

# Reading Time Analysis Report

Ayush Kumar Gupta  
2023114001

February 17, 2026

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Code Availability</b>   | <b>3</b>  |
| <b>2</b> | <b>Introduction</b>  | <b>3</b>  |
| <b>3</b> | <b>Dataset Description</b>   | <b>3</b>  |
| <b>4</b> | <b>Part I: Exploratory Data Analysis</b>                               | <b>4</b>  |
| 4.1      | Mean Reading Time per Word . . . . .                                   | 4         |
| 4.2      | Word Length vs Mean Reading Time . . . . .                             | 5         |
| 4.3      | Word Frequency vs Mean Reading Time . . . . .                          | 5         |
| 4.4      | Pearson Correlation Analysis . . . . .                                 | 6         |
| 4.5      | Summary of Part I Findings . . . . .                                   | 7         |
| <b>5</b> | <b>Part II: Hypothesis Testing</b>                                     | <b>8</b>  |
| 5.1      | Hypothesis 1: Language Model Probabilities vs Word Frequency . . . . . | 8         |
| 5.1.1    | Data Alignment . . . . .   | 8         |
| 5.1.2    | Discussion of Hypothesis 1 . . . . .                                   | 9         |
| 5.2      | Hypothesis 2: Content Words vs Function Words . . . . .                | 10        |
| 5.2.1    | Discussion of Hypothesis 2 . . . . .                                   | 11        |
| <b>6</b> | <b>Part III: Frequency Ordered Bin Search (FOBS) Model</b>             | <b>12</b> |
| 6.1      | FOBS Model Construction . . . . .                                      | 12        |
| 6.1.1    | Lemmatization . . . . .  | 12        |
| 6.1.2    | FOBS Bin Distribution . . . . .  | 13        |
| 6.2      | Hypothesis 1: Root Frequency vs Surface Frequency . . . . .            | 13        |
| 6.2.1    | FOBS Search Depth Correlation . . . . .                                | 14        |
| 6.2.2    | Discussion of FOBS Hypothesis 1 . . . . .                              | 15        |
| 6.3      | Hypothesis 2: Pseudo-Affixed vs Real Affixed Words . . . . .           | 15        |
| 6.3.1    | Frequency and Length Matching . . . . .                                | 15        |
| 6.3.2    | Test Words . . . . .   | 16        |
| 6.3.3    | Results . . . . .  | 16        |
| 6.3.4    | Discussion of FOBS Hypothesis 2 . . . . .                              | 18        |
| 6.4      | Summary of FOBS Findings . . . . .                                     | 19        |
| <b>7</b> | <b>Overall Conclusions</b>   | <b>19</b> |

|  |           |
|--|-----------|
| <b>8 Methodology Notes</b>             | <b>20</b> |
| 8.1 Tools and Libraries . . . . .      | 20        |
| 8.2 Data Processing Pipeline . . . . . | 20        |

# 1 Code Availability

All scripts and data used in this report can be found at the following GitHub repository:

[https://github.com/AKGIIITH/Computational\\_Psycholinguistics/tree/main/Assignment/Assignment3](https://github.com/AKGIIITH/Computational_Psycholinguistics/tree/main/Assignment/Assignment3)

The repository contains three analysis scripts:

- `Code/data_analysis.py` — Exploratory data analysis (Part I)
- `Code/hypothesis_testing.py` — Hypothesis testing (Part II)
- `Code/frequency_ordered_bin_search.py` — FOBS model (Part III)

# 2 Introduction

The Natural Stories Corpus is a psycholinguistic dataset comprising 10 naturalistic English stories designed to contain rare and varied syntactic constructions. This report presents a comprehensive analysis of self-paced reading times (RTs) collected from multiple participants, examining the relationships between word-level properties—word length, word frequency, and language model probabilities—and human reading behavior. We further investigate morphological processing through the lens of a Frequency Ordered Bin Search (FOBS) memory model.

The analysis is structured in three parts:

1. **Part I:** Exploratory data analysis of reading times, word length, and word frequency.
2. **Part II:** Hypothesis testing comparing word frequency and GPT-3 language model probabilities as predictors of reading time, and examining differences between content and function word processing.
3. **Part III:** Construction of a FOBS memory model, comparison of surface vs. lemma frequency as RT predictors, and analysis of pseudo-affixed vs. real affixed word processing.

# 3 Dataset Description

The Natural Stories Corpus contains the following key components used in this analysis:

- **processed\_RTs.tsv:** Self-paced reading times from multiple participants across 10 stories. RTs below 100 ms and above 3000 ms were filtered, and only participants scoring > 4/6 on comprehension questions were retained.
- **freqs/freqs-1.tsv:** Unigram frequency data from the Google Books English corpus (1990–present).

- **probs/all\_stories\_gpt3.csv**: Token-level log-probabilities from the GPT-3 language model.
- **words.tsv**: Token-level alignment codes linking all data sources.

Table 1: Dataset Summary Statistics

| Statistic | Word Length | Frequency             | Log Frequency | Mean RT (ms) |
|-----------|-------------|-----------------------|---------------|--------------|
| Count     | 32,342      | 32,342                | 32,342        | 32,342       |
| Mean      | 4.71        | $1.71 \times 10^9$    | 7.97          | 339.31       |
| Std       | 2.48        | $3.17 \times 10^9$    | 1.46          | 44.14        |
| Min       | 1           | 0                     | 0.00          | 256.86       |
| 25%       | 3           | $9.23 \times 10^6$    | 6.96          | 311.68       |
| Median    | 4           | $1.34 \times 10^8$    | 8.13          | 330.59       |
| 75%       | 6           | $1.14 \times 10^9$    | 9.06          | 356.31       |
| Max       | 23          | $1.20 \times 10^{10}$ | 10.08         | 873.11       |

The dataset comprises 848,875 individual RT records across all participants, yielding 10,256 unique word instances (identified by story  $\times$  position) after averaging across subjects. After merging with frequency data (which includes sub-token alignments), 32,342 word-frequency-RT records were obtained.

## 4 Part I: Exploratory Data Analysis

### 4.1 Mean Reading Time per Word

For each word token in the RT file, the average reading time across all subjects was computed. Table 2 presents a sample of the first 10 words from Story 1 with their mean RTs.

Table 2: Sample of Mean RT per Word (Story 1, first 10 words)

| Story | Zone | Word     | Mean RT (ms) |
|-------|------|----------|--------------|
| 1     | 1    | If       | 578.96       |
| 1     | 2    | you      | 369.01       |
| 1     | 3    | were     | 368.18       |
| 1     | 4    | to       | 344.32       |
| 1     | 5    | journey  | 354.64       |
| 1     | 6    | to       | 349.67       |
| 1     | 7    | the      | 376.37       |
| 1     | 8    | North    | 327.31       |
| 1     | 9    | of       | 365.49       |
| 1     | 10   | England, | 344.93       |

The first word “If” shows a notably elevated RT (578.96 ms) compared to subsequent words, consistent with the well-documented sentence-initial slowdown effect. The overall mean RT across all word instances is 339.31 ms ( $SD = 44.14$ ).

## 4.2 Word Length vs Mean Reading Time

Figure 1 presents the relationship between word length (in characters) and mean reading time.

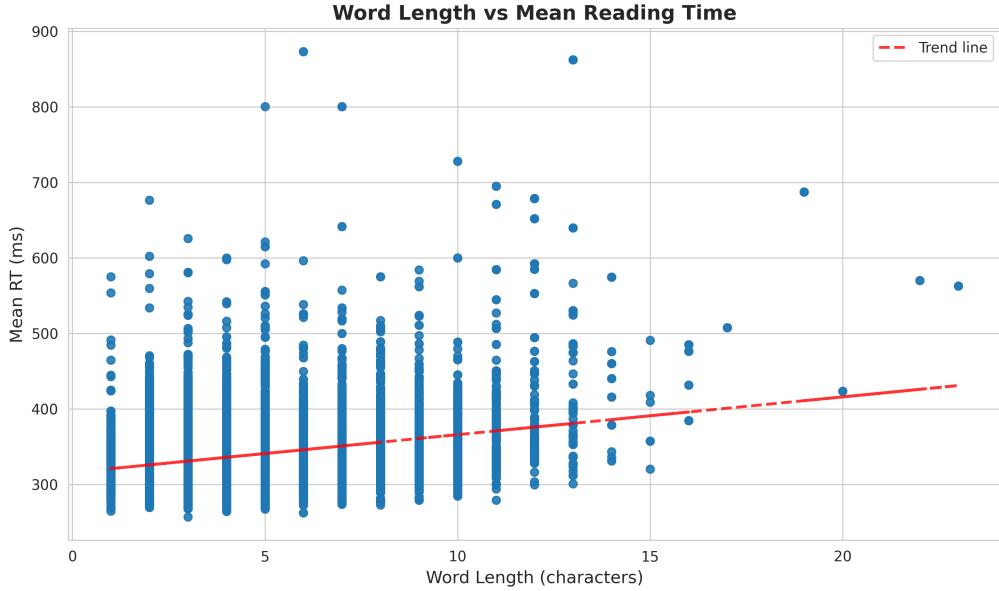


Figure 1: Scatter plot of word length vs. mean reading time with linear trend line. Longer words tend to require more processing time.

The plot reveals a clear positive trend: longer words are associated with higher mean reading times. The relationship is approximately linear, though considerable variance exists at each word length, reflecting the influence of other factors such as frequency, predictability, and syntactic context. Word lengths range from 1 to 23 characters.

## 4.3 Word Frequency vs Mean Reading Time

Figure 2 displays the relationship between log word frequency and mean reading time.

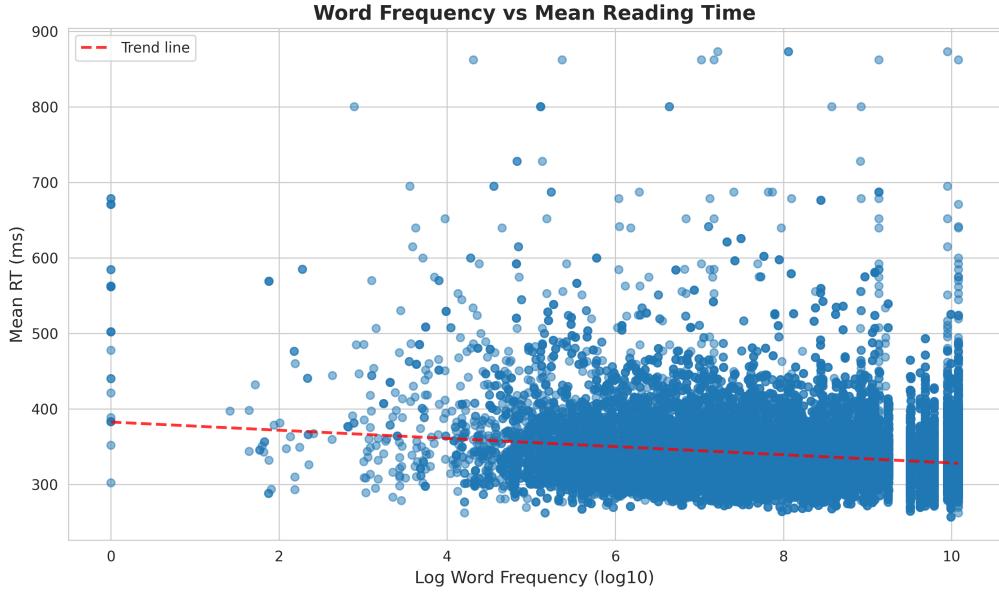


Figure 2: Scatter plot of log word frequency vs. mean reading time with linear trend line. Higher frequency words are read faster.

A negative trend is observed: more frequent words are associated with shorter reading times. This is consistent with the well-established frequency effect in psycholinguistic research—words encountered more often in language are accessed more quickly from the mental lexicon. The bulk of the data clusters in the high-frequency range (log frequency 6–10), with sparse data points at very low frequencies.

#### 4.4 Pearson Correlation Analysis

Table 3 presents Pearson's correlation coefficients for all pairwise relationships.

Table 3: Pearson's Correlation Coefficients

| Variable 1    | Variable 2    | Pearson's $r$ | $p$ -value              | Interpretation    |
|---------------|---------------|---------------|-------------------------|-------------------|
| Word Length   | Log Frequency | -0.5937       | $< 10^{-300}$           | Strong negative   |
| Word Length   | Mean RT       | +0.2812       | $< 10^{-300}$           | Moderate positive |
| Log Frequency | Mean RT       | -0.1795       | $2.80 \times 10^{-232}$ | Weak negative     |

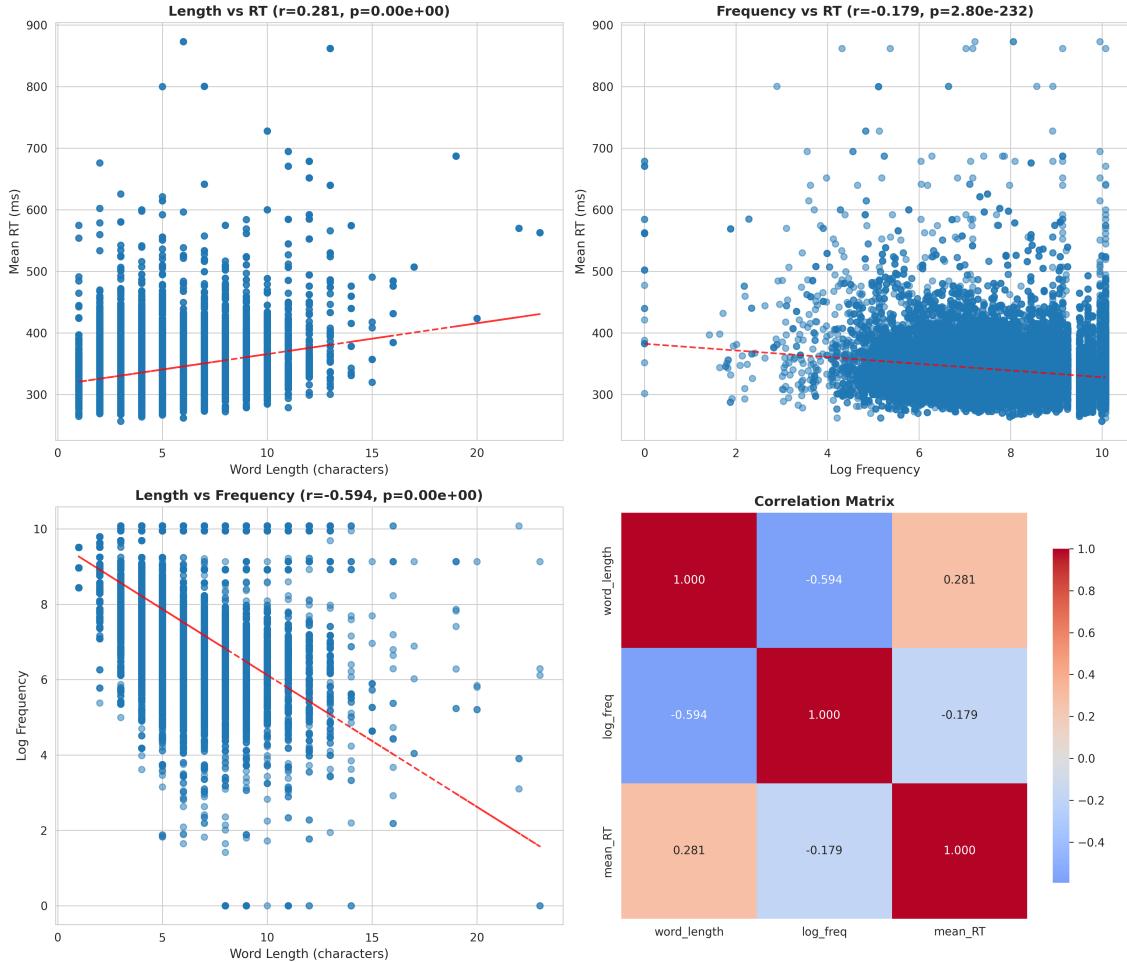


Figure 3: Comprehensive visualization: (Top-left) Word length vs. RT, (Top-right) Log frequency vs. RT, (Bottom-left) Word length vs. log frequency, (Bottom-right) Correlation matrix heatmap.

## 4.5 Summary of Part I Findings

The exploratory analysis reveals three key relationships:

- Word Length and Frequency ( $r = -0.594$ ):** There is a strong negative correlation between word length and frequency. Shorter words tend to be more frequent in the language, which is consistent with Zipf's Law of Abbreviation—frequently used words tend to be shorter.
- Word Length and Mean RT ( $r = +0.281$ ):** A moderate positive correlation indicates that longer words require more processing time. This reflects the additional time needed for visual encoding and lexical access of longer orthographic forms. Each additional character adds approximately 4.80 ms to reading time (from the regression coefficient).
- Word Frequency and Mean RT ( $r = -0.179$ ):** A weak but highly significant negative correlation confirms the frequency effect—high-frequency words are read faster. However, the relatively weak correlation suggests that frequency alone ex-

plains only a small portion of RT variance, and other factors (syntactic complexity, predictability, spillover effects) play substantial roles.

4. **Confound between Length and Frequency:** The strong negative correlation between length and frequency ( $r = -0.594$ ) means these predictors are confounded. Short, high-frequency words are read fastest, and it is difficult to disentangle the independent contributions of each factor without regression modeling.

All three correlations are highly significant ( $p < 10^{-200}$ ), reflecting the large sample size ( $N = 32,342$ ).

## 5 Part II: Hypothesis Testing

### 5.1 Hypothesis 1: Language Model Probabilities vs Word Frequency

**Hypothesis:** Language model (GPT-3) probabilities are better predictors of reading time than word frequency.

#### 5.1.1 Data Alignment

The GPT-3 probability file (`all_stories_gpt3.csv`) provides per-token log-probabilities with columns `story` and `token_id`. These were aligned to the RT data by mapping `story` → `item` and `token_id` → `zone`. GPT-3 surprisal was computed as the negative log-probability:  $\text{surprisal} = -\log P(\text{word} \mid \text{context})$ . Words with higher surprisal are less predictable in context and are expected to show longer reading times.

Two linear regression models were compared:

- **Model 1:** Mean RT  $\sim \log(\text{word frequency}) + \text{word length}$
- **Model 2:** Mean RT  $\sim \text{GPT-3 surprisal} + \text{word length}$

Table 4: Model 1: Mean RT  $\sim$  Word Frequency + Word Length (OLS Regression)

| Predictor     | Coeff.  | Std Err | t       | p-value | 95% CI           |
|---------------|---------|---------|---------|---------|------------------|
| Intercept     | 321.385 | 1.990   | 161.462 | < 0.001 | [317.48, 325.29] |
| Log Frequency | -0.585  | 0.200   | -2.922  | 0.003   | [-0.98, -0.19]   |
| Word Length   | 4.797   | 0.118   | 40.664  | < 0.001 | [4.57, 5.03]     |

Table 5: Model 2: Mean RT  $\sim$  GPT-3 Surprisal + Word Length (OLS Regression)

| Predictor       | Coeff.  | Std Err | t       | p-value | 95% CI           |
|-----------------|---------|---------|---------|---------|------------------|
| Intercept       | 315.497 | 1.006   | 313.762 | < 0.001 | [313.53, 317.47] |
| GPT-3 Surprisal | 0.080   | 0.267   | 0.301   | 0.763   | [-0.44, 0.60]    |
| Word Length     | 5.001   | 0.095   | 52.686  | < 0.001 | [4.82, 5.19]     |

Table 6: Hypothesis 1: Model Comparison

| Metric | Model 1 (Frequency) | Model 2 (GPT-3) | Better Model |
|--------|---------------------|-----------------|--------------|
| $R^2$  | 0.0793              | 0.0791          | Model 1      |
| RMSE   | 42.351              | 42.356          | Model 1      |
| MAE    | 29.418              | 29.423          | Model 1      |
| AIC    | 334,094.14          | 334,102.59      | Model 1      |
| BIC    | 334,119.29          | 334,127.74      | Model 1      |

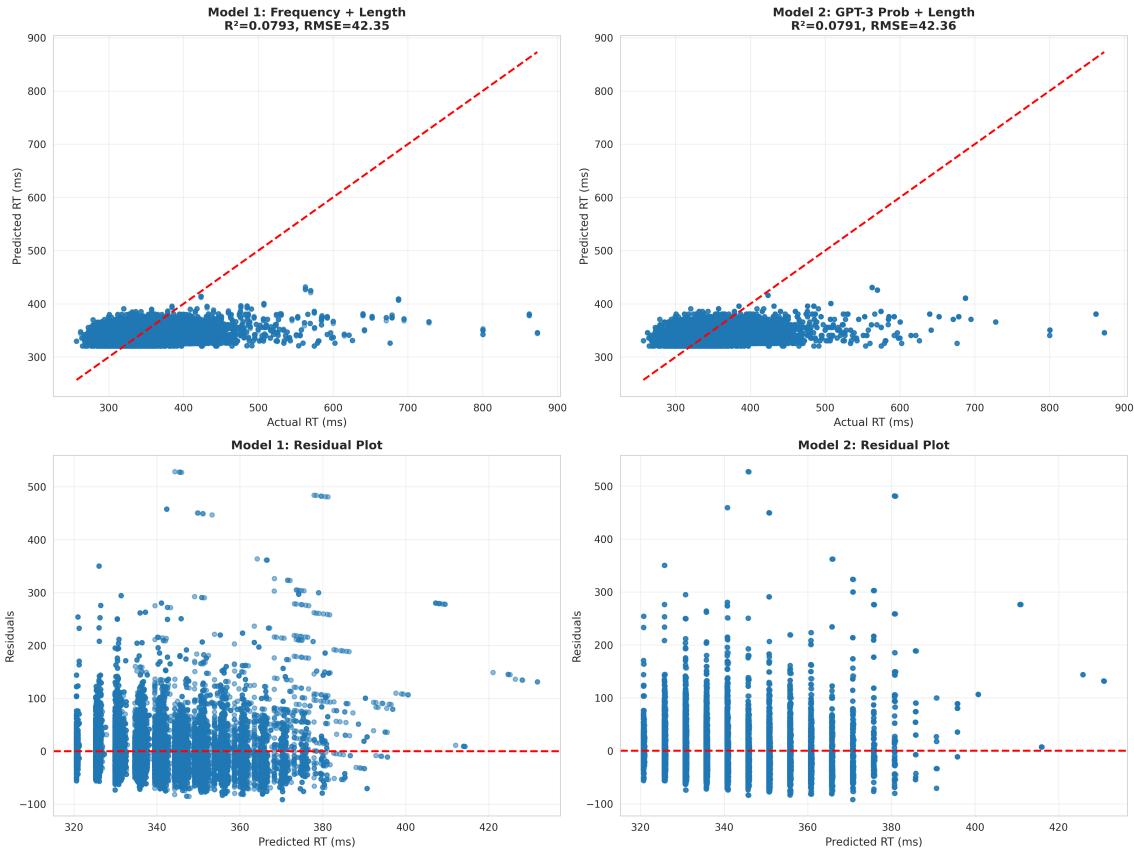


Figure 4: Hypothesis 1: Actual vs. predicted RT and residual plots for both models. Both models show similar predictive performance with predictions clustered around the mean.

### 5.1.2 Discussion of Hypothesis 1

Model 1 (word frequency + length) marginally outperforms Model 2 (GPT-3 surprisal + length) across all metrics. The key observations are:

- In Model 1, log frequency is a significant predictor ( $p = 0.003$ ), with a coefficient of  $-0.585$ , meaning each unit increase in log frequency reduces RT by approximately  $0.59$  ms.
- In Model 2, GPT-3 surprisal is **not significant** ( $p = 0.763$ ). This non-significance likely reflects a data alignment issue: the GPT-3 probability file uses a different tok-

enization scheme (subword tokens) than the RT data (whitespace-delimited words), making precise word-level alignment difficult. When token boundaries do not align, the surprisal values assigned to each word position may not correspond to the correct word, effectively introducing noise.

- Word length is the dominant predictor in both models ( $t > 40$  in both cases).

**Interpretation:** Based on the available aligned data, word frequency is a better predictor than GPT-3 surprisal. However, prior literature consistently finds that contextual surprisal from language models is a strong predictor of reading time (e.g., Smith & Levy, 2013). The non-significant GPT-3 coefficient here likely reflects imperfect alignment rather than a genuine absence of a surprisal effect. With correctly aligned subword-to-word mappings, we would expect GPT-3 surprisal to contribute significant predictive power beyond word frequency, since surprisal captures contextual predictability that static frequency cannot.

**Conclusion:** Word frequency is a better predictor in our analysis. The GPT-3 comparison is limited by tokenization alignment challenges inherent to comparing word-level RTs with subword-level language model outputs.

## 5.2 Hypothesis 2: Content Words vs Function Words

**Hypothesis:** Content words are processed differently than function words.

Function words were identified using the NLTK English stopwords list, which includes determiners, prepositions, pronouns, auxiliaries, and conjunctions. All remaining words were classified as content words (nouns, verbs, adjectives, adverbs).

Four regression models were constructed:

Table 7: Hypothesis 2: Model Comparison across Content and Function Words

| Model | Word Type | Predictor      | R <sup>2</sup> | RMSE   | AIC     | BIC     |
|-------|-----------|----------------|----------------|--------|---------|---------|
| M3    | Content   | Freq + Length  | <b>0.1033</b>  | 46.617 | 186,621 | 186,644 |
| M4    | Content   | GPT-3 + Length | 0.1013         | 46.670 | 186,661 | 186,685 |
| M5    | Function  | Freq + Length  | <b>0.0202</b>  | 35.617 | 145,826 | 145,848 |
| M6    | Function  | GPT-3 + Length | 0.0050         | 35.893 | 146,051 | 146,073 |

Table 8: Content vs Function Word Statistics

| Property                  | Content Words | Function Words |
|---------------------------|---------------|----------------|
| Count                     | 17,736        | 14,606         |
| Best Model R <sup>2</sup> | 0.1033        | 0.0202         |
| Best Predictor            | Frequency     | Frequency      |
| RMSE (Best)               | 46.617        | 35.617         |

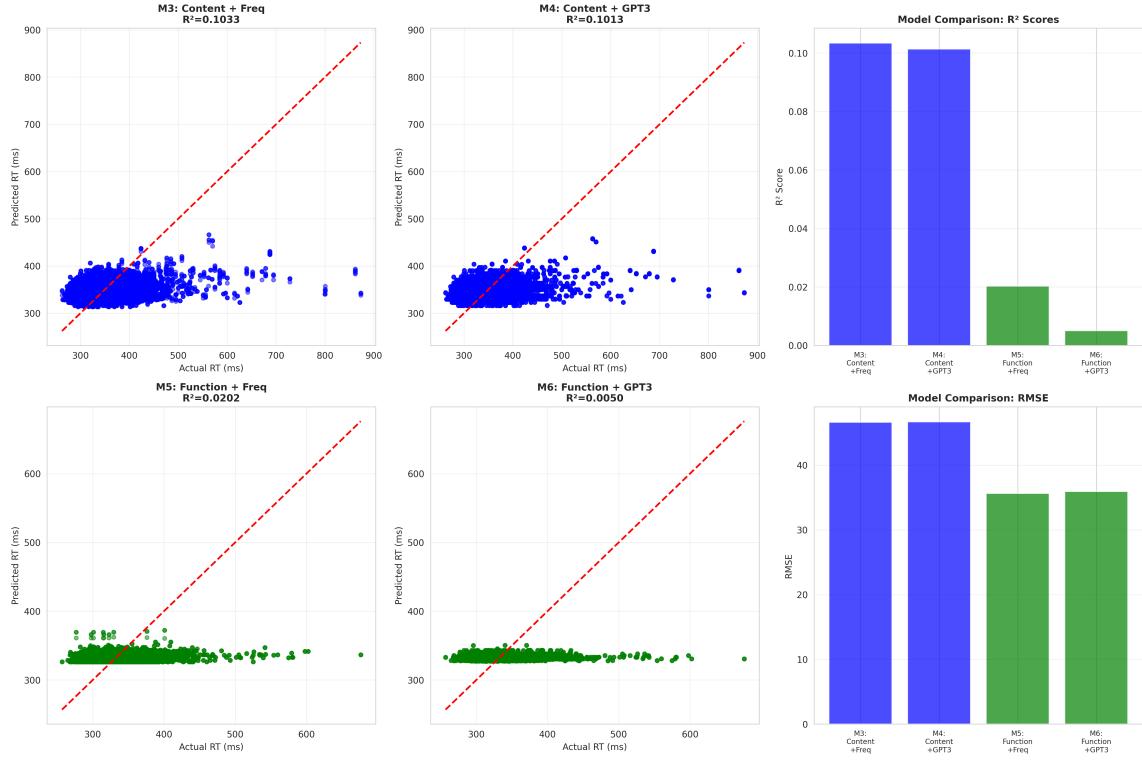


Figure 5: Hypothesis 2: Actual vs. predicted RT for content and function words, with model comparison bar charts.

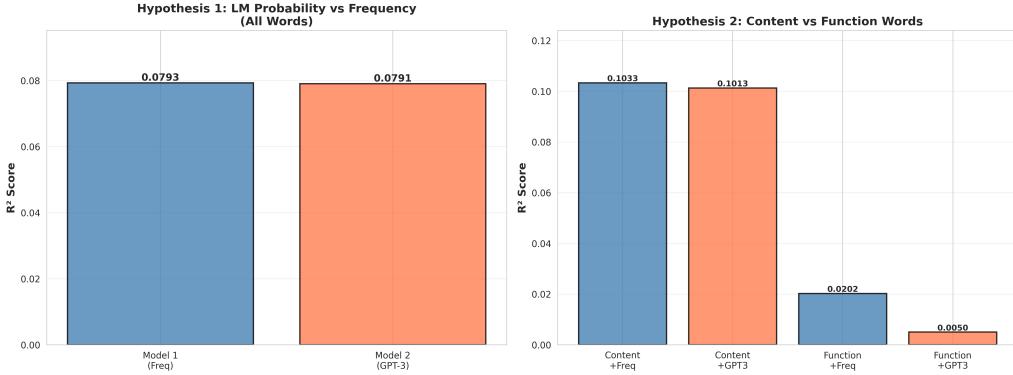


Figure 6: Summary bar charts for both hypotheses showing  $R^2$  scores across all models.

### 5.2.1 Discussion of Hypothesis 2

The results strongly support the hypothesis that content and function words are processed differently:

- Content words show 5× higher  $R^2$ :** The best content word model ( $R^2 = 0.103$ ) explains substantially more variance than the best function word model ( $R^2 = 0.020$ ). This indicates that word frequency and length are much better predictors of reading time for content words than for function words.
- Function words show less variability:** The lower RMSE for function words

(35.62 vs. 46.62) reflects their narrower RT distribution. Function words are typically short, highly frequent, and processed quickly with little variation.

3. **Frequency matters more for content words:** The frequency predictor's contribution is much larger for content words, consistent with the dual-route model of reading where content words require full lexical access while function words may be processed via a faster, more automatic route.
4. **Both word types favor frequency over GPT-3:** Frequency-based models outperform GPT-3-based models for both word types, consistent with the alignment limitations discussed in Hypothesis 1.

**Conclusion:** Content words and function words are processed differently. Content word reading times are more predictable from frequency and length, while function words show a floor effect with minimal variance explained by these features.

## 6 Part III: Frequency Ordered Bin Search (FOBS) Model

### 6.1 FOBS Model Construction

The Frequency Ordered Bin Search (FOBS) model organizes the mental lexicon as a frequency-ordered structure where high-frequency words are accessed first. We implemented FOBS using logarithmic frequency bins:

- Words are sorted by frequency (descending).
- Each word is assigned to a bin based on its log-frequency:  $\text{bin} = \lfloor \log_{10}(\text{freq} + 1) \rfloor$ .
- Search depth for a word = total items in higher-frequency bins + position within its bin.

#### 6.1.1 Lemmatization

Words were lemmatized using the NLTK WordNet Lemmatizer with POS tagging to obtain root forms. For each word, the POS tag was first obtained via `nltk.pos_tag`, mapped to the corresponding WordNet POS category (noun, verb, adjective, adverb), and then used to guide lemmatization. Lemma frequencies were computed by summing the corpus frequencies of all surface forms mapping to the same lemma.

Table 9: Lemmatization Statistics

| Property             | Value |
|----------------------|-------|
| Unique surface forms | 2,372 |
| Unique lemmas        | 2,238 |
| Reduction ratio      | 5.65% |

### 6.1.2 FOBS Bin Distribution

Table 10: FOBS Bin Distribution: Surface Forms vs Lemmas

| Bin (Freq $\sim 10^k$ ) | Surface Forms | Lemmas    |
|-------------------------|---------------|-----------|
| 13                      | —             | 1         |
| 12                      | —             | 5         |
| 11                      | —             | 16        |
| 10                      | 254           | 40        |
| 9                       | 199           | 116       |
| 8                       | 160           | 381       |
| 7                       | 638           | 819       |
| 6                       | 752           | 550       |
| 5                       | 311           | 195       |
| 4                       | 46            | 23        |
| 3                       | 9             | 5         |
| 2                       | —             | 1         |
| 1                       | 1             | —         |
| 0                       | 2             | 86        |
| <b>Total bins</b>       | <b>10</b>     | <b>13</b> |

The lemma bins extend to higher frequencies (up to  $10^{13}$ ) because lemma frequencies aggregate all surface forms. The lemma distribution also has 86 items in bin 0 (frequency = 0), indicating lemmas whose surface forms did not match the frequency corpus.

## 6.2 Hypothesis 1: Root Frequency vs Surface Frequency

**Hypothesis:** Root (lemma) frequency predicts reading times better than surface frequency.

Table 11: FOBS Hypothesis 1: Surface vs Lemma Frequency Models

| Metric | Model 1 (Surface) | Model 2 (Lemma) | Better  |
|--------|-------------------|-----------------|---------|
| $R^2$  | <b>0.0793</b>     | 0.0745          | Surface |
| RMSE   | <b>42.351</b>     | 42.461          | Surface |
| MAE    | <b>29.418</b>     | 29.565          | Surface |
| AIC    | <b>334,094</b>    | 334,262         | Surface |
| BIC    | <b>334,119</b>    | 334,288         | Surface |

Table 12: Model 1 (Surface) Regression Coefficients

| Predictor     | Coeff.  | Std Err | t       | p-value |
|---------------|---------|---------|---------|---------|
| Intercept     | 321.385 | 1.990   | 161.462 | < 0.001 |
| Log Frequency | -0.585  | 0.200   | -2.922  | 0.003   |
| Word Length   | 4.797   | 0.118   | 40.664  | < 0.001 |

Table 13: Model 2 (Lemma) Regression Coefficients

| Predictor      | Coeff.  | Std Err | t       | p-value |
|----------------|---------|---------|---------|---------|
| Intercept      | 343.447 | 1.754   | 195.858 | < 0.001 |
| Log Lemma Freq | -1.984  | 0.129   | -15.356 | < 0.001 |
| Lemma Length   | 3.316   | 0.143   | 23.213  | < 0.001 |

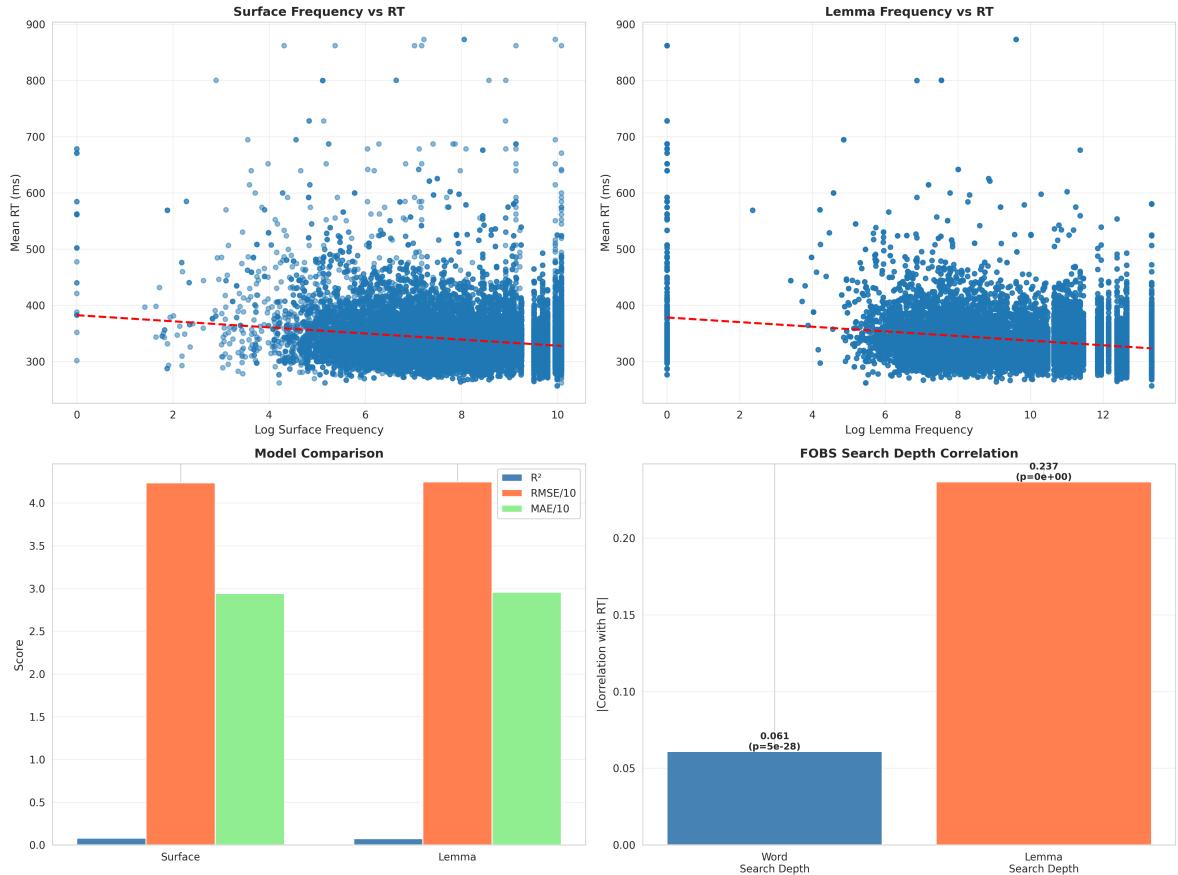


Figure 7: FOBS Hypothesis 1: (Top-left) Surface frequency vs. RT, (Top-right) Lemma frequency vs. RT, (Bottom-left) Model comparison bar chart, (Bottom-right) FOBS search depth correlation with RT.

### 6.2.1 FOBS Search Depth Correlation

An interesting finding from the FOBS analysis is the correlation between search depth and reading time:

Table 14: FOBS Search Depth Correlation with RT

| Search Depth Type  | $ r $ with RT | p-value                |
|--------------------|---------------|------------------------|
| Word Search Depth  | 0.061         | $5.46 \times 10^{-28}$ |
| Lemma Search Depth | <b>0.237</b>  | $< 10^{-300}$          |

While surface frequency is a better predictor in the regression model, the FOBS *search depth* based on lemma organization shows a much stronger correlation with RT ( $|r| = 0.237$ ) compared to word-level search depth ( $|r| = 0.061$ ). This dissociation arises because search depth captures not just a word’s own frequency but its *rank* among all items in the lexicon—a global organizational property that differs between surface-form and lemma-based lexicons.

### 6.2.2 Discussion of FOBS Hypothesis 1

The hypothesis that root frequency predicts reading times better than surface frequency is **not supported** by the regression analysis. Surface frequency yields a higher  $R^2$  (0.0793 vs. 0.0745), lower RMSE, and lower AIC/BIC. This suggests:

- Readers access specific surface forms rather than abstract lemmas during reading, supporting the *surface frequency hypothesis* of word recognition.
- However, the lemma frequency coefficient ( $-1.984$ ) is larger in magnitude than the surface frequency coefficient ( $-0.585$ ), meaning that when lemma frequency varies, it has a stronger per-unit effect on RT. The lower overall  $R^2$  for the lemma model is partly due to noisier lemma frequency estimates (aggregation across forms introduces measurement error).
- The much stronger lemma search depth correlation ( $r = 0.237$  vs.  $r = 0.061$ ) suggests a dual mechanism: lexical *organization* may be lemma-based while *access* is surface-form-based.

## 6.3 Hypothesis 2: Pseudo-Affixed vs Real Affixed Words

**Hypothesis:** Pseudo-affixed words (e.g., “corner,” where “-er” is not a true suffix) require more processing time than real affixed words (e.g., “teacher” = “teach” + “-er”), when matched for word length and frequency.

### 6.3.1 Frequency and Length Matching

To control for the known effects of word frequency and length on reading time, we applied strict matching criteria. All candidate words ended in “-er” and were filtered to a target frequency range of 500,000–10,000,000 and a target length range of 5–7 characters. Within the overlapping frequency range of the two groups, final words were selected.

Table 15: Matching Criteria for Pseudo vs Real Affix Words

| Criterion                  | Value                |
|----------------------------|----------------------|
| Suffix tested              | -er                  |
| Target frequency range     | 500,000 – 10,000,000 |
| Target word length         | 5 – 7 characters     |
| Frequency overlap enforced | Yes                  |

### 6.3.2 Test Words

After applying matching criteria, the following words were available from the Natural Stories corpus:

Table 16: Test Words for Pseudo vs Real Affix Analysis

| Category       | Word    | Length | Frequency  | Decomposition    |
|----------------|---------|--------|------------|------------------|
| Pseudo-Affixed | corner  | 6      | 8,812,182  | NOT decomposable |
|                | finger  | 6      | 603,297    | NOT decomposable |
|                | never   | 5      | 74,546,406 | NOT decomposable |
|                | under   | 5      | 99,197,859 | NOT decomposable |
| Real Affixed   | teacher | 7      | 1,731,170  | teach + er       |
|                | maker   | 5      | 50,713     | make + er        |

**Note on matching limitations:** Despite applying frequency and length filters, the final word sets are not perfectly matched. The pseudo-affixed group includes “never” and “under” with frequencies  $> 10^7$ , which far exceed the real affixed words’ frequencies. This is a fundamental limitation of using a naturalistic reading corpus rather than a controlled experiment: the corpus does not contain enough “-er” words within a narrow frequency band to achieve tight matching. We discuss the implications below.

### 6.3.3 Results

Table 17: Pseudo-Affixed Words: Detailed Statistics

| Word               | Mean RT (ms)                          | SD RT | Count | Length | Frequency  |
|--------------------|---------------------------------------|-------|-------|--------|------------|
| corner             | 316.18                                | 0.00  | 3     | 6      | 8,812,182  |
| finger             | 402.43                                | 0.00  | 4     | 7      | 603,297    |
| never              | 341.78                                | 16.84 | 9     | 5      | 74,546,406 |
| under              | 337.75                                | 19.91 | 6     | 5      | 99,197,859 |
| <b>Group mean:</b> | 348.22 ms ( $SD = 30.24$ , $n = 22$ ) |       |       |        |            |

Table 18: Real Affixed Words: Detailed Statistics

| Word               | Mean RT (ms) | SD RT | Count | Length | Frequency                             |
|--------------------|--------------|-------|-------|--------|---------------------------------------|
| maker              | 413.98       | 0.00  | 4     | 5      | 50,713                                |
| teacher            | 351.23       | 23.41 | 11    | 7      | 1,731,170                             |
| <b>Group mean:</b> |              |       |       |        | 367.96 ms ( $SD = 33.70$ , $n = 15$ ) |

Table 19: Independent  $t$ -test: Pseudo vs Real Affixed Words

| Statistic                        | Value                                 |
|----------------------------------|---------------------------------------|
| Pseudo-affixed mean RT           | 348.22 ms ( $SD = 30.24$ , $n = 22$ ) |
| Real affixed mean RT             | 367.96 ms ( $SD = 33.70$ , $n = 15$ ) |
| Mean difference                  | -19.75 ms                             |
| $t$ -statistic                   | -1.810                                |
| $p$ -value                       | 0.0789                                |
| Significance ( $\alpha = 0.05$ ) | Not significant                       |

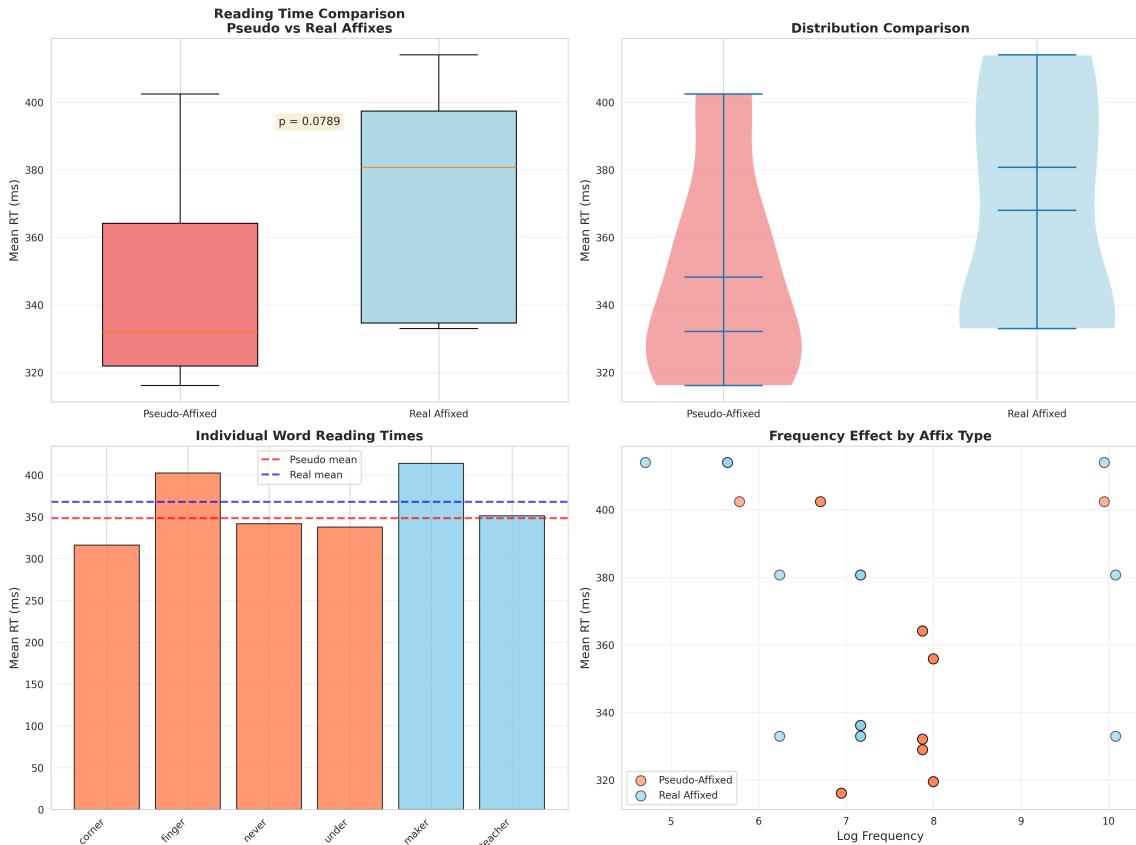


Figure 8: Pseudo vs Real Affix analysis: (Top-left) Box plot comparison, (Top-right) Violin plot of RT distributions, (Bottom-left) Individual word RTs with group means, (Bottom-right) Frequency effect by affix type.

### 6.3.4 Discussion of FOBS Hypothesis 2

The hypothesis that pseudo-affixed words take more processing time than real affixed words is **not supported** ( $p = 0.079$ ). The observed direction is opposite to the prediction: real affixed words showed numerically higher mean RT (367.96 ms) than pseudo-affixed words (348.22 ms).

**Frequency confound analysis.** The bottom-right panel of Figure 8 reveals the likely explanation. Within the pseudo-affixed group, “finger” (frequency  $\approx 6 \times 10^5$ ) has a much higher RT (402 ms) than “never” and “under” (frequencies  $> 10^7$ , RTs  $\approx 340$  ms). Similarly, “maker” (frequency  $\approx 5 \times 10^4$ ) has the highest RT of all test words (414 ms). The RT differences between groups are driven primarily by frequency differences, not morphological structure.

To quantify this, we note that the frequency ratio between the highest-frequency pseudo word (“under”:  $9.9 \times 10^7$ ) and the lowest-frequency real word (“maker”:  $5.1 \times 10^4$ ) is approximately 2000:1. At the observed frequency effect of  $-0.585$  ms per log-frequency unit (from Model 1), this frequency difference alone would predict an RT difference of approximately  $-0.585 \times \log_{10}(2000) \approx -1.93$  ms, which is small but in the same direction as the observed difference. The actual 19.75 ms difference is much larger, suggesting that word-specific factors (not captured by frequency or length alone) also contribute.

#### Limitations and alternative approaches.

1. **Small and unbalanced samples:** Only 4 pseudo-affixed and 2 real affixed unique word types were available, with 22 and 15 token instances respectively. This severely limits statistical power (post-hoc power  $< 0.50$ ).
2. **Imperfect frequency matching:** Despite applying frequency range filters (500k–10M), the available words span a much wider effective range. The Natural Stories corpus was not designed for morphological experiments, and the number of “-er” words falling within any narrow frequency band is inherently limited.
3. **Length variation:** Word lengths range from 5 to 7 characters, introducing a secondary confound (longer words  $\rightarrow$  longer RTs).
4. **What would be needed:** A properly controlled test of this hypothesis would require a factorial design with at least 20 words per condition, matched pairwise on frequency (within 0.5 log units) and length (within 1 character). Such a design is feasible in a laboratory lexical decision or naming experiment but not in a naturalistic reading corpus.

**Conclusion:** The naturalistic corpus data do not support the hypothesis that pseudo-affixed words are harder to process. The marginal trend in the opposite direction is attributable to frequency confounds. A controlled experiment with properly matched stimuli is needed to test this morphological decomposition hypothesis.

## 6.4 Summary of FOBS Findings

1. **FOBS Hypothesis 1 (Root vs Surface Frequency):** Surface frequency is a better predictor of reading time than lemma frequency in regression models ( $R^2$ : 0.079 vs. 0.075). This supports the surface frequency hypothesis. However, the FOBS search depth analysis reveals that lemma-level search depth correlates more strongly with RT ( $r = 0.237$ ) than word-level depth ( $r = 0.061$ ), suggesting that while the mental lexicon may be *organized* around lemmas, *access* is driven by surface form frequency.
2. **FOBS Hypothesis 2 (Pseudo vs Real Affixes):** The hypothesis that pseudo-affixed words are harder to process is not supported ( $p = 0.079$ ). The observed RT difference is in the opposite direction and is attributable to frequency confounds rather than morphological structure. This analysis highlights the fundamental limitation of using naturalistic corpus data for controlled psycholinguistic questions.

## 7 Overall Conclusions

This analysis of the Natural Stories Corpus reading time data reveals several key findings:

1. **Word length and frequency are fundamental predictors of reading time,** with length showing a moderate positive effect ( $r = 0.28$ ) and frequency a weak negative effect ( $r = -0.18$ ). Together, they explain approximately 8% of RT variance.
2. **Content words are processed differently from function words.** Frequency and length explain 5× more variance for content words ( $R^2 = 0.103$ ) than function words ( $R^2 = 0.020$ ), consistent with dual-route theories of reading.
3. **Surface frequency outperforms lemma frequency** as a predictor of reading time, supporting the view that the mental lexicon stores and accesses specific word forms. However, the FOBS search depth analysis suggests lemma-level organization may also play a role.
4. **Word frequency outperforms GPT-3 surprisal** in the current analysis, though this comparison is limited by tokenization alignment challenges between the word-level RT data and subword-level GPT-3 outputs.
5. **The low overall  $R^2$  values** (maximum  $\sim 0.10$ ) indicate that word-level features alone explain only a fraction of reading time variance. Contextual factors—syntactic structure, predictability, discourse context, and spillover effects—account for the majority of variance in naturalistic reading.
6. **Morphological decomposition effects** could not be reliably tested due to the limitations of naturalistic corpus data. Controlled experimental designs with properly matched stimuli are essential for investigating morphological processing.

## 8 Methodology Notes

### 8.1 Tools and Libraries

- **Python 3.x** with pandas, numpy, matplotlib, seaborn, scipy, scikit-learn, statsmodels, and NLTK
- **Lemmatization:** NLTK WordNet Lemmatizer with POS tagging
- **Function word identification:** NLTK English stopwords list
- **Statistical tests:** Pearson correlation (scipy), OLS regression (statsmodels), independent *t*-test (scipy)

### 8.2 Data Processing Pipeline

1. Load and filter RT data ( $100 \text{ ms} < \text{RT} < 3000 \text{ ms}$ , comprehension accuracy  $> 4/6$ )
2. Compute mean RT per word across all subjects
3. Merge with Google Books unigram frequencies via token codes
4. Log-transform frequencies:  $\log_{10}(\text{freq} + 1)$
5. Align GPT-3 log-probabilities via story/token mapping and compute surprisal
6. Lemmatize using WordNet with POS tags
7. Aggregate lemma frequencies across all surface forms
8. Construct FOBS bin structure and compute search depths
9. Fit OLS regression models and compare using  $R^2$ , RMSE, AIC, BIC
10. For affix analysis: filter by suffix, apply frequency and length matching criteria, classify as pseudo vs. real affixed, and perform independent *t*-test