

Machine-Learned Computational Models can Enhance the Study of Text & Discourse: A Case
Study Using Eye Tracking to Model Reading Comprehension

Sidney K. D'Mello, Rosy Southwell, & Julie Gregg

University of Colorado Boulder

Author Note

Corresponding Author

Sidney D'Mello, 594 UCB, Boulder, CO 80309, USA

sidney.dmello@colorado.edu

This research was supported by the National Science Foundation (NSF) (DRL 1235958, IIS 1523091, and DRL 1920510). Any opinions, findings and conclusions, or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of NSF.

Abstract

We propose that machine-learned computational models (MLCM) – where the model parameters and perhaps even structure are learned from data – can complement extant approaches to the study of text and discourse. Such models are particularly useful when theoretical understanding is insufficient, when the data is rife with nonlinearities and interactivity, and when researchers aspire to take advantage of “big data”. Being fully-instantiated computer programs, MLCM can also be used for autonomous assessment and real-time intervention. We illustrate these ideas in the context of an eye-movement-based MLCM of text-base comprehension during reading a long connected text. Using a dataset where 104 participants read a 6500-word text, we trained Random Forests models to predict comprehension scores from six eye movement features. The models were highly accurate (AUROC = .902; $r = .661$), robust, and generalized across participants, suggesting possible use in future studies. We conclude by arguing for an increased role of MLCMs in the future of discourse research.

Keywords: reading comprehension, eye gaze, machine learning, computational model

Machine-Learned Computational Models can Enhance Discourse Analysis: A Case Study
Using Eye Tracking to Model Reading Comprehension

Camilla: *"You, sir, should unmask."*

Stranger: *"Indeed?"*

Cassilda: *"Indeed it's time. We have all laid aside disguise but you."*

Stranger: *"I wear no mask."*

Camilla: *(Terrified, aside to Cassilda.) "No mask? NO MASK! [emphasis added]*

(Chambers, 1985)

This short dialog from *"The King in Yellow and Other Horror Stories"* by Robert W. Chambers illustrates the power of discourse in the hands of a gifted writer. In just 33 words, Chambers presents a complex narrative involving three characters embedded in a rather macabre interaction with surprise and impending terror. His brevity teases our fascination, tempting us to imagine what came before this strange meeting and what horror might occur next. Indeed, the power of discourse is not in the words themselves but in what lies beneath.

Given the complexity of discourse, which increases by orders of magnitude when moving from written prose to spoken dialogs, and even further for multiparty conversations, it might seem foolish to suggest that computational methods can contribute anything of value beyond crunching data. But this is exactly what we are suggesting. More so, it is precisely in contexts of immense complexity where their merits truly lie – in a computer's ability to efficiently sift through and identify nonobvious patterns in vast quantities of data. We argue that computational methods are a complementary, and in some cases, an essential companion, to existing approaches

to studying discourse, including qualitative analyses, code-and-count methods, reaction time studies, eye tracking, brain imaging, experimental methods, and the like. The secret is in the type of computational method advocated: machine-learned computational models (MLCM). We illustrate these points in the context of a case study involving the use of an eye-gaze-based MLCM of text comprehension during reading.

What is a Machine-Learned Computational Model (MLCM)?

A model is a representation of a thing (Frigg & Hartmann, 2018) – be it a phenomenon (e.g., reading comprehension; McNamara & Magliano, 2009), data (e.g., reading times; Graesser, Hoffman, & Clark, 1980), or a theory (e.g., landscape model of reading; Van den Broek, Young, Tzeng, & Linderholm, 1999). Models can be physical (e.g., robotic model of eye movements; Villgrattner & Ulbrich, 2010), symbolic (e.g., equations governing saccadic control; Tatler, Brockmole, & Carpenter, 2017), or a fictional entity like the phonological loop in Baddeley’s working memory model (Baddeley, 1992). A computational model is a specific type of model whose representations are *in silico* – i.e., performed on a computer or simulated by a computational device. A machine-learned computational model is a computer model *learned* from data/experience. Simply put, it is a program learned from data (Domingos, 2012). As we elaborate below, the degree and type of learning involved distinguishes MLCMs from “traditional” models and computer programs.

An MLCM has three main components: (1) a structure, such as an equation, a decision tree, an artificial neural network, or a graph; (2) feature representations, which pertain to higher order abstractions of data (e.g., fixation durations extracted from raw gaze points in eye tracking); and (3) parameters (or coefficients or weights); and (optionally) (4) hyperparameters, which control the learning process itself. For example, in the following linear regression model,

$comprehension = 5 \times total\ reading\ time + 2$, the equation is the model, comprehension is the outcome (predicted variable), total reading time is the feature, 5 and 2 are the parameters, and there are no hyperparameters. Generally, learning an MLCM from a dataset consists of adjusting the model parameters until the discrepancy between predicted and target values is minimized.

Table 1 provides a coarse-grained comparison of different computational models along a number of dimensions. The key distinguishing feature between the MLCMs (last three rows) versus traditional computational models (first row) is that the latter are more or less mathematical realization of a theory. They have a fixed structure, fixed feature representations, set parameters, and no learning. The parameters might, and probably should, be obtained from prior data, but they seldom change. Models of eye movements during reading, such as E-Z Reader (Reichle, Rayner, & Pollatsek, 2003) and SWIFT (Engbert, Nuthmann, Richter, & Kliegl, 2005), are pertinent examples of such models. Traditional computational models are theory-heavy and data-light, whereas MLCMs are data-heavy and their theoretical commitments vary, but are not as extensive as traditional models.

The simplest MLCMs are standard regression models, such as linear and logistic regressions, and their variants (e.g., ridge regression). Multicollinearity and model fitting concerns with these models often preclude the use of too many features for a given size of training dataset; they typically have a handful (usually under 10) of pre-specified features, whose coefficients are learned from data. Given the small number of parameters to be learned, this approach requires some, but not a substantial amount of data (a few hundred cases). The advantage of regression models lies in their simplicity and interpretability, but they are limited in their ability to model more complex data (such as nonlinear interactions among features). In contrast, standard machine learning methods, such as neural networks, support vector machines,

(Cortes & Vapnik, 1995) and Random Forest (Breiman, 2001) can model nonlinearities and interactions in the data. Describing each of these types of model in detail is beyond the scope of the present paper, but see Mitchell (1997); Witten and Frank (2005) for a primer on machine learning. These models typically have tens to a few hundred features, so they require more data (several hundred to a few thousand cases) to reliably estimate the parameters. Standard machine learning models can be judiciously constructed to be consistent with theory, but to a lesser extent than standard regression models. This is because there are many more free parameters to fit, and these parameters themselves determine the nature of the interaction between features and/or the nature of the function mapping features to outcomes.

Deep neural learning or deep learning (Goodfellow, Bengio, & Courville, 2016; Le Cun, Bengio, & Hinton, 2015) is a different class of models which have gained prominence over the past decade. These models are constructed by combining multiple ‘layers’ of artificial neural networks, consisting of an input layer, one or more ‘hidden’ (intermediate) layers, and an output layer; where hidden layers learn useful intermediate representations of the data by combining input features. These deep neural networks can model extremely complex phenomena. They are also capable of representation learning in that they can learn features themselves by extracting patterns from raw data instead of requiring pre-specification of features like in the other modeling approaches. Deep learning models are extremely complex with the number of free parameters in the tens to hundreds of thousands, so they require copious amounts of training data and their interpretability is low. They also have very few theoretical commitments and are basically very powerful prediction machines. As Table 1 illustrates, there is a tradeoff between theoretical commitments, explanation vs. prediction, and the amount of data needed to train

viable models. In most cases, we recommend experimenting with standard regression modeling and standard machine learning as these two approaches appropriately balance these tradeoffs.

A curious reader might ask what distinguishes standard regression models, which are extensively used in virtually all areas of science, from an MLCM regression model. In the traditional case, the focus is on the significance of the model coefficients, whereas the MLCM approach focuses on the accuracy of the model predictions to “new” data. Thus, the former approach favors explanation and description of the entire dataset whereas the latter prediction; see Yarkoni and Westfall (2017) for a detailed discussion on this issue and an enthusiastic call for psychology to engage in more predictive modeling. From a methodological standpoint, instead of building a model on the entire data set and examining p values of the coefficients and perhaps the goodness of fit, the MLCM approach will construct the model from a subset of the data and then compute fit statistics on the held-out data; a process called cross-validation. This process lessens the extent to which the model is influenced by the idiosyncrasies of individual data points (i.e., overfitting), ensuring that the model is capturing general patterns across observations (where ‘observations’ in the context of a discourse MLCM might be readers, texts, or utterances).

How Machine-Learned Computational Models can Enhance the Study of Discourse

Before delving into the potential benefits of MLCMs, let us address a common criticism that machine learning is an atheoretical fishing expedition, which produces spurious results. Although it is easy to find examples where this criticism applies, rejecting the entire field on these grounds is no different than rejecting all experimental approaches to the study of discourse based on the existence of confounded experiments (which are abundant) and careless (or unethical) data analysis methods (e.g., p -hacking - Head, Holman, Lanfear, Kahn, & Jennions,

2015). It is similar to rejecting qualitative approaches as being insufficiently rigorous or thought experiments as they are not scientific.

We think it is more productive to develop well-designed MLCMs that are scientifically rigorous *and* are useful computational tools. In our view, well-designed MLCMs of discourse should be adequately grounded in theory, but should not be overly constrained by theory. This is because discourse is complex: whilst our understanding of it is growing, it is still limited. Many of our theories apply only to particular contexts, and computational instantiations of theory are likely to fail when taken out of controlled experimental paradigms into the messiness of the real world. An MLCM is probably not needed if a phenomenon is sufficiently understood that it can be computationally instantiated with high fidelity under realistic conditions. Instead, an MLCM is most beneficial when there is some theoretical understanding, but not enough to instantiate a mechanistic model of the theory.

What can an MLCM do? It can guide theory by determining whether the right ingredients are in place or if something fundamental is missing. All things considered, a model that fails to generalize or generate accurate predictions, might suggest some missing components. For example, Bartlett, Littlewort, Frank, and Lee (2014) aimed to model expressions of pain from facial movements automatically extracted from video. They found that the temporal dynamics of the facial expressions were critical in discriminating real from posed expressions of pain. Similarly, an MLCM can ascertain which features are more important than others and what minimalistic feature set is sufficient to model the phenomenon of interest. Using the same example, a single facial movement –mouth opening – provided the most information in that the duration and variance of mouth openings as well as the interval between consecutive mouth openings was lower for faked vs. genuine pain expressions. The structure of the model itself can

provide insights into how the various components interact, for example, when one component moderates, often nonlinearly, the influence of another on the outcome prediction.

Finally, an MLCM is a computational tool that can be used for measurement, to provide feedback, to drive reflection, and for intervention. For example, Jensen, et al. (2020) developed an MLCM that automatically assesses the quality of teacher discourse in real-world classrooms. They used the model to provide feedback to teachers to help them improve their discourse. D'Mello, Mills, Bixler, and Bosch (2017) and Mills, Gregg, Bixler, and D'Mello (in press) used a previously-developed eye-gaze-based MLCM of mind wandering during reading (Faber, Bixler, & D'Mello, 2018) to deliver real-time interventions consisting of comprehension questions and self-explanations aimed at re-engaging attention and correcting any comprehension deficits associated with mind wandering.

We have developed multiple MLCMs that use linguistic, paralinguistic, behavioral, and physiological signals with the goal of understanding and/or facilitating cognitive, noncognitive, socio-affective-cognitive, and life outcomes. Such work includes a range of discourse scenarios: rhetorical, expository, pedagogical, dialogic, and collaborative discourse collected in individual, small group, multi-party, and human-computer interactions in the lab and in the wild (e.g., Bosch & D'Mello, in press; Bosch, D'Mello, Baker, Ocumpaugh, & Shute, 2016; D'Mello & Graesser, 2010; Faber, et al., 2018; Grafsgaard, Duran, Randall, Tao, & D'Mello, 2018; Hutt, et al., 2019; Kelly, Olney, Donnelly, Nystrand, & D'Mello, 2018; Stewart, et al., 2019; Stone, et al., 2019). We have also used these models for assessment (Faber, et al., 2018; Jensen, et al., 2020) and real-time intervention (Aslan, et al., 2019; D'Mello, Mills, et al., 2017; Mills, et al., in press). We have provided descriptions and tutorials of the MLCM approach along with examples in different research areas, specifically measurement of emotion (D'Mello, Kappas, & Gratch, 2018) and

engagement (D'Mello, Dieterle, & Duckworth, 2017). In the remainder of this paper, we illustrate the use of MCLMs to the study of discourse by presenting an unpublished study where we developed an MLCM of reading comprehension from eye movements.

Illustrative Example: Modeling Reading Comprehension from Eye Movements

Reading for understanding is a complex process which requires low-level text processing, active construction and maintenance of representations, retrieval and integration of information from long-term memories, and generation of predictions and inferences (Graesser, Singer, & Trabasso, 1994; Kintsch, 1988; Lesgold & Perfetti, 1978; Rayner, Pollatsek, Ashby, & Clifton Jr, 2012). With this in mind, accurately measuring reading comprehension is critical to understanding the real-time dynamics of text processing, such as whether and when readers generate elaborative inferences about the text (e.g., Graesser, et al., 1994; McKoon & Ratcliff, 1992), or whether readers are attending to text at all (e.g., Feng, D'Mello, & Graesser, 2013).

Reading comprehension is typically assessed using comprehension questions presented alongside or after the text. These questions can take various forms (multiple choice, short response, self-explanation) and can assess comprehension at different levels, for example, probing factual content from the reading (textbase-level) or deeper, inference-level understanding of the text (McNamara & Magliano, 2009). Here, we examine the question of whether an MLCM of eye movements can generate accurate, real-time and generalizable predictions of comprehension during reading. This knowledge, in turn, would contribute to theories of eye movements during reading, and the model itself can be used to assess reading comprehension as it unfolds or to trigger interventions when signs of comprehension difficulty emerge.

Background and Research Linking Eye Movements and Reading Comprehension

Given that reading requires processing of fine-grained visual stimuli (i.e., letters and words), eye movements are fundamentally linked to the cognitive processes underlying reading. Decades of research has capitalized on this insight – termed the eye-mind link (Just & Carpenter, 1976) – by using eye-tracking to investigate how readers extract coherent and even rich representations of meaning from these abstract visual stimuli. This research has demonstrated that eye movements are sensitive to text properties from the word to text levels, including word frequency (Inhoff & Rayner, 1986), lexical and syntactic ambiguity (Duffy, Morris, & Rayner, 1988; Frazier & Rayner, 1987), and text difficulty (Rayner, Chace, Slattery, & Ashby, 2006), and has made great strides toward characterizing how eyes move during reading generally (see Rayner, 2009; Rayner & Reichle, 2010 for reviews). However, in spite of decades of progress in understanding how the eyes move during reading, limited research has leveraged these insights to measure comprehension in real time.

Why might this be the case? One reason is that existing work has struggled to establish consistent links between comprehension and eye movement features. For instance, 10-25% of eye movements are regressive to an earlier part of the text (Rayner, et al., 2012). These regressions have long been interpreted as a corrective response when the reader has difficulty integrating the current word with prior context (Frazier & Rayner, 1987; Meseguer, Carreiras, & Clifton, 2002). However, some studies positively link regressions with accurate comprehension (Inhoff, Gregg, & Radach, 2016; Metzner, von der Malsburg, Vasishth, & Rösler, 2016; Schotter, Tran, & Rayner, 2014) while others show null (Christianson, Luke, Hussey, & Wochna, 2016; Wallot, Brien, Coey, & Kelty-Stephen, 2015), or even negative (Kemper, Crow, & Kemtes, 2004) associations. Relatedly, longer fixation durations have been linked to both effortful reading

(Rayner, et al., 2006) and its opposite – mind wandering (Faber, et al., 2018). On this point, recent research also suggests that people self-report re-reading the previous one or two lines of text after a mind-wandering episode (Varao-Sousa, Solman, & Kingstone, 2017), likely reflecting re-engagement with the text in an attempt to repair comprehension. Thus, although both re-engagement and mental-model repair may involve regressing to earlier parts of the text, repair might also occur covertly (i.e., resolved in working memory), typically resulting in longer fixation durations but not necessarily a regression (Meseguer, et al., 2002).

Why is establishing consistent links between eye movements and comprehension so challenging? One possibility is that eye movements primarily reflect local text processing (e.g., word identification, syntactic parsing) which is almost always successful in skilled readers, and thus does not predict later comprehension. Another possibility, however, is that mappings between eye movements and comprehension may not be consistent because they are influenced by reader- and text-specific factors. Thus, the same eye movement features may reflect different cognitive processes in different contexts. For instance, longer fixations may reflect mind-wandering (Faber, et al., 2018; Foulsham, Farley, & Kingstone, 2013), which is a negative predictor of comprehension (D'Mello, 2019; Randall, Oswald, & Beier, 2014), but may alternatively signal efforts to repair inaccurate or poor-quality text representations, presumably leading to better comprehension outcomes (Frazier & Rayner, 1982). There are also inconsistencies in the literature. For example, some studies find fewer fixations associated with mind-wandering (Bixler & D'Mello, 2016; Faber, et al., 2018; Faber, Krasich, Bixler, Brockmole, & D'Mello, in press; Smilek, Carriere, & Cheyne, 2010), whereas others find the opposite, which might be attributable to methodological differences (Faber, et al., in press).

Can an MLCM help resolve the lack of consistency between eye movements and comprehension outcomes? The answer relies on the observation that unique eye movement signatures may emerge when multiple features are considered in conjunction, whereas consistent mappings might not be evident when eye movement features are considered individually as in most studies. For instance, although longer fixations may indicate either mental model repair or mind-wandering, the two could be differentiated using other features (e.g., number of fixations and regressions). Further, processes underlying successful comprehension (e.g., motivation, attention) fluctuate over time, and may alter corresponding eye movements. When skimming or mind-wandering, for instance, eye movements exhibit less systematic correspondence with the underlying text compared to focused reading (Foulsham, et al., 2013; Reichle, Reineberg, & Schooler, 2010). Thus, examining multiple features in context of one another, in particular those which capture alignment between eye movements and the text, could illuminate systematic relationships between eye movements and comprehension. Here, we test whether an MLCM can capture these complex relationships.

Previous Work on MLCMs of Comprehension during Reading. Recent research in the human-computer interaction domain has developed MLCMs of reading comprehension, assessed alongside or immediately after reading short passages, from eye-movement features. For example, Copeland, Gedeon, and Mendis (2014); (also see Copeland & Gedeon, 2013) recorded participants' eye movements while they read a nine-slide (~400 words per slide) tutorial on conducting a web search. They trained artificial neural networks to predict performance on quiz (comprehension) questions presented alongside or immediately after each slide from slide-level eye movement features (number of fixations, mean fixation duration, total text fixation duration, number of regressions, regression fixation proportion, mean forward saccade length) as well as

eye-movement-derived features (these included the ratios of number/duration of fixations to words, and answer-seeking behavior, which was defined by saccades between the texts and the questions when the two were side-by-side). This approach yielded accurate predictions of performance on quiz questions, though comparisons with chance were not reported. See Copeland (2016) and Copeland, Gedeon, and Caldwell (2015) for additional studies which further investigate factors that influence prediction accuracy, e.g., text difficulty and whether or not the text was in the reader's first language.

Other studies have examined whether MLCMs of eye movements can be used to predict general language skill. For example, Martinez-Gomez and Aizawa (2014) examined whether participants' level of understanding, as well as English language skill, could be predicted based on their eye movements during reading. Participants read two short (~450 word) educational texts while their eye movements were recorded. They answered 8 questions to assess their understanding after each text. Random forest models with leave-one-participant-out cross validation (see below) yielded significantly above-baseline ($p < 0.001$; ~50% error reduction) predictions of binarized (high vs. low) text understanding, where the baseline consisted of simply always predicting the majority class. Notably, eye movement features were more discriminative than linguistic features in these models. However, models predicting continuous comprehension scores performed at chance levels. The researchers also obtained above-baseline ($p = 0.015$) performance when training an MLCM to predict participants' English skill as measured by their scores on standardized English tests (Test of English for International Communication [TOEIC] or Test of English as a Foreign Language [TOEFL]). Similarly, using support vector machines (SVMs), Lou, Liu, Kaakinen, and Li (2017) could discriminate high vs. low literacy participants with 80.3% accuracy (as assessed by a Chinese standardized test similar to the SAT; baseline

accuracy was not reported for comparison) using text-mapped eye-movement features (e.g., fixation times on section headers).

Though these studies provide encouraging evidence that MLCMs of eye movements could be plausibly used for diagnosis and intervention in the face of reading difficulties, there are some important limitations to consider. If eye-movement-based models are to serve as a viable mechanism for monitoring comprehension, they must be able to generate accurate predictions for previously *unseen* individuals. However, only a few studies have examined generalizability of models to new readers by using person-independent cross-validation, which entails testing models on data from participants not used for model training. Specifically, the model reported in Copeland and Gedeon (2013) showed poor performance on held-out participants. Remaining work used data-point-level cross-validation (Copeland, et al., 2015; Copeland, et al., 2014), as opposed to participant-level cross-validation, which has been acknowledged as a limitation (Copeland, 2016). Data-point-level cross-validation means that features at the level of individual observations (e.g. single pages) are randomly held out of the training data, but as a result, the same participant's data, albeit from different pages, will recur in the training and test partitions, potentially jeopardizing the generalizability of the resulting model. Participant-level cross-validation, where model performance is assessed using data from participants who did not appear in the training data, is required to demonstrate that the model generalizes to unseen participants.

Further, while the model in Martinez-Gomez and Aizawa (2014) successfully predicted the comprehension levels for the best and worst performers, they excluded data from participants with intermediate comprehension levels, which are possibly the more difficult cases. Their validation method also did not ensure generalizability to new participants. Whereas, Lou, et al. (2017) did utilize appropriate participant-level cross-validation, they focused on predicting

binarized (high vs. low) literacy skills not comprehension outcomes. Thus, existing models have only been able to make very coarse-grain predictions about new participants using eye movements. This could be useful for diagnostic purposes, but may not provide sufficient granularity to inform theory or to deploy real-time interventions during reading.

Design Considerations for a Gaze-based MLCM of Reading Comprehension. Our aim was to examine whether eye movements can be used to train an MLCM model of reading comprehension in a way that is generalizable across readers. The choice of our modeling approach was guided by theory and empirical research on eye movements during reading (cited above) and was based on a number of design considerations. First, we know that relationships between eye movement features and comprehension can be interactive, that is, features may have different relationships with comprehension depending on reader- and text-specific factors. There is also the conflicting goal of balancing prediction accuracy with model explainability (Molnar, 2019), and generalizability. With these considerations in mind, we selected Random Forest (Breiman, 2001) for our classifier. This is a classifier architecture based on decision trees. A decision tree can be thought of as a flow-chart describing possible paths to a decision based on binary decisions determined by the values of particular features, where the outcome of each decision defines which ‘branch’ to progress to next (e.g., if number of fixations < 5 , then look at number of regressions). Decision trees capture nonlinearity and interactivity between features in an interpretable fashion. For instance, a decision tree could “branch” based on number of regressions being greater or lesser than some threshold amount, and make different predictions about comprehension accuracy based on other features within each branch (e.g., few regressions could predict accurate comprehension with short vs. long reading times – i.e., the branching factor). Decision trees are relatively interpretable because their structure consists of a sequence

of readable “if, then” rules. A random forest consists of an ensemble of different decision trees, each using random subsets of training examples selected with replacement for each tree (called bagging); and within each tree, random subsets of features at each branch point. Due to these properties, random forests are more likely to generalize: by randomly leaving out fractions of the full dataset, both in terms of features and training examples, the final model comprising the full ‘forest’ is less prone to overfit to the data.

Second, to assure that the learned relationships generalize to “new” readers, we apply a person-independent cross-validation technique by testing the model using data from held-out participants that are not used in model training. If the learned relationships are too specific to the training participants (i.e., the model is overfit), then the model will perform poorly when generating predictions for previously-unseen readers. However, if performance on the held-out participants is high, this would indicate generalization to new readers albeit with data collected in similar contexts.

Third, we restricted our choice of eye gaze features to six that we largely based on prior literature. This was an important design consideration to balance the tradeoff between alignment with theory with allowing room for new discovery. We intentionally selected a smaller feature set so that the resultant models could be interpreted and to avoid the criticism of engaging in unbridled exploration. Further, in order for this approach to be applicable to real-world applications (i.e., eye-movement-based interventions), it would need to be easily applied to new texts and robust to routine eye-tracking errors. To this end, we focused on global eye movement features (e.g., number of fixations and mean fixation duration on a page) which are less affected by calibration errors, less reliant on positional information (e.g., which word is fixated), and do not need to be mapped to local text properties (e.g., the frequency of a fixated word). Our final

set of features included the number of fixations, mean fixation duration, regression fixation proportion, mean saccade length, horizontal saccade proportion, and fixation dispersion.

Fourth, we aim for the model to be fine-grained in that it can predict comprehension on individual items. Thus, rather than predicting passage- or person-level differences in comprehension (e.g., Copeland & Gedeon, 2013; Lou, et al., 2017; Martinez-Gomez & Aizawa, 2014), we predict comprehension at the page-level (where a page refers to the text presented on a computer screen). Generating page-level predictions would allow for real-time automated assessment of comprehension, which could be used to deploy interventions and to study comprehension processes on-line.

In what follows, we discuss the steps towards building the aforementioned MLCM, beginning with collection of training data for machine learning. This entailed interrupting the reader with online comprehension assessments, a required step in order to collect training data to build the model. If successful, the model can then be used to generate the assessments for new readers without interruptions as summarized below:

Training (eye gaze features + **comprehension scores**) → *computational model*
 Deployment (eye gaze features + *computational model*) → **comprehension scores**

Data

We leveraged data from a previous eye-tracking study that collected assessments of reading comprehension during a computerized reading task (D'Mello, Mills, et al., 2017). At the time of writing, the eye tracking data collected in this study has not been previously published.

Participants. Participants were 104 students at a private Midwestern university in the U.S. who participated in exchange for course credit. Participants signed a written informed consent form prior to participating, and the study was approved by the university's Institutional Review Board (IRB).

Materials and Procedure. Participants read a 6500-word excerpt from a book about the surface tension of liquids, *Soap Bubbles and the Forces which Mould Them* (Boys, 1895). The excerpt was taken from the first 35 pages of the book and was modified to remove images and associated references in text, which were not necessary for comprehension. The text was divided into 57 pages (screens; 115 words average per screen) and presented on a computer screen in 35-point Courier New font. Sentences, but not words, could be split across page boundaries. Left and right eye movements were recorded using the Tobii TX300 remote eye tracker sampling at 120 Hz. Head position was unrestrained, so participants could select a comfortable position for reading. Reading was self-paced, and participants advanced through the text one-page-at-a-time via a key press. However, they could not return to a previously-read page.

Comprehension was assessed during reading using four-option multiple-choice questions that tapped page-specific, textbase-level (i.e., factual) content of the text. Below is an example of a sentence of text and associated comprehension assessment.

Text: “Plateau in his famous work, *Statique des Liquides*, quotes a passage from a book by Henry Berthoud, to the effect that there is an Etruscan vase in the Louvre in Paris in which children are represented blowing bubbles from a pipe.”

Question: “The suggestion that there is an Etruscan vase in the Louvre that depicts children blowing bubbles from a pipe was put forth by:”

Choices: (a) Lord Rayleigh; (b) Van der Mensbrugghe; (c) Millais; (d) Plateau (**correct answer**).

Questions could occur after reading any of the 57 pages (apart from the first two); but the number of questions ($M = 15$, $SD = 4$) and exact pages on which questions appeared differed by participant. There were two groups of participants in this study, and whether a question was

asked on a given page was determined as follows. For the experimental group, the computer interface presented comprehension questions as determined by another eye-movement-based MLCM model of mind wandering (Faber, et al., 2018). This mind-wandering MLCM uses eye movement features to generate a probability that the participant is currently mind wandering. If this probability exceeded a threshold, this triggered the reading interface to display a comprehension question. A different set of yoked-control participants received identical interventions to the experimental participants – i.e. comprehension questions occurring on the same pages, irrespective of mind-wandering. Participants were allowed to advance if they answered the first of two questions correctly, or after the second question regardless. Participants were given the option to re-read the preceding page prior to answering the second question (only data from the first read of a page were included in analyses). Our present focus is training an MLCM to predict accuracy on the first comprehension question on a page from eye gaze recorded during the first read of that page.

Participants completed a posttest with 38 questions from the same pool immediately after reading, which averaged 34 minutes after reading began. These questions are not analyzed here as the current work focuses on modeling comprehension during reading.

Modeling Approach. Figure 1 depicts an overview of the modeling approach, which focused on training and validating supervised machine learning models which predict comprehension from eye movements. We examined differences in eye movement features and comprehension across conditions (intervention vs. yoked control). No significant differences were observed, so we merged across conditions to obtain the maximum amount of data for machine learning and to improve model generalizability. We also modeled these conditions separately to ascertain whether the mind-wandering intervention confounded the link between

eye movements and comprehension, i.e., whether the comprehension MLCM performance varied by experimental condition.

Eye movement features. The eye data from both eyes were averaged and then fixations and saccades were estimated using a dispersion-based fixation filter from Open Gaze and Mouse Analyzer (OGAMA; Voßkühler, Nordmeier, Kuchinke, & Jacobs, 2008). Fixations were defined as consecutive eye movement samples within a range of 57 pixels (approximately 1 degree of visual angle), and saccades were computed from the fixations. We examined first-read eye movements and accuracy for the first question (during reading) from the 1618 pages with accompanying questions.

Eye movement features were chosen based on prior literature as well as the principle that page-level fluctuations in comprehension would be best captured by features which may index alignment between eye movements and the text. Four of the eye movement features were literature-based, including *number of fixations* on the page and their *mean fixation duration* in ms. Further, we selected *regression fixation proportion*, which is the proportion of all fixations on a page that were preceded by a regression (defined as any fixations on a word with an index lower than that of the previous word that was fixated on), and *mean saccade length*, which is the average number of pixels between two subsequent fixations. These features have been used in similar modeling efforts (Copeland, 2016; Copeland & Gedeon, 2013; Copeland, et al., 2015; Copeland, et al., 2014; Martinez-Gomez & Aizawa, 2014) and have been empirically linked to factors that influence text comprehension, including mind wandering (Reichle, et al., 2010; Uzzaman & Joordens, 2011) and text difficulty (Rayner, et al., 2006). This literature-driven choice of features illustrates how MLCMs can be at least in part constrained by theory.

The remaining two eye movement features, horizontal saccade proportion and fixation dispersion, have yet to be examined in the literature, but were included for their potential to capture overall alignment of eye movements with the text, which may be linked to comprehension (Wallot, et al., 2015). Specifically, *horizontal saccade proportion* was calculated as the proportion of saccades with an angle no more than 30 degrees above or below the x-axis, and *fixation dispersion* as the root-mean-square of the distance of each fixation to the average fixation. Thus, these two features may capture especially sparse or erratic eye movements that could signal a disconnect between eye movement patterns and the text on the page. The use of this highly constrained feature set also minimized multicollinearity, which is a concern for some machine learning models (see Table 2 for between-feature correlations).

Because head movement was not restrained to allow for naturalistic reading, participants could shift position leading to eye-tracking disruptions. To compensate, we excluded pages that were clearly unread (reading time < 1 sec; 7 pages) and those without recorded eye movement data (i.e., one fixation or fewer on the page; 81 pages). Only 5% of pages (88 of 1618) were discarded as a result of these criteria, leaving 1530 pages for modeling. To address outliers, we replaced eye gaze feature values greater than 2.5 median absolute deviations (MAD) with the highest observed value of that particular feature within these bounds. As elaborated above, we further minimized the impact of eye tracking errors by focusing on global eye movement features, which are based on relative rather than absolute eye position (i.e., word-specific features like gaze duration and dwell time were not included).

Supervised classification and validation. Random Forest classification models were implemented in R with the caret package (Kuhn, 2008). We used the default parameters, which were (a) varying the number of features to split at each node, then selecting the value that

optimizes model fit; and (b) building 500 decision trees. We compared random forests, which capture interactivity and nonlinearity, to logistic regression, which is a linear additive model. This also illustrates how an MLCM can be used to answer a pertinent question on the nature of the relationship between eye movements and comprehension. Lastly, we examined whether model performance relied on systematic correspondence between eye movements and comprehension by comparing the random forest model to shuffled surrogates, created by shuffling the comprehension scores.

We used a participant-level four-fold cross-validation procedure to ensure generalizability to new participants. Specifically, data were split into four subsets at the participant level, trained on three of the subsets, and tested on the remaining subset. This process was repeated four times so that predictions were generated for all four test subsets. To assess stability of results, the entire process was repeated 100 times. There was very little variability across runs (see Table 3 for stability of model performance across the 100 runs) so we focused on the median-performing model (based on the correlation metric – see below) in our analyses.

Model evaluation. The model output is a probability that the question on a given page was answered correctly, i.e. a continuous variable between 0 and 1, termed the probability of the correct class. We computed page-level accuracy as the area under the receiver operating characteristic curve (AUROC; chance = 0.5) between the class probability of the correct class (between 0 and 1) and observed scores (1 or 0). We also averaged both predicted (0 or 1, after applying a threshold at 0.5 to the class probabilities) and observed (0 or 1) comprehension scores over all pages for each subject, and then computed the correlation between the predicted and observed performance at the participant level. Note that the model was *trained* based on its performance at the page-level for individual participants, but through averaging its predictions at

the page level we use the same MLCM to make predictions at the participant level. Whereas page-level accuracy evaluates the models' ability to make fine-grained predictions on individual comprehension items based on eye movements from corresponding pages, participant-level accuracy is a coarser measure focusing on between-subject variability. AUROCs were compared using the `roc.test` function in the `pROC` package (Robin, et al., 2011) and correlations were compared using tests of dependent correlations as implemented in the `cocor` package (Diedenhofen & Musch, 2015).

Modeling Results

Model Accuracy. The results are summarized in Table 3 and ROC curves for all models are depicted in Figure 2. The random forest model generated highly accurate predictions of both page- and participant-level comprehension (page-level AUROC = 0.902; participant-level $r = 0.661$, $p < 0.001$). As illustrated in Figure 3, the predicted distribution of comprehension scores closely aligned with the observed distribution (upper left panel), and the model accurately captured participant- (upper right panel) and page-level (bottom panel) variation in comprehension. This indicates that eye movements can be reliably linked to comprehension during reading.

To examine whether interactivity of features improved model performance, we compared the random forest model to a logistic regression (an additive) model trained on the same features. The logistic regression model also generated above-chance predictions of page- and participant-level comprehension (page-level AUROC = 0.879, participant-level $r = 0.594$, $p < 0.001$). Interactivity significantly improved page-level model performance ($z = 3.07$, $p = 0.002$), but less so for participant-level performance ($z = 1.82$, $p = 0.07$), suggesting a slight improvement.

Robustness Checks. We examined whether above-chance predictions resulted from systematic correspondence between eye movements and comprehension. The shuffled model performed at or near chance on the page- and participant-levels (AUROC = 0.48; $r = -0.02$, $p = 0.85$) and reliably worse than its non-shuffled counterpart (page-level: $z = 20.70$, $p < 0.001$; participant-level: $z = 5.81$, $p < 0.001$).

Next, to examine the impact of outlier treatment on model performance, we compared the performance of the random forest model (with outliers addressed as described above) to an identical model in which outliers were not addressed. Results were equivalent at both the page-level (model with outliers: AUROC = 0.904; comparison: $z = -0.60$, $p = 0.55$) or participant-levels (model with outliers: $r = 0.663$, $p < 0.001$; comparison: $z = -0.15$, $p = 0.87$), suggesting that similar model performance may be achieved without outlier treatment. Results were also similar and not significantly different ($z = -1.35$, $p = 0.177$) when no pages were removed and missing feature values were replaced with zeroes (zero imputation; AUROC = 0.864; $r = 0.633$, $p < 0.001$).

Checks for Confounds. It was critical to determine whether model performance could be attributed to the contingency between eye movements and interventions that were triggered by mind wandering (in the experimental condition). Thus, we analyzed whether predictions were dependent on experimental condition (intervention vs. yoked control) by modeling each condition separately. Both models performed above-chance on the page- and participant-levels (intervention: AUROC = 0.882, $r = 0.652$, $p < 0.001$; yoked-control: AUROC = 0.892, $r = 0.664$, $p < 0.001$), and there was no reliable difference in performance across the models (participant-level $z = -0.100$, $p = 0.921$; page-level $z = 0.522$, $p = 0.601$). Thus, model performance did not rely on the mind-wandering-contingent interventions.

Relatedly, to what extent do associations between eye movement features and comprehension reflect lower-level processes, such as whether participants were attending to the text? To examine this, we computed correlations between participants' average mind-wandering likelihood (derived from the eye-movement-based MLCM model used to trigger the interventions) for 100 participants with available data and their predicted comprehension scores from the median random forest model. As expected, we observed a significant negative correlation between mind-wandering and predicted comprehension ($r = -0.21$, $p = 0.04$), though mind-wandering accounted for only 4.4% of the variance in predictions ($R^2 = 0.044$). This suggests that the model captured higher-level reading processes in addition to whether or not the reader was attending to the text.

Predictive Features. How do eye movement features relate to comprehension? To investigate this question, we analyzed the models further. We opted to focus on the logistic regression model in lieu of the Random Forest model because the former, though slightly less accurate, is more interpretable; this is because feature values are linearly combined in logistic regression such that the coefficients are readily interpretable as the contribution of each feature to model performance. One way to quantify coefficient importance is to examine the model coefficients by computing the mean and variability across the cross-validated folds. The alternate approach, which we utilized here, is to focus on significance of the coefficients. Accordingly, we fit mixed effects logistic regression models in order to quantify coefficient significance using the lme4 library in R (Bates & Maechler, 2010) with participant as an intercept-only random effect. To address outliers, feature values were MAD-scaled prior to modeling (scaled values were truncated at ± 2.5 MAD).

Model coefficients are presented in Table 4. We found that more fixations, shorter saccades, fewer horizontal saccades, and lower fixation dispersion were associated with better comprehension scores. Why might this be the case? To illustrate, Figure 4 depicts example eye movements from two participants corresponding to correct (left) and incorrect (right) responses to a comprehension question immediately after reading a page. Longer saccades and fewer fixations are indicative of skimming (Rayner, et al., 2012), which is the inverse of what predicted accurate comprehension in our model. Further, higher horizontal saccade proportion may result from eye movements which too rigidly align with the text, which could also reflect skimming (as in the right compared to the left panel of Figure 4). Lastly, higher fixation dispersion may correspond with sparse and erratic eye movements with poor correspondence to the text. Thus, lower horizontal saccade proportion and fixation dispersion, in combination with more fixations and shorter saccades, may indicate that the reader was sufficiently engaged to construct high-quality representations of the text.

Note that the marginal R^2 for the mixed-effects model was 0.566 and the conditional R^2 was 0.680. This corresponds to marginal $r = 0.752$ and conditional $r = 0.824$. This is indeed a highly accurate model, with correlation coefficients higher than for both the MLCMs (random forest and logistic regression). However this can be attributed to the difference in the way these models were fit. The mixed model is an explanatory model, fit to the entire dataset as a whole. Therefore, the marginal and conditional R^2 do not reflect the predictive accuracy of the model. In contrast, the same model fit with the MLCM approach, has a correlation of .627 when fit with cross-validation. Put differently, this reduction in correlation provides an index of the extent to which the mixed modeling approach overfits to the data. This model, if it were used to predict

comprehension on unseen participants, would likely underperform as the model will have overfit to the training data.

Discussion of Modeling Approach and Results

Our goal was to demonstrate how eye movements could be used to develop an MLCM of reading comprehension. Our results showed that our models produced highly accurate predictions of page- and participant-level comprehension, suggesting that this approach yields predictions of sufficient accuracy and granularity (AUROC of .90; $r = .66$) to be used for research as well as intervention and diagnostic purposes. Further, we demonstrated that the success of these models relies on the page-level correspondence between eye movements and comprehension, specifically, as illustrated by comparisons with a random surrogate model.

This current model offers several crucial improvements over prior work. First, it extends previous research predicting overall comprehension (Copeland, 2016; Copeland & Gedeon, 2013; Copeland, et al., 2015; Copeland, et al., 2014; Martinez-Gomez & Aizawa, 2014) to real-time, page-level comprehension during reading. This is critical because real-time assessment of comprehension is needed to better our theoretical understanding of short-lived, dynamic phenomena during reading. For example, as suggested by the bottom panel of Figure 3, comprehension may fluctuate throughout a reading as attention waxes and wanes, or if particular text regions are especially difficult to integrate. By leveraging such a real-time, eye-movement-based measure of comprehension, we anticipate that researchers will be better-poised to study text comprehension dynamics without interrupting the reader.

This work also provides encouraging evidence for the generalizability and scalability of this approach. Namely, using cross-validated models, we achieved highly-accurate predictions on held-out participants, suggesting that the models can be used to predict comprehension for new

readers in a similar context. In further support of this, the current modeling approach generalized across two experimental conditions (intervention vs. yoked-control) and were highly stable across 100 runs (see the lower half of Table 3). In addition, by using global eye movement features (e.g., number and duration of fixations on a page), the current approach minimizes the need for human coding of the text (e.g., defining text regions of interest and defining linguistic properties for those regions) and minimizes the impact of routine eye-tracking and calibration errors that would be expected in real-world settings (e.g., classrooms).

In addition, our results provide further insights on how eye movements are linked to comprehension. We hypothesized based on previous research that modeling interactivity among eye movements is key towards capturing the link between eye movements and comprehension. In line with this, when we compared random forest (interactive) and regression (additive) algorithms, we observed that modeling interactivity slightly improved page-level model performance, but the improvement was quite minor, suggesting that an additive model would suffice. Indeed a benefit of the regression model over the random forest model is that the contribution of individual features to model performance was more easily quantified. The trade-off between the ability to model more complex data and the ability to interpret the structure of the resulting model in terms of its individual features is indeed an important consideration in the choice of specific model to use with an MLCM approach. Accordingly, we found that more fixations, along with shorter saccades (the inverse of eye movements typically observed during skimming) was predictive of accurate comprehension. Further, features which indexed spatial alignment between eye movements and the text (horizontal saccade proportion and fixation dispersion) also predicted comprehension. Collectively, these findings suggest that eye

movements during reading capture higher-level processes (e.g., motivation) which can be used to accurately diagnose comprehension in real-time.

Notably, models of eye movements during reading assume the reader is attending, and that local text processing is generally successful (e.g., E-Z Reader model: (Reichle, Pollatsek, Fisher, & Rayner, 1998). The current research suggests that such assumptions might be reconsidered when the goal is to model naturalistic reading. Instead, a comprehensive model of eye movements during reading of longer, connected texts must integrate attentional and motivational processes (e.g., mind-wandering) that are neglected in current models. How could this be accomplished? One approach could be to develop models which detect the unique eye movement signatures of reading behaviors and associated mental states (e.g., mind-wandering, skimming, and repair) that interfere with or facilitate comprehension.

Encouragingly, an MLCM modeling approach has been applied to detect mind-wandering (Bixler & D'Mello, 2016; Faber, et al., 2018) and skimming (Biedert, Hees, Dengel, & Buscher, 2012). However, much work still needs to be done to model what occurs during motivated, attentive reading, some of which might be covert (e.g., developing models to predict covert comprehension repair and inferencing). These models of the reading process could also offer improvement over the current results by integrating predictions from several sub-models (e.g., of skimming, mind-wandering, and repair) in an ensemble-like fashion to generate highly-robust predictions of comprehension accuracy.

Despite the success of the present MLCM that models eye movements from eye gaze, many open questions remain. For example, the current approach models only one dataset with a single text and only textbase-level assessments of comprehension. Thus, the extent to which the current findings generalize to other contexts (e.g., texts) and whether eye movement measures

capture deeper text processing remains to be investigated. A critical test of the current MLCM would be to test its performance on gaze data from uninterrupted reading. If this model successfully predicts multiple-choice item performance *after* reading, then it could be usable as a fine-grained, page-level comprehension measure without interrupting the participant. Further, we did not explicitly measure the underlying reading behaviors and processes involved in reading (such as skimming, inferencing and error-monitoring), which would be a critical step towards building more comprehensive models of real-time reading comprehension, which would make more explicit the link between the model structure and theory-derived constructs. An additional avenue for future work using this MLCM would be to assess whether it generalizes to reading different texts and under different contexts.

Keith Rayner – perhaps the most influential scholar of eye movements and reading of the past 40 years – and his colleagues (2006) imagined that once eye-tracking technology became more portable and affordable, eye movements could be used as a diagnostic and intervention tool for comprehension difficulties. With consumer-off-the-shelf (COTS) eye-tracking technology becoming more cost-effective and accurate (Gibaldi, Vanegas, Bex, & Maiello, 2017) that day has seemingly arrived from a technological standpoint. Whereas we use an expensive research-grade eye tracker in this work, other studies have successfully developed MLCM using consumer-off-the-shelf (COTS) eye trackers in authentic environments (Hutt, et al., 2019). For example, Hutt, et al. (2019) successfully developed MLCMs of mind wandering while interacting with a learning technology using data collected with COTS eye-trackers in a classroom context. The researchers also found the model to be comparably accurate to models based on data collected in a laboratory setting with COTS eye-trackers (Hutt, Mills, White, Donnelly, & D’Mello, 2016). Thus, recent advances in eye-tracking technology coupled with the

MLCM approach advocated here take us a step towards bringing Keith Rayner's predictions about the future of eye-movement-based diagnostic and intervention tools from the realm of science fiction to plausible reality.

Concluding Remarks

It is generally accepted that computational analyses of discourse can complement other methods including think-alouds, code and count, experimental methods, and the like. Here we suggest that machine-learned computational models (MLCMs) provide a unique, yet complementary, approach to the study of discourse. To make our case, we distinguished MLCMs from traditional computational models, highlighted different types of MLCMs, discussed their potential benefits, and demonstrated the overall idea by developing and validating an MLCM model of reading comprehension from eye movements. Despite the encouraging success of these models, there is still a long way to go. Even the most sophisticated MLCM is unlikely to deeply understand, for example, the excerpt from *The King in Yellow* reproduced in the opening lines of this article. We do not know the answer yet, except to advocate for a future of discourse research that incorporates computational models in its existing arsenal of theoretical, observational, and experimental methods. In this future, MLCMs are rapidly instantiated to analyze experimental data while also serving as measurement and intervention tools. Qualitative inspection of their parameters and behavior provide insights into the underlying theories, which in turn, can be tested via different instantiations of the models. Thus, computational models and theory development go hand in hand.

References

- Aslan, S., Alyuz, N., Tanriover, C., Mete, S. E., Okur, E., D'Mello, S. K., & Esme, A. A. (2019). Investigating the Impact of a Real-time, Multimodal Student Engagement Analytics Technology in Authentic Classrooms *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI 2019)*. . New York: ACM.
- Baddeley, A. (1992). Working memory. *Science*, 255(5044), 556-559.
- Bartlett, M. S., Littlewort, G. C., Frank, M. G., & Lee, K. (2014). Automatic decoding of facial movements reveals deceptive pain expressions. *Current Biology*, 24(7), 738-743.
- Bates, D. M., & Maechler, M. (2010). lme4: Linear mixed-effects models using S4 classes. Retrieved from <http://CRAN.R-project.org/package=lme4>
- Biedert, R., Hees, J. o., rn, Dengel, A., & Buscher, G. (2012). *A robust realtime reading-skimming classifier*. Paper presented at the Proceedings of the Symposium on Eye Tracking Research and Applications.
- Bixler, R., & D'Mello, S. K. (2016). Automatic gaze-based user-independent detection of mind wandering during computerized reading. *User Modeling & User-Adapted Interaction*, 26, 33-68.
- Bosch, N., & D'Mello, S. K. (in press). Automatic Detection of Mind Wandering from Video in the Lab and in the Classroom. *IEEE Transactions on Affective Computing*.
- Bosch, N., D'Mello, S., Baker, R., Ocumpaugh, J., & Shute, V. (2016). Using video to automatically detect learner affect in computer-enabled classrooms. *ACM Transactions on Interactive Intelligent Systems*, 6(2), 17.11-17.31.
- Boys, C. V. (1895). *Soap bubbles, their colours and the forces which mold them*: Society for Promoting Christian Knowledge.

- Breiman, L. (2001). Random forests. *Machine learning*, 45, 5-32.
- Christianson, K., Luke, S. G., Hussey, E. K., & Wochna, K. L. (2016). Why reread? Evidence from garden-path and local coherence structures. *The Quarterly Journal of Experimental Psychology*, 0218(May), 1-51.
- Copeland, L. (2016). *Eye Tracking to Support eLearning*. Retrieved from <https://openresearch-repository.anu.edu.au/handle/1885/108880>
- Copeland, L., & Gedeon, T. (2013). Measuring reading comprehension using eye movements. *2013 IEEE 4th International Conference on Cognitive Infocommunications (CogInfoCom)*, 791-796.
- Copeland, L., Gedeon, T., & Caldwell, S. (2015). *Effects of text difficulty and readers on predicting reading comprehension from eye movements*. Paper presented at the 6th IEEE Conference on Cognitive Infocommunications, CogInfoCom 2015.
- Copeland, L., Gedeon, T., & Mendis, S. (2014). Predicting reading comprehension scores from eye movements using artificial neural networks and fuzzy output error. *Artificial Intelligence Research*, 3(3).
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273-297.
- D'Mello, S. K. (2019). What do we think about when we learn? In K. Millis, J. Magliano, D. Long & K. Wiemer (Eds.), *Deep Comprehension: Multi-Disciplinary Approaches to Understanding, Enhancing, and Measuring Comprehension* (pp. 52-67). New York, NY: Routledge.
- D'Mello, S. K., Dieterle, E., & Duckworth, A. L. (2017). Advanced, Analytic, Automated (AAA) Measurement of Engagement during Learning. *Educational Psychologist*, 52(2), 104-123.

- D'Mello, S. K., & Graesser, A. (2010). Multimodal semi-automated affect detection from conversational cues, gross body language, and facial features. *User Modeling and User-adapted Interaction*, 20(2), 147-187.
- D'Mello, S. K., Mills, C., Bixler, R., & Bosch, N. (2017). Zone out no more: Mitigating mind wandering during computerized reading. In X. Hu, T. Barnes, A. Hershkovitz & L. Paquette (Eds.), *Proceedings of the 10th International Conference on Educational Data Mining* (pp. 8-15): International Educational Data Mining Society.
- D'Mello, S. K., Kappas, A., & Gratch, J. (2018). The Affective Computing Approach to Affect Measurement. *Emotion Review*, 10(2), 174-183.
- Diedenhofen, B., & Musch, J. (2015). cocor: A comprehensive solution for the statistical comparison of correlations. *PLoS ONE*, 10(4), e0121945.
- Domingos, P. (2012). A few useful things to know about machine learning. *Communications of the ACM*, 55(10), 78-87.
- Duffy, S. A., Morris, R. K., & Rayner, K. (1988). Lexical ambiguity and fixation times in reading. *Journal of Memory and Language*, 27(4), 429-446.
- Engbert, R., Nuthmann, A., Richter, E. M., & Kliegl, R. (2005). SWIFT: a dynamical model of saccade generation during reading. *Psychological Review*, 112(4), 777.
- Faber, M., Bixler, R., & D'Mello, S. K. (2018). An automated behavioral measure of mind wandering during computerized reading. *Behavior Research Methods*, 50(1), 134-150.
- Faber, M., Krasich, K., Bixler, R., Brockmole, J., & D'Mello, S. K. (in press). The Eye-Mind Wandering Link: Identifying Gaze Indices of Mind Wandering Across Tasks. *Journal of Experimental Psychology: Human Perception and Performance*.

- Feng, S., D'Mello, S., & Graesser, A. (2013). Mindwandering while reading easy and difficult texts. *Psychonomic Bulletin & Review*, 20(1), 586-592.
- Foulsham, T., Farley, J., & Kingstone, A. (2013). Mind wandering in sentence reading: Decoupling the link between mind and eye. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 67(1), 51.
- Frazier, L., & Rayner, K. (1982). Making and correcting errors during sentence comprehension. *Cognitive Psychology*, 14(2), 178-210.
- Frazier, L., & Rayner, K. (1987). Resolution of syntactic category ambiguities: Eye movements in parsing lexically ambiguous sentences. *Journal of Memory and Language*, 26(5), 505-526.
- Frigg, R., & Hartmann, S. (2018). Models in science. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Summer 2018 Edition ed.).
- Gibaldi, A., Vanegas, M., Bex, P. J., & Maiello, G. (2017). Evaluation of the Tobii EyeX Eye tracking controller and Matlab toolkit for research. *Behavior research methods*, 49(3), 923-946.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. Cambridge, MA: MIT press.
- Graesser, A., Singer, M., & Trabasso, T. (1994). Constructing inferences during narrative text comprehension. *Psychological Review*, 101(3), 371-395.
- Graesser, A. C., Hoffman, N. L., & Clark, L. F. (1980). Structural components of reading time. *Journal of Verbal Learning and Verbal Behavior*, 19(2), 135-151.
- Grafsgaard, J. F., Duran, N. D., Randall, A. K., Tao, C., & D'Mello, S. K. (2018). Generative Multimodal Models of Nonverbal Synchrony in Close Relationships In S. K. D'Mello, L.

- Yin, L. P. Morency & M. Valstar (Eds.), *Proceedings of the 13th IEEE Conference on Automatic Face and Gesture Recognition (FG'18)* (pp. 195-202). Washington, DC: IEEE.
- Head, M. L., Holman, L., Lanfear, R., Kahn, A. T., & Jennions, M. D. (2015). The extent and consequences of p-hacking in science. *PLoS biology*, 13(3), e1002106.
- Hutt, S., Krasich, K., Mills, C., Bosch, N., White, S., Brockmole, J., & D'Mello, S. K. (2019). Automated Gaze-Based Mind Wandering Detection during Computerized Learning in Classrooms. *User Modeling & User-Adapted Interaction*.
- Hutt, S., Mills, C., White, S., Donnelly, P. J., & D'Mello, S. K. (2016). The Eyes Have It: Gaze-based Detection of Mind Wandering during Learning with an Intelligent Tutoring System. In T. Barnes, M. Chi & M. Feng (Eds.), *Proceedings of the 9th International Conference on Educational Data Mining (EDM 2016)* (pp. 86-93): International Educational Data Mining Society.
- Inhoff, A. W., Gregg, J., & Radach, R. (2016). Eye movement programming and reading accuracy. *The Quarterly Journal of Experimental Psychology*, 0(0), 1-9.
- Inhoff, A. W., & Rayner, K. (1986). Parafoveal word processing during eye fixations in reading: effects of word frequency. *Perception & psychophysics*, 40(6), 431-439.
- Jensen, E., Dale, M., Donnelly, P., Stone, C., Kelly, S., Godley, A., & D'Mello, S. K. (2020). Toward Automated Feedback on Teacher Discourse to Enhance Teacher Learning *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI 2020)*.
- Just, M. A., & Carpenter, P. A. (1976). Eye fixations and cognitive processes. *Cognitive Psychology*, 8(4), 441-480.

- Kelly, S., Olney, A. M., Donnelly, P., Nystrand, M., & D'Mello, S. K. (2018). Automatically Measuring Question Authenticity in Real-World Classrooms. *Educational Researcher*, 47(7), 451-464.
- Kemper, S., Crow, A., & Kemtes, K. (2004). Eye-fixation patterns of high- and low-span young and older adults: down the garden path and back again. *Psychology and aging*, 19(1), 157-170.
- Kintsch, W. (1988). The role of knowledge in discourse comprehension: A Construction-Integration model. *Psychological Review*, 95, 163-182.
- Kuhn, M. (2008). Building predictive models in R using the caret package. *Journal of statistical software*, 28(5), 1-26.
- Le Cun, Y., Bengio, Y., & Hinton, G. E. (2015). Deep learning. *Nature*, 521, 436-444.
- Lesgold, A. M., & Perfetti, C. A. (1978). Interactive processes in reading comprehension. *Discourse Processes*, 1(4), 323-336.
- Lou, Y., Liu, Y., Kaakinen, J. K., & Li, X. (2017). Using support vector machines to identify literacy skills: Evidence from eye movements. *Behavior research methods*, 49(3), 887-895.
- Martinez-Gomez, P., & Aizawa, A. (2014). *Recognition of Understanding Level and Language Skill using Measurements of Reading Behavior*. Paper presented at the IUI.
- McKoon, G., & Ratcliff, R. (1992). Inference during reading. *Psychological Review*, 99(3), 440-466.
- McNamara, D. S., & Magliano, J. (2009). Toward a comprehensive model of comprehension. *Psychology of learning and motivation*, 51, 297-384.

- Meseguer, E., Carreiras, M., & Clifton, C. (2002). Overt reanalysis strategies and eye movements during the reading of mild garden path sentences. *Memory & cognition*, 30(4), 551-561.
- Metzner, P., von der Malsburg, T., Vasishth, S., & Rösler, F. (2016). The Importance of Reading Naturally: Evidence From Combined Recordings of Eye Movements and Electric Brain Potentials. *Cognitive Science*, 1-32.
- Mills, C., Gregg, J., Bixler, R., & D'Mello, S. K. (in press). Eye-Mind Reader: An Intelligent Reading Interface that Promotes Long-term Comprehension by Detecting and Responding to Mind Wandering. *Human-Computer Interaction*.
- Mitchell, T. (1997). *Machine Learning*: Mc-Graw-Hill.
- Molnar, C. (2019). *Interpretable machine learning*.
- Randall, J. G., Oswald, F. L., & Beier, M. E. (2014). Mind-wandering, cognition, and performance: A theory-driven meta-analysis of attention regulation. *Psychological Bulletin*, 140(6), 1411-1431.
- Rayner, K. (2009). Eye movements in reading: Models and data. *Journal of Eye Movement Research*, 2(5), 1-10.
- Rayner, K., Chace, K. H., Slattery, T. J., & Ashby, J. (2006). Eye movements as reflections of comprehension processes in reading. *Scientific Studies of Reading*, 10(3), 241-255.
- Rayner, K., Pollatsek, A., Ashby, J., & Clifton Jr, C. (2012). *Psychology of reading*. New York: Psychology Press.
- Rayner, K., & Reichle, E. D. (2010). Models of the reading process. *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(6), 787-799.

- Reichle, E. D., Pollatsek, A., Fisher, D. L., & Rayner, K. (1998). Toward a model of eye movement control in reading. *Psychological Review*, 105(1), 125.
- Reichle, E. D., Rayner, K., & Pollatsek, A. (2003). The EZ Reader model of eye-movement control in reading: Comparisons to other models. *Behavioral and Brain Sciences*, 26(4), 445-476.
- Reichle, E. D., Reineberg, A. E., & Schooler, J. W. (2010). Eye movements during mindless reading. *Psychological Science*, 21(9), 1300-1310.
- Robin, X., Turck, N., Hainard, A., Tiberti, N., Lisacek, F., Sanchez, J.-C., & Müller, M. (2011). pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC bioinformatics*, 12:77.
- Schotter, E. R., Tran, R., & Rayner, K. (2014). Don't Believe What You Read (Only Once): Comprehension Is Supported by Regressions During Reading. *Psychological Science*, 25(6), 1218-1226.
- Smilek, D., Carriere, J. S. A., & Cheyne, J. A. (2010). Out of mind, out of sight: Eye blinking as indicator and embodiment of mind wandering. *Psychological Science*, 21(6), 786-789.
- Stewart, A., Vrzakova, H., Sun, C., Yonehiro, J., Stone, C., Duran, N., Shute, V., & D'Mello, S. K. (2019). I Say, You Say, We Say: Using Spoken Language to Model Socio-Cognitive Processes during Computer-Supported Collaborative Problem Solving. *Proceedings of the ACM on Human-Computer Interaction, Volume 3(CSCW)*, 39:31-19.
- Stone, C., Quirk, A., Gardener, M., Hutt, S., Duckworth, A. L., & D'Mello, S. K. (2019). Language as Thought: Using Natural Language Processing to Model Noncognitive Traits that Predict College Success *Proceedings of the 9th International Learning Analytics and Knowledge Conference (LAK'19)*.

- Tatler, B. W., Brockmole, J. R., & Carpenter, R. H. (2017). LATEST: A model of saccadic decisions in space and time. *Psychological Review*, 124(3), 267-300.
- Uzzaman, S., & Joordens, S. (2011). The eyes know what you are thinking: Eye movements as an objective measure of mind wandering. *Consciousness and Cognition*, 20(4), 1882-1886.
- Van den Broek, P., Young, M., Tzeng, Y., & Linderholm, T. (1999). The landscape model of reading: Inferences and the online construction of a memory representation. *The construction of mental representations during reading*, 71-98.
- Varao-Sousa, T. L., Solman, G. J., & Kingstone, A. (2017). Re-reading after mind wandering. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 71(3), 203-211.
- Villgrattner, T., & Ulbrich, H. (2010). Optimization and dynamic simulation of a parallel three degree-of-freedom camera orientation system *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems* (pp. 2829-2836). Washington DC: IEEE.
- Voßkühler, A., Nordmeier, V., Kuchinke, L., & Jacobs, A. M. (2008). OGAMA (Open Gaze and Mouse Analyzer): open-source software designed to analyze eye and mouse movements in slideshow study designs. *Behavior Research Methods*, 40(4), 1150-1162.
- Wallot, S., Brien, B. A. O., Coey, C. A., & Kelty-Stephen, D. G. (2015). *Power-law fluctuations in eye movements predict text comprehension during connected text reading*. Paper presented at the CogSci.
- Witten, I., & Frank, E. (2005). *Data Mining: Practical machine learning tools and techniques* (2nd ed.). San Francisco, CA: Morgan Kaufmann.

Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science*, 12(6), 1100-1122.

Table 1. Comparison of different types of computational models.

Structure	Feature Representations	Parameters/ Coefficients	Functions Learned	Theoretical Commitments	Data Required	Model Type	Examples
Fixed	Fixed	Fixed	None	Most	Fewest	Mathematical models	EZ Reader; SWIFT
Fixed	Fixed	Learned	Linear classifiers	More	Fewer	Standard regression modeling	Generalized linear models (e.g., Linear & Logistic regression)
Fixed	Fixed	Learned	Linear and nonlinear	Fewer	More	Standard machine learning	Random Forest; Support vector machines; Shallow neural networks
Fixed	Fixed/Learned	Learned	Very complex nonlinear functions	Fewest	Most	Deep neural learning	Convolutional neural networks; Long short-term memory neural networks ¹

Note. ¹ Although subsumed under deep learning methods, these models can be considered shallow if only a single hidden layer is used.

Table 2. Descriptive and correlation table of eye gaze features and comprehension accuracy. Means and standard deviations in column 1 are computed over participants.

	Mean (SD)	Mean Fix Dur	N Fix	Reg Fix Prop	Mean Sacc Length	Horiz Sacc Prop	Fix Disp
Mean Fixation Duration (ms)	267 (33)	-					
Number of Fixations	100 (36)	0.323	-				
Regression Fixation Proportion	0.136 (0.037)	0.180	0.269	-			
Mean Saccade Length (pixels)	237 (29)	-0.429	-0.419	-0.102	-		
Horizontal Saccade Proportion	0.891 (0.073)	-0.194	0.170	-0.210	0.030	-	
Fixation Dispersion	0.397 (0.050)	-0.276	-0.440	-0.274	0.245	0.211	
Comprehension Accuracy	0.684 (0.169)	0.228	0.372	0.200	-0.260	-0.430	-0.421

Table 3. Summary of classification model performance.

Model	Page-Level (AUROC)	Participant-Level (Correlation)
Median Models	AUROC [95% CI]	Pearson's r [95% CI]
Random Forest	0.902* [0.885, 0.919]	0.661*** [0.537, 0.757]
Logistic Regression	0.879* [0.860, 0.899]	0.594*** [0.453, 0.706]
Shuffled Random Forest	0.475 [0.494, 0.556]	-0.019 [-0.211, 0.174]
Across 100 Iterations	Mean [min, max]	Mean [min, max]
Random Forest	0.902 [0.89, 0.91]	0.663 [0.602, 0.705]
Logistic Regression	0.878 [0.866, 0.883]	0.592 [0.548, 0.617]
Shuffled Random Forest	0.463 [0.437, 0.486]	-0.024 [-0.193, 0.125]

Note. Significant AUROCs (bootstrapped 95% CIs non-overlapping with chance [0.5]) and correlations are marked with asterisks (*). For correlations, * = $p < .05$, ** = $p < .01$, *** = $p < .001$.

Table 4. Coefficients from a generalized linear mixed model with eye movement features as fixed effects and participants as an intercept-only random effect.

Predictors	Comprehension Accuracy		
	<i>Coefficient</i>	<i>95% CI</i>	<i>p</i>
(Intercept)	1.34	1.06 – 1.63	<0.001
Mean Fixation Duration	-0.05	-0.27 – 0.17	0.642
Number of Fixations	1.22	0.97 – 1.47	<0.001
Regression Fixation Proportion	-0.01	-0.21 – 0.19	0.948
Mean Saccade Length	-0.47	-0.68 – -0.27	<0.001
Horizontal Saccade Proportion	-1.57	-1.78 – -1.36	<0.001
Fixation Dispersion	-0.49	-0.68 – -0.30	<0.001

Figure 1. A schematic of the machine learning approach. Models were trained to predict page-level multiple-choice comprehension question performance (correct vs. incorrect) from the eye movement features of the associated page using data from computerized reading studies.

“OGAMA” = Open Gaze and Mouse Analyzer; “AUROC” = area under the receiver operating characteristic curve.

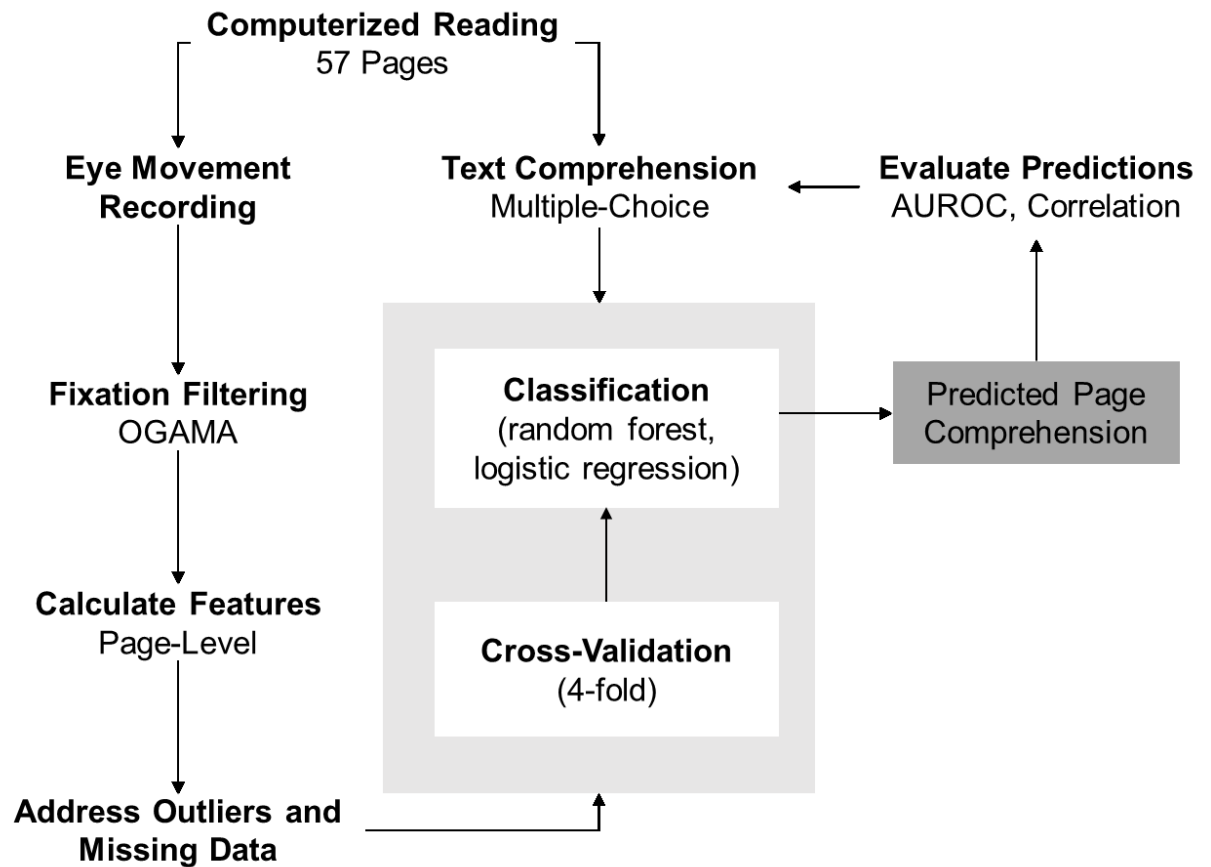


Figure 2. Receiver Operating Characteristic (ROC) curves for the median models.

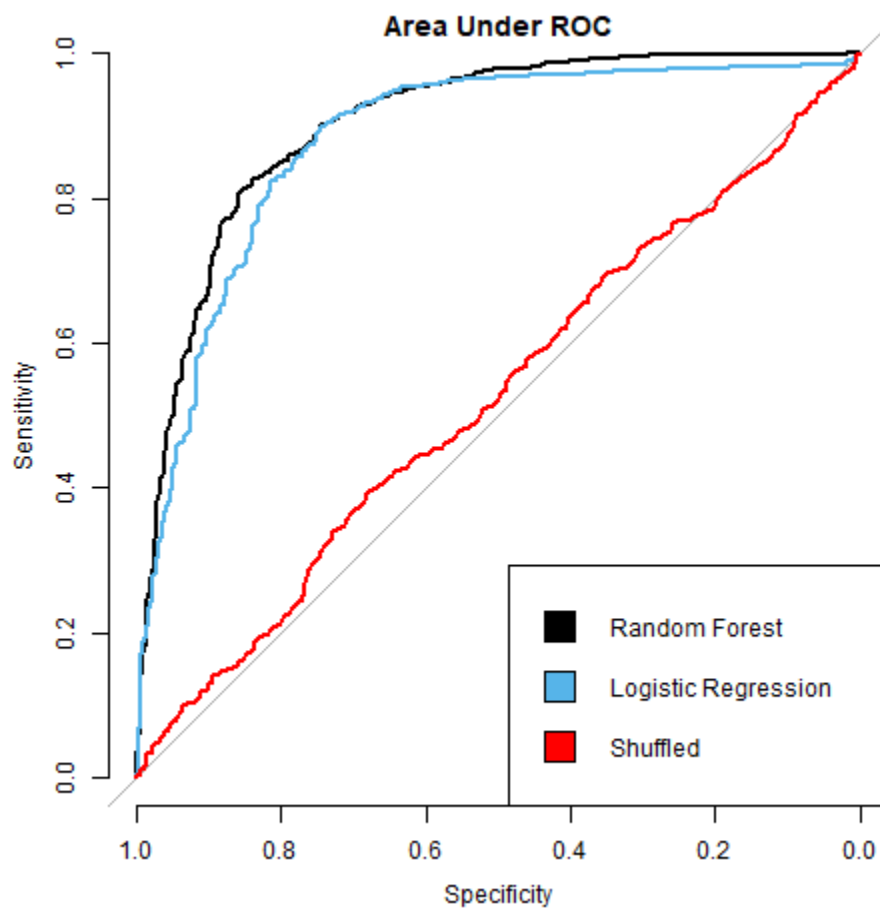


Figure 3. Visualization of page- and participant-level performance for the median random forest model predicting comprehension from eye movements. Upper left: predicted and observed participant-level mean comprehension accuracy. Upper right: correlation between predicted and observed participant-level mean comprehension accuracy. Lower: mean page-level accuracy; predicted and observed. Error bars are 95% confidence intervals computed across participants.

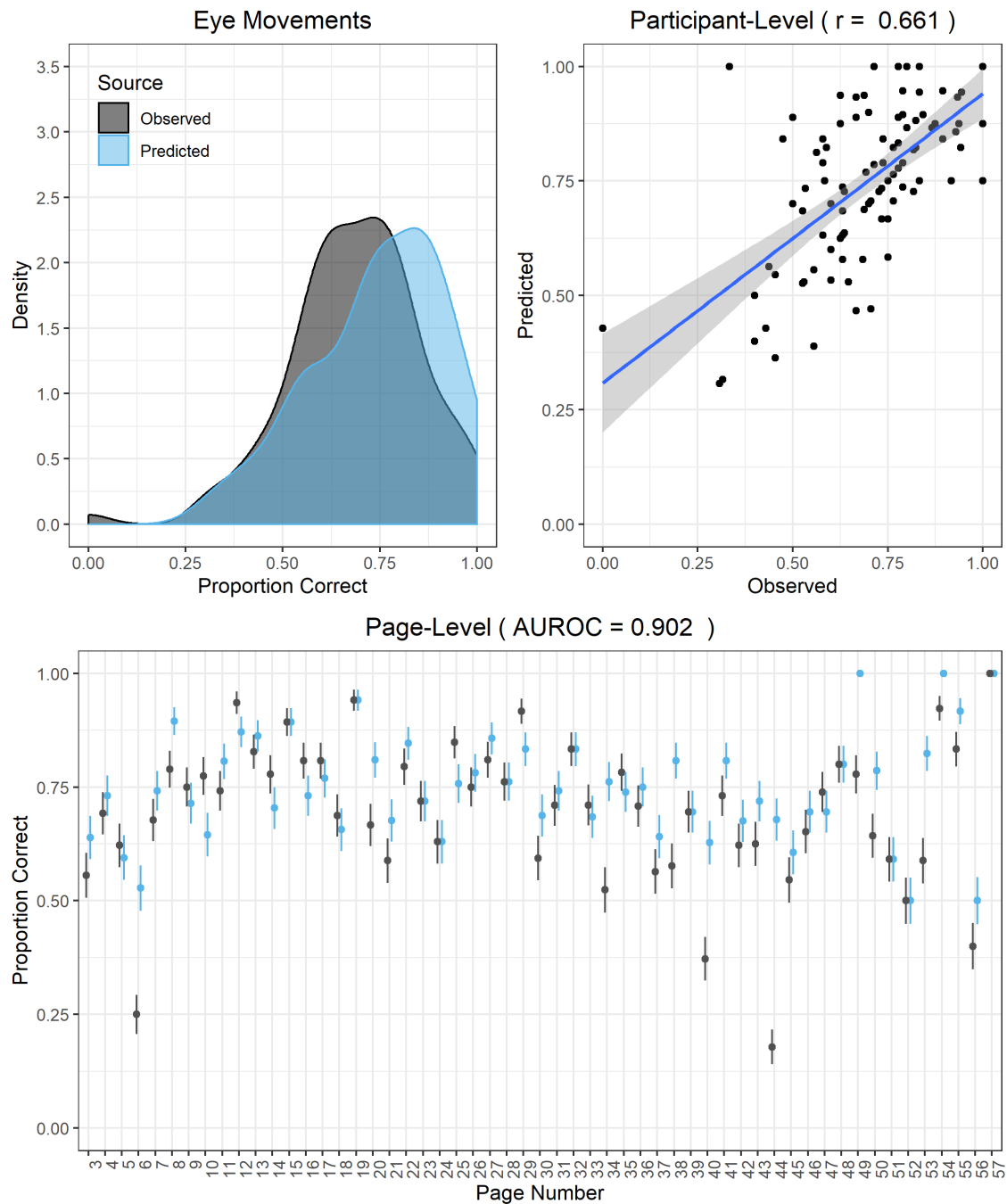


Figure 4. Comparison of eye movements from two participants reading the same page. Circles represent fixations and are scaled by duration, and lines represent saccades. The eye movements in the left and right panels preceded correct and incorrect responses to a subsequent comprehension question, respectively.

