this reason, they are commonly  
<sup>168</sup> the primary visual cortex (e.g., Seeliger et al., 2018); for this reason, they are commonly  
<sup>174</sup> implemented for image classification. To test the possibility that the image data are better  
<sup>169</sup> implemented for image classification. To test the possibility that the image data are better  
<sup>175</sup> suited to the CNN classifier, the data were also transformed from raw timelines into simple  
<sup>170</sup> suited to the CNN classifier, the data were also transformed from raw timelines into simple  
<sup>176</sup> image representations. To our knowledge, no study has attempted to address the inverse  
<sup>171</sup> image representations. To our knowledge, no study has attempted to address the inverse  
<sup>177</sup> Yarbus problem using any combination of the following methods: (1) Non-aggregated data,  
<sup>172</sup> Yarbus problem using any combination of the following methods: (1) Non-aggregated data,

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

## 182 Methods

### 183 Participants

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

### 192 Materials and Procedures

193 Each participant viewed a series of indoor and outdoor scene images while carrying out  
194 a search, memorization, or rating task. For the memorization task, participants were  
195 instructed to memorize the image for a forced choice recognition test. At the end of each  
196 trial, the participants were prompted to indicate which of two images was just  
197 analyzed but were included in the experiment design to encourage searching behavior on  
other Search trials. Trials containing the target were excluded because search behavior was

193 likely to stop if the target was found padding considerable phoserto. They movement data.  
194 For consistency between trial types, participants were prompted to indicate if they found a  
195 “Z” or “N” at the end each Search trial. For the memorization task, participants were  
196 instructed to memorize the image for a forced choice recognition test. At the end of each to  
197 Memorize trial, the participants were prompted to indicate which of two images was just  
198 presented. Before the rating task, participants were asked to think about how they would rate  
199 the image on a scale from 1e (very unpleasant) to 7c (very pleasant). At the end of the trial,  
200 the participants were prompted to provide a rating dimm “Z”mediately N” after viewing of the image.  
201 The same materials were used in both experiments with a minor variation in the procedures.  
202 In the Confirmatory experiment, participants were directed as to where search targets might  
203 appear in the image (e.g., on flat surfaces). No such instructions were provided in the  
204 procedures. In the Confirmatory experiment, participants were directed as to where search  
205 Exploratory experiment.  
206 targets might appear in the image (e.g., on flat surfaces). No such instructions were provided  
207 in the Exploratory experiment.  
208 In both experiments, participants completed three mixed or uniform blocks of 40 trials  
209 ( $n = 120$  trials). When the blocks were mixed, the trial types were randomly intermixed.  
210 In both experiments, participants completed three mixed or uniform blocks of 40 trials  
211 within the block. For uniform blocks, each block consisted entirely of one of the three  
212 ( $n = 120$  trials). Block type was assigned in counterbalanced order. When the blocks were  
213 conditions (Search, Memorize, Rate) presented in random order. Each stimulus image was  
214 mixed, the trial types were randomly intermixed within the block. For uniform blocks, each  
215 presented for 8 seconds. The pictures were presented in color, with a size of 1024 x 768  
216 block consisted entirely of one of the three conditions (Search, Memorize, Rate), with block  
217 pixels, subtending a visual angle of 23.8° x 18.0°.  
218 types presented in random order. Each stimulus image was presented for 8 seconds. The  
219 pictures were presented in color, with a size of 1024 x 768 pixels, subtending a visual angle of  
220 Eye movements were recorded using an SR Research EyeLink 1000 eye tracker with a  
221 23.8° x 18.0°  
222 sampling rate of 1000Hz. Only the right eye was recorded. The system was calibrated using  
223 a nine-point accuracy and validity test. Errors greater than 1° or averaging greater than 0.5°  
224 Eye movements were recorded using an SR Research EyeLink 1000 eye tracker with a  
225 in total were re-calibrated.  
226 sampling rate of 1000Hz. Only the right eye was recorded. The system was calibrated using  
227 a nine-point accuracy and validity test. Errors greater than 1° or averaging greater than 0.5°  
228 **Datasets**  
229 in total were re-calibrated.

221 On some trials, a probe was presented on the screen six seconds after the onset of the  
222 trial. To avoid confounds resulting from the probe, only the first six seconds of the data for

223 Each trial was analyzed. Trials that contained fewer than 6000 samples within the first six  
 224 seconds of the trial were excluded before analysis. For both datasets, the trials were pooled  
 225 On some trials, a probe was presented on the screen six seconds after the onset of the  
 226 across participants. After excluding trials, the Exploratory dataset consisted of 12,177 of the  
 227 total trials, and the Confirmatory dataset consisted of 9,301 of the 10,395 total trials.  
 228 each trial was analyzed. Trials that contained fewer than 6000 samples within the first six  
 229 seconds of the trial were excluded before analysis.

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

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

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

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

representing the x- and y- coordinates, and pupil size for each data point collected in the first six seconds of each trial. To systematically assess the predictive value of each XYP (i.e., x-coordinates, y-coordinates, pupil size) component of the data, the timeline and image datasets were batched into subsets that excluded one of the components (i.e., XYØ, XØP, ØYP), or contained only one of the components (i.e., XØØ, ØYØ, ØØP). For the timeline

*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.

The data were replaced with zeros because removing the columns would change the structure of the data. The same systematic batching process was carried out for the image datasets. See Figure 2 for an example of each of these image data subsets. The data were replaced with zeros. The data were replaced with zeros because removing the columns would change the structure of the data. The same systematic batching process was carried out for the image dataset. See Figure 2 for an example of each of these image data subsets.

ØYP

XØP

XYØ

**Classification**

Deep CNN model architectures were implemented to classify the trials into Search, Memorize, or Rate categories. Because CNNs act as a digital filter sensitive to the number of features in the data, the differences in the structure of the timeline and image data formats

necessitated separate CNN model architectures. The model architectures were developed with the intent of establishing a generalizable approach to classifying cognitive processes from eye movement data.

XØØ

ØYØ

ØØP

The development of these models was not guided by any formal theoretical assumptions regarding the patterns or features likely to be extracted by the classifier. Like many HCI models, the development of these models followed general intuitions concerned

with building a model architecture capable of transforming the data inputs into an interpretable feature set that would not overfit the dataset. The models were developed

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.

## 251 Classification

XoP

XYo

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

oYø

øoP

258 The development of these models was not guided by any formal theoretical assumptions  
 259 regarding the patterns or features likely to be extracted by the classifier. Like many HCI  
 260 models, the development of these models followed general intuitions concerned with building  
 261 a model architecture capable of transforming the data inputs into an interpretable feature

262 set that would not overfit the dataset. The models were developed using version 0.3b of the  
 263 *Figure 1: PRO images were used to represent each type of data subset. As with the trials in the AVE dataset, the time*  
*span of each trial was represented by the shading of the dot. The first sample of each trial was white, and the last*  
*sample was black.*

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

275 The models were developed and tested on the Exploratory dataset. Model  
 276 hyperparameters were adjusted until the classification accuracies appeared to peak. The  
 277 model architecture with the highest classification accuracy on the Exploratory dataset was  
 278 trained, validated, and tested independently on the Confirmatory dataset. This means that  
 279 the model that was used to analyze the Confirmatory dataset was not trained on the  
 280 Exploratory dataset. The model architectures used for the timeline and plot image datasets

277 Are shown in Figure 3. The model architectures used for the timeline and plot image datasets  
 284 are shown in Figure 3.

## 278 Analysis

### 285 Analysis

279 Results for the CNN architecture that resulted in the highest accuracy on the

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

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

## 297 Results

### 304 Timeline Data Classification

### Results

#### 305 Timeline Data Classification

299 **Exploratory.** Classification accuracies for the XYP timeline dataset were well above

306 chance ( $M = .526$ ,  $SD = .018$ ;  $t(9) = 34.565$ ,  $p < .001$ ). Accuracies for all  
 307 classifications of the 33 batched data subsets were all better than chance (See Figure 4). As

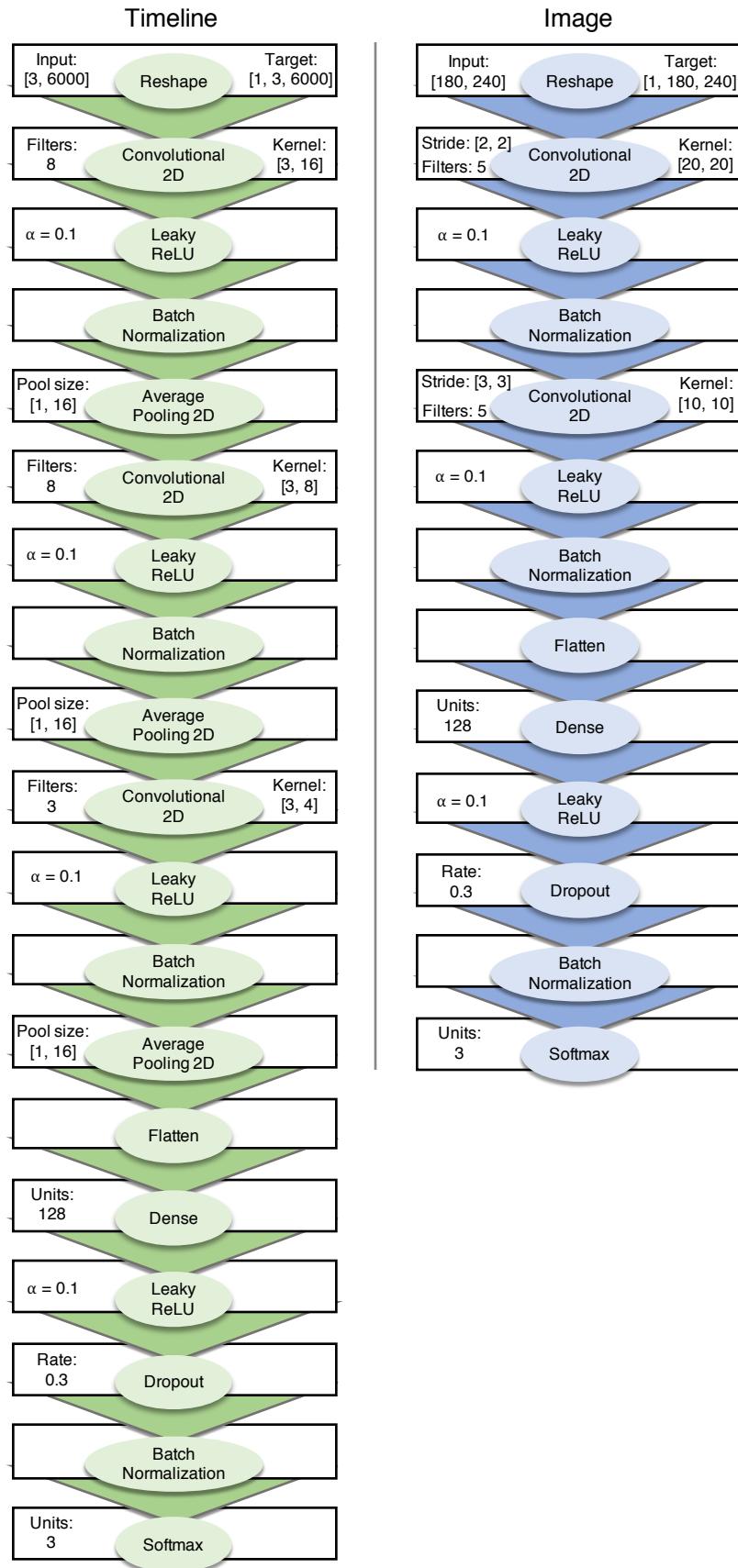


Figure 3. Two different model architectures were used to classify the timeline and image data. Both models were compiled using a categorical crossentropy loss function, and optimized with the Adam algorithm.

302 shown in the confusion matrices displayed in Figure 5, the data subsets (switched) overall  
 303 classification accuracies almost always classified the Memorize condition at or below chance  
 304 levels of accuracy. Misclassifications of the Memorize condition were split relatively evenly  
 305 between the Search and Rate conditions. The Memorize condition were split relatively evenly  
 312 between the Search and Rate conditions.

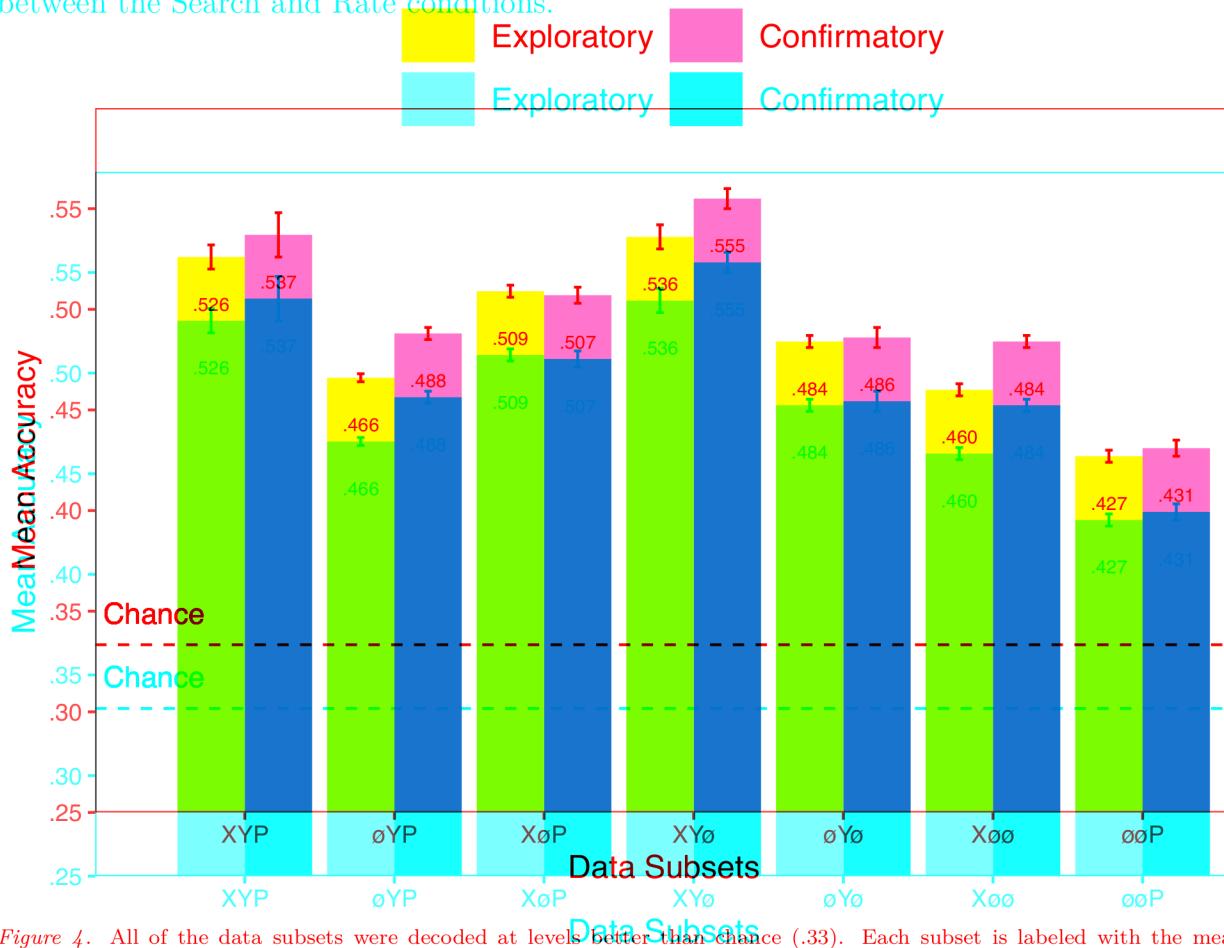


Figure 4. All of the data subsets were decoded at levels better than chance (.33). Each subset is labeled with the mean accuracy. The error bars represent standard errors.

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.

306 There was a difference in classification accuracy for the XYP dataset and the subsets  
 307 that had the pupil size,  $x$ -coordinate, and  $y$ -coordinate data systematically removed (Esets)  
 308 (t = 4.71471,  $p < 0.001$ ,  $\text{size} = \text{x}0.798$ ). In Post-hoc comparisons against the XYP dataset showed that  
 309 classification accuracies were not affected by the removal of pupil size or  $x$ -coordinated data  
 310 (see Table 2). A thermal effect present when pupil size was removed suggests that the pupil  
 311 size data were not contributing unique information that was not otherwise provided by the  $x$ -  
 312 and  $y$ -coordinates. A strict significance threshold of  $t = w0.05$  implies this same conclusion for

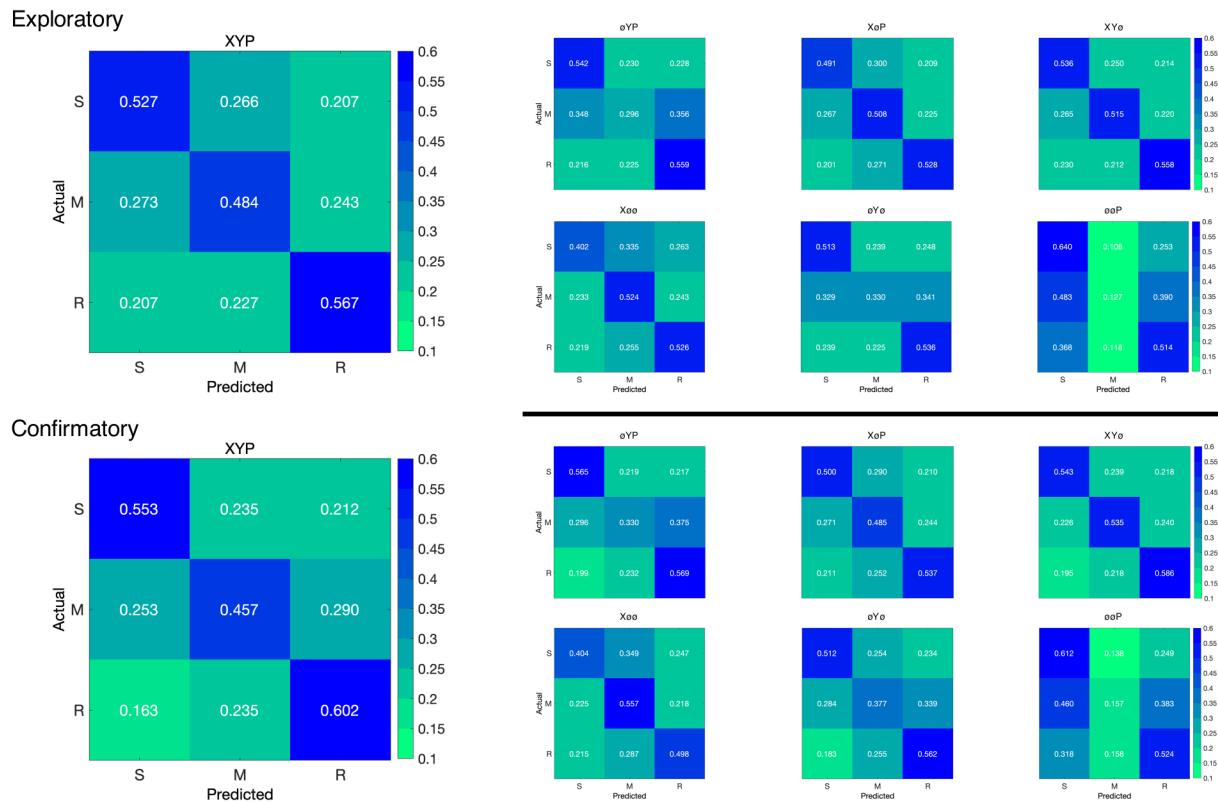


Figure 5. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents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

313 the  $y$ -coordinate data, but the relatively low degrees of freedom ( $df=18$ ) and the borderline  
 314 observed  $p$  value ( $p=1.056$ ) afford the possibility that there exists a small effect. However,  
 315 classification for the  $\emptyset Y\emptyset$  subset was significantly lower than the XYP dataset, showing that  
 316 the  $x$ -coordinate data were uniquely informative to the classification. XYP dataset, showing that

323 Table 2 coordinate data were uniquely informative to the classification.

#### Timeline Subset Comparisons

Table 2

#### Timeline Subset Comparisons

Comparison	Exploratory		Confirmatory	
	t	p	t	p
XYP vs. $\emptyset Y\emptyset$	9.420	< .001	5.210	< .001
XYP vs. XoP	2.645	.056	3.165	.016
XYP vs. XYo	9.339	< .001	5.319	< .001
XYP vs. XoO	2.645	.056	3.165	.016
XoP vs. $\emptyset Y\emptyset$	5.187	< .001	0.495	.874
XoP vs. XYo	1.635	.372	1.805	.288
XoP vs. XoO	12.213	< .001	10.178	< .001
$\emptyset Y\emptyset$ vs. $\emptyset Y\emptyset$	7.026	< .001	9.683	< .001
XoO vs. $\emptyset Y\emptyset$	12.213	< .001	10.178	< .001
$\emptyset Y\emptyset$ vs. $\emptyset Y\emptyset$	7.026	< .001	9.683	< .001

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

318 subsets $F_{(2,27)} = 75.145, p < .001, \eta^2 = 0.848$ ) in Post-hoc comparisons. No significant differences were found between the XOP and YOP subsets.

319 Classification accuracy was 75.14% for the XOP subset, which was lower than the XOP and YOP subsets.

320 Classification accuracy for the XOP subset was higher than the YOP subset. Altogether,

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

322 decisions, while the x-coordinates data was the most uniquely informative to classification

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

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

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

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

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

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

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

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

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

332 and Rate tasks (see Figure 5).

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

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

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

336 XYP Exploratory and Confirmatory timeline datasets were normally distributed, but Levene's test indicated that

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

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

339 0.406. These findings indicate that the learning model decoded the Exploratory and

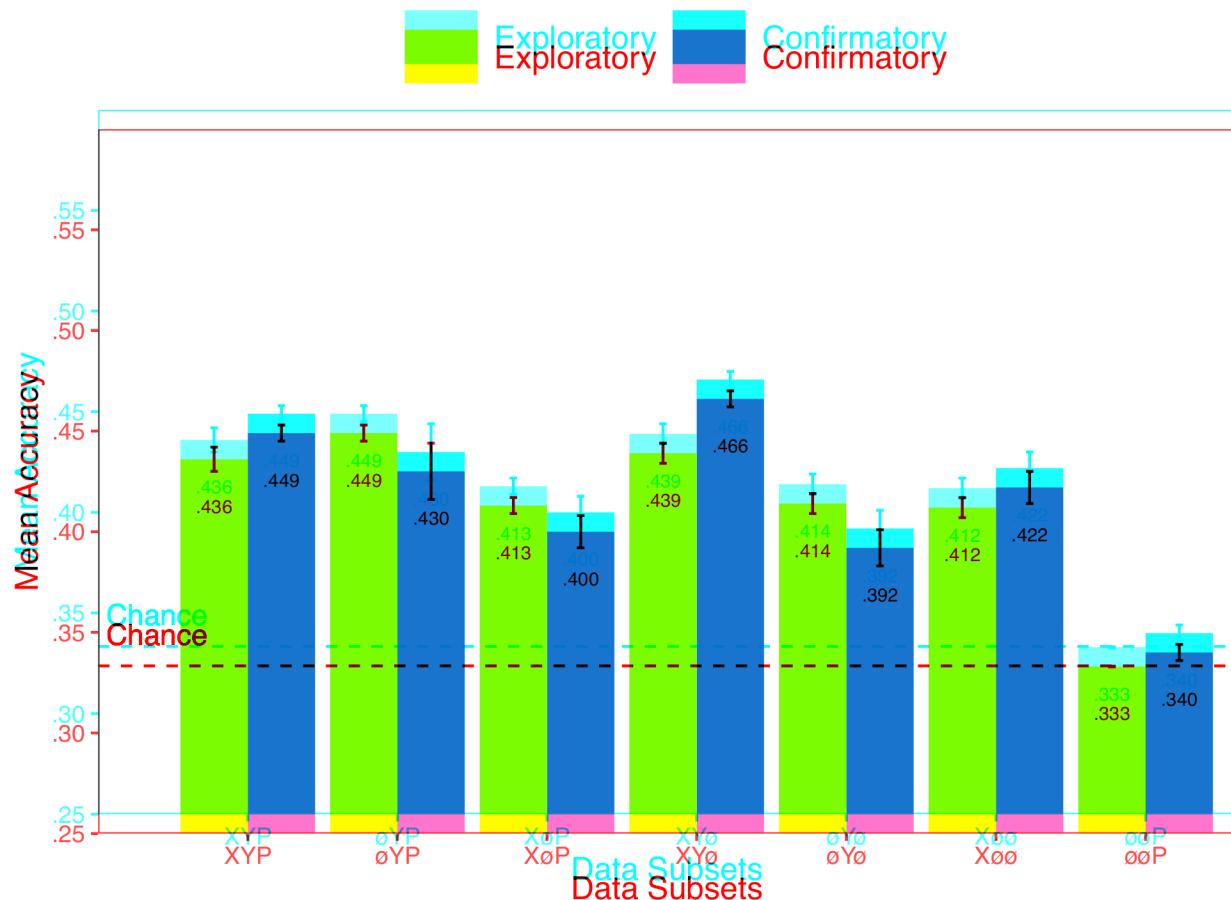
340 Confirmatory timeline datasets equally well, but the Confirmatory dataset classifications

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

342 deviation).

<sup>342</sup> **Plot Image Classification**

<sup>343</sup> **Exploratory.** Classification accuracies for the XYP plot image data were better  
<sup>344</sup> than chance ( $M \equiv .436$ ,  $SD \equiv .020$ ,  $p \leq .001$ ), but were less accurate than the classifications  
<sup>345</sup> for the XYP Exploratory timeline data ( $t_{(18)} \equiv 10.813$ ,  $p \leq .001$ ). Accuracies for the  
<sup>346</sup> classifications for all subsets of the plot image data except the ØØP subset were better than  
<sup>347</sup> chance (see Figure 6). Following the pattern expressed by the timeline dataset, the confusion  
<sup>348</sup> matrices showed that the Memorize condition was misclassified more often than the other  
<sup>349</sup> conditions, and appeared to be equally mis-identified as a Search or Rate condition (see  
<sup>350</sup> Figure 7).



*Figure 6.* All of the data subsets except for the Exploratory ØØP dataset were decoded at levels better than chance (.33). Each subset is labeled with the mean accuracy. The error bars represent standard errors.

<sup>351</sup> There was a difference in classification accuracy between the XYP dataset and the data  
<sup>351</sup> There was a difference in classification accuracy between the XYP dataset and the data

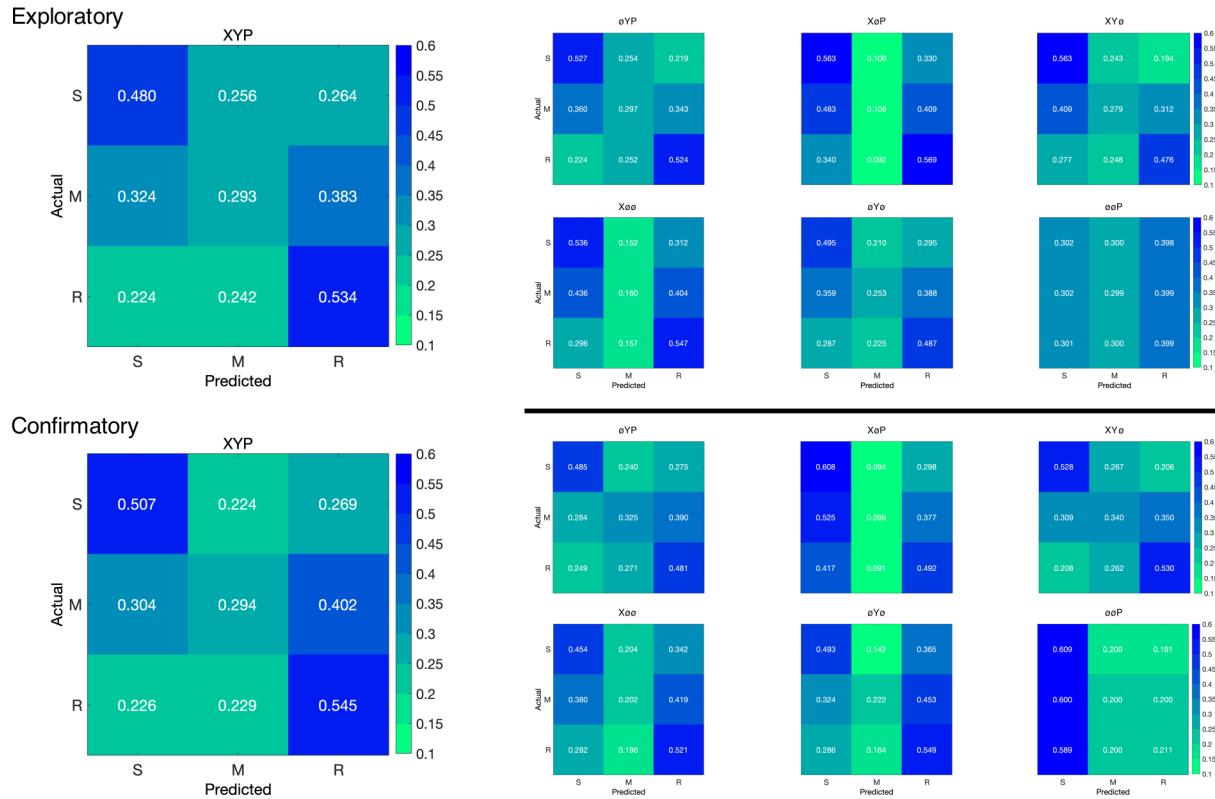


Figure 7. The confusion matrices represent the average classification accuracies for each condition of the image data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents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

352 subsets ( $F_{(4,45)} = 7.093, p < .001, \eta^2 = .387$ ). Post-hoc comparisons showed that compared  
 353 to the XYP dataset, there was no effect of removing pupil size or the x-coordinates, but  
 354 classification accuracy was worse when the y-coordinates were removed (see Table 3).

Table 3  
*Image Subset Comparisons*

Comparison	Exploratory		Confirmatory	
	t	p	t	p
XYP vs. ØYP	1.792	.391	1.623	.491
XYP vs. XØP	2.939	.039	4.375	< .001
XYP vs. XYØ	0.474	.989	1.557	.532
XØØ vs. ØYØ	0.423	.906	2.807	.204
XØØ vs. ØØP	13.569	< .001	5.070	< .001
ØYØ vs. ØØP	13.235	< .001	7.877	< .001

355 There was also a difference in classification accuracies between the XØØ, ØYØ, and  
 356 ØØP subsets (Levene's test:  $F_{(2,27)} = 3.815, p = .035$ ; Welch correction for lack of

<sup>357</sup> homogeneity of variances:  $F_{(2,17.993)} = 228.137, p < .001, \eta^2 = .899$ ). Post-hoc comparisons  
<sup>358</sup> showed that there was no difference in classification accuracies for the XØØ and ØYØ  
<sup>359</sup> subsets, but classification for the ØØP subset were less accurate than the XØØ and ØYØ  
<sup>360</sup> subsets.

<sup>361</sup> **Confirmatory.** Classification accuracies for the XYP confirmatory image dataset  
<sup>362</sup> were well above chance ( $M = .449, SD = 0.012, t_{(9)} = 31.061, p < .001$ ), but were less  
<sup>363</sup> accurate than the classifications of the confirmatory timeline dataset ( $t_{(18)} = 11.167, p <$   
<sup>364</sup> .001). Accuracies for classifications of the data subsets were also all better than chance (see  
<sup>365</sup> Figure 6). The confusion matrices followed the pattern showing that the Memorize condition  
<sup>366</sup> was confused most often, and was relatively evenly misidentified as a Search or Detect trial  
<sup>367</sup> (see Figure 7). As with the timeline data, the general trend showing that pupil size data was  
<sup>368</sup> the least informative to the model was replicated in the Confirmatory dataset (see Table 3).

<sup>369</sup> To test the generalizability of the model architecture, the classification accuracies for  
<sup>370</sup> the XYP Exploratory and Confirmatory plot image datasets were compared. The  
<sup>371</sup> independent samples *t*-test comparing the classification accuracies for the Exploratory and  
<sup>372</sup> Confirmatory plot image datasets did not show a significant difference,  $t_{(18)} = 1.777, p =$   
<sup>373</sup> .092, Cohen's *d* = 0.795.

## <sup>381</sup> Discussion

<sup>382</sup> The present study aimed to produce a practical and reliable example of a black box  
<sup>383</sup> solution to the inverse Yarbus problem. To implement this solution, we classified raw  
<sup>384</sup> timeline and minimally processed plot image data using a CNN model architecture. To our  
<sup>385</sup> knowledge, this study was the first to provide a solution to determining mental state from  
<sup>386</sup> eye movement data using each of the following: (1) Non-aggregated eye tracking data (i.e.,  
<sup>387</sup> raw x-coordinates, y-coordinates, pupil size); (2) timeline and image data formats (see  
<sup>388</sup> Figure 2); and (3) a black box CNN architecture. This study probed the relative predictive  
<sup>389</sup> value of the x-coordinate, y-coordinate, and pupil size components of the eye movement data

383 using a CNN. The CNN was able to decode the timeline and plot image data better than  
384 chance, although only the timeline datasets were decoded with accuracies comparable to  
385 other state-of-the-art approaches. Datasets with lower classification accuracies were not able  
386 to differentiate the cognitive processes underlying the Memorize task from the cognitive  
387 processes underlying the Search and Rate tasks. Decoding subsets of the data revealed that  
388 pupil size was the least uniquely informative component of the eye movement data. This  
389 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

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

409 timeline and plot image formats, there was good replicability between the Exploratory and  
410 Confirmatory datasets, demonstrating that these architectures performed similarly from  
411 experiment to experiment. An appropriately tuned CNN should be capable of learning any  
412 arbitrary function, but given that the upper bound for decodability of these datasets is  
413 unknown, there is the possibility that a model architecture exists that is capable of  
414 classifying the plot image data format more accurately than the model used to classify the  
415 timeline data. Despite this possibility, the convergence of these findings with other studies  
416 (see Table 1) suggests that the results of this study are approaching a ceiling for the  
417 potential to solve the inverse Yarbus problem with eye movement data. Although the true  
418 capacity to predict mental state from eye movement data is unknown, standardizing datasets  
419 in the future could provide a point for comparison that can more effectively indicate which  
420 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  
To more clearly understand how the different tasks influenced the decodability of the  
eye movement data, additional analyses were conducted on the Exploratory and  
eye movement data, additional analyses were conducted on the Exploratory and  
Confirmatory timeline datasets (see Appendix). For the main supplementary analysis, the  
Confirmatory timeline datasets (see Appendix). For the main supplementary analysis, the  
data subsets were re-submitted to the CNN and re-classified as 2-category task sets. In  
data subsets were submitted to the model in 2-category task sets. In addition to the

435 supplementary analysis, the results from the primary analysis were pre-calculated from mere  
436 3-category task sets to 2-category task sets. These analyses showed a tendency for the model  
437 to mis-classify the Search and Rate trials as Memorize. In the primary analyses, one of the  
438 Memorized condition was predicted with the lowest accuracy, but misclassifications of the of  
439 Search and Rate trials were most often categorized as Memorizer. As a whole, this pattern of  
440 results indicated a general bias for certain trials to be categorized as Memorize. Overall,  
441 the findings from this supplemental analysis show that conclusions drawn from comparisons  
442 between approaches that do not use the same task sets for the same number of tasks could  
443 be potentially uninterpretable because the features underlying the task categories are  
444 interpreted differently by the neural network algorithm. same number of tasks, could be  
445 potentially uninterpretable because the features underlying the task categories are  
446 When determining the relative contributions of the eye movement features used in  
447 this study (x-coordinates, y-coordinates, pupil size), the pupil size data was consistently the  
448 least uniquely informative. When pupil size was removed from the Exploratory and  
449 Confirmatory timeline and plot image datasets, classification accuracy remained stable (vs.  
450 XYR uniquely informative). Furthermore, classification accuracy of the OGP subset was the lowest of all  
451 of the data subsets, and in one instance, was no better than chance. Although these findings  
452 indicate that, in this case, pupil size was a relatively uninformative component of the eye  
453 movement data, previous research has associated changes in pupil size as indicators of  
454 working memory load (Kahneman & Beatty, 1966; Karatekin, Couperus, & Marcus, 2004),  
455 arousal (Wang et al., 2018), and cognitive effort (Porter, Truscott, & Giedd, 2007). The  
456 results of the current study indicate that the changes in pupil size associated with these  
457 underlying processes were not useful in delineating the tasks being classified (i.e., Search,  
458 Memorize, Rate), potentially because these tasks did not evoke a reliable pattern of changes  
459 in pupil size. Additionally, properties of the stimuli known to influence pupil size, such as  
460 Memorize and contrast, were not controlled in these datasets. Given that stimuli were  
461 randomly assigned, there is the potential that uncontrolled stimulus properties known to  
462 affect pupil size made it difficult for the CNN's capacity to detect patterns in the pupil size data.

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. XYR uniquely informative). Furthermore, classification accuracy of the OGP 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, Truscott, & Giedd, 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 Memorize and contrast, were not controlled in these datasets. Given that stimuli were randomly assigned, there is the potential that uncontrolled stimulus properties known to affect pupil size made it difficult for the CNN's capacity to detect patterns in the pupil size data.

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

496 task-diagnostic features in the data and the level of distinctiveness required to classify the  
497 tasks, future research could more definitively conceptualize the cognitive processes  
498 underlying the task-at-hand.

499 Classifying mental state from eye movement data is often carried out in an effort to  
500 advance technology to improve educational outcomes, strengthen the independence of  
501 physically and mentally handicapped individuals, or improve HCI's (Koochaki &  
502 Classifying mental state from eye movement data is often carried out in an effort to  
503 Najafizadeh, 2018). Given the previous questions raised regarding the reliability and  
504 advance technology to improve educational outcomes, strengthen the independence of  
505 generalizability of black-box CNN classification, the current study first tested models on an  
506 physically and mentally handicapped individuals, or improve HCI's (Koochaki &  
507 exploratory dataset, then confirmed the outcome using a second independent dataset.  
508 Najafizadeh, 2018). Given the previous questions raised regarding the reliability and  
509 Overall, the findings of this study indicate that this black-box approach is capable of  
510 generalizability of black-box CNN classification, the current study first tested models on an  
511 producing a stable and generalizable outcome. Additionally, the supplementary analyses  
512 exploratory dataset, then confirmed the outcome using a second independent dataset.  
513 showed that different task sets, or a different number of tasks, could lead the algorithm to

489 interpret features differently, which should be taken into account when comparing task  
490 classification approaches. Future studies that incorporate features from the stimulus might  
491 have the potential to surpass current state-of-the-art classification. According to Bulling,  
492 Weichelt, and Gellersen (2013), incorporating stimulus features into the dataset  
493 may provide improved accuracy relative to decoding gaze location data and pupil size. might  
494 Alternatively, Borji and Itti (2014) suggested that accounting for salient features in the  
495 stimuli might leave little room for theoretically defined classifiers to consider mental  
496 state. Future research should examine the potential for the inclusion of stimulus features in  
497 addition to the eye movement data to boost black-box CNN classification  
498 accuracy of image data beyond that of timeline data. In addition, it is important to consider mental state. Future  
500 research should examine the potential for the inclusion of stimulus feature information in  
501 addition to the eye movement data to boost black-box CNN classification accuracy of image  
502 data beyond that of timeline data.

## References

- 499 Boisvert, J. F. G., & Bruce, N. D. B. (2016). Predicting task from eye movements: On the  
 500 importance of spatial distribution, dynamics, and image features. *Neurocomputing*,  
 501 207, 653–668. <https://doi.org/10.1016/j.neucom.2016.05.047>
- 502  
 503 Borji, A.; & Itti, L. (2014). Defending Yarbus: Eye movements reveal observers' task.  
 504 *Journal of Vision*, 14(3), 1–21. <https://doi.org/10.1167/14.3.29>
- 505 Bulling, A., Weichel, C., & Gellersen, H. (2013). EyeContext: Recognition of high-level  
 506 Bulling, A., Weichel, C., & Gellersen, H. (2013). EyeContext: Recognition of high-level  
 507 contextual cues from human visual behaviour. In *Proceedings of the SIGCHI  
 Conference on Human Factors in Computing Systems - CHI '13* (p. 305). Paris,  
 508 France: ACM Press. <https://doi.org/10.1145/2470654.2470697>  
 France: ACM Press. <https://doi.org/10.1145/2470654.2470697>
- 509 Castelhano, M. S., Mack, M. L., & Henderson, J. M. (2009). Viewing task influences eye  
 510 movement control during active scene perception. *Journal of Vision*, 9(3), 1–15.  
 511 movement control during active scene perception. *Journal of Vision*, 9(3), 6–6.  
 512 <https://doi.org/10.1167/9.3.6>  
 513 <https://doi.org/10.1167/9.3.6>
- 514 Coco, M. I., & Keller, F. (2014). Classification of visual and linguistic tasks using  
 515 Coco, M. I., & Keller, F. (2014). Classification of visual and linguistic tasks using  
 516 eye-movement features. *Journal of Vision*, 14(3), 1–18.  
 517 eye-movement features. *Journal of Vision*, 14(3), 11–11.  
 518 <https://doi.org/10.1167/14.3.11>  
 519 <https://doi.org/10.1167/14.3.11>
- 520 DeAngelus, M., & Pelz, J. B. (2009). Top-down control of eye movements: Yarbus revisited.  
 521 DeAngelus, M., & Pelz, J. B. (2009). Top-down control of eye movements: Yarbus revisited.  
 522 *Visual Cognition*, 17(6-7), 790–811. <https://doi.org/10.1080/13506280902793843>  
 523 *Visual Cognition*, 17(6-7), 790–811. <https://doi.org/10.1080/13506280902793843>
- 524 Greene, M. R., Liu, T., & Wolfe, J. M. (2012). Reconsidering Yarbus: A failure to predict  
 525 Greene, M. R., Liu, T., & Wolfe, J. M. (2012). Reconsidering Yarbus: A failure to predict  
 526 observers' task from eye movement patterns. *Vision Research*, 62, 1–8.  
 527 observers' task from eye movement patterns. *Vision Res*, 62, 1–8.  
 528 <https://doi.org/10.1016/j.visres.2012.03.019>  
 529 <https://doi.org/10.1016/j.visres.2012.03.019>
- 530 Haji-Abolhassani, A., & Clark, J. J. (2014). An inverse Yarbus process: Predicting observers'  
 531 Haji-Abolhassani, A., & Clark, J. J. (2014). An inverse Yarbus process: Predicting observers'  
 532 task from eye movement patterns. *Vision Research*, 103, 127–142.  
 533 task from eye movement patterns. *Vision Research*, 103, 127–142.

- 522 https://doi.org/10.1016/j.visres.2014.08.014
- 523 Henderson, J. M., Shinkareva, S. V., Wang, J., Luke, S. G., & Olejarczyk, J. (2013).  
524 Predicting Cognitive State from Eye Movements. *PLoS ONE*, 8(5), e64937.  
525 https://doi.org/10.1371/journal.pone.0064937
- 526 Kahneman, D., & Beatty, J. (1966). Pupil Diameter and Load on Memory. *Science*,  
527 Kahneman, D., & Beatty, J. (1966). Pupil Diameter and Load on Memory. *Science*,  
528 154(3756), 1583–1585. Retrieved from <http://www.jstor.org/stable/1720478>  
529 154 (3756), 1583–1585. Retrieved from <https://www.jstor.org/stable/1720478>
- 530 Kanan, C., Ray, N. A., Bseiso, D. N. F., Hsiao, J. H., & Cottrell, G. W. (2014). Predicting  
531 Kanan, C., Ray, N. A., Bseiso, D. N. F., Hsiao, J. H., & Cottrell, G. W. (2014). Predicting  
532 an observer's task using multi-fixation pattern analysis. In *Proceedings of the*  
533 an observer's task using multi-fixation pattern analysis. In *Proceedings of the*  
534 *Symposium on Eye Tracking Research and Applications - ETRA '14* (pp. 287–290).  
535 *Symposium on Eye Tracking Research and Applications - ETRA '14* (pp. 287–290).  
536 Safety Harbor, Florida: ACM Press. <https://doi.org/10.1145/2578153.2578208>  
537 Safety Harbor, Florida: ACM Press. <https://doi.org/10.1145/2578153.2578208>
- 538 Karatekin, C., Couperus, J. W., & Marcus, D. J. (2004). Attention allocation in the  
539 Karatekin, C., Couperus, J. W., & Marcus, D. J. (2004). Attention allocation in the  
540 dual-task paradigm as measured through behavioral and psychophysiological  
541 dual-task paradigm as measured through behavioral and psychophysiological  
542 responses. *Psychophysiology*, 41(2), 175–185.  
543 responses. *Psychophysiology*, 41(2), 175–185.  
544 https://doi.org/10.1111/j.1469-8986.2004.00147.x  
545 https://doi.org/10.1111/j.1469-8986.2004.00147.x
- 546 Koochaki, F., & Najafizadeh, L. (2018). Predicting Intention Through Eye Gaze Patterns.  
547 Koochaki, F., & Najafizadeh, L. (2018). Predicting Intention Through Eye Gaze Patterns.  
548 In *2018 IEEE Biomedical Circuits and Systems Conference (BioCAS)* (pp. 1–4).  
549 In *2018 IEEE Biomedical Circuits and Systems Conference (BioCAS)* (pp. 1–4).  
550 https://doi.org/10.1109/BIOCAS.2018.8584665
- 551 Król, M. E., & Król, M. (2018). The right look for the job: Decoding cognitive processes  
552 involved in the task from spatial eye movement patterns. *Psychological Research*, 84,  
553 involved in the task from spatial eye movement patterns. *Psychological Research*.  
554 https://doi.org/10.1007/s00426-018-0996-5
- 555 Lukander, K., Toivanen, M., & Puolamäki, K. (2017). Inferring Intent and Action from Gaze  
556 in Naturalistic Behavior: A Review. *International Journal of Mobile Human-Computer Interaction*, 9(4), 41–57. <https://doi.org/10.4018/IJMHCI.2017100104>
- 557 Lukander, K., Toivanen, M., & Puolamäki, K. (2017). Inferring Intent and Action from Gaze  
558 in Naturalistic Behavior: A Review. *International Journal of Mobile Human-Computer Interaction*, 9(4), 41–57. <https://doi.org/10.4018/IJMHCI.2017100104>
- 559 Lukander, K., Toivanen, M., & Puolamäki, K. (2017). Inferring Intent and Action from Gaze  
560 in Naturalistic Behavior: A Review. *International Journal of Mobile Human-Computer Interaction*, 9(4), 41–57. <https://doi.org/10.4018/IJMHCI.2017100104>

- 545 MacInnes, W., Joseph, Hunt, A. R., Clarke, A. D. F., & Dodd, M. D. (2018). A Generative  
546 Model of Cognitive State from Task and Eye Movements. *Cognitive Computation*,  
547 10(5), 703–717. <https://doi.org/10.1007/s12559-018-9558-9>
- 548 Mills, M., Hollingworth, A., Van der Stigchel, S., Hoffman, L., & Dodd, M. D. (2011).  
549 Examining the influence of task set on eye movements and fixations. *Journal of*  
550 *Vision*, 11(8), 17417. <https://doi.org/10.1167/118.17>
- 551 Porter, G., Troscianko, T., & Gilchrist, I. D. (2007). Effort during visual search and  
552 counting: Insights from pupillometry. *Quarterly Journal of Experimental Psychology*  
553 (2006), 60(2), 211–229. <https://doi.org/10.1080/17470210600673818>
- 554 Seeliger, K., Fritzsche, M., Güçlü, Ü., Schoenmakers, S., Schoffelen, J.-M., Bosch, S. E., &  
555 van Gerven, M. A. J. (2018). Convolutional neural network-based encoding and  
556 decoding of visual object recognition in space and time. *NeuroImage*, 180, 253–266.  
557 <https://doi.org/10.1016/j.neuroimage.2017.07.018>
- 558 Tatler, B. W., Wade, N. J., Kwan, H., Findlay, J. M., & Velichkovsky, B. M. (2010). Yarbus,  
559 Eye Movements, and Vision. *I-Perception*, 1(1), 7–27. <https://doi.org/10.1068/i0382>
- 560 Wang, C.-A., Baird, T., Huang, J., Coutinho, J. D., Brien, D. C., & Munoz, D. P. (2018).  
561 Arousal Effects on Pupil Size, Heart Rate, and Skin Conductance in an Emotional  
562 Face Task. *Frontiers in Neurology*, 9. 1029. <https://doi.org/10.3389/fnme.2018.01029>
- 563 Yarbus, A. (1967). *Eye Movements and Vision*. New York, NY: Plenum Press, from  
564 [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)
- 565 Zhou, B., Bau, D., Oliva, A., & Torralba, A. (2019). Comparing the Interpretability of Deep  
566 Networks via Network Dissection. In W. Samek, G. Montavon, A. Vedaldi, L. K.  
567 Hansen, & K. R. Müller (Eds.), *Explaining, Interpreting, Explaining and  
568 Visualizing Deep Learning* (pp. 243–252). Cham: Springer International Publishing.  
569 [https://doi.org/10.1007/978-3-030-28954-6\\_12](https://doi.org/10.1007/978-3-030-28954-6_12): Interpreting, Explaining and

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

570        [https://doi.org/10.1007/978-3-030-28954-6\\_12](https://doi.org/10.1007/978-3-030-28954-6_12)  
 579        Additional analyses were conducted in an attempt to clarify the effect of task on

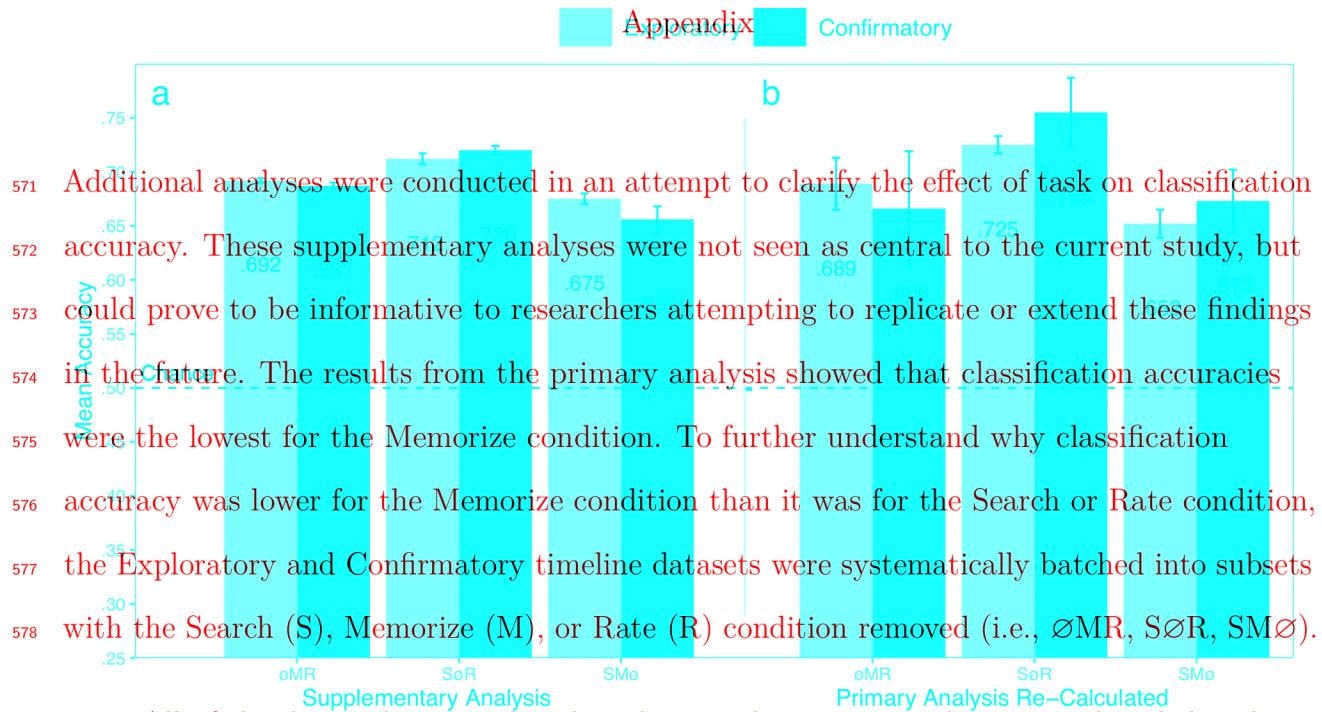
580        classification accuracy. These supplementary analyses were not seen as central to the current  
 581        study, but could prove to be informative to researchers attempting to replicate or extend  
 582        these findings in the future. The results from the primary analysis showed that 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	
	t	p	t	p
$\emptyset$ 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



All of the data subsets analyzed in this supplementary analysis were decoded with

*Figure 8.* The graph represents the average accuracy reported for each subset of the Exploratory and Confirmatory timeline better (i.e., supplementary analysis,  $\text{SoR}$ ,  $\text{oMR}$ ) than chance (i.e.,  $.50$ ). The results for all three conditions were decoded at levels better than chance ( $.50$ ). The error bars represent standard errors.

both the Exploratory and Confirmatory datasets. When the Memorize condition was removed, classification accuracy improved (see Table A1; see Figure A1a). When the Rate condition was removed, classification was the worst. When the Memorize condition was drawn from a comparison of the results of analyses could be misleading. For this reason, we included (i.e.,  $\text{SM}\emptyset$  and  $\text{oMR}$ ), mis-classifications were biased toward Memorize, and the Memorize condition was more accurately predicted than the Search and Rate conditions (see equivalent to a 50% chance threshold. Because the cross-validation scheme implemented by Figure A2).

*Table A1*

*Supplementary Subset Comparisons*

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  $\text{SoR}$  vs.  $\text{oMR}$  analysis (see Figure 3.248). The predictions were  $\text{oMR}$  vs.  $\text{SM}\emptyset$  3.248,  $t = .008$ ,  $p = .994$ ;  $\text{SoR}$  vs.  $\text{oMR}$  3.248,  $t = .008$ ,  $p = .994$ ;  $\text{oMR}$  vs.  $\text{SoR}$  3.248,  $t = .008$ ,  $p = .994$ ;  $\text{SoR}$  vs.  $\text{SM}\emptyset$  3.248,  $t = .008$ ,  $p = .994$ ;  $\text{SM}\emptyset$  vs.  $\text{oMR}$  3.248,  $t = .008$ ,  $p = .994$ ;  $\text{SM}\emptyset$  vs.  $\text{SoR}$  3.248,  $t = .008$ ,  $p = .994$ .

re-calculated using the following formula:  $\text{Prediction}_{(A,A,\text{ABC})} = \text{Prediction}_{(A,A,\text{ABC})} / (\text{Prediction}_{(A,A,\text{ABC})} + \text{Prediction}_{(A,C,\text{ABC})})$ . For example, accuracy for the Search

The accuracies for all of the data subsets observed in the supplementary analysis were calculated with the following formula:  $\text{Prediction}_{(\text{S},\text{S},\text{S}\emptyset\text{R})} = \text{Prediction}_{(\text{S},\text{S},\text{S}\emptyset\text{R})} / (\text{Prediction}_{(\text{S},\text{S},\text{S}\emptyset\text{R})} + \text{Prediction}_{(\text{S},\text{R},\text{S}\emptyset\text{R})})$ .

higher than the accuracies observed in the main analysis. Although there is a clear difference in accuracy, the primary analysis was classifying three categories (chance = .33) and the

inaccuracy, the primary analysis was classifying three categories (chance = .33) and the

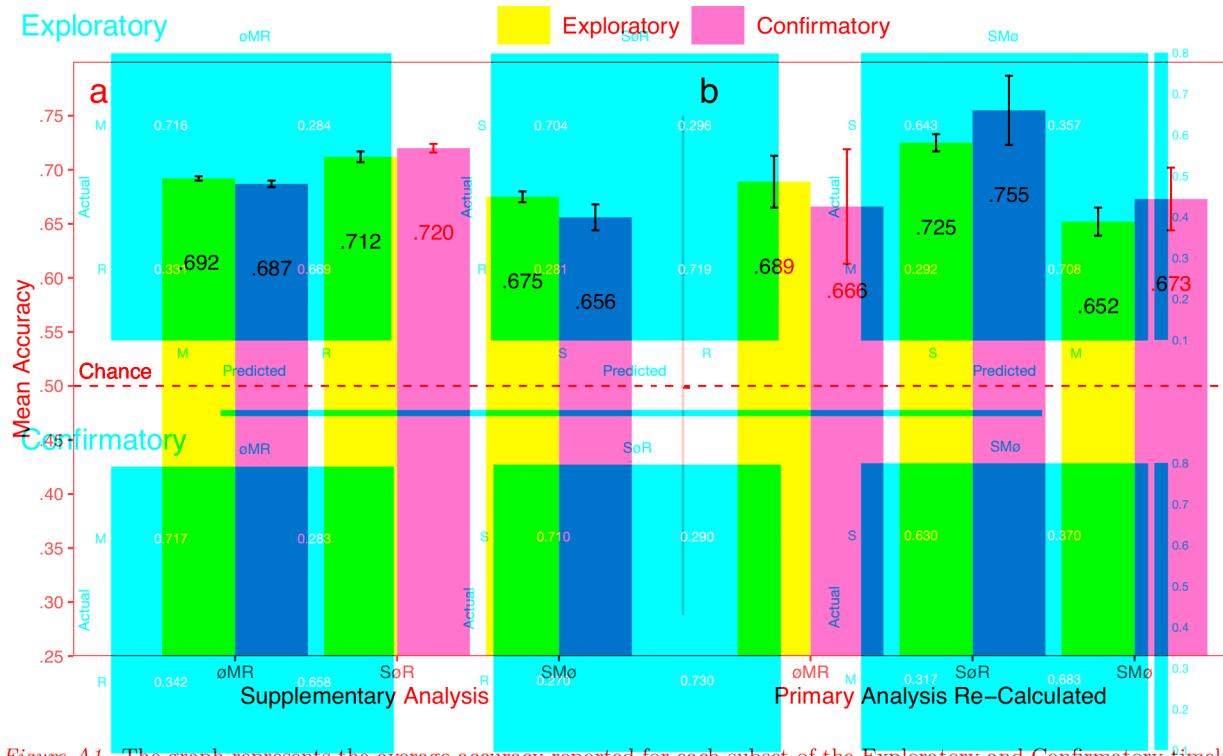


Figure A1. The graph represents the average accuracy reported for each subset of the Exploratory and Confirmatory timeline data for (a) the re-calculated accuracies from the primary analysis; and (b) the supplementary analysis. All of the data subsets were decoded at levels better than chance (.50).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

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.

supplementary analysis was classifying two categories (chance = .50). Because the baseline the ratio of Search trials that were misclassified as Rate. 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 classifications predicted the Memorize conditions with the lowest accuracy (c.f., Search and Rate conditions), and mis-classifications of the Search and Rate conditions were most often categorized as Memorize (see Figure 5). Because the Memorize condition was mis-classified an equal number of trials in the test set are assigned to each condition for each dataset, we more often than the other conditions in the primary analysis, the removal of the third class were able to re-calculate 2-category predictions from the 3-category predictions presented in in the re-calculated SMø and ØMR subsets resulted in a disproportionate amount of the confusion matrices from the primary analysis (see Figure 5). The predictions were mis-classified Memorize trials being removed from the data subset somewhat eliminating the re-calculated using the following formula:  $\text{Prediction}_{(A,A,A \ominus C)} = \text{Prediction}_{(A,A,ABC)} / (\text{Prediction}_{(A,A,ABC)} + \text{Prediction}_{(A,C,ABC)})$ . For example, accuracy for the Search the re-calculated SMø and ØMR subsets were classified less accurately than SøR. classification for SøR would be calculated with the following:  $\text{Prediction}_{(S,S,S \ominus R)} = \text{Prediction}_{(S,S,S \ominus R)} / (\text{Prediction}_{(S,S,S \ominus R)} + \text{Prediction}_{(S,R,S \ominus R)})$ , where  $\text{Prediction}_{(S,R,S \ominus R)}$  is

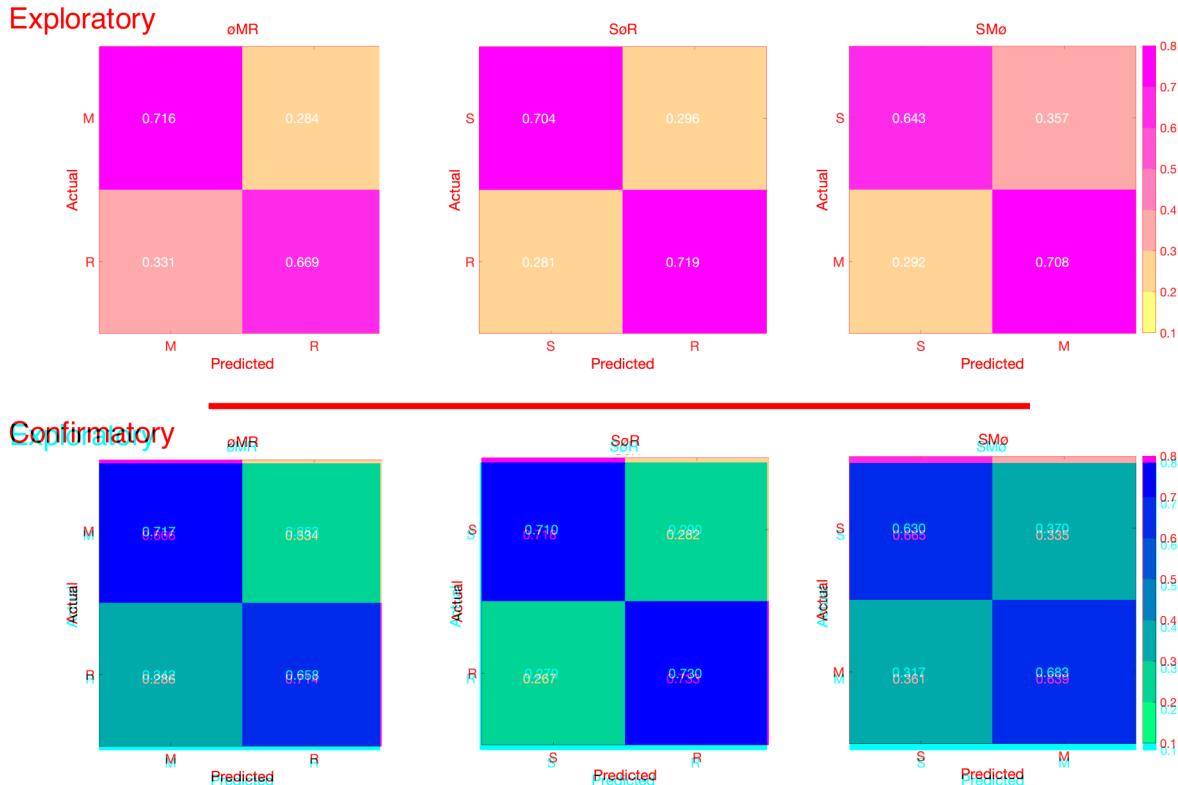


Figure A2. The confusion matrices represent the average classification accuracies for each condition of the timeline data (S = Search, M = Memorize, R = Rate). The vertical axis of the confusion matrices represents the actual condition for the trial. The horizontal axis of the confusion matrices represents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

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

604 The results for the re-calculated predictions followed a pattern similar to the

605 supplementary analysis (see Figure A1b). This is supported by the persisting tendency of

606 the algorithm to mis-classify Search and Rate trials in the SM $\varnothing$  and  $\varnothing$ MR subsets as

607 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).

608 classifications predicted the Memorize conditions with the lowest accuracy (c.f., Search and

609 Rate conditions), mis-classifications of the Search and Rate conditions were most often

610 categorized as Memorize (see Figure 5). This overall pattern points toward a general bias to

611 categorize uncertain trials as Memorize.

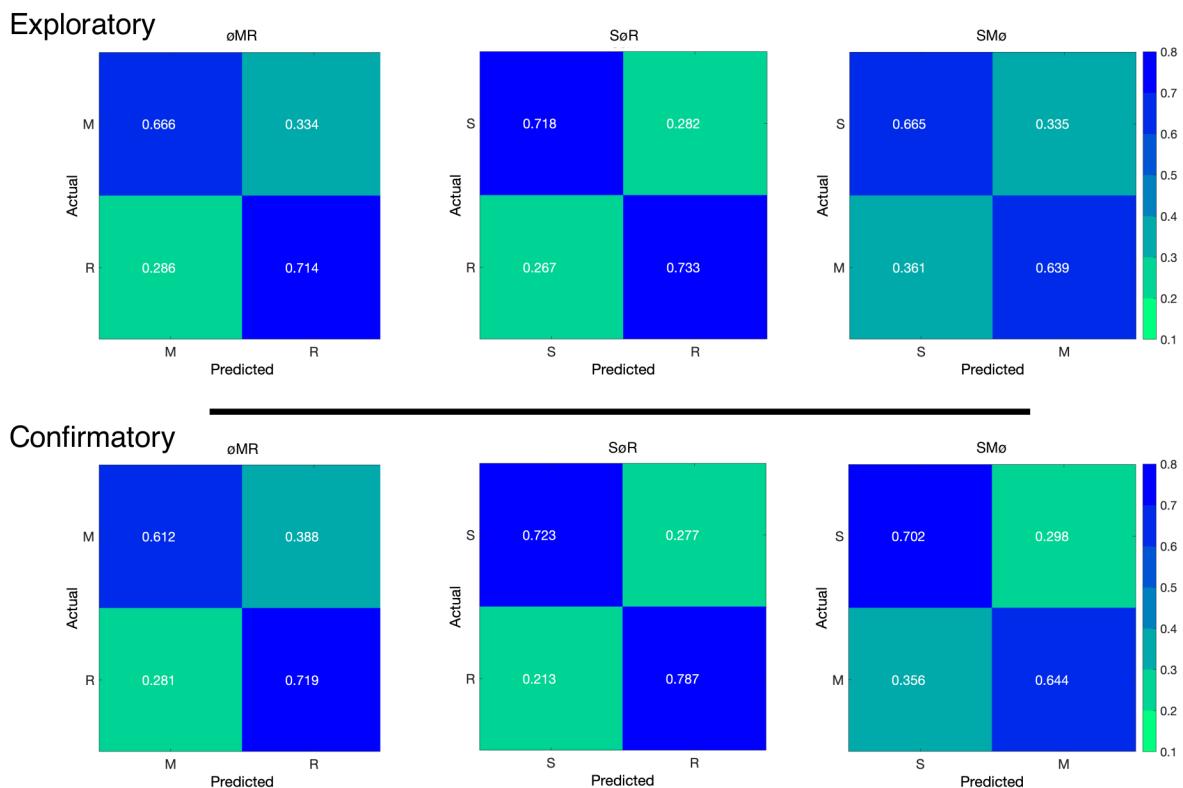


Figure A3. The confusion matrices represent a re-calculation of the classification accuracies for each category from the primary analysis. This re-calculation is meant to make the accuracies presented in the primary analysis (chance = .33) equivalent to the classification accuracies presented in the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