

Prediction and preview strongly affect reading times but not skipping during natural reading

Micha Heilbron^{1,2}, Jorie van Haren¹, Peter Hagoort^{1,2}, Floris P. de Lange¹

¹Donders Institute, Radboud University, Nijmegen, the Netherlands

²Max Planck Institute for Psycholinguistics, Nijmegen, the Netherlands

micha.heilbron@donders.ru.nl

Abstract

In a typical text, readers look much longer at some words than at others and fixate some words multiple times, while skipping others altogether. Historically, researchers explained this variation via low-level visual or oculomotor factors, but today it is primarily explained via cognitive factors, such as how well words can be predicted from context or discerned from parafoveal preview. While the existence of these effects has been well established in experiments, the relative importance of prediction, preview and low-level factors for eye movement variation in natural reading is unclear. Here, we address this question in three large datasets (n=104, 1.5 million words), using a deep neural network and Bayesian ideal observer to model linguistic prediction and parafoveal preview from moment to moment in natural reading. Strikingly, neither prediction nor preview was important for explaining word skipping – the vast majority of skipping was explained by a simple oculomotor model. For reading times, by contrast, we found strong but independent contributions of both prediction and preview, and effect sizes matching those from controlled experiments. Together, these results challenge dominant models of eye movements in reading by showing that linguistic prediction and parafoveal preview are not important determinants of word skipping.

INTRODUCTION

When reading a text, readers move their eyes across the page to bring new information to the centre of the visual field, where perceptual sensitivity is highest. While it may subjectively feel as if the eyes smoothly slide along the text, they in fact traverse the words with rapid jerky movements called *saccades*, followed by brief stationary periods called *fixations*. Across a text, saccades and fixations are highly variable and seemingly erratic: Some fixations last less than 100 ms, others more than 400; and while some words are fixated multiple times, many other words are skipped altogether [1, 2]. What explains this striking variation?

Historically, researchers have pointed to low-level non-linguistic factors like word length, oculomotor noise, or the relative position where the eyes happen

to land [2–5]. Such explanations were motivated by the idea that oculomotor control was largely *autonomous*. In this view, readers can adjust saccade lengths and fixation durations to global characteristics like text difficulty or reading strategy, but not to subtle word-by-word differences in language processing [2–4, 6].

As reading was studied in more detail, however, it became clear that the link between eye movements and cognition was more direct. For instance, it was found that fixation durations were shorter for words with higher frequency [7, 8]. Eye movements were even shown to depend on how well a word's identity could be inferred *before* fixation. Specifically, researchers found that words are read faster and skipped more often if they are *predictable* from linguistic context [9, 10] or if they are identifiable from a *parafoveal preview* [11–13].

These demonstrations of a direct link between eye movements and language processing overturned the autonomous view, replacing it by cognitive accounts describing eye movements during reading as largely, if not entirely, controlled by linguistic processing [14, 15]. Today, many studies still build on the powerful techniques like gaze-contingent displays that helped overturn the autonomous view, but now ask much more detailed questions, like whether word identification is a distributed or sequential process [16, 17]; how many words can be processed in the parafovea [18]; at which level they are analysed [19], and how this might differ between writing systems or orthographies [20, 21].

Here, we ask a different, perhaps more elemental question: how much of the variation in eye movements do linguistic prediction, parafoveal preview, and non-linguistic factors each explain? That is, how important are these factors for determining how the eyes move during reading? Dominant, cognitive models explain eye movement variation primarily as a function of ongoing processing. Skipping, for instance, is modelled as the probability that a word is identified before fixation [14, 22, 23]. Some, however, have questioned this purely cognitive view, suggesting that low-level features like word eccentricity or length might be more important [24–26]. Similarly, one may ask what drives next-word identification: is identifying the next word mostly driven by linguistic predictions [27] or by parafoveal perception? Remarkably, while it is well-established that both linguistic and oculomotor, and both predictive and parafoveal processing, all affect eye-movements [13, 24, 28, 29], a comprehensive picture of their relative explanatory power is currently missing, perhaps because they are seldom studied all at the same time.

To arrive at such a comprehensive picture we focus on natural reading, analysing three large datasets of participants reading passages, long articles, and even an entire novel – together encompassing 1.5 million (un)fixated words, across 108 individuals [30–32]. Instead of manipulating word predictability or perturb-

ing parafoveal perceptibility, we combine deep neural language modelling [33] and Bayesian ideal observer analysis [34] to quantify how much information is conveyed by both factors, on moment-by-moment basis. This way, we can probe the effect of both prediction and preview on *each* word during natural reading. Such a broad-coverage approach has been applied to the effects of predictability on reading before [29, 35–38], but either without considering preview or only through coarse heuristics such as using frequency as a proxy for parafoveal identifiability [17, 39, 40] (cf. [34]). By contrast, here we explicitly model both, in addition to low-level explanations like autonomous oculomotor control. To assess explanatory power, we use set theory to derive the unique and shared variation in eye movements explained by each model.

To preview the results, this revealed a striking dissociation between skipping and reading times. For word skipping, the overwhelming majority of variation could be explained – mostly *uniquely* explained – by a non-linguistic oculomotor model. For reading times, by contrast, we found strong effects of both prediction and preview, tightly matching effect sizes from controlled designs. Interestingly, linguistic prediction and parafoveal preview seem to operate independently: we found strong evidence against Bayes-optimal integration of the two. Together, these results challenge dominant cognitive models of reading, and show that skipping (or the decision of *where* to fixate) and reading times (i.e. *how long* to fixate) are governed by different principles.

RESULTS

We analysed eye movements from three large datasets of participants reading texts ranging from isolated paragraphs to an entire novel. Specifically, we considered three datasets: Dundee [32] (N=10, 51.502 words per participant), Geco [31] (N=14, 54.364 words per participant) and Provo [30] (N=84, 2.689 words per participant). In each corpus, we analysed both skipping

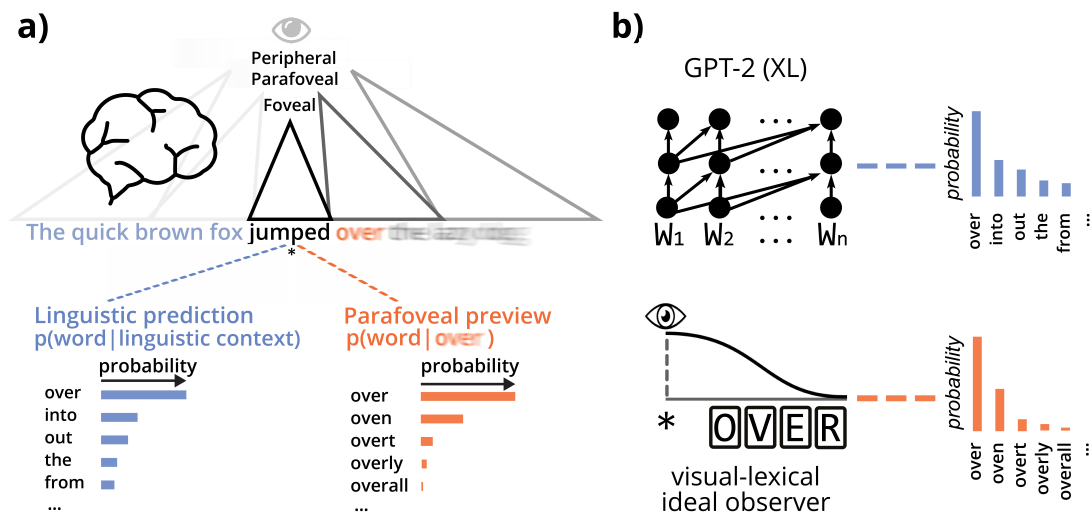


Figure 1: Quantifying two types of context during natural reading

a) Readers can infer the identity of the next word before fixation either by predicting it from context or by discerning it from the parafovea. Both can be cast as a probabilistic inference about the next word, either given the preceding words (prediction, blue) or given a parafoveal percept (preview, orange). **b)** To model prediction, we use GPT-2, one of the most powerful publicly available language models [33]. For preview, we use an ideal observer [34] based on well-established ‘Bayesian Reader’ models [41–43]. Importantly, we do not use either model as a cognitive model *per se*, but rather as a tool to quantify how much information is *in principle* available from prediction or preview on a moment-by-moment basis.

and reading times (indexed by gaze duration), as they are thought to reflect separate processes: the decision of *where* vs *how long* to fixate, respectively [14, 24].

To estimate the effect of linguistic prediction and parafoveal preview, we quantified the amount of information conveyed by both factors for each word in the corpus (for preview, this was tailored to each individual participant, since each word was previewed at a different eccentricity by each participant). To this end, we formalised both processes as a probabilistic belief about the identity of the next word, given either the preceding words (prediction) or a noisy parafoveal percept (preview; see Figure 1a). As such, we could describe these disparate cognitive processes using a common information-theoretic currency. To compute the probability distributions, we used GPT-2 for prediction [33] and a Bayesian ideal observer for preview [34] (see Figure 1b and *Methods*).

Prediction and preview increase skipping rates and reduce reading times

We first asked whether our formalisations allowed us to observe the expected effects of prediction and preview, while statistically controlling for oculomotor and lexical variables in a multiple regression model. Because the decisions of whether to skip and how long to fixate a word are made at different moments, we modeled each separately with a different set of explanatory variables; but for both, we considered the full model (detailed below).

As expected, we found in all datasets that words were more likely to be skipped if there was more information available from the linguistic prediction (Bootstrap: Dundee, $p = 0.023$; GECO, $p = 0.034$; Provo $p < 10^{-5}$) and/or the parafoveal preview (Bootstrap: Dundee, $p = 4 \times 10^{-5}$; GECO, $p < 10^{-5}$; Provo $p < 10^{-5}$). Similarly, reading times were reduced for words that

were more predictable (all $p's < 3.2 \times 10^{-4}$) or more identifiable from the parafovea (all $p's < 4 \times 10^{-5}$).

Together this confirms that our model-based approach can capture the expected effects of both prediction [15] and preview [13] in natural reading, while statistically controlling for other variables.

Skipping can be largely explained by non-linguistic oculomotor factors

After confirming that prediction and preview had a statistically significant influence on word skipping and reading times, we went on to assess their relative explanatory power. After confirming the effects of prediction and preview, we then further examined their relative explanatory power. That is, we asked the question how important these factors were, by examining how much variance was explained by each. To this end, we grouped the variables from the full regression model into different types of explanations, and assessed how well each type accounted for the data.

For skipping, we considered three explanations. First, a word might be skipped *purely* because it could be predicted from context – i.e. purely as a function of the amount of information conveyed by the prediction. Secondly, a word might be skipped because its identity could be gleaned from a parafoveal preview – that is, purely as a function of the informativeness of the preview. Finally, a word might be skipped simply because it is so short or so close to the prior fixation location that a saccade of average length will overshoot it, irrespective of its linguistic properties – in other words, purely as a function of length and eccentricity. Note that we did not include often-used lexical attributes like frequency to predict skipping, because using attributes of word_{*n*+1} already pre-supposes parafoveal identification. Moreover, to the extent that a lexical attribute like frequency might influence a word's parafoveal identifiability, this should already be captured by the parafoveal entropy (see Figure S3 and *Methods* for more details).

For each word, we thus modelled the probability of

skipping either as a function of prediction, preview, or oculomotor information, or by any combination of the three. Then we partitioned the unique and shared cross-validated variation explained by each account. Strikingly, this revealed that the overwhelming majority of explained skipping variation (94 %) could be accounted for by the non-linguistic baseline (Figure 2). Moreover, the majority of the variation was *only* explained by the baseline, which explained 10 times more unique variation than prediction and preview combined. There was a large degree of overlap between preview and the oculomotor baseline, which is unsurprising since a word's identifiability decreases as a function of its eccentricity and length. Interestingly, there was even more overlap between the prediction and baseline model: almost all skipping variation that could be explained by contextual constraint could be equally well explained by the oculomotor baseline factors.

Importantly, while the contribution of prediction and preview was small, it was significant both for prediction (bootstrap: Dundee, $p = 0.015$; Geco, $p = 0.0001$; Provo, $p < 10^{-5}$) and preview (all $p's < 5 \times 10^{-5}$), confirming that both factors do affect skipping. Crucially however, the vast majority of skipping that could be explained by either prediction or preview was equally well explained by the more parsimonious oculomotor model – which also explained much more of the skipping data overall.

Reading times are strongly modulated by prediction and preview

For reading times (operationalised through gaze durations, so considering foveal 'reading' only), we also considered three explanatory factors. First, a word might be read faster because it was predictable from the preceding context, which we formalised via linguistic surprise. Second, a word might be read faster if it could already be partly identified from the parafoveal preview (before fixation). This informativeness of the preview was again formalised via the parafoveal preview entropy. Finally, a word might be read faster due

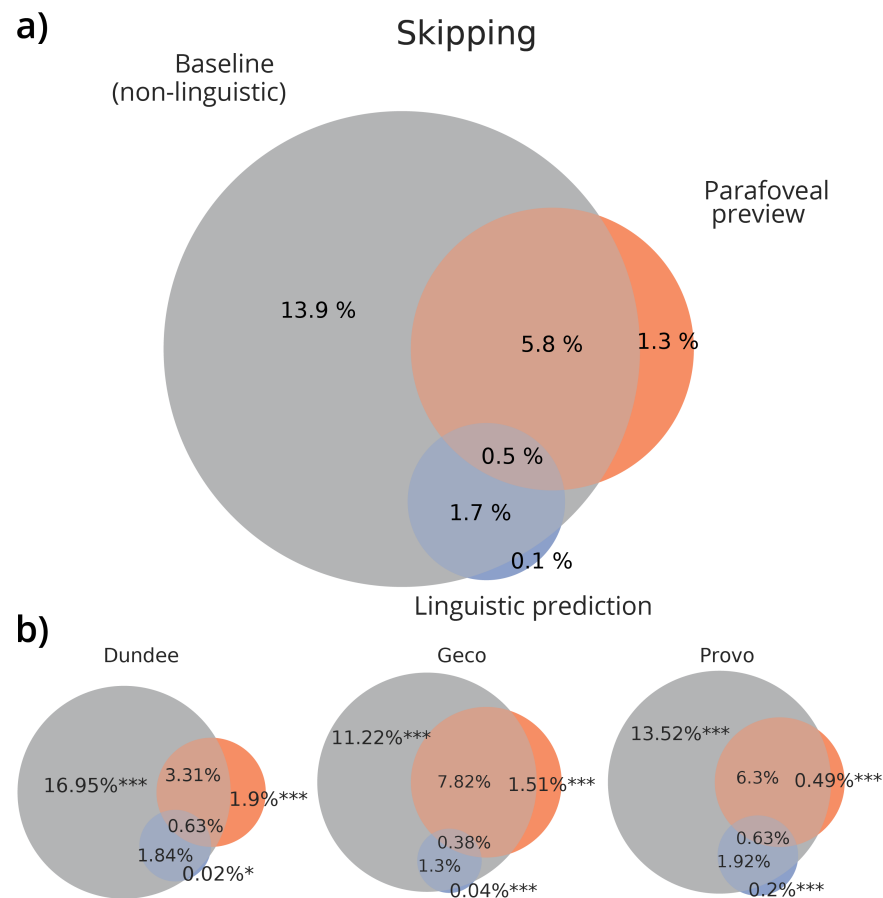


Figure 2: Variation in skipping explained by predictive, parafoveal and autonomous oculomotor processing

a) Proportions of cross-validated variation explained by prediction (blue), preview (orange) oculomotor baseline (grey) and their overlap; averaged across datasets (each dataset weighted equally). **b)** Variation partitions for each individual dataset, including statistical significance of variation uniquely explained by predictive, parafoveal or oculomotor processing. Stars indicate significance-levels of the cross-validated unique variation explained (bootstrap t-test against zero): $p < 0.05$ (*), $p < 0.05$ (**), $p < 0.001$ (***) For results of individual participants, and their consistency, see Figure S5.

to attributes of the word itself, such as lexical frequency. This last explanatory factor functioned as an aggregate baseline model that captured key non-contextual word attributes, both linguistic and non-linguistic (see Methods).

In all datasets, prediction (all $p's < 7.7 \times 10^{-3}$), preview (all $p's < 1.2 \times 10^{-4}$) and non-contextual word attributes (all $p's < 1.8 \times 10^{-4}$) again all explained significant unique variation. The non-contextual baseline explained the most variance, which shows – perhaps

unsurprisingly – that properties of the word itself are more important than contextual factors in determining how long a word is fixated. Critically however, compared to skipping the *unique* contribution of prediction and preview was more than three times higher (see Fig 3). Specifically, while prediction and preview could only uniquely account for 6% of explained word skipping variation, they uniquely accounted for more than 18 % of explained variation in reading times. Importantly, the *non-contextual* baseline used to predict reading times

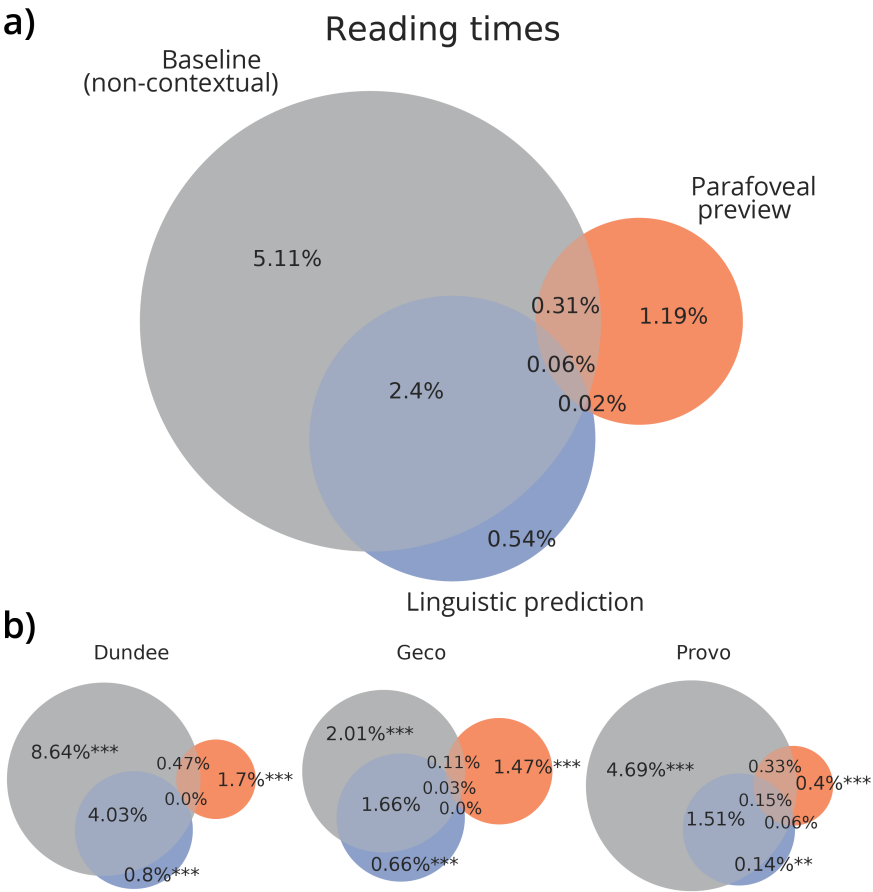


Figure 3: Variation in reading times explained by predictive, parafoveal and non-contextual information

a) Grand average of partitions of cross-validated variance in reading times (indexed by gaze durations) across datasets (each dataset weighted equally) explained by non-contextual factors (grey), parafoveal preview (orange), and linguistic prediction (blue). **b)** Variance partitions for each individual dataset, including statistical significance of the cross-validated variance explained uniquely by the predictive, parafoveal or non-contextual explanatory variables. Stars indicate significance-levels of the cross-validated unique variance explained (bootstrap t-test against zero): $p < 0.05$ (**), $p < 0.001$ (***). Note that the non-contextual model here included both lexical attributes (e.g. frequency) and oculomotor factors (relative viewing or landing position); assessing these separately reveals that reading time variation uniquely explained by oculomotor factors was small (see Fig S7). For results of individual participants, see Figure S6.

included both linguistic (e.g. lexical frequency) and non-linguistic information (viewing position) of the current word. When we analysed these separately, we found that the unique contribution of non-linguistic factors was small (see S7). This shows that contrary to skipping, variation in reading time is heavily influenced by online linguistic processing.

Naturalistic prediction and preview benefit effect match experimental effect sizes

The previous result confirms that reading times (indexed via gaze durations) are highly sensitive to linguistic and parafoveal context, in line with decades of research on eye movements in reading [44]. But how well do our results compare exactly to established find-

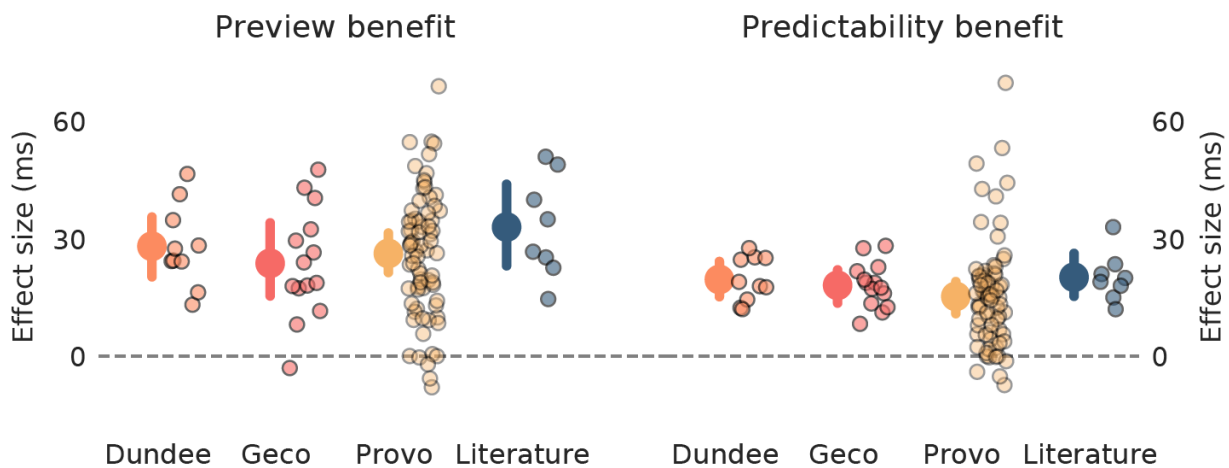


Figure 4: Simulated preview and predictability benefits match those reported in experimental literature

Preview (left) and predictability benefits (right) inferred from our analysis of each dataset, and observed in a sample of studies (see Table S1). In this analysis, preview benefit was simulated as the expected difference in gaze duration after a preview of average informativeness versus after no preview at all. Predictability benefit was defined as the difference in gaze duration for high versus low probability words; 'high' and 'low' were defined by subdividing the cloze probabilities from provo into equal thirds of 'low', 'medium' and 'high' probability (see Methods). In each plot, small dots with dark edges represent either individual subjects within one dataset or individual studies in the sample of the literature; larger dots with error bars represent the mean effect across individuals or studies, plus the bootstrapped 99% confidence interval.

ings from the experimental literature?

To directly address this question, we simulated, for each participant the effect size of two well-established effects that would be expected to be obtained if we would conduct a well-controlled factorial experiment. Specifically, because we estimated how much additional information from either prediction or preview (in bits) reduced reading times (in milliseconds) we could predict reading times for words that are expected vs unexpected (predictability benefit [28, 45]) or have valid vs invalid preview (i.e. preview benefit [13]).

The simulated effects tightly corresponded to those from experimental studies (see Fig 4). This shows that our analysis does not strongly underfit or otherwise underestimate the effect of either prediction or preview. Moreover, it shows that the effect sizes, which are well-established in controlled designs, generalise to natural reading. This is especially interesting for the preview benefit, because it implies that this effect can

be largely explained through parafoveal lexical identification, rather than visual preprocessing or interference effects (see Discussion).

No integration of prediction and preview

So far, we have treated prediction and preview as being independent. However, it might be that these processes, while using different information, are integrated – such that a word is parafoveally more identifiable when it is *also* more predictable in context. Bayesian probability theory proposes an elegant and mathematically optimal way to integrate these sources of information: the prediction of the next word could be incorporated as a prior in perceptual inference. Such a contextual prior fits into hierarchical Bayesian models of vision [46], and has been observed in speech perception, where a contextual prior guides the recognition of words from a partial sequence of phonemes [47]. Does such a prior also guide word recognition in reading,

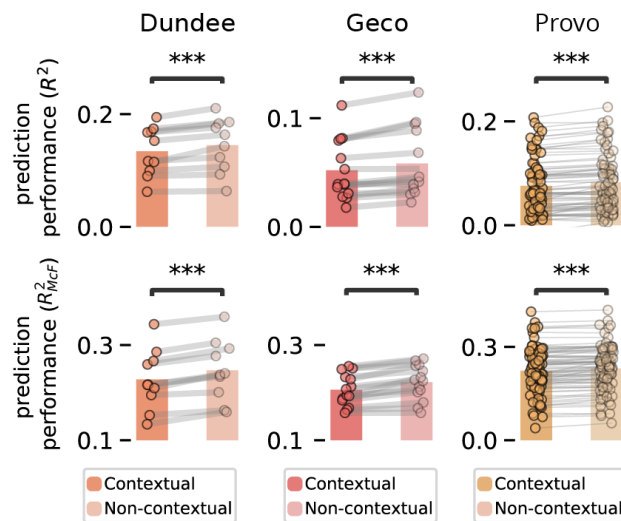


Figure 5: Evidence against bayesian integration of linguistic prediction and parafoveal preview

Cross-validated prediction performance of the full reading times (top) and skipping (bottom) model (including all variables), equipped with parafoveal preview information either from the contextual observer or from the non-contextual observer. Dots with connecting lines indicate participants; stars indicate significance: $p < 0.001$ (***).

based on a partial parafoveal percept?

To test this, we recomputed the parafoveal identifiability of each word for each participant, but now with an ideal observer using the prediction from GPT-2 as a prior. As expected, bayesian integration enhanced perceptual inference: on average, the observer using linguistic prediction as a prior extracted more information from the preview (± 6.25 bits) than the observer not taking the prediction into account (± 4.30 bits; $T_{1.39 \times 10^6} = 1.35 \times 10^{11}$, $p \approx 0$). Interestingly however, it provided a worse fit to the human reading data. This was established by comparing two versions of the full regression model: one with parafoveal entropy from the (theoretically superior) contextual ideal observer and one from the non-contextual ideal observer. In all datasets both skipping and reading times were better explained by a model including parafoveal identifi-

bility from the non-contextual observer (skipping: all $p's < 10^{-5}$; reading times: $p's < 10^{-5}$; see Figure 5).

Together, this suggests linguistic prediction and parafoveal preview are not integrated, but instead operate independently – thereby highlighting a remarkable sub-optimality in reading, and potentially an intriguing difference between visual and auditory word recognition.

DISCUSSION

Eye movements during reading are highly variable. Across three large datasets, we have assessed the relative importance of two major cognitive explanations for this variability – linguistic prediction and parafoveal preview – compared to alternative non-linguistic and non-contextual explanations. This revealed a stark dissociation between skipping and reading times. For word skipping, neither prediction nor preview were especially important, as the overwhelming majority of variation could be explained – mostly *uniquely* explained – by an oculomotor baseline model using just word length and eccentricity. For reading times, by contrast, we observed clear contributions of both prediction and preview – in addition to non-contextual features like frequency – and effect sizes matching those obtained in controlled experiments. Interestingly, preview effects were best captured by a non-contextual observer, suggesting that while readers use both linguistic prediction and preview, these do not appear to be integrated online. Together, the results underscore the dissociation between skipping and reading times, and show that for word skipping, the link between eye movements and cognition is less direct than commonly thought.

Our results on skipping align well with earlier findings by Drieghe and colleagues [24]. They analysed effect sizes from studies on skipping and found a disproportionately large effect of length, compared to proxies of processing-difficulty like frequency and predictability. We significantly extend their findings by modelling skipping itself (rather than effect sizes from studies) and

making a direct link to processing mechanisms. For instance, based on their analysis it was unclear how much of the length effect could be attributed to the decreasing visibility of longer words – i.e. how much of the length effect may be an identifiability effect [24, p. 19]. We show that length and eccentricity alone explained three times as much variation as parafoveal identifiability – and that most of the variation explained by identifiability was equally well explained by length and eccentricity. This demonstrates that length and eccentricity themselves – not just to the extent they reduce parafoveal identifiability – are key drivers of skipping.

This conclusion challenges dominant, cognitive models of eye movements, which describe lexical identification as the primary driver behind skipping [14, 22, 23]. However, it does not challenge the notion of predictive or parafoveal word identification itself. In fact, we believe this happens routinely – after all, most skips are not followed by regressions. Rather, our results challenge the notion that moment-to-moment decisions of whether to skip individual words are primarily determined by the identification of those words. Instead, they support a much simpler strategy, which is primarily sensitive to a word’s length and eccentricity.

One such simpler strategy would be a ‘blind’ random walk: making saccades of some average length, plus oculomotor noise. However, we do not think this is likely, since landing positions are distributed with preferred positions with respect to word boundaries [24, 48]. Instead, we suggest an alternative view, in which the decision of where to look next is based on an analysis of the parafovea – but at a very low level, aimed to discern mostly the next word’s length and eccentricity (see also [24, 25]). This is not the whole story, since preview and prediction explain some unique skipping variation that cannot be reduced to low-level variables (or other variables [34]). Our results may thus support a hybrid account, in which most skipping decisions are made by a low-level ‘autopilot’, whereas in some limited cases skipping is influenced by high-level contextual

information. How the brain arbitrates between these strategies is an interesting question for future research.

A distinctive feature of our approach is that we focus on a limited number of computationally explicit functional explanations, rather than using lexical attributes as proxies for functional explanations (e.g. a word’s frequency as a proxy for its identifiability). For instance, we model parafoveal identifiability using a single variable that should in principle capture all important effects such as that of frequency and orthography (see Figure S3 and *Methods*). A limitation of this approach is that an imperfection in the model can cause an underestimation of preview importance. However, a key advantage of using explicit modelling rather than proxies is that it can avoid confound-related misinterpretations. For instance, word frequency is strongly correlated with length, so when using frequency as a proxy for identifiability (e.g. to predict skipping), one may find apparent identifiability effects which are in fact length effects, and strongly overestimate the importance of preview [49]. Therefore, we have only used explanatory variables that explicitly relate to the dependent variable (*Methods*). After all, our goal was not to measure as many effects as possible, but to gain a clear picture of the importance of two cognitive explanations for eye movement variation. Based on the effect sizes for gaze duration (Fig 4) we do not believe that we strongly underestimate either prediction or preview, and we are optimistic the results provide the comprehensive, interpretable picture we aimed for.

When comparing Figures 2, 3 and 5, the numerical R^2 values of the reading times regression may seem rather small, potentially indicating a poor fit. However, our (cross-validated) R^2 s for gaze durations are not lower than R^2 s reported by other regression analyses of gaze durations in natural reading [e.g. 17]; moreover we find effect sizes in line with the experimental literature (Fig 4). Therefore, we do not believe we either overfit or underfit the gaze durations. Instead, what the relatively low R^2 values indicate, we suggest, is that gaze

433 durations are inherently noisy, and that only a limited
434 amount of the word-by-word variation is *systematic* vari-
435 ation, due to e.g. preview or frequency effects. While
436 this is interesting in itself, it is not of primary interest in
437 this study, which focusses on the relative importance
438 of different *explanations*, and hence only on systematic
439 variation. Therefore, what matters is not as much the
440 absolute R^2 values, but rather the relative importance
441 of different explanations – in other words, the relative
442 size of the circles in Figures 2, 3 and S7, their overlap,
443 and the explanations each circle represents. It is on
444 this level of analysis that we find the stark dissociation
445 – that for skipping (but not for reading times) a simple
446 low-level heuristic can account for almost all of the ex-
447 plained variation – and not on the level of numerical
448 values for variation explained.

449 A remarkable result is that we found preview bene-
450 fits comparable to effect sizes from controlled designs
451 (Figure 4), despite major methodological differences.
452 In controlled designs preview benefits are defined as
453 the difference in reading time for words with preview,
454 versus words where the preview was masked or made
455 invalid (i.e. where a different word was previewed than
456 subsequently perceived at fixation). As such, it seemed
457 *a priori* plausible that a significant portion of preview
458 benefits observed in controlled studies might reflect
459 interference or mismatch between the preview (or the
460 mask) and the subsequent percept, rather than reflect-
461 ing purely the lack of parafoveal identification of a word.
462 Our analysis modelled the preview benefit purely in
463 terms of lexical identification, and yielded only slightly
464 smaller effect sizes (Fig. 4). This may suggest that pre-
465 view benefits may largely reflect lexical parafoveal iden-
466 tification, and that interference or visual ‘preprocessing’
467 may only play a minor role [compare 13, 14].

468 Another notable finding is that preview was best ex-
469 plained by a non-contextual observer – a model which
470 only takes word frequency (and not contextual pre-
471 dictability) into account. This replicates and extends the
472 only study that explicitly compared contextual and non-

473 contextual accounts of parafoveal preview [34]. That
474 study only analysed skipping (in the Dundee corpus);
475 the fact that we find the same for reading times (where
476 preview and prediction effects are stronger) and repli-
477 cate the result in other corpora, considerably strength-
478 ens the conclusion that parafoveal word recognition
479 is not informed by linguistic context. This conclusion
480 seems to contradict experiments finding an interaction
481 between linguistic context and preview, which was inter-
482 preted as context constraining preview [e.g. 9, 50–52].
483 One explanation for this discrepancy stems from how
484 the effect is measured. Experimental studies did not
485 explicitly model contextual and non-contextual recog-
486 nition, but looked at the effect of context on the dif-
487 ference in reading time after valid versus invalid pre-
488 view [51, 52]. This may reveal a context effect not
489 on recognition, but at a later stage (e.g. priming be-
490 tween the context, preview and fixated word). Arguably,
491 these scenarios yield different predictions: if context af-
492 fects recognition it may allow identification of otherwise
493 unidentifiable words. However, if the interaction occurs
494 later it might only *amplify* processing of recognisable
495 words. Constructing a model that formally reconciles
496 this discrepancy – and predicts the context-preview in-
497 teraction using a non-contextual prior – is an interesting
498 challenge for future work.

499 The lack of influence of contextual constraint on
500 parafoveal preview might be specific to parafoveal pre-
501 view, perhaps due to time-constraints imposed by eye
502 movements. Given that readers on average only look
503 some 250 ms at a word in which they have to both
504 recognise the foveal word and process the parafoveal
505 percept, this perhaps leaves too little time to fully let
506 the foveal word and context inform parafoveal percep-
507 tion. On the other hand, word recognition based on
508 partial information also occurs in speech perception,
509 where it also occurs under significant time-constraints.
510 And yet in auditory word recognition, contextual effects
511 are found [53, 54], and a formally highly similar analysis
512 of word recognition based on partial phonemic infor-

mation recently showed clear support for a contextual prior; i.e. the exact opposite of what we find here [47]. An alternative, more speculative explanation for the lack of context effect in reading but not speech perception is that this may reflect a difference between visual and auditory word recognition. This could be related to the fact that contrary to auditory word recognition, visual word recognition is an acquired skill and occurs throughout areas in the visual system repurposed for reading [55, 56], where perhaps the dynamic sentence context cannot exert as much of an influence as rapidly, allowing for facilitation by lexical or orthographic context [57–60], but not as much of sentence context.

Given that readers use both prediction and preview, why would they strongly affect reading times but hardly word skipping? To understand this dissociation, it is important to consider that they reflect different decisions, namely *where* versus *how long* to fixate, which are made at different moments. Specifically, the decision of where to fixate – and hence whether to skip the next word – is made early in saccade programming, which can take 100–150 ms [24, 44, 61]. Although the exact sequence of operations leading to a saccade remains debated, given that readers on average only look some 250 ms at a word, it is clear that skipping decisions are made under strong time constraints, especially given the lower processing rate of parafoveal information. Our results suggest that the brain meets this constraint by resorting to a computationally frugal policy, largely based on low-level characteristics such as length and eccentricity. *How long* to fixate, by contrast, mostly depends on foveal information, which is processed more rapidly and may thus directly influence the decision to either keep dwelling and accumulate more information or initiate a saccade (and/or an attention shift).

In conclusion, we have found that two important contextual sources of information in reading, linguistic prediction and parafoveal preview, strongly drive variation in reading times, but hardly affect word skipping, which is largely based on low-level factors. Our

results show that as readers, we do not always use all information available to us; and that we are, in a sense, of two minds: consulting complex inferences to decide how long to look at a word, while employing semi-mindless scanning routines to decide where to look next. It is striking that these disparate strategies operate mostly in harmony. Only occasionally they go out of step – then we notice that our eyes have moved too far and we have to look back, back to where our eyes left cognition behind.

METHODS

We analysed eye-tracking data from three, big, naturalistic reading corpora, in which native English speakers read texts while eye-movement data was recorded [31, 32, 38].

Eye-tracking data and stimulus materials

We considered the English-native portions of the Dundee, Geco and Provo corpora. The Dundee corpus comprises eye-movements from 10 native speakers from the UK ([32]), who read a total of 56,212 words across 20 long articles from The Independent newspaper. Secondly, the English portion of the Ghent Eye-tracking Corpus (Geco) [31] is a collection of eye movement data from 14 UK English speakers who each read Agathe Christie's *The Mysterious Affair at Styles* in full (54,364 words per participant). Lastly, the Provo corpus ([30]) is a collection of eye movement data from 84 US English speakers, who each read a total of 55 paragraphs (extracted from diverse sources) for a total of 2,689 words.

Language model

Contextual predictions were formalised using a language model – a model computing the probability of each word given the preceding words. Here, we used GPT-2 (XL) – currently among the best publicly released English language models. GPT-2 is a transformer-based model, that in a single pass turns a sequence of tokens (representing either whole words or word-parts) $U = (u_1, \dots, u_k)$ into a sequence of conditional probabilities, $(p(u_1), p(u_2|u_1), \dots, p(u_k | u_1, \dots, u_{k-1}))$.

Roughly, this happens in three steps: first, an embedding step encodes the sequence of symbolic tokens as a sequence

of vectors, which can be seen as the first hidden state h_o . Then, a stack of n transformer blocks each apply a series of operations resulting in a new set of hidden states h_l , for each block l . Finally, a (log-)softmax layer is applied to compute (log-)probabilities over target tokens. In other words, the model can be summarised as follows:

$$h_0 = UW_e + W_p \quad (1)$$

$$h_l = \text{transformer_block}(h_{l-1}) \forall l \in [1, n] \quad (2)$$

$$P(u) = \text{softmax}(h_n W_e^T), \quad (3)$$

where W_e is the token embedding and W_p is the position embedding.

The key component of the transformer-block is *masked multi-headed self-attention* (Fig S1). This transforms a sequence of input vectors $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$ into a sequence of output vectors $(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k)$. Fundamentally, each output vector \mathbf{y}_i is simply a weighted average of the input vectors: $\mathbf{y}_i = \sum_{j=1}^k w_{ij} \mathbf{x}_j$. Critically, the weight w_{ij} is not a parameter, but is *derived* from a dot product between the input vectors $\mathbf{x}_i^T \mathbf{x}_j$, passed through a softmax and scaled by a constant determined by the dimensionality d_k : $w_{ij} = (\exp \mathbf{x}_i^T \mathbf{x}_j / \sum_{j=1}^k \exp \mathbf{x}_i^T \mathbf{x}_j) \frac{1}{\sqrt{d_k}}$. Because this is done for each position, each input vector \mathbf{x}_i is used in three ways: first, to derive the weights for its own output, \mathbf{y}_i (as the *query*); second, to derive the weight for any other output \mathbf{y}_j (as the *key*); finally, in it used in the weighted sum (as the *value*). Different linear transformations are applied to the vectors in each cases, resulting in Query, Key and Value matrices (Q, K, V). Putting this all together, we obtain:

$$\text{self_attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V. \quad (4)$$

To be used as a language model, two elements are added. First, to make the operation position-sensitive, a position embedding W_p is added during the embedding step – see Equation (1). Second, to enforce that the model only uses information from the past, attention from future vectors is masked out. To give the model more flexibility, each transformer block contains multiple instances ('heads') of the self-attention mechanisms from Equation (4).

In total, GPT-2 (XL) contains $n = 48$ blocks, with 12 heads each; a dimensionality of $d = 1600$ and a context window of $k = 1024$, yielding a total 1.5×10^9 parameters. We used the PyTorch implementation of GPT-2 provided by HuggingFace's *Transformers* package [62]. For words spanning multiple tokens, we computed their joint probability.

Ideal observer

To compute parafoveal identifiability, we implemented an ideal observer based on the formalism by Duan & Bicknell [34]. This model formalises parafoveal word identification using Bayesian inference and builds on previous well-established 'Bayesian Reader' models [41–43]. It computes the probability of the next word given a noisy percept by combining a prior over possible words with a likelihood of the noisy percept, given a word identity:

$$p(w | \mathcal{I}) \propto p(w)p(\mathcal{I}|w), \quad (5)$$

where \mathcal{I} represents the noisy visual input, and w represents a word identity. We considered two priors (see Fig 5): a non-contextual prior (the overall probability of words in English based on their frequency in Subtlex ([63]), and a contextual prior based on GPT2 (see below). Below we describe how visual information is represented and perceptual inference is performed. For a graphical schematic of the model, see Fig S2; for some distinctive simulations showing how the model captures key effects of linguistic and visual characteristics on word recognition, see Fig S3.

Sampling visual information

Like in other Bayesian Readers [41–43], noisy visual input is accumulated by sampling from a multivariate Gaussian which is centered on a one-hot 'true' letter vector – here represented in an uncased 26-dimensional encoding – with a diagonal covariance matrix $\Sigma(\epsilon) = \lambda(\epsilon)^{-1/2} I$. The shape of Σ is thus scaled by the sensory quality $\lambda(\epsilon)$ for a letter at eccentricity ϵ . Sensory quality is computed as a function of the perceptual span: this uses using a Gaussian integral based follows the perceptual span or processing rate function from the SWIFT model [22]. Specifically, for a letter at eccentricity ϵ , λ is given by the integral within the bounding box of the letter:

$$\lambda(\epsilon) = \int_{\epsilon-.5}^{\epsilon+.5} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx, \quad (6)$$

which, following [34, 43], is scaled by a scaling factor Λ . Unlike SWIFT, the Gaussian in Equation 6 is symmetric, since we only perform inference on information about the next word. By using one-hot encoding and a diagonal covariance matrix, the ideal observer ignores similarity structure between letters. This is clearly a simplification, but one with significant computational benefits; moreover, it is a simplification shared by all Bayesian Reader-like models [34, 41, 43], which can nonetheless capture many important aspects of visual word recognition and reading. To determine parameters Λ and σ , we performed a grid search on a subset of Dundee and Geco (see Fig S4), resulting in $\Lambda = 1$ and $\sigma = 3$. Note that this σ value is close to the average σ value of SWIFT and (3.075) and corresponds well to prior literature on the size of the perceptual span (± 15 characters; [13, 22, 43]).

Perceptual inference

Inference is performed over the full vocabulary. This is represented as a matrix which can be seen as a stack of word vectors, $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_v$, obtained by concatenating the letter vectors. The vocabulary is thus a $V \times d$ matrix, with V the number of words in the vocabulary and d the dimensionality of the word vectors (determined by the length of the longest word: $d = 26 \times l_{max}$).

To perform inference, we use the belief-updating scheme from [34], in which the posterior at sample t is expressed as a $(V - 1)$ dimensional log-odds vector $\mathbf{x}^{(t)}$, in which each entry $\mathbf{x}_i^{(t)}$ represents the log-odds of \mathbf{y}_i relative to the final word \mathbf{y}_v . In this formulation, the initial value of \mathbf{x} is thus simply the prior log odds, $\mathbf{x}_i^{(0)} = \log p(w_i) - \log p(w_v)$, and updating is done by summing prior log-odds and the log-odds likelihood. This procedure is repeated for T samples, each time taking the posterior of the previous timestep as the prior in the current timestep. Note that using log odds in this way avoids renormalization:

$$\begin{aligned} \mathbf{x}_i^{(t)} &= \log \frac{p(w_i | \mathcal{I}^{(0, \dots, t)})}{p(w_v | \mathcal{I}^{(0, \dots, t)})} \\ &= \log \frac{p(w_i | \mathcal{I}^{(0, \dots, t-1)}) p(\mathcal{I}^{(t)} | w_i)}{p(w_v | \mathcal{I}^{(0, \dots, t-1)}) p(\mathcal{I}^{(t)} | w_v)} \\ &= \log \frac{p(w_i | \mathcal{I}^{(0, \dots, t-1)})}{p(w_v | \mathcal{I}^{(0, \dots, t-1)})} + \log \frac{p(\mathcal{I}^{(t)} | w_i)}{p(\mathcal{I}^{(t)} | w_v)} \\ &= \mathbf{x}_i^{(t-1)} + \Delta \mathbf{x}_i^{(t)}. \end{aligned} \quad (7)$$

In other words, as visual sample $\mathcal{I}^{(t)}$ comes in, beliefs are updated by summing the prior log odds $\mathbf{x}^{(t-1)}$ and the log-odds likelihood of the new information $\mathbf{x}^{(t)}$.

For a given word w_i , the log-odds likelihood of each new sample is the difference of two multivariate Gaussian log likelihoods, one centred on \mathbf{y}_i and one on the last vector \mathbf{y}_v . This can be formulated as a linear transformation of \mathcal{I} :

$$\begin{aligned} \Delta \mathbf{x}_i &= \log p(\mathcal{I} | w_i) - \log p(\mathcal{I} | w_v) \\ &= \log p(\mathcal{I} | \mathcal{N}(\mathbf{y}_i, \Sigma)) - \log p(\mathcal{I} | \mathcal{N}(\mathbf{y}_v, \Sigma)) \\ &= \left[-\frac{1}{2} (\mathcal{I} - \mathbf{y}_i)^T \Sigma^{-1} (\mathcal{I} - \mathbf{y}_i) \right] - \\ &\quad \left[-\frac{1}{2} (\mathcal{I} - \mathbf{y}_v)^T \Sigma^{-1} (\mathcal{I} - \mathbf{y}_v) \right] \\ &= \frac{\mathbf{y}_v^T \Sigma^{-1} \mathbf{y}_v - \mathbf{y}_i^T \Sigma^{-1} \mathbf{y}_i}{2} + (\mathbf{y}_i - \mathbf{y}_v)^T \Sigma^{-1} \mathcal{I}, \end{aligned} \quad (8)$$

which implies that updating can be implemented by sampling from a multivariate normal. To perform inference on a given word, we performed this sampling scheme until convergence (using $T = 50$), and then transformed the posterior log-odds into the log posterior, from which we computed the Shannon entropy as a metric of parafoveal identifiability.

To compute the parafoveal entropy for each word in the corpus, we make the simplifying assumption that parafoveal preview only occurs during the last fixation prior to a saccade, thus computing the entropy as a function of the word itself and its distance to the last fixation location within the previously fixated word (which is not always the previous word). Because this distance is different for each participant, it was computed separately for each word, for each participant. Moreover, because the inference scheme is based on sampling, we repeated it 3 times, and averaged these to compute the posterior entropy of the word. The amount of information obtained from the preview is then simply the difference between prior and posterior entropy.

The ideal observer was implemented in custom Python code, and can be found in the data sharing collection (see below).

Contextual vs non-contextual prior

We considered two observers: one with a non-contextual prior capturing the overall probability of a word in a language,

and with a contextual prior, capturing the contextual probability of a word in a specific context. For the non-contextual prior, we simply used lexical frequencies from which we computed the (log)-odds prior used in equation (7). For the contextual prior, we derived the contextual prior from log-probabilities from GPT-2. This effectively involves constructing a new Bayesian model for each word, for each participant, in each dataset. To simplify this process, we did not take the full predicted distribution of GPT-2, but only the ‘nucleus’ of the top k predicted words with a cumulative probability of 0.95, and truncated the (less reliable) tail of the distribution. Further, we simply assumed that the rest of the tail was ‘flat’ and had a uniform probability. Since the prior odds can be derived from relative frequencies, we can think of the probabilities in the flat tail as having a ‘pseudocount’ of 1. If we similarly express the prior probabilities in the nucleus as implied ‘pseudofrequencies’, the cumulative implied nucleus frequency is then complementary to the length of the tail, which is simply the difference between the vocabulary size and nucleus size ($V - k$). As such, for word i in the text, we can express the nucleus as implied frequencies as follows:

$$\text{freqs}_\psi = P_{tr}(w^{(i)}|\text{context}) \frac{V - k}{1 - \sum_{j=1}^k P(w_j^{(i)}|\text{context})}, \quad (9)$$

where $P_{tr}(w^{(i)}|\text{context})$ is the truncated lexical prediction, and $P(w_j^{(i)}|\text{context})$ is predicted probability that word i in the text is word j in the sorted vocabulary. Note that using this flat tail not only simplifies the computation, but also deals with the fact that the vocabulary of GPT-2 is smaller than of the ideal observer – using this tail we can still use the full vocabulary (e.g. to capture orthographic uniqueness effects), while using 95% of the density from GPT-2.

Data selection

In all our analyses, we focus strictly on first-pass reading, analysing only those fixations or skips when none of the subsequent words have been fixated before. We extensively preprocessed the corpora so that we could include as many words as possible. However, we had to impose some additional restrictions. Specifically we did not include words if they a) contained non-alphabetic characters; b) if they were adjacent to blinks; c) if the distance to the prior fixation location was more than 24 characters ($\pm 8^\circ$); moreover, for the gaze duration we excluded d) words with implausibly short

(< 70ms) or long (> 900ms) gaze durations. Criterion c) was chosen because some participants occasionally skipped long sequences of words, up to entire lines or more. Such ‘skipping’ – indicated by saccades much larger than the perceptual span – is clearly different from the skipping of words during normal reading, and was therefore excluded. Note that these criteria are comparatively mild (cf. [34, 35]), and leave approximately 1.1 million observations for the skipping analysis, and 593,000 reading times observations.

Regression models: skipping

Skipping was modelled via logistic regression in scikit-learn [64], with three sets of explanatory variables (or ‘models’) each formalising a different explanation for why a word might be skipped.

First, a word might be skipped because it could be confidently predicted from context. We formalise this via *linguistic entropy*, quantifying the information conveyed by the prediction from GPT-2. We used entropy, not (log) probability, because using the next word’s probability directly would presuppose that the word is identified, undermining the dissociation of prediction and preview. By contrast, prior entropy specifically probes the information available from prediction only.

Secondly, a word might be skipped because it could be identified from a parafoveal preview. This was formalised via parafoveal entropy, which quantifies the parafoveal preview uncertainty (or, inversely, the amount of information conveyed by the preview). This is a complex function integrating low-level visual (e.g. decreasing visibility as a function of eccentricity) and higher-level information (e.g. frequency or orthographic effects) and their interaction (see Fig S3). Here, too we did not use lexical features (e.g. frequency) of the next word to model skipping directly, as this presupposes that the word is identified; and to the extent that these factors are expected to influence identifiability, this is already captured by the parafoveal entropy (Fig S3).

Finally, a word might be skipped simply because it is too short and/or too close to the prior fixation location, such that a fixation of average length would overshoot the word. This autonomous oculomotor account was formalised by modelling skipping probability purely as a function of a word’s length and its distance to the previous fixation location.

Note that these explanations are not mutually exclusive, so we also evaluated their combinations (see below).

Regression models: reading time

As an index of reading time, we analysed first-pass *gaze duration*, the sum of a word's first-pass fixation durations. We analyse gaze durations as they arguably most comprehensively reflect how long a word is looked at, and are the focus of similar model-based analyses of contextual effects in reading [35, 37]. For reading times, we used linear regression, and again considered three sets of explanatory variables, each formalising a different kind of explanation.

First, a word may be read more slowly because it is unexpected in context. We formalised this using surprisal – $\log(p)$, a metric of a word's unexpectedness – or how much information is conveyed by a word's identity in light of a prior expectation about the identity. To capture spillover (R; regpaper; smith) we included not just the surprisal of the current word, but also that of the previous two words.

Secondly, a word might be read more slowly because it was difficult to discern from the parafoveal preview. This was formalised using the parafoveal entropy (see above).

Finally, a word might be read more slowly because of non-contextual factors of the word itself. This is an aggregate baseline explanation, aimed to capture all relevant non-contextual word attributes, which we contrast to the two major contextual sources of information about a word identity that might affect reading times (prediction and preview). We included word class, length, log-frequency, and the relative landing position (quantified as the distance to word centre in characters. For log-frequency we used the UK or US version of SUBTLEX depending on the corpus and included the log-frequency of the past two words to capture spillover effects.

Note that, while for skipping, we used a *non-linguistic* baseline, for reading times we use a *non-contextual* baseline. This is because for skipping the most interesting contrast is between the role of non-linguistic oculomotor control vs an account that explains skipping via ease of next-word identification (either through prediction or preview). For reading times, by contrast, the most interesting comparison is between properties of the word itself versus contextual cues, as a purely non-linguistic account for gaze duration variation

seemed less plausible (indeed, see Figure S7 for a supplementary analysis confirming that the limited relative importance of a purely non-linguistic account for reading time variation).

Model evaluation

We compared the ability of each model to account for the variation in the data by probing prediction performance in a 10-fold cross-validation scheme, in which we quantified how much of the observed variation in skipping rates and gaze durations could be explained.

For reading times, we did this using the coefficient of determination, defined via the ratio of residual and total sum of squares: $R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$. The ratio $\frac{SS_{res}}{SS_{tot}}$ relates the error of the model (SS_{res}) to the error of a 'null' model predicting just the mean (SS_{tot}), and gives the variance explained. For skipping, we use a tightly related metric, the McFadden R^2 . Like the R^2 it is computed by comparing the error of the model to the error of a null model with only an intercept: $R_{McF}^2 = 1 - \frac{LM}{L_{null}}$, where L indicates the loss.

While R^2 and R_{McF}^2 are not identical, they are formally tightly related – critically, both are zero when the prediction is constant (no variation explained) and go towards one proportionally as the error decreases to zero (i.e. towards all variation explained). Note that in a cross-validated setting, both metrics can become negative when prediction of the model is worse than the prediction of a constant null-model.

Variation partitioning

To assess relative explanatory power, we used variation partitioning to estimate how much of the explained variation could be attributed to each set of explanatory variables. This is also known as *variance* partitioning, as it is originally based on partitioning sums of squares in regression analysis; here we use the more general term 'variation' following [65].

Variation partitioning builds on the insight that when two (groups of) explanatory variables (A and B) both explain some variation in the data y , and A and B are independent, then variation explained by combining A and B will be approximately additive. By contrast, when A and B are fully redundant – e.g. when B only has an *apparent* effect on y through its correlation with A – then a model combining A and B will not explain more than the two alone.

Following [66], we generalise this logic to three (groups of) explanatory variables, by testing each individually and all combinations, and use set theory notation and graphical representation for its simplicity and clarity. For three groups of explanatory variables (A , B and C), we first evaluate each separately and all combinations, resulting in 7 models:

$$A, B, C, A \cup B, A \cup C, B \cup C, A \cup B \cup C.$$

From these 7 models we obtain 7 ‘empirical’ scores (expressing variation explained), from which we derive the 7 ‘theoretical’ variation partitions: 4 overlap partitions and 3 unique partitions. The first overlap partition is the variation explained by all models, which we can derive as:

$$A \cap B \cap C = A \cup B \cup C + A + B + C - A \cup B - A \cup C - B \cup C. \quad (10)$$

The next three overlap partitions contain all pairwise intersections of models that did not include the other model:

$$\begin{aligned} (A \cap B) \setminus C &= A + B - A \cup B - A \cap B \cap C \\ (A \cap C) \setminus B &= A + C - A \cup C - A \cap B \cap C \\ (B \cap C) \setminus A &= B + C - B \cup C - A \cap B \cap C. \end{aligned} \quad (11)$$

The last three partitions are those explained exclusively by each model. This is the relative complement: the partition unique to A is the relative complement of BC: BC^{RC} . For simplicity we also use a star notation, indicating the unique partition of A as A^* . These are derived as follows:

$$\begin{aligned} A^* &= BC^{RC} = A \cup B \cup C - B \cup C \\ B^* &= AC^{RC} = A \cup B \cup C - A \cup C \\ C^* &= AB^{RC} = A \cup B \cup C - A \cup B. \end{aligned} \quad (12)$$

Note that, in the cross-validated setting, the results can become paradoxical and depart from what is possible in classical statistical theory, such as partitioning sums of squares. For instance, due to over-fitting, a model that combines multiple EVs could explain *less* variance than all of the EVs alone, in which case some partitions would become negative. However, following [66], we believe that the advantages of using cross-validation outweigh the risk of potentially paradoxical results in some subjects. Partitioning was carried out for each subject, allowing to statistically assess whether the additional variation explained by a given model was significant. On average, none of the partitions were paradoxical.

Simulating effect sizes

Preview benefits were simulated as the expected difference in gaze duration after a preview of average informativeness versus after no preview at all (...). This this best corresponds to an experiment in which the preceding preview was masked (e.g. XXXX) rather than invalid (see Discussion). To compute this we compared the took the difference in parafoveal entropy between an average preview and the prior entropy. Because we standardised our explanatory variables, this was transformed to subject-specific z-scores and then multiplied by the regression weights to obtain an expected effect size.

For the predictability benefit, we computed the expected difference in gaze duration between ‘high’ and ‘low’ probability words. ‘High’ and ‘low’ was empirically defined based on the human-normed cloze probabilities in *provo*, which we divided into thirds using percentiles. The resulting cutoff points (low < 0.02; high > 0.25) were log-transformed, applied to the surprisal values from GPT-2, and multiplied by the weights to predict effect sizes. Note that these definitions of ‘low’ and ‘high’ may appear low compared to those in literature – however, most studies collect cloze only for specific ‘target’ words in relatively predictable contexts, which biases the definition of ‘low’ vs ‘high’ probability. By contrast, we analysed cloze probabilities for *all* words, yielding these values.

Statistical testing

Statistical testing was performed across participants within each dataset. Because two of the three corpora had a low number of participants (10 and 14 respectively) we used non-parametric bootstrap t-tests, by creating resampling a null-distribution with zero mean counting how likely a t-value at least as extreme as the true t-value was to occur. Each test used at least 10^4 bootstraps; p-values were computed without assuming symmetry (equal-tail bootstrap).

Data and code availability

The *Provo* and *Geco* corpora are freely available ([30, 31]). All additional data and code needed to reproduce these results will be made public at the Donders Data Repository.

Contributions

Conceptualisation: MH. Data wrangling and preprocessing: JvH. Formal analysis: MH, JvH. Statistical analysis and visu-

962 alisations: MH, JvH. Supervision: FPdL, PH. Initial draft: MH.
963 Final draft: MH, JvH, PH, FPdL.

964 Acknowledgements

965 We thank Maria Barrett, Yunyan Duan, and Benedikt
966 Ehinger for useful input and inspiring discussions during vari-
967 ous stages of this project, and Ashley Lewis for helpful com-
968 ments on an earlier version of this manuscript. This work was
969 supported by The Netherlands Organisation for Scientific Re-
970 search (NWO Research Talent grant to M.H.; NWO Vidi grant
971 to F.P.d.L.; Gravitation Program Grant Language in Interaction
972 no. 024.001.006 to P.H.) and the European Union Horizon
973 2020 Program (ERC Starting Grant 678286 to F.P.d.L).

974 REFERENCES

- 975 1. Rayner, K. & Pollatsek, A. in *Attention and per-*
976 *formance 12: The psychology of reading* 327–362
977 (Lawrence Erlbaum Associates, Inc, Hillsdale, NJ,
978 US, 1987). ISBN: 978-0-86377-083-8 978-0-86377-
979 084-5.
- 980 2. Dearborn, W. F. *The psychology of reading: an exper-*
981 *imental study of the reading pauses and movements*
982 *of the eye ...* <<https://catalog.hathitrust.org/Record/000359636>> (visited on 05/21/2021)
983 (Archives of philosophy, psychology and scientific
984 methods, no. 4., n. p., 1906).
- 985 3. Bouma, H. & Voogd, A. H. d. On the control of
986 eye saccades in reading. English. *Vision Research*
987 **14**. Publisher: Elsevier, 273–284. ISSN: 0042-6989
988 (1974).
- 989 4. Buswell, G. T. *An experimental study of the eye-voice*
990 *span in reading* **17** (University of Chicago, 1920).
- 991 5. O'Regan, J. K. The control of saccade size and fix-
992 ation duration in reading: The limits of linguistic
993 control. en. *Perception & Psychophysics* **28**, 112–
994 117. ISSN: 1532-5962 (March 1980).
- 995 6. Morton, J. The Effects of Context upon Speed of
996 Reading, Eye Movements and Eye-voice Span. en.
997 *Quarterly Journal of Experimental Psychology* **16**.
998 Publisher: SAGE Publications, 340–354. ISSN: 0033-
999 555X (December 1964).
- 1000 7. Rayner, K. Visual attention in reading: Eye move-
1001 ments reflect cognitive processes. en. *Memory &*
1002 *Cognition* **5**, 443–448. ISSN: 1532-5946 (July 1977).
- 1003 8. Inhoff, A. W. Two stages of word processing during
1004 eye fixations in the reading of prose. en. *Journal of*
1005 *Verbal Learning and Verbal Behavior* **23**, 612–624.
1006 ISSN: 0022-5371 (October 1984).
- 1007 9. Balota, D. A., Pollatsek, A. & Rayner, K. The inter-
1008 action of contextual constraints and parafoveal
1009 visual information in reading. en. *Cognitive Psychol-*
1010 *ogy* **17**, 364–390. ISSN: 0010-0285 (July 1985).
- 1011 10. Ehrlich, S. F. & Rayner, K. Contextual effects on
1012 word perception and eye movements during read-
1013 ing. en. *Journal of Verbal Learning and Verbal Be-*
1014 *havior* **20**, 641–655. ISSN: 0022-5371 (December
1015 1981).
- 1016 11. Rayner, K. The perceptual span and peripheral
1017 cues in reading. en. *Cognitive Psychology* **7**, 65–81.
1018 ISSN: 0010-0285 (January 1975).
- 1019 12. McConkie, G. W. & Rayner, K. The span of the
1020 effective stimulus during a fixation in reading. *Per-*
1021 *ception & Psychophysics* **17**. Publisher: Springer,
1022 578–586 (1975).
- 1023 13. Schotter, E. R., Angele, B. & Rayner, K. Parafoveal
1024 processing in reading. en. *Attention, Perception, &*
1025 *Psychophysics* **74**, 5–35. ISSN: 1943-393X (January
1026 2012).
- 1027 14. Reichle, E. D., Rayner, K. & Pollatsek, A. The E-Z
1028 reader model of eye-movement control in read-
1029 ing: comparisons to other models. eng. *The Be-*
1030 *havioral and Brain Sciences* **26**, 445–476, 445–476.
1031 ISSN: 0140-525X (August 2003).
- 1032

15. Clifton, C. *et al.* Eye movements in reading and information processing: Keith Rayner's 40 year legacy. en. *Journal of Memory and Language* **86**, 1–19. ISSN: 0749-596X (January 2016).
16. Kliegl, R., Risse, S. & Laubrock, J. Preview benefit and parafoveal-on-foveal effects from word n + 2. *Journal of Experimental Psychology: Human Perception and Performance* **33**. Place: US Publisher: American Psychological Association, 1250–1255. ISSN: 1939-1277(Electronic),0096-1523(Print) (2007).
17. Kliegl, R., Nuthmann, A. & Engbert, R. Tracking the mind during reading: The influence of past, present, and future words on fixation durations. *Journal of Experimental Psychology: General* **135**. Place: US Publisher: American Psychological Association, 12–35. ISSN: 1939-2222(Electronic),0096-3445(Print) (2006).
18. Rayner, K., Juhasz, B. J. & Brown, S. J. Do readers obtain preview benefit from word N + 2? A test of serial attention shift versus distributed lexical processing models of eye movement control in reading. eng. *Journal of Experimental Psychology. Human Perception and Performance* **33**, 230–245. ISSN: 0096-1523 (February 2007).
19. Hohenstein, S. & Kliegl, R. Semantic preview benefit during reading. eng. *Journal of Experimental Psychology. Learning, Memory, and Cognition* **40**, 166–190. ISSN: 1939-1285 (January 2014).
20. Yan, M., Kliegl, R., Shu, H., Pan, J. & Zhou, X. Parafoveal load of word N+1 modulates preprocessing effectiveness of word N+2 in Chinese reading. eng. *Journal of Experimental Psychology. Human Perception and Performance* **36**, 1669–1676. ISSN: 1939-1277 (December 2010).
21. Tiffin-Richards, S. P. & Schroeder, S. Children's and adults' parafoveal processes in German: Phonological and orthographic effects. *Journal of Cognitive Psychology* **27**. Publisher: Routledge _eprint: <https://doi.org/10.1080/20445911.2014.999076>, 531–548. ISSN: 2044-5911 (July 2015).
22. Engbert, R., Nuthmann, A., Richter, E. M. & Kliegl, R. SWIFT: a dynamical model of saccade generation during reading. eng. *Psychological Review* **112**, 777–813. ISSN: 0033-295X (October 2005).
23. Engbert, R. & Kliegl, R. The game of word skipping: Who are the competitors? *Behavioral and Brain Sciences* **26**. Publisher: Cambridge University Press, 481 (2003).
24. Drieghe, D., Brysbaert, M., Desmet, T. & De Baecke, C. Word skipping in reading: On the interplay of linguistic and visual factors. en. *European Journal of Cognitive Psychology* **16**, 79–103. ISSN: 0954-1446, 1464-0635 (January 2004).
25. Reilly, R. G. & O'Regan, J. K. Eye movement control during reading: A simulation of some word-targeting strategies. en. *Vision Research* **38**, 303–317. ISSN: 0042-6989 (January 1998).
26. Vitu, F., O'Regan, J. K., Inhoff, A. W. & Topolski, R. Mindless reading: Eye-movement characteristics are similar in scanning letter strings and reading texts. en. *Perception & Psychophysics* **57**, 352–364. ISSN: 1532-5962 (April 1995).
27. Goodman, K. S. Reading: A psycholinguistic guessing game. *Journal of the Reading Specialist* **6**. Place: United Kingdom Publisher: Taylor & Francis, 126–135 (1967).
28. Staub, A. The Effect of Lexical Predictability on Eye Movements in Reading: Critical Review and Theoretical Interpretation. en. *Language and Linguistics Compass* **9**. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/lnc3.12151>, 311–327. ISSN: 1749-818X (2015).
29. Kliegl, R., Grabner, E., Rolfs, M. & Engbert, R. Length, frequency, and predictability effects of words on eye movements in reading. en. *European Journal of Cognitive Psychology* **16**, 262–284. ISSN: 0954-1446, 1464-0635 (January 2004).

- 1111 30. Luke, S. G. & Christianson, K. The Provo Corpus:
1112 A large eye-tracking corpus with predictability
1113 norms. en. *Behavior Research Methods* **50**, 826–
1114 833. ISSN: 1554-3528 (April 2018).
- 1115 31. Cop, U., Dirix, N., Drieghe, D. & Duyck, W. Present-
1116 ing GECO: An eyetracking corpus of monolingual
1117 and bilingual sentence reading. en. *Behavior Re-
1118 search Methods* **49**, 602–615. ISSN: 1554-3528 (April
1119 2017).
- 1120 32. Kennedy, A. *The dundee corpus [cd-rom]* 2003.
- 1121 33. Radford, A. *et al.* Language models are unsuper-
1122 vised multitask learners. *OpenAI Blog* **1**, 8 (2019).
- 1123 34. Duan, Y. & Bicknell, K. A Rational Model of Word
1124 Skipping in Reading: Ideal Integration of Visual
1125 and Linguistic Information. en. *Topics in Cognitive
1126 Science* **12**, 387–401. ISSN: 1756-8757, 1756-8765
1127 (January 2020).
- 1128 35. Smith, N. J. & Levy, R. The effect of word pre-
1129 dictability on reading time is logarithmic. en. *Cog-
1130 nition* **128**, 302–319. ISSN: 0010-0277 (September
1131 2013).
- 1132 36. Frank, S. L., Fernandez Monsalve, I., Thompson,
1133 R. L. & Vigliocco, G. Reading time data for eval-
1134 uating broad-coverage models of English sen-
1135 tence processing. en. *Behavior Research Methods*
1136 **45**, 1182–1190. ISSN: 1554-3528 (December 2013).
- 1137 37. Goodkind, A. & Bicknell, K. *Predictive power of
1138 word surprisal for reading times is a linear func-
1139 tion of language model quality* in *Proceedings of the
1140 8th Workshop on Cognitive Modeling and Computa-
1141 tional Linguistics (CMCL 2018)* (Association for Com-
1142 putational Linguistics, Salt Lake City, Utah, Janu-
1143 ary 2018), 10–18. doi:10.18653/v1/W18-0102.
1144 <[https://www.aclweb.org/anthology/W18-
1145 0102](https://www.aclweb.org/anthology/W18-0102)> (visited on 05/22/2020).
- 1146 38. Luke, S. G. & Christianson, K. Limits on lexical pre-
1147 diction during reading. en. *Cognitive Psychology* **88**,
1148 22–60. ISSN: 0010-0285 (August 2016).
- 1149 39. Pynte, J. & Kennedy, A. An influence over eye move-
1150 ments in reading exerted from beyond the level
1151 of the word: Evidence from reading English and
1152 French. en. *Vision Research* **46**, 3786–3801. ISSN:
1153 0042-6989 (October 2006).
- 1154 40. Kennedy, A., Pynte, J., Murray, W. S. & Paul, S.-A.
1155 Frequency and predictability effects in the Dundee
1156 Corpus: an eye movement analysis. eng. *Quarterly
1157 Journal of Experimental Psychology (2006)* **66**, 601–
1158 618. ISSN: 1747-0226 (2013).
- 1159 41. Norris, D. The Bayesian reader: explaining word
1160 recognition as an optimal Bayesian decision pro-
1161 cess. eng. *Psychological Review* **113**, 327–357. ISSN:
1162 0033-295X (April 2006).
- 1163 42. Norris, D. Putting it all together: a unified account
1164 of word recognition and reaction-time distribu-
1165 tions. eng. *Psychological Review* **116**, 207–219. ISSN:
1166 0033-295X (January 2009).
- 1167 43. Bicknell, K. & Levy, R. *A Rational Model of Eye Move-
1168 ment Control in Reading* in *Proceedings of the 48th
1169 Annual Meeting of the Association for Computa-
1170 tional Linguistics* (Association for Computational
1171 Linguistics, Uppsala, Sweden, July 2010), 1168–
1172 1178. <[https://www.aclweb.org/anthology/
1173 P10-1119](https://www.aclweb.org/anthology/P10-1119)> (visited on 05/15/2020).
- 1174 44. Rayner, K. Eye movements and attention in read-
1175 ing, scene perception, and visual search. eng.
1176 *Quarterly Journal of Experimental Psychology (2006)*
1177 **62**, 1457–1506. ISSN: 1747-0226 (August 2009).
- 1178 45. Rayner, K. & Well, A. D. Effects of contextual con-
1179 straint on eye movements in reading: A further
1180 examination. en. *Psychonomic Bulletin & Review* **3**,
1181 504–509. ISSN: 1531-5320 (December 1996).
- 1182 46. Lee, T. S. & Mumford, D. Hierarchical Bayesian
1183 inference in the visual cortex. eng. *Journal of the
1184 Optical Society of America. A, Optics, Image Science,
1185 and Vision* **20**, 1434–1448. ISSN: 1084-7529 (July
1186 2003).

47. Heilbron, M., Armeni, K., Schoffelen, J.-M., Hagoort, P. & Lange, F. P. d. A hierarchy of linguistic predictions during natural language comprehension. en. *bioRxiv*, 2020.12.03.410399 (2021).
48. O'Regan, J. K. en. in *Eye Movements and Visual Cognition: Scene Perception and Reading* (ed Rayner, K.) 333–354 (Springer, New York, NY, 1992). ISBN: 978-1-4612-2852-3. doi:10.1007/978-1-4612-2852-3_20. <https://doi.org/10.1007/978-1-4612-2852-3_20> (visited on 06/18/2021).
49. Brysbaert, M. & Drieghe, D. Please stop using word frequency data that are likely to be word length effects in disguise. *Behavioral and Brain Sciences* **26**. Publisher: [New York]: Cambridge University Press, 1978-, 479 (2003).
50. McClelland, J. L. & O'Regan, J. K. Expectations increase the benefit derived from parafoveal visual information in reading words aloud. *Journal of Experimental Psychology: Human Perception and Performance* **7**. Place: US Publisher: American Psychological Association, 634–644. ISSN: 1939-1277(Electronic),0096-1523(Print) (1981).
51. Schotter, E. R., Lee, M., Reiderman, M. & Rayner, K. The effect of contextual constraint on parafoveal processing in reading. en. *Journal of Memory and Language* **83**, 118–139. ISSN: 0749-596X (August 2015).
52. Veldre, A. & Andrews, S. Parafoveal preview effects depend on both preview plausibility and target predictability. *Quarterly Journal of Experimental Psychology* **71**. Publisher: SAGE Publications, 64–74. ISSN: 1747-0218 (January 2018).
53. Zwitserlood, P. The locus of the effects of sentential-semantic context in spoken-word processing. en. *Cognition* **32**, 25–64. ISSN: 0010-0277 (June 1989).
54. McClelland, J. L. & Elman, J. L. The TRACE model of speech perception. en. *Cognitive Psychology* **18**, 1–86. ISSN: 0010-0285 (January 1986).
55. Dehaene, S. *Reading in the brain: The new science of how we read* (Penguin, 2009).
56. Yeatman, J. D. & White, A. L. Reading: The Confluence of Vision and Language. *Annual Review of Vision Science* **7**. _eprint: <https://doi.org/10.1146/annurev-vision-093019-113509>, null (2021).
57. Reicher, G. M. Perceptual recognition as a function of meaningfulness of stimulus material. eng. *Journal of Experimental Psychology* **81**, 275–280. ISSN: 0022-1015 (August 1969).
58. Wheeler, D. D. Processes in word recognition. en. *Cognitive Psychology* **1**, 59–85. ISSN: 00100285 (January 1970).
59. Lupyan, G. Objective effects of knowledge on visual perception. eng. *Journal of Experimental Psychology. Human Perception and Performance* **43**, 794–806. ISSN: 1939-1277 (2017).
60. Heilbron, M., Richter, D., Ekman, M., Hagoort, P. & de Lange, F. P. Word contexts enhance the neural representation of individual letters in early visual cortex. en. *Nature Communications* **11**. Number: 1 Publisher: Nature Publishing Group, 321. ISSN: 2041-1723 (January 2020).
61. Deubel, H., O'Regan, J. K. & Radach, R. en. in *Reading as a Perceptual Process* (eds Kennedy, A., Radach, R., Heller, D. & Pynte, J.) 355–374 (North-Holland, Oxford, January 2000). ISBN: 978-0-08-043642-5. doi:10.1016/B978-008043642-5/50017-6. <<https://www.sciencedirect.com/science/article/pii/B9780080436425500176>> (visited on 06/29/2021).
62. Wolf, T. et al. HuggingFace's Transformers: State-of-the-art Natural Language Processing. *arXiv:1910.03771 [cs]*. arXiv: 1910.03771. <<http://arxiv.org/abs/1910.03771>> (visited on 07/17/2020) (July 2020).

-
- 1263 63. Brysbaert, M. & New, B. Moving beyond Kucera
1264 and Francis: A critical evaluation of current word
1265 frequency norms and the introduction of a new
1266 and improved word frequency measure for Amer-
1267 ican English. en. *Behavior Research Methods* **41**,
1268 977–990. ISSN: 1554-3528 (November 2009).
- 1269 64. Pedregosa, F. *et al.* Scikit-learn Machine Learning
1270 in Python. *Journal of Machine Learning Research* **12**,
1271 2825–2830. ISSN: ISSN 1533-7928 (2011).
- 1272 65. Legendre, P. Studying beta diversity: ecological
1273 variation partitioning by multiple regression and
1274 canonical analysis. *Journal of Plant Ecology* **1**, 3–8.
1275 ISSN: 1752-9921 (March 2008).
- 1276 66. De Heer, W. A., Huth, A. G., Griffiths, T. L., Gallant,
1277 J. L. & Theunissen, F. E. The Hierarchical Cortical
1278 Organization of Human Speech Processing. *Jour-
1279 nal of Neuroscience* **37**, 6539–6557 (July 2017).

Supplementary materials

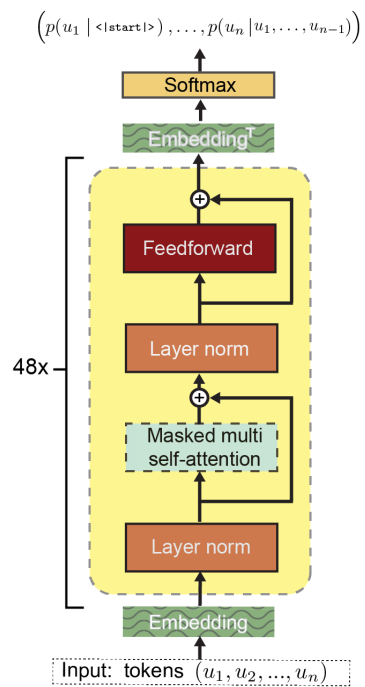


Figure S1 – GPT-2 Architecture. Note that this panel is based on the original GPT schematic, with some operations modified and re-arranged to reflect the slightly different architecture of GPT-2. The most important and distinctive step of each transformer block is masked multi-headed self-attention (see Methods). Not visualised here is the initial tokenisation, mapping a sequence of characters into a sequence of tokens.

Table S1 – Literature sample for effect size ranges

Effect type	Publication	Effect size
preview benefit	Inhoff, A. W. (1989). Lexical access during eye fixations in reading: Are word access codes used to integrate lexical information across interword fixations?. <i>Journal of Memory and Language</i> , 28(4), 444-461.	51
preview benefit	Veldre, A., & Andrews, S. (2018). Parafoveal preview effects depend on both preview plausibility and target predictability. Lexical access during eye fixations in reading: <i>Quarterly Journal of Experimental Psychology</i> , 71(1), 64-74.	49

Continued on next page

Table S1 – *Continued from previous page*

Effect type	Publication	Effect size
preview benefit	Inhoff, A. W., & Rayner, K. (1986). Parafoveal word processing during eye fixations in reading: Effects of word frequency. <i>Perception & psychophysics</i> , 40(6), 431-439.	40
preview benefit	McDonald, S. A. (2006). Parafoveal preview benefit in reading is only obtained from the saccade goal. <i>Vision Research</i> , 46(26), 4416-4424.	35
preview benefit	Williams, C. C., Perea, M., Pollatsek, A., & Rayner, K. (2006). Previewing the neighborhood: The role of orthographic neighbors as parafoveal previews in reading. <i>Journal of Experimental Psychology: Human Perception and Performance</i> , 32(4), 1072.	26.7
preview benefit	Kennison, S. M., & Clifton, C. (1995). Determinants of parafoveal preview benefit in high and low working memory capacity readers: Implications for eye movement control. <i>Journal of Experimental Psychology: Learning, Memory, and Cognition</i> , 21(1), 68.	25.25
preview benefit	Blanchard, Harry E., Alexander Pollatsek, and Keith Rayner. "The acquisition of parafoveal word information in reading." <i>Perception & Psychophysics</i> 46.1 (1989): 85-94.	22.6
preview benefit	Schroyens, W., Vitu, F., Brysbaert, M., & d'Ydewalle, G. (1999). Eye movement control during reading: Foveal load and parafoveal processing. <i>The Quarterly Journal of Experimental Psychology Section A</i> , 52(4), 1021-1046.	14.6
prediction benefit	Ehrlich, S. F., & Rayner, K. (1981). Contextual effects on word perception and eye movements during reading. <i>Journal of verbal learning and verbal behavior</i> , 20(6), 641-655.	33
prediction benefit	Rayner, K., & Well, A. D. (1996). Effects of contextual constraint on eye movements in reading: A further examination. <i>Psychonomic Bulletin & Review</i> , 3(4), 504-509.	20
prediction benefit	R.J. Altarriba, J. Kroll, A. Sholl, K. Rayner. (1996) The influence of lexical and conceptual constraints on reading mixed-language sentences: Evidence from eye fixations and naming times <i>Memory & Cognition</i> , 24 (1996), pp. 477-492.	21

Continued on next page

Table S1 – *Continued from previous page*

Effect type	Publication	Effect size
prediction benefit	Ashby, J., Rayner, K., & Clifton Jr, C. (2005). Eye movements of highly skilled and average readers: Differential effects of frequency and predictability. <i>The Quarterly Journal of Experimental Psychology Section A</i> , 58(6), 1065-1086.	23.5
prediction benefit	Rayner, K., Ashby, J., Pollatsek, A., & Reichle, E. D. (2004). The effects of frequency and predictability on eye fixations in reading: implications for the EZ Reader model. <i>Journal of Experimental Psychology: Human Perception and Performance</i> , 30(4), 72	19
prediction benefit	Rayner, K., Binder, K. S., Ashby, J., & Pollatsek, A. (2001). Eye movement control in reading: Word predictability has little influence on initial landing positions in words. <i>Vision Research</i> , 41(7), 943-954.	15
prediction benefit	Rayner, K., Slattery, T. J., Drieghe, D., & Liversedge, S. P. (2011). Eye movements and word skipping during reading: effects of word length and predictability. <i>Vision Research</i> , 41(7), 943-954.	18
prediction benefit	Hand, C. J., Miellet, S., O'Donnell, P. J., & Sereno, S. C. (2010). The frequency-predictability interaction in reading: It depends where you're coming from. <i>Journal of Experimental Psychology: Human Perception and Performance</i> , 36(5), 1294-1313.	12

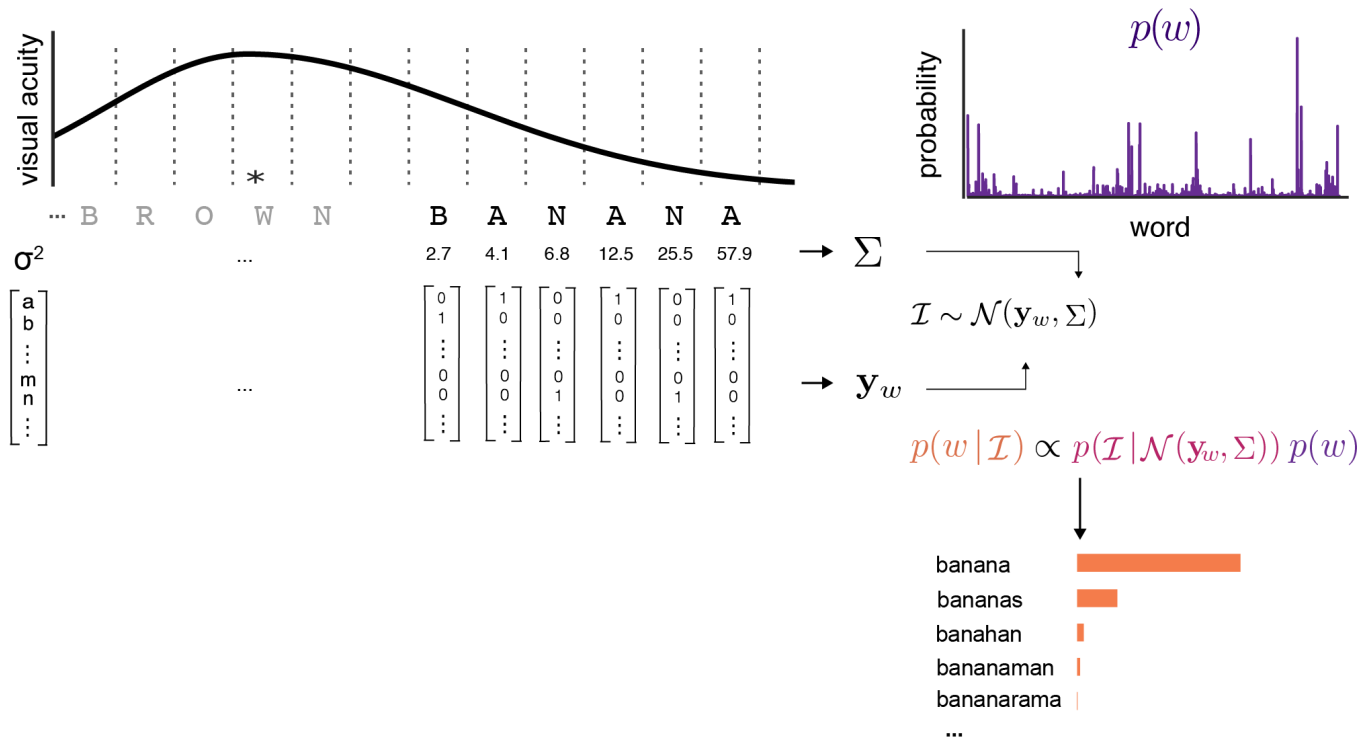


Figure S2 – Encoding and inference scheme of the ideal observer analysis. A word at a given eccentricity is converted into a noisy visual percept, after which a posterior probability of the identity of the word given the noisy percept was computed using Bayesian inference. The uncertainty of this posterior (expressed in terms of Shannon entropy) was then used to quantify the expected uncertainty in the parafoveal percept – or, inversely, a word's *parafoveal identifiability*.

In this scheme, words are represented as a concatenation of one-hot encoded letter vectors. Visual information (\mathcal{I}) is sampled from a multivariate Gaussian centred on the word vector \mathbf{y}_w with a diagonal covariance matrix Σ , the values of which (σ^2) are inversely related to the integral under the visual acuity function around each letter. The posterior is then computed by combining the likelihood of the visual information \mathcal{I} given a particular word, with a prior probability of that word $p(w)$ (e.g. derived from lexical frequency). This computation was performed using a log-odds formulation that exploits the proportionality in Bayes' rule to perform belief-updating without renormalisation (see Methods).

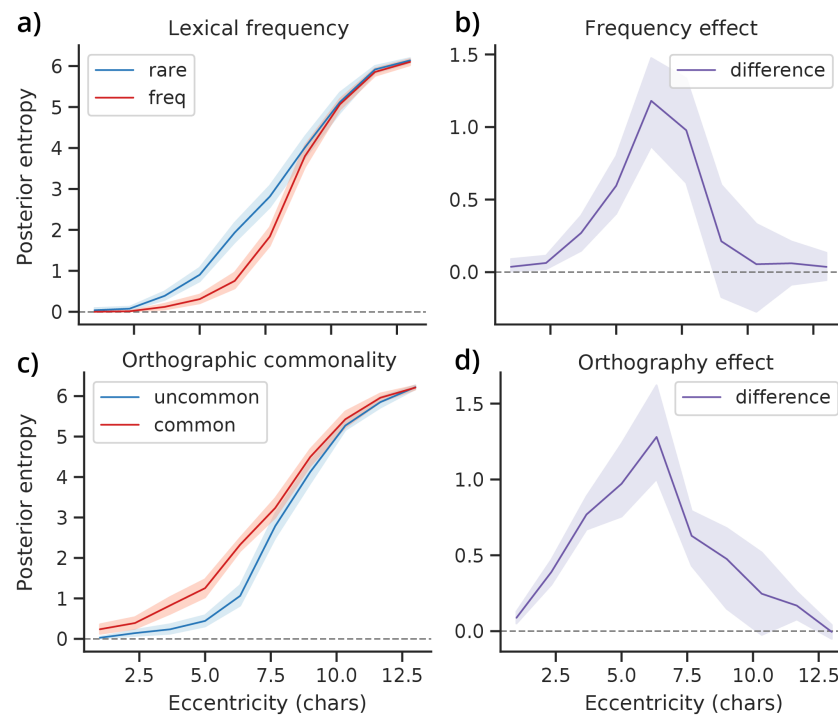


Figure S3 – Modulation of parafoveal identifiability by visual and linguistic features, and their interaction. The parafoveal entropy for a given word (Fig S2) is a complex function that integrates linguistic and visual characteristics, and which can account for various known effects, such as the effect of lexical frequency and orthographic neighbourhood on visual word recognition. To illustrate this, we simulated some characteristic effects of eccentricity, frequency (a,b) and orthographic distinctiveness (c,d). For frequency (a), we randomly sampled 20 ‘rare’ and ‘frequent’ 5-letter words (based on a quartile split), and computed the parafoveal identifiability (quantified via posterior entropy) at increasing eccentricities. As can be seen, the percept becomes uncertain at increasing eccentricities more quickly for low-frequency words, showing that lexical frequency boosts parafoveal identifiability. For orthography (c), we similarly sampled 20 7-letter words that were classified as orthographically common or uncommon based on the first three letters. Here, commonality was again defined using a quartile split but now on the number of alternative words starting with the same three letters. For instance, the letters ‘awk’ in the word ‘awkward’ are highly uncommon and allow to identify the entire word with high confidence based on just those three letters. As can be seen, the model predicts that orthographic uniqueness boosts parafoveal identifiability – as observed in experiments (see [13]). Notably, when we consider the difference between the two classes of words (b,d), an inverted U shape is apparent: the effects are strongest at intermediate visibility. This demonstrates the well-established fact that the effects of prior (linguistic) knowledge is strongest at intermediate levels of perceptual uncertainty (see [41] for discussion). (Note that, while both the orthography and frequency effects are effects of "prior linguistic knowledge", only the frequency effect is technically an effect of the *prior*, since the orthography effect is driven by the generative model.) In all plots, thick lines represent the mean entropy across words; shaded regions indicate bootstrapped 95% confidence intervals.

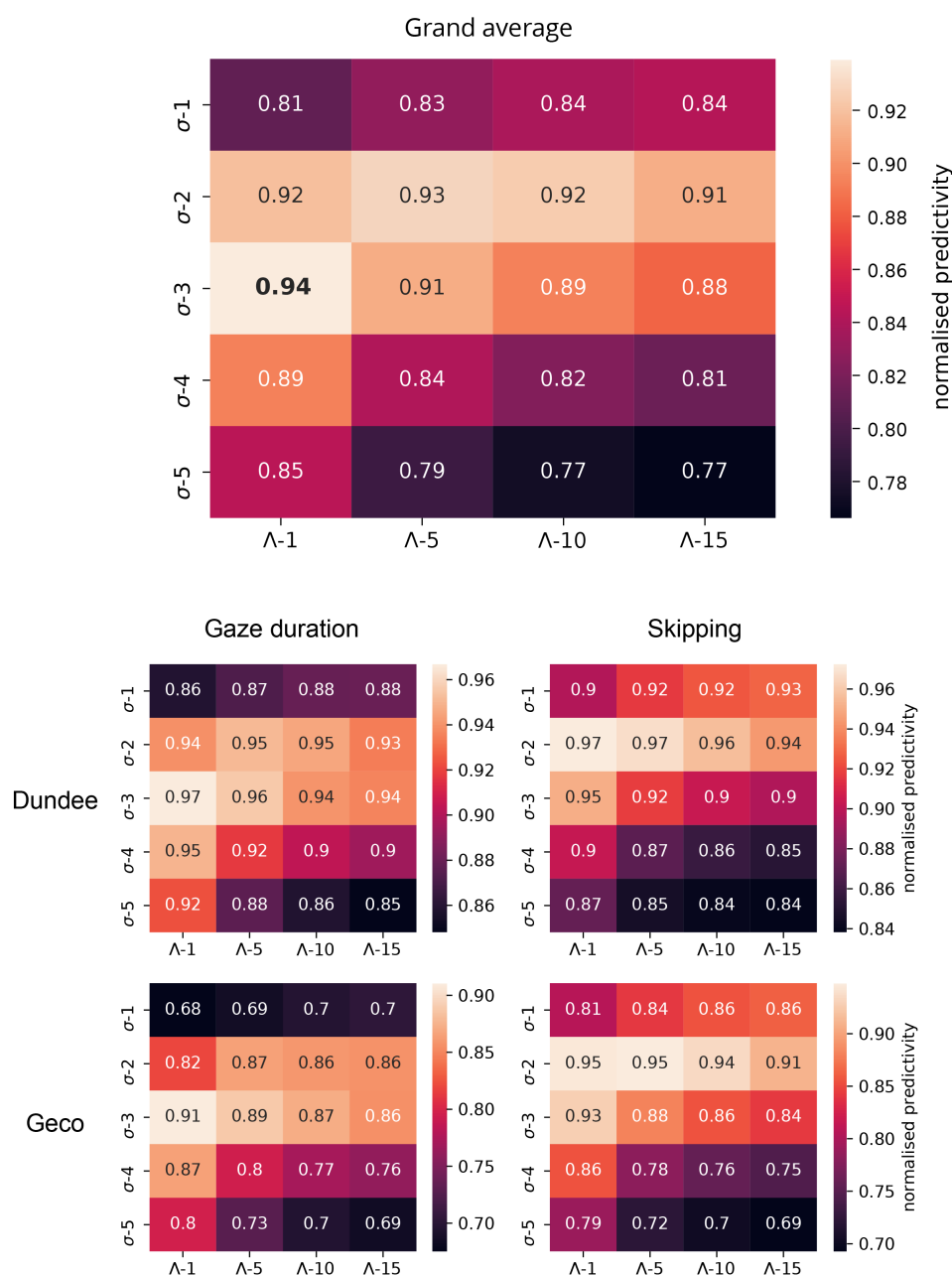


Figure S4 – Grid search to establish ideal observer parameters. Grid search result grand average (top) and individual results for different corpora and analyses (bottom). To decide on the values for σ and Λ , a grid search was performed on a random subset of 25% of the Dundee and Geco corpus; we did not apply it to PROVO because there was not enough data per participant. In both skipping and reading times, we performed a 10-fold cross-validation with the full model, using parafoveal entropy as computed with different visual acuity parameters σ and Λ (Equation 6). To avoid biasing the contextual vs non-contextual model comparison (Figure 5), we used both the contextual and non-contextual prior and averaged the results to obtain the results for each analysis in each corpus. To ensure that different analyses and corpora are weighted equally in the grand average, the prediction scores (R^2 or R^2_{McF}) were normalised by dividing the prediction score of each parameter combination by the highest score (i.e. score of the best parameter combination) for each subject, for each analysis. This resulted in $\sigma = 3$ and $\Lambda = 1$, which we have used in all analyses. Note that σ determines the perceptual span (see Figure S2) and that $\sigma = 3$ corresponds well to what is known about the size of the perceptual span and is close to default parameters in other models (see Methods).

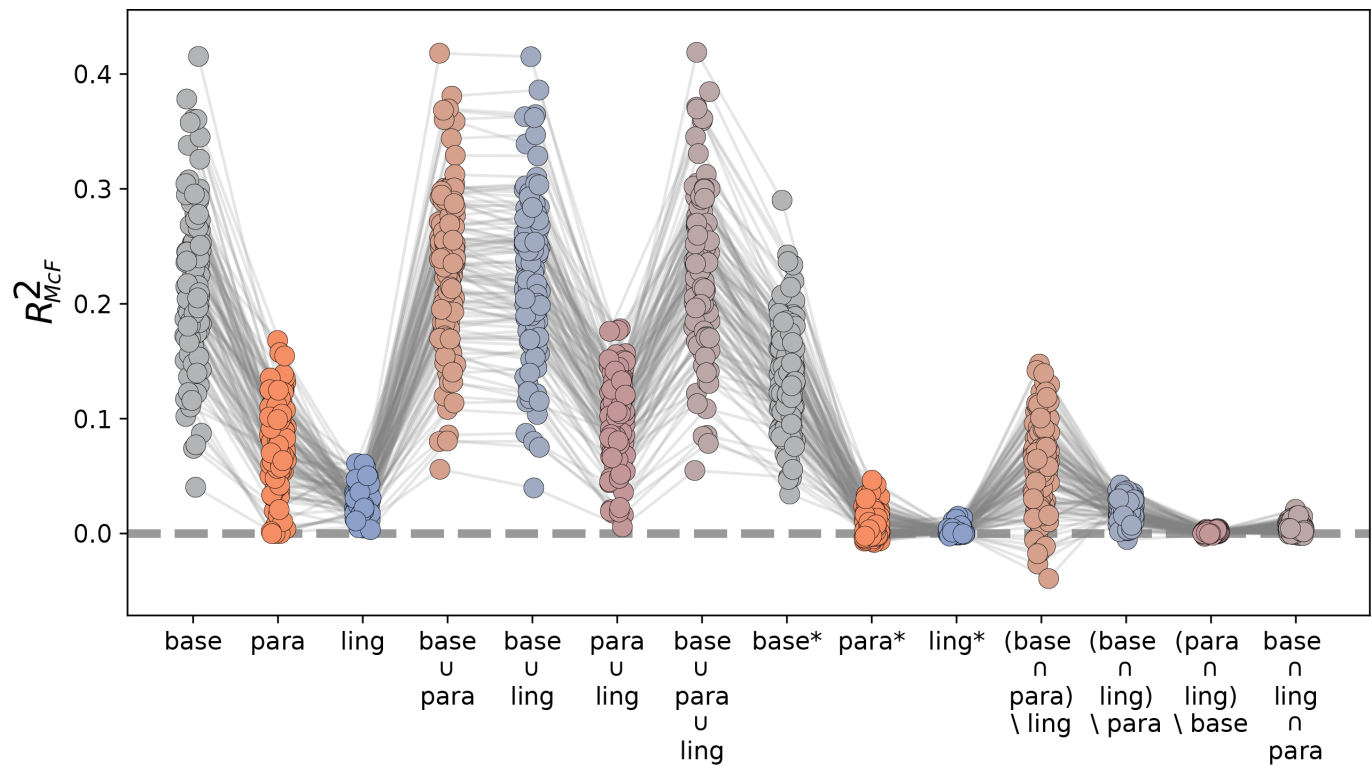


Figure S5 – Skipping variation partitioning for all participants. Explained cross-validated variation partition for skipping (see Fig 2) of each partition, for each participant, for the skipping analysis. Models for the baseline, parafoveal preview and linguistic prediction are indicated by ‘base’, ‘para’, and ‘ling’, respectively. Unions are indicated by \cup , intersections by \cap ; for the relative complement we use the asterisk-notation: e.g. ‘para*’ indicates variation explained uniquely by parafoveal preview. Note that due to cross-validation, the amount of variation explained can become negative in some partitions for individual participants (see Methods).

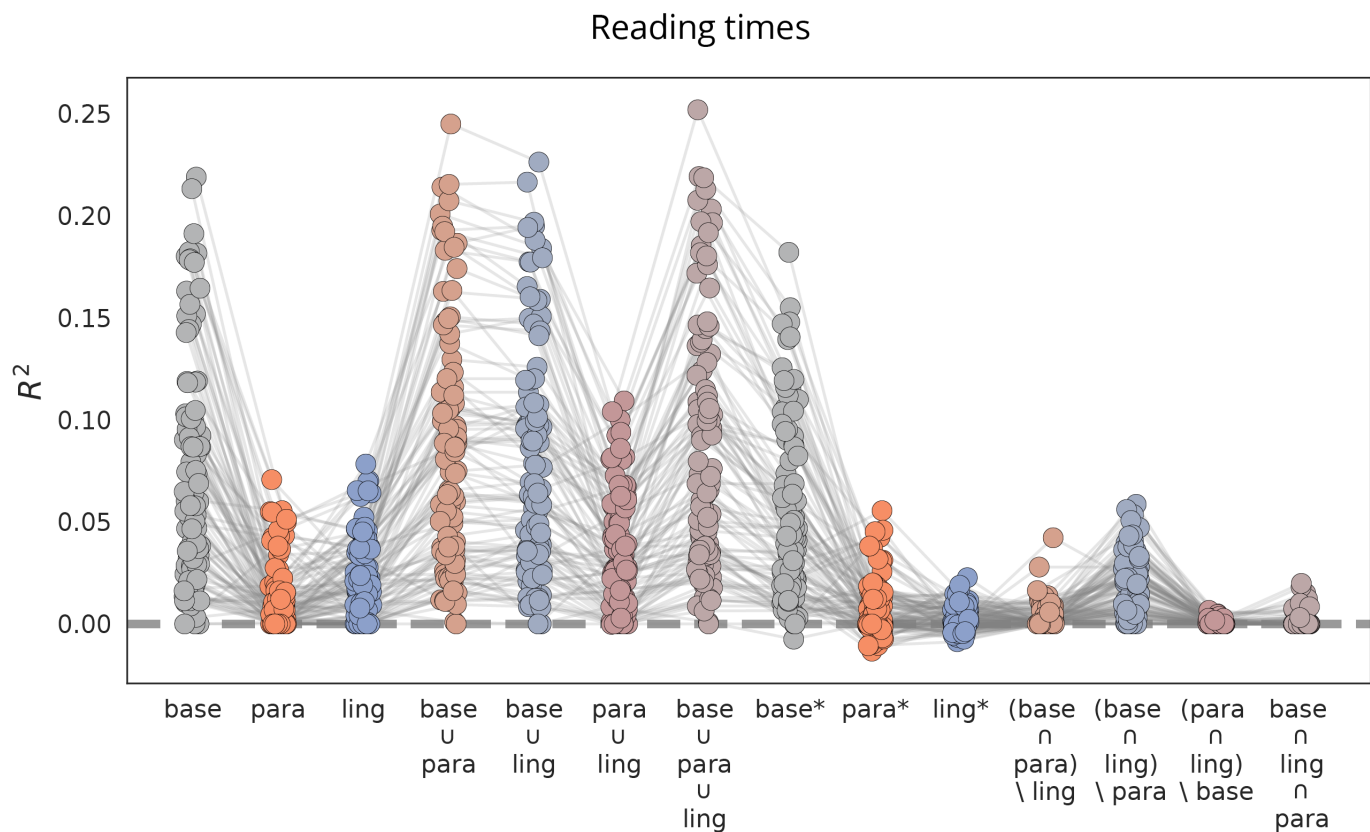


Figure S6 – Reading times variance partitioning. Explained cross-validated variation partition for skipping (see Fig 3) of each partition, for each participant, for the skipping analysis. Models for the baseline, parafoveal preview and linguistic prediction are indicated by ‘base’, ‘para’, and ‘ling’, respectively. Unions are indicated by \cup , intersections by \cap ; for the relative complement we use the asterisk-notation: e.g. ‘para*’ indicates variation explained uniquely by parafoveal preview (see Methods). Note that due to cross-validation, the amount of variation explained can become negative in individual participants (see Methods).

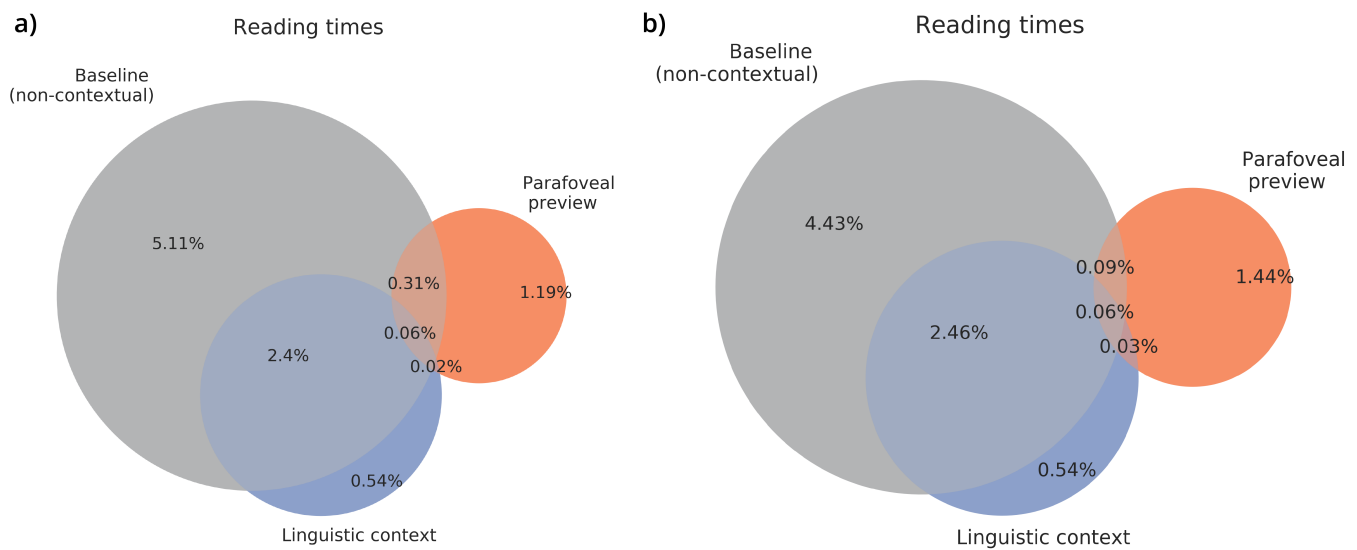


Figure S7 – Reading times variance partitioning with and without non-linguistic factors Same as in Fig 3, but comparing the baseline with (a) or without (b) the primary non-linguistic explanatory factor for reading time variation – viewing position [48]. Including the viewing position adds 0.7% additional variance explained. This demonstrates while that viewing position affect reading times, the amount of variance uniquely explained by non-linguistic factors is much lower for reading times than for skipping.