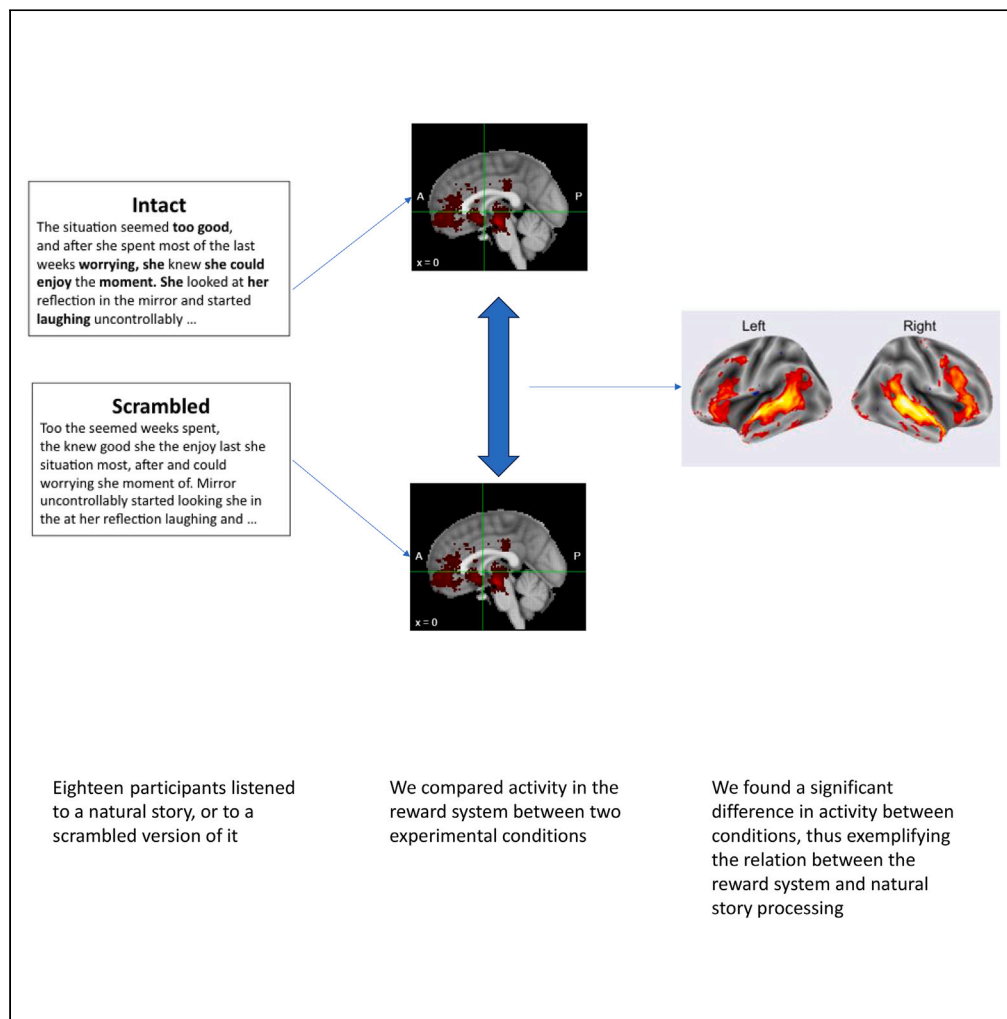


Article

Reward-related regions play a role in natural story comprehension



Oren Kobo, Yaara Yeshurun, Tom Schonberg

orenkobo@mail.tau.ac.il (O.K.)
schonberg@tauex.tau.ac.il (T.S.)

Highlights

The involvement of the reward system in language processing is assessed

We used fMRI data to compare activity in the reward system during story comprehension

We compared an intact version versus a scrambled one

We show that activity in the reward system involved in processing of natural stories

Kobo et al., iScience 27, 109844
June 21, 2024 © 2024 The Authors. Published by Elsevier Inc.
<https://doi.org/10.1016/j.isci.2024.109844>



Article

Reward-related regions play a role in natural story comprehension

Oren Kobo,^{1,5,*} Yaara Yeshurun,^{2,4} and Tom Schonberg^{3,4,*}

SUMMARY

The reward system was shown to be involved in a wide array of processes. Nevertheless, the exploration of the involvement of the reward system during language processing has not yet been directly tested. We investigated the role of reward-processing regions while listening to a natural story. We utilized a published dataset in which half of the participants listened to a natural story and the others listened to a scrambled version of it to compare the functional MRI signals in the reward system between these conditions and discovered a distinct pattern between conditions. This suggests that the reward system is activated during the comprehension of natural stories. We also show evidence that the fMRI signals in reward-related areas might potentially correlate with the predictability level of processed sentences. Further research is needed to determine the nature of the involvement and the way the activity interacts with various aspects of the sentences.

INTRODUCTION

The reward system has been shown to be involved in a wide variety of processes in the brain.¹ Prediction is a hallmark of the reward system.² Accordingly, several theories propose that the brain's primary objective is to reduce surprise given sensory input.^{3,4} The reward prediction error signals encodes the difference between received and predicted rewards using the phasic activity of dopamine neurons.⁵

In psycholinguistics, prediction is a key explanation of the human ability to comprehend language efficiently.^{6,7} The negative correlation between N400 evoked response potential component and word predictability provides neural evidence to the relationship between predictability and speed of processing and exemplifies the major role statistics plays in prediction^{8,9} and language comprehension.^{10,11} Importantly, several recent studies^{11–14} managed to isolate the semantic prediction signal even prior to the appearance of the critical word, providing additional evidence of the profound role prediction plays in language processing. Several fMRI studies focused on brain activity evoked by a surprising ending or semantic plausibility¹⁵ and showed that the left inferior frontal cortex is consistently activated in such cases. Others have studied the interaction between unexpected input and evoked neural activity, e.g., the effect of sentence type,¹⁶ word knowledge,¹⁷ or type of anomaly.¹⁸ It was also demonstrated that sentence comprehension specifically is attributed to the middle temporal gyrus.¹⁹ However, all those studies concentrated on traditional linguistic areas in the brain and on mapping different aspects of the input and their influence on evoked neuronal activity. Here, in contrary to these previous studies, we aimed to directly test the involvement of the reward system in language comprehension of natural stories and test the effect of predictability during processing of natural stories.

Scrambling of sentences is a tool often utilized to distill cognitive semantic processes when contrasted with corresponding intact sentences. It was used to explore properties of timescale hierarchy and temporal accumulation of information during processing across the cortex. It has been shown that the longest timescales were for default-mode networks, shorter for intermediate areas along the superior temporal gyrus, and the shortest to early sensory regions.^{20–22} Recently, natural language processing (NLP) computational modeling has become more commonly used in the investigation of processing during comprehension of a naturalistic story. Several studies correlated a predictability-related metric with various facets of neural activity while listening to natural stories.^{23,24} Using stories rather than controlled stimuli has many advantages. For example, it can be used to explore various questions in a single dataset and often reveal more widespread responses to language than controlled stimuli.^{25–28} Computational tools are now used to formalize the attributes of the process in question and correlate it with elicited neural activity to gain a better understanding of underlying cognitive processes. The predictability of the input is often used as that regard, with the measures of semantic distance²³ entropy & surprise²⁴ or perplexity²⁹ as various ways to formalize it.

Here we used sentence perplexity as a metric, in order to test the role of the reward system during comprehension of sentences. Perplexity is a common metric in NLP used to evaluate text³⁰ and was also used in cognition and neuroimaging.^{22,29,31} We pre-registered this as our target measure.

¹Sagol School of Neuroscience, Tel Aviv University, Tel Aviv, Israel

²School of Neurobiology, Biochemistry and Biophysics, The George S. Wise Faculty of Life Sciences and Sagol School of Neuroscience, Tel Aviv University, Tel Aviv, Israel

³School of Psychological Science and Sagol School of Neuroscience, Tel Aviv, Israel

⁴These authors contributed equally

⁵Lead contact

*Correspondence: orenkobo@mail.tau.ac.il (O.K.), schonberg@tauex.tau.ac.il (T.S.)

<https://doi.org/10.1016/j.isci.2024.109844>



Differentiating Conditions

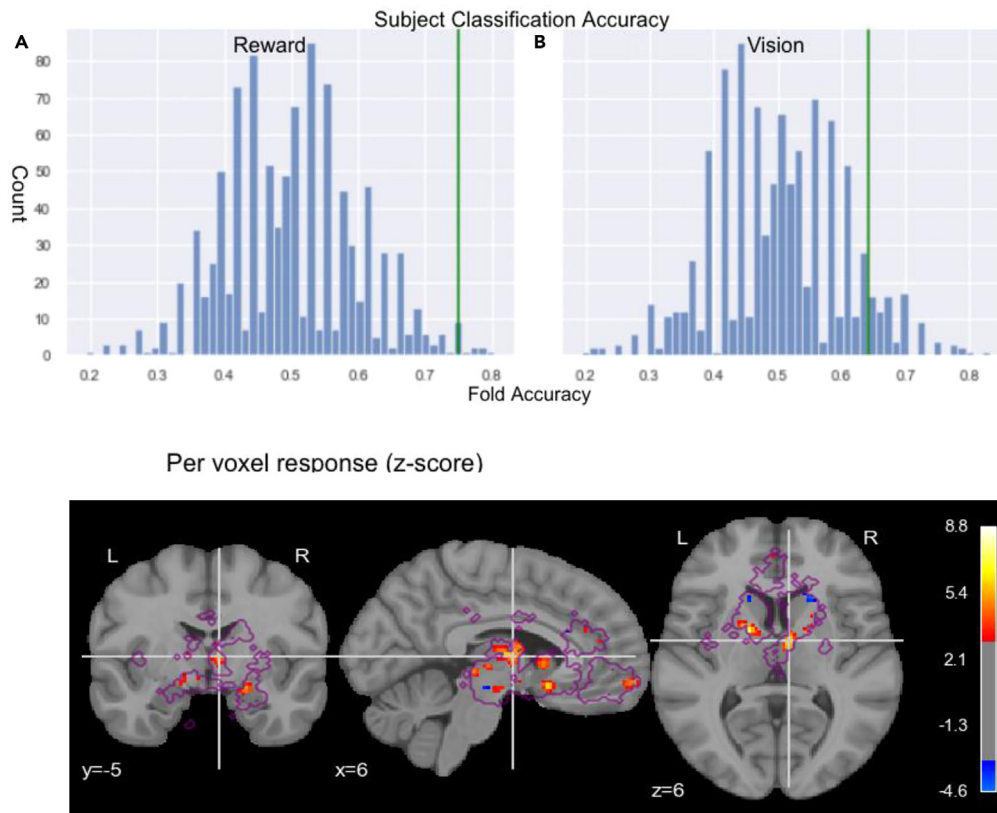


Figure 1. Differentiating between activity in intact vs. scrambled conditions

Up: Results of condition classification per participant: rank of accuracy score in the null distribution results (permutation test). The null distribution was generated by shuffling of the labels, to obtain p values. (A) Reward system, (B) vision system (control). The green vertical line indicated the mean accuracy across folds. Down: per-voxel difference (Z score) between neural response to the intact vs. the scrambled story at $(x = 6.44, y = -5.63, z = 6.71)$. The map was thresholded at $z = 3$. The contours of the mask are plotted in purple and the activations are in red/yellow scale.

We analyzed an openly shared fMRI dataset where participants listened to a story and a scrambled version of the same story. We surmise that given the prominent role of the reward system in generating and calibrating predictions, and the role predictions play in language processing, the reward system is likely to take part in sentence comprehension. We hypothesized we would be able to distinguish between conditions based on activity in the reward system and exemplify that the nature of this distinction can be attributed to sentence predictability. Accordingly, we tested the involvement of the reward system in natural stories processing, and whether the predictability of the linguistic input interacts with elicited neural activity in the reward system, which to date has not been directly demonstrated.

RESULTS

Differentiating between activity in intact vs. scrambled conditions

We used a support vector machine classifier to predict the condition of a participant, in order to assess the differentiation between activity in the response to the intact and scrambled stories. We achieved an accuracy of 74% in predicting the condition. This accuracy value was significant according to the permutation test ($p = 0.012$, see Figure 1A). As a control, we also ran the same analysis based on activity in the visual system, for which we received an insignificant result in a permutation test ($p = 0.085$). Although a p value of 0.085 is approaching significance, note that no voxel in the visual system survived the subsequent analysis we conducted to measure differentiation between the conditions, unlike in the reward system. We replicated this analysis also for the ventral striatum (defined by the nucleus accumbens) and the dorsal striatum (defined by the unification of the caudate and the putamen), to check if any of these sub-regions showed significantly better accuracy. Both sub-regions yielded similarly significant accuracy levels at the classification task (See Figure S1 in the supplemental information).

To further explore the differentiation between conditions, we calculated the Euclidean distance between the averaged time courses for the intact and scrambled conditions in each voxel across the reward system (12,031 voxels). For each voxel, we obtained the Z score of the intact-vs-scrambled Euclidean distance, compared to the null distribution, and calculated the adjusted p (with false discovery rate [FDR]). This

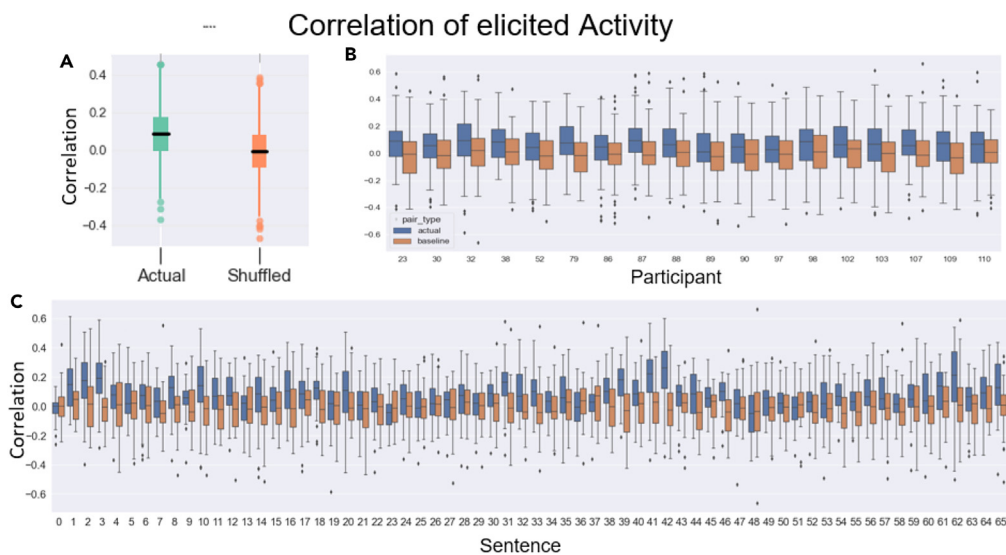


Figure 2. Breakdown of the correlation in the elicited activity during the story comprehension

(A) Correlation of elicited activity per sentence is higher for actual data compared to baseline.

(B) Results per participants.

(C) Results per sentence.

analysis revealed 268 voxels in the reward system in which there was a significant difference in the response to the intact vs. the scrambled story (See Figure 1B). In the control visual system, this analysis revealed no such voxels. See Supplementary for plotting of the reward and vision masks we used (Figure S2). We also conducted post hoc Euclidean analysis between the dorsal and the ventral striatum, hoping to probe for differences in activity in the reward centers: At each TR (repetition time), for each of the selected sub-regions (dorsal striatum – defined by the nucleus accumbens and ventral striatum – defined by the putamen and the caudate), we calculated the mean activity (mean across all participants) in that TR for the sub-region (denoted as TRm) and then calculated the sum of distances D_i , such that D_i is the distance between TRm and the activity at the same TR of participant i . To account for the different number of voxels in each region, we first applied principal component analysis ($n = 10$) at the participant level. We ran a t test to compare the distribution of distances and obtained a significant difference $t(268) = 40.3$ ($p < 0.0001$, see Figure S6 in the supplementary for plotting of the per-TR total distance distribution).

Encoding of sentence-level information in the reward system

We conducted pattern similarity analysis – i.e., we correlated activity between sentences in the reward mask to detect a between-participants correlation of evoked activity per sentence in the reward system. This was done to test for shared aspects of processing of sentences across participants in the reward system. We obtained a weak but significant correlation of 0.09 for the actual data compared to 0.01 for the baseline (randomly chosen pairs) (see Figure 2). For the scrambled condition, the correlation was 0.03. We did not find a significant correlation between perplexity and pattern similarity per sentence, meaning we couldn't identify a relation between values pattern similarity of sentences to perplexity values. We also tested pattern similarity per participant and per sentence to assess the robustness of such a correlation. At the participant level, we found that for each of the participants there is a significant difference between the spatial correlation of elicited activity in the real compared to randomly chosen (baseline) pairs ($p = 0.04$ for a single participant and $p < 0.01$ for all others). For the scrambled condition, only two participants had $p < 0.01$ and 5 others had $p < 0.05$. All the others had a p value of above 0.05. We performed a t test between the p value in the scrambled and the intact to ensure it is significant and obtained $t(17) = -3.42$, $p = 0.003$. At the sentence level, we found a significant correlation only for 34 out of the 66 sentences. For scrambled, we found a significant correlation only for 16 out of the 66 sentences. We performed a t test between the p -score in the scrambled and the intact to ensure it is significant and obtained $t(65) = -3.63$, $p = 0.0005$. This may imply that activity in the reward system is moderated or activated by certain sentences (or types of sentences) to a higher extent than by others. Guided by this finding, we also trained a multi-sentence classifier. The classifier was trained to predict the processed sentences based on the activity in the reward area. We obtained an accuracy of 0.053 (while chance level is 0.01515 since there are 66 sentences). This was significantly above chance level ($p = 0.016$) in the permutation test, whereas for the control visual system and the control scrambled condition we obtained non-significant results ($p = 0.28$ and $p = 0.29$, respectively) (see Figure 3). We did not find correlations between sentence perplexity and classification accuracy. Meaning, level of perplexity did not influence the performance of the model. However, there was a strong correlation between the accuracy of sentence classification and the pattern similarity of a sentence – (Spearman $r = 0.35$, $p = 0.003$), meaning that sentences that got higher accuracy rates also tended to have higher pattern similarity. This could further indicate that certain sentences elicit activity in the reward regions in a different manner than other sentences, depending on their specific attributes. However, it is worth

Multi-Sentence Classifier results

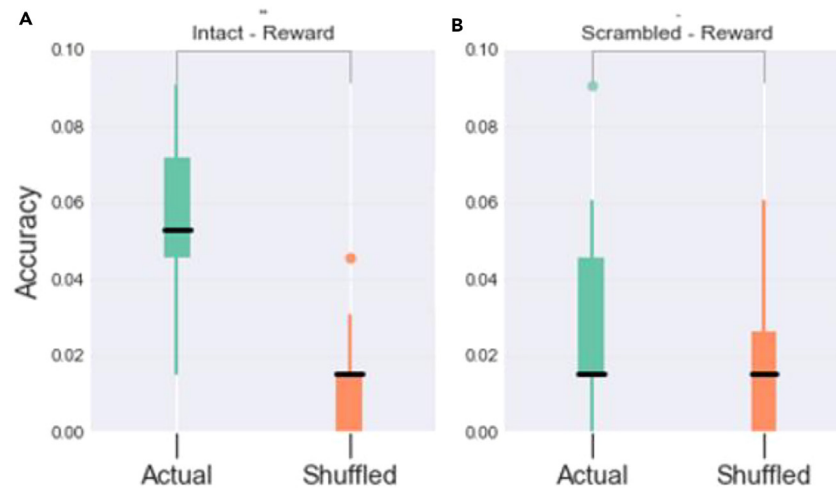


Figure 3. Multi-sentences classifier: Distribution of per fold accuracy: actual vs. shuffled labels (permutations)

(A) For the intact condition.

(B) For the scrambled condition. Models are based on neural data acquired in the reward system.

mentioning an alternative explanation for this correlation: since both analyses rely on between-subject similarity in sentence representations, a certain sentence might inherently be successful in both types of analysis. See Supplementary for a more detailed analysis (Figures S3 and S4 - results with all possible pairs as a null distribution; Figure S5 - classification results in the vision control region).

The relation between activity and predictability

In these analyses, our goal was to examine if the relation between activity in the reward system could be attributed to predictability. Meaning, in these analyses, we aimed to probe whether the differences that were detected in previous ones could be attributed to differences in predictability between conditions.

Ordinal classifier for sentence perplexity

We trained an ordinal classifier for sentence perplexity, based on activity in the reward system in the intact condition and in the scrambled condition as control. The difference between the intact results and the results with shuffled labels was insignificant (as well as in scrambled and vision) (See Figure 4 Left).

Binary perplexity classifier. In this pre-registered analysis, we trained a binary classifier to predict the perplexity of a sentence (above or below the median), given elicited neural data from the reward system in the intact condition, and compared the results to a baseline obtained by using shuffled labels. We obtained insignificant ($p = 0.13$) results in permutation test (rank of the mean in the null distribution) but significant results in KS (Kolmogorov–Smirnov) test between the distributions ($KS = 0.5$, $p = 0.001$). The same comparison in the scrambled condition was insignificant (0.27 in permutation test. $KS = 0.27$, $p = 0.5$ in KS test). Results in the control visual region were also insignificant ($p = 0.37$ in permutation test, $KS = 0.33$, $p = 0.27$ in KS test) (see Figure 4 Middle).

Regression of sentence perplexity. We trained a regression model for sentence perplexity and in the vision system as control. We used R-square as an evaluation metric. We obtained significant ($p < 0.001$) results in the permutation test (rank of the mean in the null distribution) and in the KS test between the distributions ($KS = 0.61$, $p = 0.001$). The same comparison in the scrambled condition was insignificant (0.28 in permutation test. $KS = 0.22$, $p = 0.78$ in KS test). It was also insignificant in the visual region ($p = 0.166$ in permutation test, $KS = 0.38$, $p = 0.13$ in KS test) (See Figure 4 Right).

Within-participant perplexity classifier

We trained a separate classifier per participant, to predict the perplexity level of a given sentence, from the evoked activity while it is being processed. We did not obtain any significant result in the permutation test, meaning when trained a separate classifier per participant (instead of the aggregated data of all the participants), we could not successfully predict whether a certain sentence had above median perplexity based on acquired data in the reward system on this specific participant.

Perplexity Analysis Results



Figure 4. Results of various types of perplexity analysis models

Left: Ordinal classifier (distribution of per fold MSE).

Middle: Binary classifier (distribution of per fold accuracy).

Right: Regression (distribution of per fold R-squared).

Whole-brain searchlight analysis

We applied a voxel-level Euclidean distance whole-brain analysis to test for regions that had a significant difference in their response to the intact and scrambled stories (Figure 5). As expected, based on previous studies and the language-comprehension nature of the difference between the two stories, this exploratory analysis revealed mainly regions involved in language and comprehension processes (Table 1). This demonstrates the dominance of language-related regions in such tasks and explains why our analysis of the role of the reward region should be done within a specified mask and not the whole brain.

DISCUSSION

In this study, we explored the role of the reward system in language processing during narrative comprehension. We reanalyzed a published dataset, composed of two conditions – a natural story and a scrambled version of it.³² We found that the reward system was involved in processing the intact (but not the scrambled) story. Moreover, using a sentences classifier and pattern similarity analysis, we found that the reward system encoded information at the sentence level. We provide an indication that this can be driven by sentence-level information, possibly the predictability of the language input, by assessing the dynamics through which fMRI responses were affected by sentence-level information and specifically by the predictability (formalized by perplexity) of the processed sentence, although more research is required to specifically inspect the role of prediction in driving the reward system response during narrative comprehension.

The reward system has been shown to be activated during language processing, mostly in the context of a presence of an external reinforcer,^{33,34} humor,^{35,36} or a pleasing content.³⁷ The reward system, and specifically the ventral striatum, has been shown to have an increased activation when learning new words.^{38,39} In the current study, we aimed to focus on the role of the reward system in the processing of natural language as opposed to individual words. The core difference between learning-related and language-related findings highlights the contribution of the current study. Semantic processing and the reward system have also been linked through the finding of the spread of semantic priming activation, mainly using patients with Parkinson's disease that display abnormal activation.^{40,41} However, it has yet to be exemplified that the reward system is activated during normal comprehension of a natural story. The ability to consistently differentiate test participants at an above-chance level between conditions based on activity in the reward areas suggests that the processing of a natural story is to the very least supported by these regions. Prediction of upcoming input is a core strategy in the brain, activated across many cognitive domains, including perception, sensory-motor processing, and learning (for reviews see Clark, Friston^{42,43}). The reward system takes part in predictions using prediction error signals,⁵ which also seem to be crucial in language comprehension due to the fast processing required.⁴⁴ Thus, a successful prediction leads the system to improve and facilitate faster reaction times.

In this manuscript, we first demonstrated that the activity in the reward system is moderated by the experimental condition – a natural story (intact), or a scrambled version of the same story: we successfully trained a classifier to detect the condition (type of story that is being processed by a participant), and we showed that the activity of a large number of voxels in the reward system operates significantly different between conditions. In both cases, we verified that the finding does not replicate in the control region of interest – the visual system. We further tested the hypothesis that the activation in the reward system is related to the actual processed sentence at a certain time. We found an increased pattern similarity between sentences and above chance accuracy of a classifier that predicted the processed sentence from neural activity, suggesting there was a distinct sentence-level representation. While a correlation between accuracy in sentence classification prediction and the predictability of a sentence was not found, we did find a correlation between accuracy in sentence classification and levels of pattern similarity (meaning, sentences that tend to have higher pattern similarity tend to have higher classification accuracy). This raises the possibility that specific types of sentences are more represented or processed in the reward system than others. Given these findings, we concluded that there might be some extent of sentence-level information encoded in the reward system, and that the reward system encoding might be related to specific attributes of the sentence. As far as we know, this is the first time that a relation between sentence-related attributes and representation in the reward system has been exemplified. However, more research with controlled stimuli is required to further assess this.

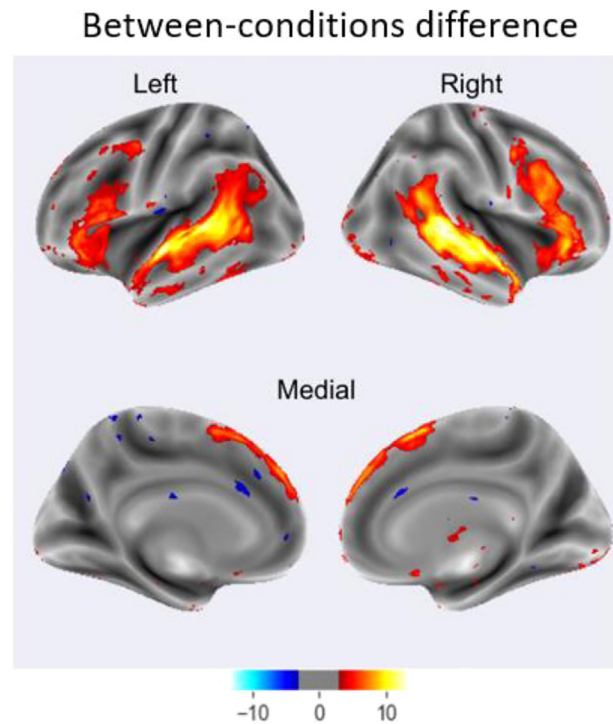


Figure 5. Voxels activity difference between scrambled and milky conditions (Z score)

As the two main factors that were modified between the two conditions (intact and scrambled) were the predictability of the input and the existence of a structured narrative in the intact condition, we then aimed to disentangle them by focusing on demonstrating the influence of predictability on neural activity. We relied on perplexity estimation as an indicator of sequence predictability. While there have been studies that selected another measure from information theory (mainly entropy) to gauge predictability, each of them captures a specific sense of word predictability. Perplexity is a measure taken from the field of NLP that should arguably model sentence probability and was also utilized in cognition.²⁹ In these analyses, we trained four types of models to predict perplexity level based on neural activity elicited by a sentence: an ordinal classifier, a regression model, a binary classifier, and a within-participant classifier. While it is important to note that the results of these analyses (aimed to probe specifically relation to predictability) were less decisive (results of the ordinal classifier and within-participant classifier were not significant compared to chance-level control), the accumulation of the results certainly suggests that this is an area worth further exploration. Unfortunately, the within-participant predictability analysis, aimed to provide evidence of the relation of the effect to predictability, yielded null results. This may be reasonable with the low amount of data we have per participant (only 66 sentences). Considering that the perplexity binary classifier and perplexity regression (two analyses that had the same purpose) yielded significant results, a further investigation that specifically focuses on this perspective is needed. This could be accomplished in a new study that utilizes a within-participant design. It is worth mentioning that each condition has a completely different set of participants. Therefore, every analysis conducted on this dataset and compares the two different conditions ought to be between participants.

Overall, we provide supporting evidence that the difference in the reward system between conditions is not to be associated merely with changes in attention, arousal, engagement, the existence of narrative, or any other factor that inherently changes when scrambling sentences, by virtue of showing that (1) activity in the reward system is related to sentence-level information (sentences classifier, pattern similarity analysis) and (2) there might be a relation between elicited activity and predictability level (regression and binary classification of perplexity level). While the per-participant classifier and ordinal classification failed to predict the level of perplexity, thus hurting the claim that activity is related to prediction, it is reasonable to assume that this was due to the low amount of data and the design that was sub-optimal for this specific question.

The human brain constantly predicts what stimulus is expected based on current information.^{42,45} Prediction was shown to be a key component in sentence comprehension, mainly using reading times⁴⁶ and fMRI.^{23,29} The reward system is known to take part in prediction and the prediction error signal.⁵ However, a link between the reward system and predictability during language comprehension has yet to be found. Accordingly, we hypothesized a link between prediction and the activity in the reward system, and probed it here by trying to predict levels of statistical predictability of sentences in the intact condition based on activity in the reward areas. While using natural stimuli may increase complexity, add possible artifacts, or decrease effect size, it entails other advantages. Natural language studies often reveal much more widespread (and less left-lateralized) responses to language than studies with controlled stimuli.²⁶ Natural language studies also seem to be more sensitive as they could be used for exploring many different questions using a single dataset.^{26,46,47} showed that

Table 1. Pearson correlations between searchlight classification maps and Neurosynth term-based reverse inference activation maps obtained by uploading the map to Neurosynth

Term	Correlation (r)
Linguistic	0.424
Sentences	0.421
Language	0.413
Spoken	0.397
Comprehension	0.39

The 5 most highly correlated functional terms are listed.

different aspects of story processing were encoded in different brain networks. Therefore, although many previous studies used tightly controlled stimuli, such as a set of single isolated words, sentences or curated passages, and explicit lexical semantic tasks.^{47–49} We chose to use natural stories to provide initial evidence to our claim. However, this approach also has limitations in that the data are often not optimized to answer in-depth questions about specific attributes of the data. Accordingly, a study with a within-participant design could serve as a subsequent step to isolate and explore specific aspects that might moderate activity in the reward system. Such a study should be very carefully designed in order to control all other aspects and avoid semantic priming. This future work may use controlled stimuli set to better account for artifacts and other alternative explanations while precisely isolating and exploring the effect of predictability on the evoked activity and the extent to which semantic data are represented in reward areas. Such design could be fruitful in effectively identifying the exact role of the reward system in comprehension and the various aspects of the sentences that cause the evoked activity in the reward system. To precisely map which attributes of a sentence implicate higher involvement of the reward system, such design should control additional aspects of the sentences, such as sentiment, length, plausibility, etc. A dedicated design with forward inference manipulating specifically reward in the context of language will allow future studies to use a more precise functional mapping of the involved regions, obtained using a task localizer, thus enabling a direct masking procedure. Future work could also concentrate on the specific attributes of the reward-system involvement during processing, specifically the proposed representation of sentence-related information in the reward system, aiming to define what types of sentences are represented and how, and the relation to language representation and connectivity with this area during processing. As noted, a within-participant design, in which the same participant processes sentences from both conditions, should be very beneficial for answering this question, while also utilizing more standard univariate analysis. Specifically, the within-participant predictability analysis, which yielded null results in our paper, should benefit from such a design and should be of high importance.

In order to gain a more fine-grained anatomical understanding, it might be useful for future work with pre-registered hypotheses to study the involvement of specific substriatal regions and the differences between them. It will also enable us to better interpret the precise characteristics the reward system plays in language processing. In this regard, we performed post hoc analyses (See S5, S6) using one sub-region separation of the system (dorsal vs. ventral based on maps retrieved from the Harvard-Oxford atlas, such that ventral striatum is defined as the nucleus accumbens, and the dorsal striatum is the caudate and the putamen) and did not find a notable difference in the ability to classify between conditions but an indication for a difference in the way the pattern of activity changes between participants (S6).

Overall, we demonstrated for the first time the involvement of the reward system during natural language comprehension, showing that the reward system plays a role in story processing. We found evidence indicating that this role is related to predictability, but this hypothesis should be further explored. We propose an interpretation by which the more predictable nature of the intact condition caused this effect and that its response is related to the predictability level of the input, implying that the reward network activity might reflect a rewarding nature of language input with higher predictability. This interpretation should be further explored with more suitable experimental design. Future work should be done to better interpret the precise characteristics of the role the reward system plays in language processing, its specific mechanism, and the role of each of the sub-regions of the reward system (specifically the ventral and dorsal regions). Importantly, to better characterize the relation of evoked activity to predictability and characterize the possible influence of other attributes of the processed sentence, it might be necessary to create a new design (ideally within-subject) to ensure that any effect found between experimental conditions is to be attributed to the predictability of the input. Using a more controlled design can be fruitful in effectively identifying the role of the reward system in comprehension and the various aspects of the sentences that cause the evoked activity in the reward system. Such a study will also allow us to use a more precise functional map of the reward region, obtained using a task localizer, thus enabling more anatomically accurate data analysis.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- RESOURCE AVAILABILITY
 - Lead contact
 - Materials availability

- Data and code availability
- **METHOD DETAILS**
 - Stories and experimental design
 - Sentence predictability
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
 - Data reduction
 - Rationale
 - Whole-brain searchlight analysis

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2024.109844>.

ACKNOWLEDGMENTS

This work was supported by the ISRAEL SCIENCE FOUNDATION (ISF) grant number 1996/20.

AUTHOR CONTRIBUTIONS

O.K., Y.Y., and T.S. conceptualized, performed required investigation and research, and designed the required analyses. O.K. implemented the formal analysis under the guidance and supervision of Y.Y. and T.S. O.K., Y.Y., and T.S. wrote and reviewed the paper.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: July 13, 2023

Revised: February 18, 2024

Accepted: April 25, 2024

Published: April 29, 2024

REFERENCES

1. Smillie, L.D., and Wacker, J. (2014). Dopaminergic foundations of personality and individual differences [Editorial]. *Front. Hum. Neurosci.* *8*, 874. <https://doi.org/10.3389/fnhum.2014.00874>.
2. Schultz, W., Dayan, P., and Montague, P.R. (1997). A neural substrate of prediction and reward. *Science* *275*, 1593–1599.
3. Friston, K. (2010). The free-energy principle: A unified brain theory? *Nat. Rev. Neurosci.* *11*, 127–138. <https://doi.org/10.1038/nrn2787>.
4. Den Ouden, H.E.M., Kok, P., and de Lange, F.P. (2012). How prediction errors shape perception, attention, and motivation. *Front. Psychol.* *3*, 548.
5. Schultz, W. (2016). Dopamine reward prediction error coding. *Dialogues Clin. Neurosci.* *18*, 23–32. <https://doi.org/10.31887/DCNS.2016.18.1/wschultz>.
6. Crystal, T.H., and House, A.S. (1990). Articulation rate and the duration of syllables and stress groups in connected speech. *J. Acoust. Soc. Am.* *88*, 101–112. <https://doi.org/10.1121/1.399955>.
7. Liberman, A.M., Cooper, F.S., Shankweiler, D.P., and Studdert-Kennedy, M. (1967). Perception of the speech code. *Psychol. Rev.* *74*, 431–461.
8. Kuperberg, G.R., and Jaeger, T.F. (2016). What do we mean by prediction in language comprehension? *Lang. Cognit. Neurosci.* *31*, 32–59.
9. Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition* *106*, 1126–1177. <https://doi.org/10.1016/j.cognition.2007.05.006>.
10. Conway, C.M., Bauernschmidt, A., Huang, S.S., and Pisoni, D.B. (2010). Implicit statistical learning in language processing: Word predictability is the key. *Cognition* *114*, 356–371. <https://doi.org/10.1016/j.cognition.2009.10.009>.
11. Saffran, J.R. (2003). Statistical language learning: Mechanisms and constraints. *Curr. Dir. Psychol. Sci.* *12*, 110–114. <https://doi.org/10.1111/1467-8721.01243>.
12. Grisoni, L., Miller, T.M., and Pulvermüller, F. (2017). Neural correlates of semantic prediction and resolution in sentence processing. *J. Neurosci.* *37*, 4848–4858.
13. Boux, I., Tomasello, R., Grisoni, L., and Pulvermüller, F. (2021). Brain signatures predict communicative function of speech production in interaction. *Cortex* *135*, 127–145.
14. León-Cabrera, P., Flores, A., Rodríguez-Fornells, A., and Moris, J. (2019). Ahead of time: Early sentence slow cortical modulations associated to semantic prediction. *Neuroimage* *189*, 192–201.
15. Lau, E.F., Phillips, C., and Poeppel, D. (2008). A cortical network for semantics: (De)constructing the N400. *Nat. Rev. Neurosci.* *9*, 920–933. <https://doi.org/10.1038/nrn2532>.
16. Stringaris, A.K., Medford, N., Giora, R., Giampietro, V.C., Brammer, M.J., and David, A.S. (2006). How metaphors influence semantic relatedness judgments: The role of the right frontal cortex. *Neuroimage* *33*, 784–793.
17. Hagoort, P. (2005). On Broca, brain, and binding: a new framework. *Trends Cognit. Sci.* *9*, 416–423. <https://doi.org/10.1016/j.tics.2005.07.004>.
18. Kuperberg, G.R., Sitnikova, T., and Lakshmanan, B.M. (2008). Neuroanatomical distinctions within the semantic system during sentence comprehension: Evidence from functional magnetic resonance imaging. *Neuroimage* *40*, 367–388. <https://doi.org/10.1016/j.neuroimage.2007.10.009>.
19. Hickok, G., and Poeppel, D. (2007). The cortical organization of speech processing. *Nat. Rev. Neurosci.* *8*, 393–402.
20. Hasson, U., Chen, J., and Honey, C.J. (2015). Hierarchical process memory: memory as an integral component of information processing. *Trends Cognit. Sci.* *19*, 304–313.
21. Lerner, Y., Honey, C.J., Silbert, L.J., and Hasson, U. (2011). Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *J. Neurosci.* *31*, 2906–2915.
22. Yeshurun, Y., Nguyen, M., and Hasson, U. (2017). Amplification of local changes along the timescale processing hierarchy. *Proc. Natl. Acad. Sci. USA* *114*, 9475–9480.
23. Frank, S.L., and Willems, R.M. (2017). Word predictability and semantic similarity show distinct patterns of brain activity during language comprehension. *Lang. Cognit. Neurosci.* *32*, 1192–1203.
24. Willems, R.M., Frank, S.L., Nijhof, A.D., Hagoort, P., and Van Den Bosch, A. (2016). Prediction during natural language comprehension. *Cerebr. Cortex* *26*, 2506–2516. <https://doi.org/10.1093/cercor/bhw075>.
25. Tikochinski, R., Goldstein, A., Yeshurun, Y., Hasson, U., and Reichart, R. (2023).

- Perspective changes in human listeners are aligned with the contextual transformation of the word embedding space. *Cerebr. Cortex* 33, 7830–7842.
26. Hamilton, L.S., and Huth, A.G. (2020). The revolution will not be controlled: natural stimuli in speech neuroscience. *Lang. Cognit. Neurosci.* 35, 573–582.
 27. de Heer, W.A., Huth, A.G., Griffiths, T.L., Gallant, J.L., and Theunissen, F.E. (2017). The hierarchical cortical organization of human speech processing. *J. Neurosci.* 37, 6539–6557.
 28. Huth, A.G., De Heer, W.A., Griffiths, T.L., Theunissen, F.E., and Gallant, J.L. (2016). Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature* 532, 453–458.
 29. Lopotopolo, A., Frank, S.L., van den Bosch, A., and Willems, R.M. (2017). Using stochastic language models (SLM) to map lexical, syntactic, and phonological information processing in the brain. *PLoS One* 12, e0177794. <https://doi.org/10.1371/journal.pone.0177794>.
 30. Martin, J.H. (2009). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition* (Pearson/Prentice Hall).
 31. Hale, J., Kuncoro, A., Hall, K., Dyer, C., and Brennan, J. (2019, November). Text genre and training data size in human-like parsing. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5846–5852.
 32. Nastase, S.A., Liu, Y.F., Hillman, H., Zadbood, A., Hasenfratz, L., Keshavarzian, N., Chen, J., Honey, C.J., Yeshurun, Y., Regev, M., et al. (2021). The “Narratives” fMRI dataset for evaluating models of naturalistic language comprehension. *Sci. Data* 8, 250. <https://doi.org/10.1038/s41597-021-01033-3>.
 33. Kaltwasser, L., Ries, S., Sommer, W., Knight, R.T., and Willems, R.M. (2013). Independence of valence and reward in emotional word processing: electrophysiological evidence. *Front. Psychol.* 4, 168.
 34. De Loof, E., Ergo, K., Naert, L., Janssens, C., Talsma, D., Van Opstal, F., and Verguts, T. (2018). Signed reward prediction errors drive declarative learning. *PLoS One* 13, e0189212.
 35. Shibata, M., Terasawa, Y., and Umeda, S. (2014). Integration of cognitive and affective networks in humor comprehension. *Neuropsychologia* 65, 137–145.
 36. Mobbs, D., Greicius, M.D., Abdel-Azim, E., Menon, V., and Reiss, A.L. (2003). Humor modulates the mesolimbic reward centers. *Neuron* 40, 1041–1048.
 37. Bohrn, I.C., Altmann, U., Lubrich, O., Menninghaus, W., and Jacobs, A.M. (2013). When we like what we know—a parametric fMRI analysis of beauty and familiarity. *Brain Lang.* 124, 1–8.
 38. Ripollés, P., Marco-Pallarés, J., Hielscher, U., Mestres-Missé, A., Tempelmann, C., Heinze, H.J., Rodríguez-Fornells, A., Noesselt, T., and Noesselt, T. (2014). The role of reward in word learning and its implications for language acquisition. *Curr. Biol.* 24, 2606–2611.
 39. Ripollés, P., Marco-Pallares, J., Alicart, H., Tempelmann, C., Rodríguez-Fornells, A., and Noesselt, T. (2016). Intrinsic monitoring of learning success facilitates memory encoding via the activation of the SN/VTA-Hippocampal loop. *Elife* 5, e17441.
 40. Kischka, U., Kammer, T., Maier, S., Weisbrod, M., Thimm, M., and Spitzer, M. (1996). Dopaminergic modulation of semantic network activation. *Neuropsychologia* 34, 1107–1113.
 41. Tiedt, H.O., Ehlen, F., and Klostermann, F. (2022). Dopamine-Related Reduction of Semantic Spreading Activation in Patients With Parkinson’s Disease. *Front. Hum. Neurosci.* 16, 837122.
 42. Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behav. Brain Sci.* 36, 181–204.
 43. Friston, K. (2005). A theory of cortical responses. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 360, 815–836.
 44. Ehrlich, S.F., and Rayner, K. (1981). Contextual effects on word perception and eye movements during reading. *J. Verb. Learn. Verb. Behav.* 20, 641–655.
 45. Frank, S.L., Otten, L.J., Galli, G., and Vigliocco, G. (2015). The ERP response to the amount of information conveyed by words in sentences. *Brain Lang.* 140, 1–11.
 46. Goodkind, A., and Bicknell, K. (2018, January). Predictive power of word surprisal for reading times is a linear function of language model quality. In *Proceedings of the 8th workshop on cognitive modeling and computational linguistics (CMCL 2018)*, pp. 10–18.
 47. Wehbe, L., Murphy, B., Talukdar, P., Fyshe, A., Ramdas, A., and Mitchell, T. (2014). Simultaneously uncovering the patterns of brain regions involved in different story reading subprocesses. *PLoS One* 9, e112575.
 48. Chee, M.W., O’Craven, K.M., Bergida, R., Rosen, B.R., and Savoy, R.L. (1999). Auditory and visual word processing studied with fMRI. *Hum. Brain Mapp.* 7, 15–28.
 49. Buchweitz, A., Mason, R.A., Tomitch, L.M.B., and Just, M.A. (2009). Brain activation for reading and listening comprehension: An fMRI study of modality effects and individual differences in language comprehension. *Psychol. Neurosci.* 2, 111–123.
 50. Yarkoni, T., Poldrack, R.A., Nichols, T.E., Van Essen, D.C., and Wager, T.D. (2011). Large-scale automated synthesis of human functional neuroimaging data. *Nat. Methods* 8, 665–670. <https://doi.org/10.1038/nmeth.1635>.
 51. Marslen-Wilson, W.D. (1975). Sentence Perception as an Interactive Parallel Process. *Science* 189, 226–228. <https://doi.org/10.1126/science.189.4198.226>.
 52. Tipping, M.E., and Bishop, C.M. (1999). Probabilistic principal component analysis. *J. Roy. Stat. Soc. B Stat. Methodol.* 61, 611–622.
 53. Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. Roy. Stat. Soc. B* 57, 289–300.
 54. Chen, J., Leong, Y.C., Honey, C.J., Yong, C.H., Norman, K.A., and Hasson, U. (2017). Shared memories reveal shared structure in neural activity across individuals. *Nat. Neurosci.* 20, 115–125. <https://doi.org/10.1038/nn.4450>.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Software and algorithms		
Python	Python	RRID:SCR_008394
Other		
fMRI dataset	Nastase et al. 2021 ³²	https://openneuro.org/datasets/ds002345/versions/1.1.4

RESOURCE AVAILABILITY

Lead contact

Further information for resources should be directed to Oren Kobo (oren.kobo@gmail.com).

Materials availability

This study did not generate new unique reagents.

Data and code availability

- This paper uses the publicly available. The link is available in the [key resources table](#).
- All original analysis code has been uploaded to github and are publicly available as of the date of publication at https://github.com/orenpapers/Reward_Predictability_Paper.
- Any information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

METHOD DETAILS

This study was pre-registered at: <https://osf.io/mt7yd/>. Analyses outside of the pre-registration are mentioned specifically in the text. We performed all analyses within a pre-registered reward mask and a control area (within the visual network), both obtained from [Neurosynth.org](https://neurosynth.org)⁵⁰ by entering the corresponding terms in the meta-analyses tool and downloading the association map. We used classification and computational models to probe whether there was a different response in the reward network to the experimental conditions of a natural story (intact) and a scrambled version of it (scrambled). We further aimed to inspect the relatedness of these findings specifically to predictability.

Stories and experimental design

The current study reanalyzes a previously published fMRI dataset²² that was shared as part of the narrative dataset.⁵⁰ In the dataset reanalyzed in this study, 36 right-handed participants listened inside the MRI scanner to one of two narratives: 18 participants listened to the *intact* story and 18 other participants listened to the *scrambled* story. The *intact* story was about a woman that was obsessed with an American Idol judge, meets a psychic, and then becomes fixated on Vodka. The *scrambled* story had the exact same words as the intact, but the words within each sentence were randomly scrambled (The below image Edited from (Yeshurun et al., 2017)²²). The order of the sentences remained the same between the conditions as in the original study. The difference was that in the scrambled condition it had a different ordering of a few of the words within a sentence. This manipulation created an implausible input throughout the *scrambled* story. Both stories were read and recorded by the same actor. The beginning of each sentence was aligned post-recording. Each story was 6:44 minutes long and was preceded by 18 seconds of neutral music and 3 seconds of silence. Stories were followed by an additional 15 seconds of silence. These music and silence periods were discarded from all analyses. Experimental procedures were approved by the Princeton University Committee on Activities Involving Human Participants. All participants provided written informed consent.

Experimental Conditions

Intact	Scrambled
The situation seemed too good , and after she spent most of the last weeks worrying , she knew she could enjoy the moment . She looked at her reflection in the mirror and started laughing uncontrollably ...	Too the seemed weeks spent, the knew good she the enjoy last she situation most, after and could worrying she moment of. Mirror uncontrollably started looking she in the at her reflection laughing and ...

Experimental conditions

We used 2 of the 4 conditions of the original study. Intact story had the same words as scrambled, but the words were scrambled within a sentence. Thus, there was no words change, but scrambled words that resulted in unpredictable input.

Sentence predictability

The effect of prediction on comprehension and generation of expectation has been shown in the syntactic⁹ and phonetic⁵¹ expectations. In the current study, we chose to focus on semantic predictability (the probability of a certain word in context regardless of the grammatical structure). This approach was used in various cases to formalize predictability by statistical expectation, for example by using statistical measures such as transitional probability and the likelihood of two words to co-occur.

We measure sequence predictability using perplexity. The perplexity (PP) of a language model on a test set is the inverse probability of the test set, normalized by the number of words (See Equation 1). Thus, minimizing perplexity is equivalent to maximizing the test set probability according to the language model.³⁰

$$PP(W) = \sqrt[N]{\prod_i^N \frac{1}{P(w_i|w_1 \dots w_{i-1})}} \quad (\text{Equation 1})$$

QUANTIFICATION AND STATISTICAL ANALYSIS

In all the analyses, we calculated an accuracy level for the intact condition by multiple iterations of training on a training set and testing on a different set (test set). Thus, we generated a global accuracy score and a distribution of accuracy per iteration. To obtain statistical significance, we compared actual accuracies to a corresponding controlled null distribution, generated by either using the data from the scrambled condition or by shuffling the actual labels, depending on the specific question in hand. When relevant, we ran the exact same analysis on the preregistered control area (visual-related region), to ensure any effect that was detected in the reward region does not replicate in areas that are not involved. We set the significance level in our analyses to $p=0.05$ and added FDR multiple comparisons correction when relevant. All statistical tests were performed in Python.

Data reduction

In our ML-based Analyses (all besides the Euclidean distance and pattern similarity), we performed a data reduction process to decrease the number of features (original data was at the voxel level, meaning the activity at each specific voxel is an input feature for the model. There were 12031 voxels in the reward mask, which is far too many for standard ML tools for this amount of data) The method we used for dimensionality reduction was Principal Component Analysis (PCA).⁵² In all cases, the PCA was trained on the train data and then was used to transform the test data, to prevent no danger of leakage.

Rationale

The analyses performed in this manuscript can be divided into three general goals: A) Showing that the reward system is involved in processing of natural stories in some aspect by demonstrating that the activity is modulated by the experimental condition. B) Once this is established, we aimed to show that during processing, there is some extent of sentence-level information that is encoded in the reward system. C) We aimed to explore the relation between the activity and the predictability of the input.

Differentiating between activity in intact vs. scrambled conditions

We tested if there was a difference in the evoked neural response between the intact and scrambled conditions in reward-related areas. If there is indeed such a difference, this suggests that the reward system operates differently between conditions, therefore involved in processing and affected by the difference between the input of the conditions (intact and scrambled). For that, we applied the following analyses:

1. Condition classifier at the participant level

In this pre-registered analysis, we trained a support vector machine model (SVM) classifier at the participant level to predict a condition given participants' activity. The rationale of this analysis was to show that there was a difference in the activity of the reward system that was moderated by the condition. If the reward system was completely agnostic to the language processing task, we would not have been able to perform such a classification. For that purpose, a participant activity vector (sized 1×269 , as 269 is the number of TRs) was created by averaging participants' responses across sentences. Thus, this vector provides a representation of the neural activity of the participant during the story (See below image). We then ran 1000 iterations with 5-fold cross-validation. In each iteration, we trained the model on the training-set participants and evaluated it on the test-set participants, as assigned by the cross-validation. Due to the high dimensionality of the data relative to the number of samples, we first applied a PCA with 20 components to the data. Then, we obtained p-values using a permutation test (artificially generated by running 1000 iterations of the model with shuffled labels). We applied the exact same procedure for the pre-determined control vision area.

2. Euclidean Distance Measure

In this pre-registered analysis, we aimed to evaluate the extent to which activity in the reward system alternates between conditions. To do so, we used the per-voxel Euclidean distance as a measure of the difference between the two conditions.²² Like the previous analysis, the purpose of the analysis was to demonstrate a different activity between the experimental conditions in the reward system throughout the task. Therefore, it demonstrates the involvement of the reward system in the task. We measured the significance of the difference between the activity of the 2 conditions, per voxel, as follows: For each condition, we created a matrix of mean activity (across participants) per TR, per voxel – which resulted in 2 matrices with the shape of $n_TR \times n_voxels$. Then, for each voxel, we calculated the Euclidean distance between its respective time series – yielding a distance vector with the shape of $1 \times n_Voxels$, that is a representation of the overall activity of the voxel throughout the task. Thus, at each point, the vector stores the Euclidean distance (over time) of the averaged (across participants) distance. Afterward, we ran 1000 iterations of label shuffling, to obtain the null distribution of 1000 artificial distance vectors. Finally, for each voxel, we calculated the rank of the real value compared to the null distribution and calculated z-scores based on this rank, and adjusted the p-value for each voxel using False Discovery Rate method.⁵³

Sentence-Level information in the reward system

We tested whether there was any information regarding some properties of the sentences themselves that were encoded in the reward system. Such a finding would support the suggestion that the reward system was involved in the processing of stories by showing it is related to the perceived language information. This is as opposed to a possible alternative explanation - that the difference in plausibility or engagement between conditions can explain the different pattern activity across conditions (the finding that was shown in the analyses of section A). We ran the following analyses:

1. Pattern similarity between sentences

We tested if the same sentence elicited similar activity across all participants in the reward system by using an analysis that measures between-participants spatial similarity.⁵⁴ We iterated over all sentences and all participants. In each iteration, we calculated the correlation between the activity from the specific pairings of participants and the sentence to the activity elicited by the same sentence in all other participants. Additionally, to create a baseline null distribution, an additional sentence was randomly selected in every iteration and its correlation with the original sentence was calculated (see above image). We sampled only one sentence per iteration for the distributions to be equal in size as required by the parametric significance test we used. The results obtained when comparing to the distribution of all possible pairs are in the Supplementary (Figure S3).

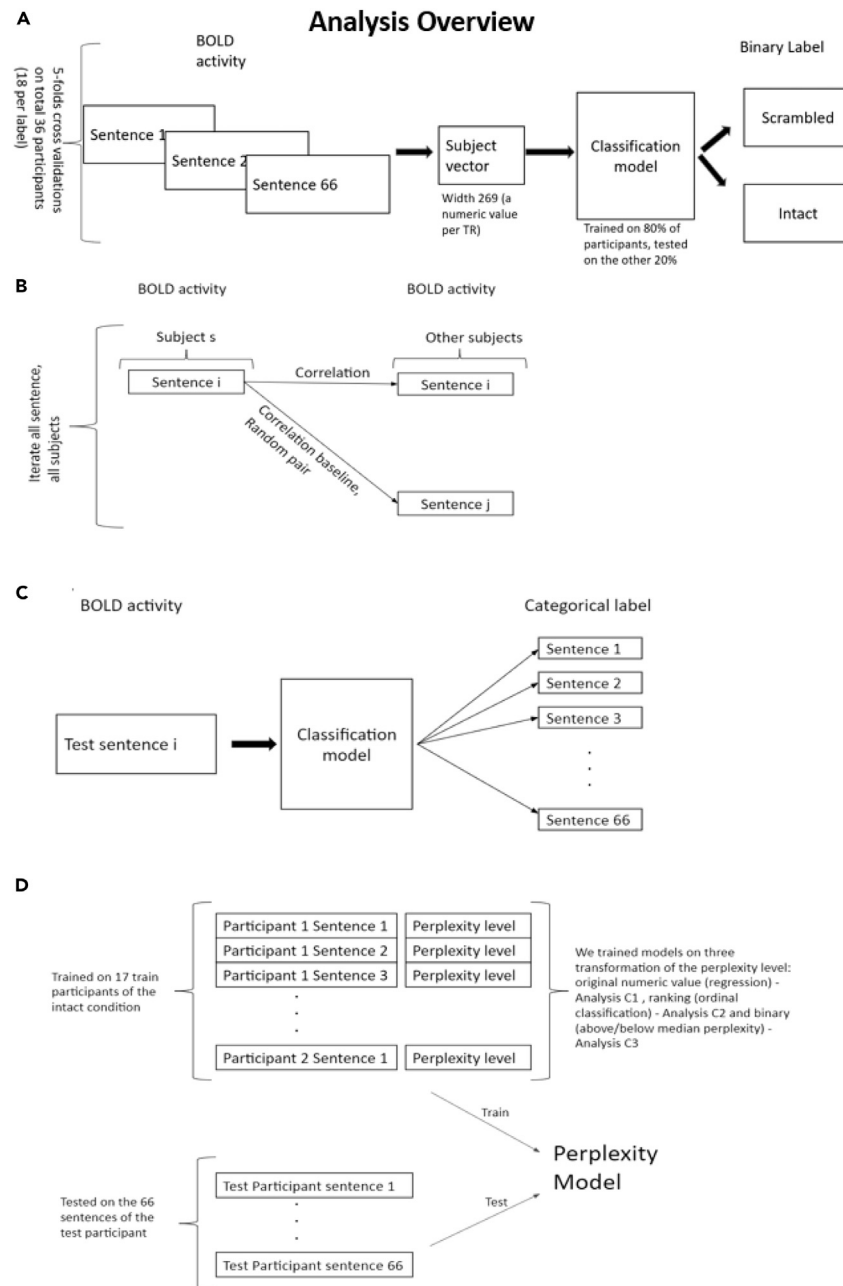
2. Multisentences classifier

We measured the extent of sentence-level information encoded in the reward system during processing by probing if we could train a classifier to predict the sentence that is processed based on the elicited neural activity. We ran iterations of Leave-One-Participant-Out and trained a multilabel classifier (66 labels) on 17 train set participants. We then used the trained model to classify the sentences of the left-out participant based on the neural activity (such that 66 sentences were predicted, each according to the neural activity elicited by a specific sentence) and obtained an accuracy score (see below image). Similarly to the condition classifier, we applied PCA with 20 components followed by an SVM and generated a null distribution with 1000 iterations of randomized labels, in order to obtain the significance level.

Probing the relation between activity and predictability

We further conducted the following analyses to accommodate the hypothesis that activity in the reward system during processing was related specifically to the predictability of the linguistic input, and not to other possible artifacts. Accordingly, these analyses were designed to characterize specifically the relation between activity in the reward system and sentence predictability (see below image). We used the following control in these analyses to ensure reliability: First, intact scores were compared to the same data with shuffled labels (null distribution). Then, to demonstrate the effect is related to perplexity and not to the words themselves, we performed the same analysis for the scrambled condition. We also repeated the process for the visual system, to verify that any effect in the reward system is not part of the whole-brain activity.

1. Ordinal Classification of sentence perplexity



Overview of analysis methods used in the manuscript

(A) The process of creating a classifier for each participant's condition. We generated a single vector per participant with a width of 269 (number of TRs). This was used as the input for the model, and the output for a binary label that was the condition.

(B) Between-participants pattern similarity of sentences. For each sentence, neural data were averaged across time within each voxel, resulting in a single representing vector per sentence. Correlations were computed between sentences across participants. The null distribution was generated by choosing a random sentence, in addition to the actual pair.

(C) Multiclass sentence classifier. In each fold, we trained the model to predict one of the 66 experimental sentences based on elicited neural activity. We then predicted the correct labels for each of the 66 sentences of the test participant.

(D) Perplexity analysis. We added several types of perplexity models. In each of the iterations, we trained a model to predict the perplexity level of 17 train participants (the analysis was done only on the intact condition) and used it to predict the perplexity level of the test participant. We tried three types of labels: numerical, ordinal, and binary.

We used ordinal classification to predict the rank of sentence predictability within the sentence vector. We chose the approach of ordinal classification as the variable can be viewed as having an arbitrary scale where only the relative ordering between different values is significant. Similarly to the Multisentences classifier described above, here we also ran iterations of Leave-One-Out participants in the intact condition. We trained a multilabel classifier (66 labels) on 17 train-set participants and tested the prediction on the remaining test participants. We reduced the dimensionality of the neural data using PCA. The ordinality was implemented as follows: We trained the model to predict if a sentence belongs to a label. Then, to obtain ordinality, we aggregated probabilities such that $\Pr(V_i) = \Pr(y > V_{i-1}) - \Pr(y > V_i)$ for each class V_i . The final prediction is then the label with the highest probability. We used mean squared error (MSE) between the predicted and actual label as the evaluation metric. Since each sentence in the scrambled condition was a direct variant of the corresponding sentence in the intact condition, they both had the exact same words. With that in mind, after creating a null distribution by 1000 iterations of labels shuffling, we used the scrambled condition as a baseline to ensure that any effect between actual data and null distribution does not replicate also for the scrambled data. Therefore, we can conclude that the effect is to be attributed to the sentence and its characteristics as a whole, and not to the appearance of the specific words in the input. We also repeated the entire process in the control area (the visual system).

2. Binary perplexity classifier

In this pre-registered analysis we trained a binary classifier to predict the perplexity of a sentence (above or below the median) given elicited neural data from the reward system in the intact condition and compared results to the baseline of shuffled labels.

For each condition (intact, scrambled) and ROI (reward, vision) we ran 18 iterations of left-one-participant-out. We stacked all the per-sentence neural response signals from all the participants together. We then divided each sentence into low and high predictability (above/below median) using perplexity scores. We used PCA for dimensionality reduction. We calculated accuracy based on predictions of the left-out participant at each iteration.

3. Regression of sentence perplexity

In this non-pre-registered analysis, we trained a regression model to predict the level of perplexity based on the neural activity this sentence elicited. We applied a principal component analysis (PCA) on the data to reduce dimensionality and fitted a regression model. At each iteration, we left one participant out and used it as the test. We ran 400 iterations of labels shuffling to create a null distribution. We repeated the same process for the visual system as control. We evaluated to model using R-squared.

4. Within participant perplexity classifier

In this pre-registered analysis, we divided the sentences into low and high predictability (above/below median) based on perplexity. Then, for each participant's data (meaning we trained a model per participant), we trained and evaluated 1000 iterations with a 20% test set (20% of the sentences were used for test) and used it to yield an accuracy level. We obtained p-values by comparing them to the null distribution accuracy. This was done separately for reward and visual areas.

Whole-brain searchlight analysis

We also applied an exploratory analysis, in which we applied the Euclidean Distance measure procedure at the whole-brain level (and not within the specific pre-registered ROIs). For each voxel, we evaluated the extent to which it was indicative of one of the two conditions by assessing the z-score of the distance between neural data of the conditions and the per-voxel distance that was generated by a null distribution. To establish the cognitive functions in which these regions were most involved, we conducted a formal reverse inference analysis using NeuroSynth,⁵⁰ correlating our searchlight map with the neural activation maps for each term in the NeuroSynth database.