

# Factors Influencing Early Childhood Comprehension of American English\*

Analyzing the impact of early childhood on English comprehension

Ziyuan Shen

November 30, 2024

This paper investigates the factors that influence language comprehension in young children by using a Bayesian linear regression model. We analyzed predictors including age, word production, caregiver education, and birth order. The analysis showed that factors such as caregiver education and word production have a significant association with comprehension levels. This study highlights the importance of both environmental and developmental factors in shaping early language abilities, providing a clearer understanding of how these elements contribute to a child's comprehension.

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data</b>	<b>2</b>
2.1	Data Overview . . . . .	2
2.2	Data Measurement . . . . .	3
2.3	Outcome variables . . . . .	4
2.3.1	Comprehension(Number of Words Understand By Children) . . . . .	4
2.4	Predictor variables . . . . .	5
2.4.1	Age . . . . .	5
2.4.2	Birth Order and Caregiver Education . . . . .	6
2.4.3	Race . . . . .	8
<b>3</b>	<b>Model</b>	<b>9</b>
3.1	Alternative Models . . . . .	9

---

\*Code and data are available at: [Factors Influencing Early Childhood Comprehension of American English](#).

3.2	Model set-up . . . . .	9
3.3	Model justification . . . . .	10
<b>4</b>	<b>Results</b>	<b>11</b>
4.1	Model results and interpretation . . . . .	11
<b>5</b>	<b>Discussion</b>	<b>14</b>
5.1	The Role of Caregiver Education in Language Development . . . . .	14
5.2	Effect of birth order on language comprehension . . . . .	14
5.3	Model Limitations and Uncertainties . . . . .	15
5.4	Recommendations for Future Research . . . . .	15
	<b>Appendix</b>	<b>16</b>
<b>A</b>	<b>CDI methodology overview and assessment</b>	<b>16</b>
A.1	Overview . . . . .	16
A.2	Target population, frame, and sample . . . . .	16
A.3	Sample recruitment . . . . .	16
A.4	Sampling approach and trade-offs . . . . .	17
A.5	Non-response handling . . . . .	17
A.6	Longitudinal Stability of CDI Measurements . . . . .	18
A.7	Psychometric Analysis and Item Response Theory (IRT) . . . . .	18
A.8	Simulation of CDI Sampling Approach . . . . .	18
A.9	Conclusion . . . . .	19
<b>B</b>	<b>Additional data details</b>	<b>19</b>
B.1	Data manipulation and cleaning . . . . .	19
B.2	Descriptions of Each Predictor Variable . . . . .	20
<b>C</b>	<b>Model details</b>	<b>21</b>
C.1	Model summary . . . . .	21
C.2	Posterior predictive check . . . . .	21
C.3	Variance Inflation Factors for Each Predictor in the Model . . . . .	21
C.4	Diagnostics . . . . .	22
	<b>References</b>	<b>24</b>

# 1 Introduction

Early childhood language comprehension is a fundamental area of cognitive development that influences future literacy and communication skills. Research has identified several key predictors of language comprehension in young children. For example, a study by Taylor et

al.(Zubrick, Taylor, and Christensen (2015)) found that socioeconomic factors, family background, and personal characteristics significantly influence early childhood receptive vocabulary development. As highlighted by Lawton et al.(Lawton, Araujo, and Kufaishi (2023)), the quality and quantity of language exposure in infancy are associated with brain development and long-term language achievement. In addition, Bergelson’s study(Zonarich (2024)) showed that babies begin to understand common nouns as early as 6 to 7 months of age, which suggests that comprehension begins earlier than previously thought. Despite these insights, I believe there is still a need for robust research that considers the range of demographic, developmental, and environmental factors that influence the early comprehension of American English learners.

In this paper, my estimand is the number of words comprehended by children learning American English. Our focus is on quantifying how developmental factors such as age, birth order, caregiver influences such as education level, and language context such as monolingual versus bilingual status affect comprehension outcomes. Using Bayesian regression modelling, I assessed the contributions of these predictors and identified the factors that have the greatest impact on early childhood comprehension.

My findings show that developmental factors such as age, production ability of language, and caregiver education are important predictors of young children’s level of comprehension. Children with more educated caregivers demonstrated better comprehension, which highlights the role of the home environment in language development. These findings are important because they provide a basis for carefully targeted early childhood interventions, which highlight the areas that need to be supported in order to promote language comprehension. By identifying the most influential factors, this study contributes to a better understanding of early language learning, which is essential for the development of effective educational practices and policies.

The structure of this paper is organized as follows: after this introduction, Section 2 details the dataset used, including the data collection and cleaning processes, and provides an overview of the key variables. Section 3 introduces the regression models employed in the analysis and discusses why these models are suitable for estimating comprehension outcomes. Section 4 presents the results, emphasizing the relationships between different predictors and comprehension. Finally, Section 5 concludes the paper by evaluating the implications of the findings, discussing their importance for educational policy, and outlining potential directions for future research.

## 2 Data

### 2.1 Data Overview

I used the statistical programming language R (R Core Team 2023) to retrieve, simulate, clean, analyze, and test early childhood language comprehension data. The dataset used for

this analysis is from the Wordbank database of children’s vocabulary growth (Frank, n.d.), specifically focusing on American English children. It contains comprehension and a variety of predictors, which include developmental, demographic, and environmental factors. Following the methodology discussed in “Telling Stories with Data” (Alexander 2023), I explore how different predictors affect language comprehension in early childhood.

The dataset used in this study provides different data on the language comprehension skills of young children among American English language learners. The dataset consists of 14,826 rows and 22 columns covering a range of early childhood developmental attributes such as comprehension, output and parental education. To ensure the validity of the analysis, I filtered the data to include children whose comprehension assessments and caregiver information were complete, which would ensure a comprehensive, unbiased representation of the target population. This cleaning allowed us to focus on high-quality, well-integrated cases, which provided reliable results in young children’s comprehension skills. For key operations, please refer to Appendix B.1 for detailed data processing steps.

## 2.2 Data Measurement

The Early Childhood Language Development Assessment Measurement Data Set was collected using a standardized instrument such as the Communicative Development Inventory (CDI). Specifically, comprehension and vocabulary data are obtained through parent reports. These CDIs are survey instruments provided to parents or caregivers that allow them to document their child’s language skills at specific developmental stages.

For example, let your child learn the meaning of different objects at home. This process naturally occurs through interaction with caregivers and is translated into quantifiable data when parents marked items on the CDI checklist that their child understood or said. The resulting scores become items in our dataset, which reflect each child’s language comprehension skills in numerical form.

In the data cleaning process, we prepared the raw election data for analysis by applying transformations, filtering, and restructuring using several R packages, including `dplyr` (Wickham et al. 2023), `arrow` (Richardson et al. 2024), and `rsample` (Frick et al. 2024). Also, the following packages were used in this study:

- `here`(Müller 2020): Helps manage file paths in a project directory, which make it easier to locate and load data files.
- `ggplot2` (Wickham 2016): Used for creating data visualizations like charts and graphs.
- `modelsummary`(Arel-Bundock 2022): Summarizes model outputs in tables that are easy to read.
- `rstanarm`(Goodrich et al. 2022): Fits Bayesian linear regression models using the **Stan** framework.
- `knitr`(Xie 2014): Converts R code and markdown text into reports, documents, or presentations.

- `kableExtra`(Zhu 2024): Beautify tables made with `knitr::kable` by adding formatting options.
- `ggridges`(Wilke 2024): Creates ridgeline plots, which are useful for visualizing distributions.
- `car`(Fox and Weisberg 2019): Provides tools for regression diagnostics and advanced linear model analysis. The detailed cleaning processes are in Appendix B.1.

## 2.3 Outcome variables

### 2.3.1 Comprehension(Number of Words Understand By Children)

The primary outcome variable was the child’s comprehension of American English, measured by the number of words the child understood. These data include responses from children of different ages (months) and allow for age-specific and general measures of language comprehension in early development.

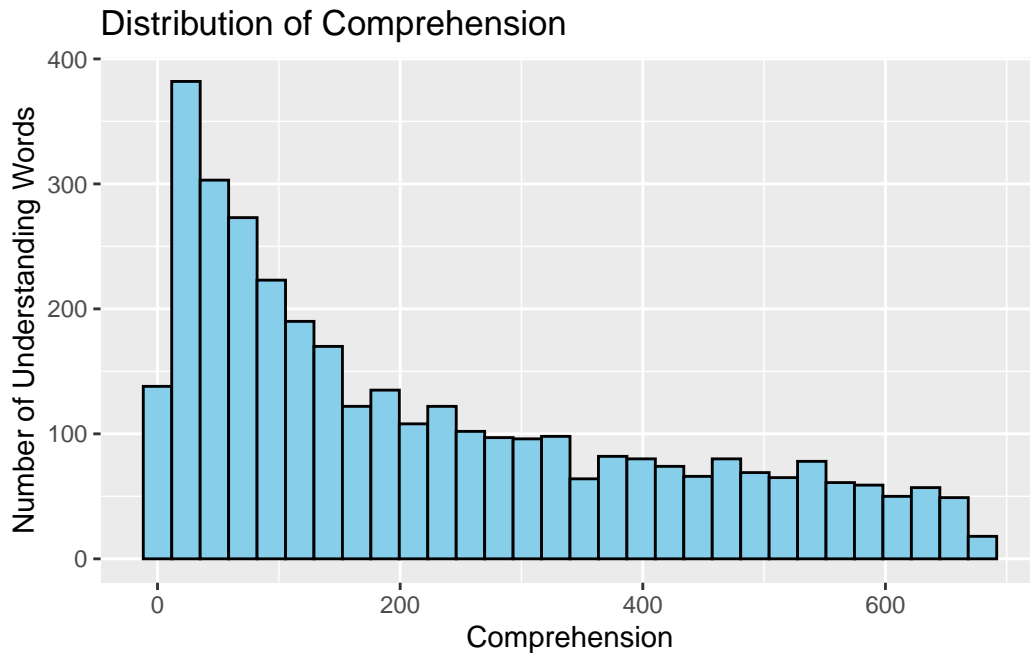


Figure 1: Distribution of comprehension among children. Each bar represents the count of children who understand a specific range of words, illustrating the varying levels of vocabulary comprehension in the dataset.

Using the package `ggplot2` (Wickham 2016), Figure 1 shows the distribution of comprehension, which is the number of words children understand in early childhood. Each bar in the histogram represents the count of children who comprehend a specific range of words.

The histogram shows a right-skewed distribution, with a higher frequency of children having fewer understood words. The count gradually decreases as the number of comprehended words increases, suggesting that most children in the dataset have relatively lower comprehension levels, while fewer children understand a larger vocabulary.

## 2.4 Predictor variables

There are eight predictor variables for my study after cleaning. The Descriptions of Each Predictor Variable is in Table 2 shows at Appendix B.2.

### 2.4.1 Age

The **age** variable indicates the age (in months) of each child. This variable is important for understanding how comprehension develops over time in early childhood. By using age in months, we can capture growth on a more granular level rather than grouping children into broader age categories, which may overlook subtle developmental differences. Age is an important factor because it is directly related to a child's cognitive and language development, which can affect their ability to understand and produce words.

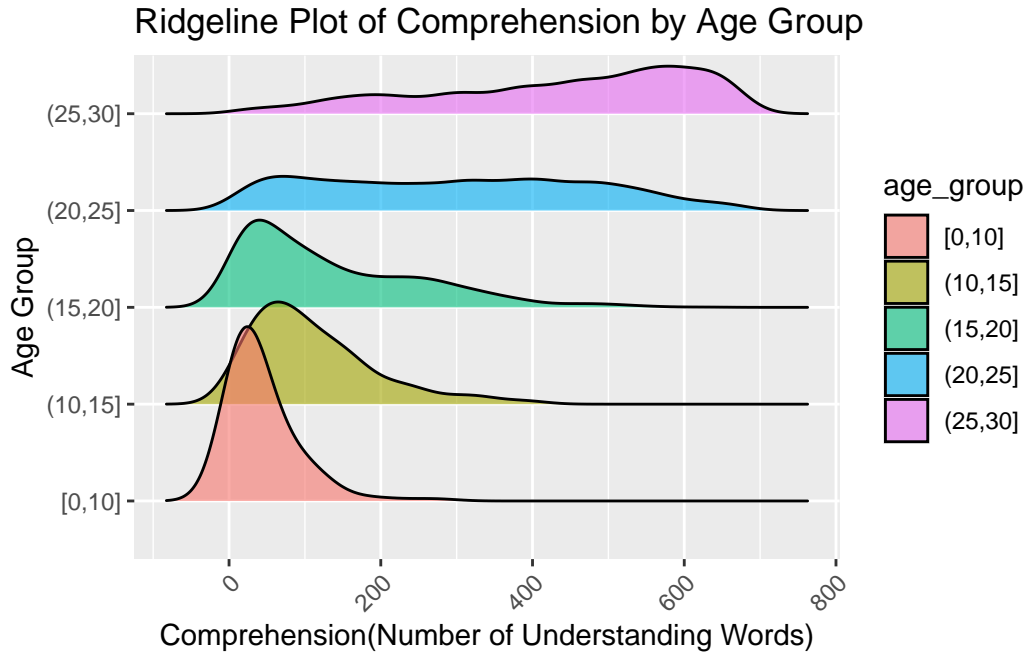


Figure 2: Ridgeline Plot of Comprehension Scores Across Different Age Groups. This plot illustrates the distribution of comprehension scores for each age group, highlighting how comprehension changes across age intervals.

Figure 2 presents a ridgeline plot that illustrates the distribution of comprehension across different age groups. The data is segmented into different age groups to analyze how comprehension levels vary across different stages of early childhood. Each ridge represents an age group, with the x-axis showing the number of words understood (comprehension) and the y-axis indicating the age group category. The density of each ridge shows how comprehension varies within each age group.

As we can see from the plot, comprehension is generally lower for children from 0 to 10 months of age at earlier ages. As children grow up, the distribution of comprehension skills is skewed toward higher levels. It is interesting to note that 15 to 20 month old have a wider distribution of comprehension skills, which indicates greater variability.

#### 2.4.2 Birth Order and Caregiver Education

The `birth_order` variable indicates the position of the child in terms of the order of births within their family. It ranges from 1 (first-born) to 8, which allow us to understand whether being a first-born, middle, or later-born child influences comprehension. Birth order can play an important role in a child's language development due to varying attention levels or interactions they receive from caregivers and siblings.

The `caregiver_education` variable reflects the highest level of education attained by the child's caregiver. This variable includes levels such as "Primary," "Secondary," "College," "Some College," "Graduate". Caregiver education is also important in the analysis because it often serves as a proxy for socioeconomic status and access to language resources, both of which can significantly affect a child's language environment and comprehension ability.

Figure 3 shown above visualizes the distribution of birth order along with caregiver education level. Each bar represents the number of children with a specific birth order, ranging from 1 (first-born) to 8. The bars are stacked by caregiver education level, with different colors indicating categories such as "College", "Graduate", "Primary", and others.

As can be seen from the graph, the majority of children are first or second births, while the number of children in higher birth orders is lower. Caregivers with a "college" or "graduate" level of education make up a large percentage of first and second births, while lower birth order caregivers with "some college" or "secondary" level of education also make up a large proportion of caregivers. The predominance of lower birth orders suggests that most families in the data set tend to have only one or two children.

Figure 4 shows the distribution of children's number of comprehension words grouped by the caregiver's educational level. Children whose caregivers had higher education, such as **Graduate** and **Secondary** tended to have higher median comprehension scores compared to children with primary education. Median comprehension scores are lowest in the primary school category, with a narrower range compared to the other categories, which suggests that there is less variation in comprehension for children with caregivers at this level of education.

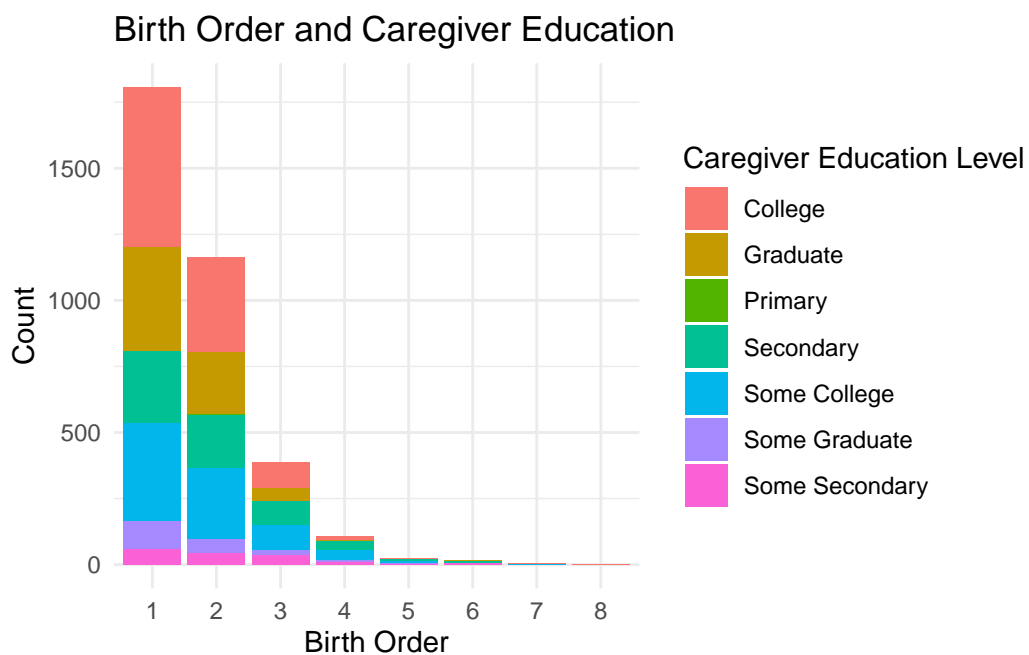


Figure 3: Distribution of Caregiver Education Levels by Birth Order. The stacked bar plot shows the count of children for each birth order, colored by the education level of their primary caregiver.



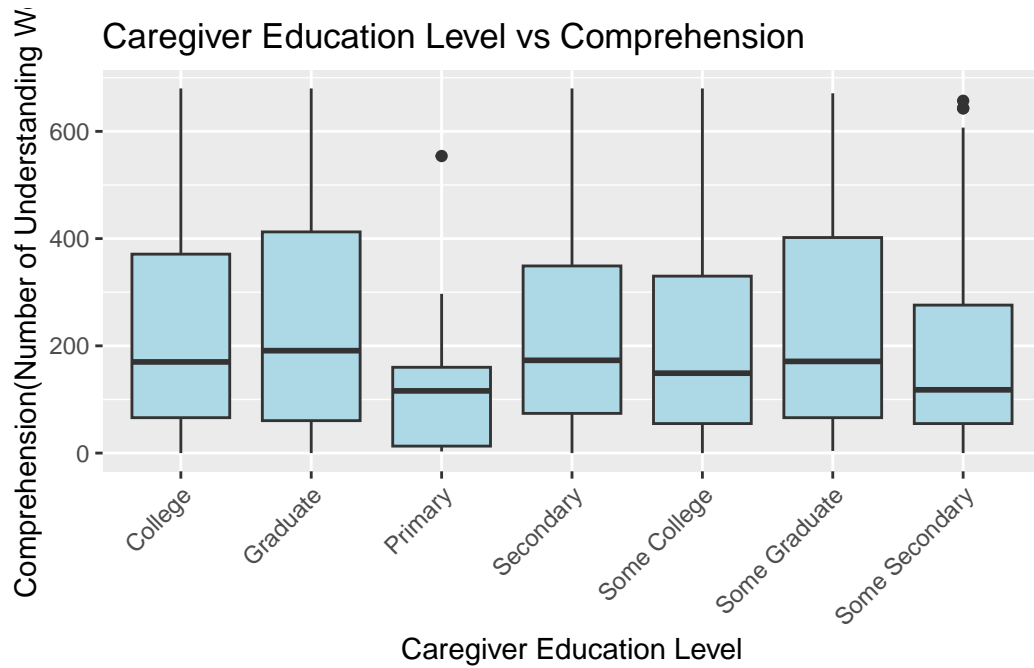


Figure 4: Boxplot of Comprehension by Caregiver Education Level. This plot shows the distribution of comprehension scores for children based on the education level of their caregiver, which highlights differences in comprehension across different educational backgrounds.

**Graduate** and **Some College** show larger IQRs, indicating more variability in comprehension scores among children.

There are several outliers present, particularly for **Some Graduate** and **Graduate** education levels, suggesting that some children scored significantly higher or lower compared to the majority.

### 2.4.3 Race

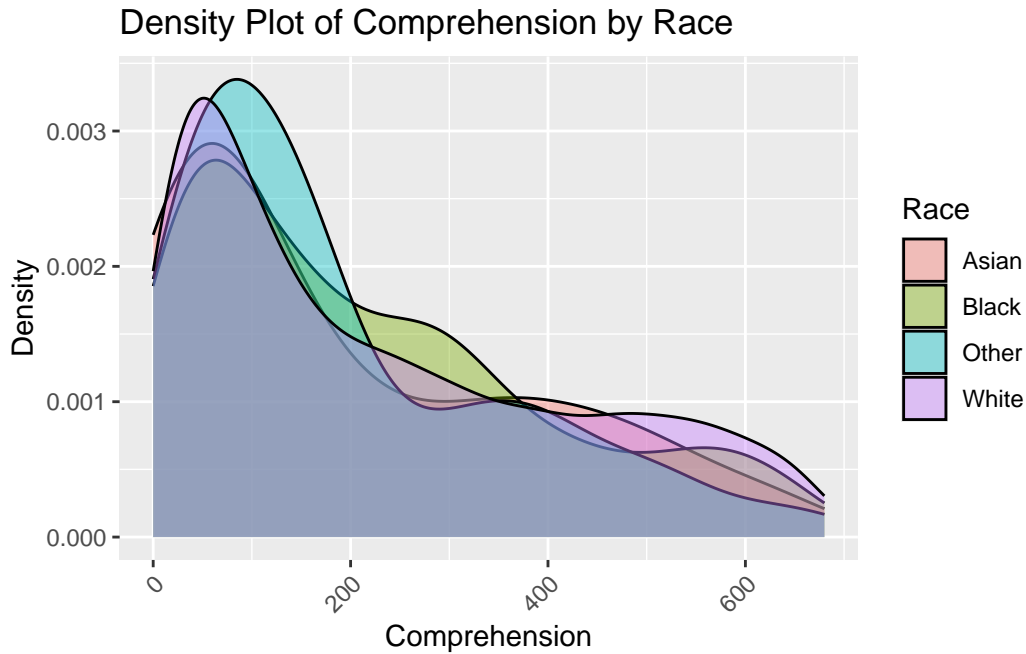


Figure 5: Comprehension by Race

Figure 5 visualizes the distribution of comprehension across different race groups, including Asian, Black, Other, and White. Each coloured density curve represents the distribution of comprehension scores for children from each race group.

From the plot, we can see that **White** children have a wider distribution of comprehension compared to other groups, with a peak density occurring at a higher number of words compared to other groups. **Asian**, **Black**, and **Other** race groups appear to have slightly lower comprehension distributions compared to White children, with peaks occurring at a similar lower range.

The overlap in the density curves shows comprehension across race groups share some similarities in their distributions, although the White group tends to show a broader range and higher scores.

## 3 Model

My modelling approach aims to quantify the relationship between various child and caregiver characteristics and comprehension in early childhood. For this analysis, I use a Bayesian linear model to examine how factors such as age, production (number of words produced), whether the child is in a norming status, birth order, caregiver education level, race, sex, and whether the child is monolingual influence comprehension. The model is implemented using the `stan_glm` function, with a Gaussian distribution to capture the variability in comprehension.

The model assumes that the distribution of comprehension follows a normal distribution when these factors are considered. This Gaussian assumption favours parameter estimation, which is a standard method for linear regression. To prevent overfitting and maintain interpretability, we assume a modest before ensuring that uncertainty is balanced between predictors. This approach allowed us to assess the effects of child and caregiver characteristics on comprehension while maintaining the stability of the findings. See in Appendix C for more background details and diagnostics.

### 3.1 Alternative Models

Alternative models, including a Bayesian logistic regression was considered. Logistic regression was rejected since the outcome (comprehension) is continuous, not binary. The Bayesian linear model was ultimately chosen for its balance between interpretability and the ability to handle uncertainty explicitly through priors.

### 3.2 Model set-up

The model predicts the comprehension of children using the following predictor variables:

- Age(**age**): Represents the child's age in months.
- Production Words(**production**): The number of words the child can produce.
- Norming Status (**is\_norming**): A binary variable, shoews if the child is being used for norming purposes (1 if true, 0 false).
- Birth Order (**birth\_order**): Represents the child's birth order, ranging from 1 to 8.
- Caregiver Education (**caregiver\_education**): The highest education level attained by the child's caregiver.
- Race (**race**): Represents the race of the child.
- Gender (**sex**): A binary variable, the gender of children (1 for female, 0 for male).

- Monolingual (`monolingual`): A binary variable, which indicates if the child is monolingual (1 if true, 0 false).

The model takes the form:

$$\begin{aligned}
y_i \mid \mu_i, \sigma &\sim \text{Normal}(\mu_i, \sigma) \\
\mu_i &= \beta_0 + \beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{production}_i + \beta_3 \cdot \text{is\_norming}_i + \beta_4 \cdot \text{birth\_order}_i \\
&\quad + \beta_5 \cdot \text{caregiver\_education}_i + \beta_6 \cdot \text{race}_i + \beta_7 \cdot \text{sex}_i + \beta_8 \cdot \text{monolingual}_i + \epsilon_i \\
\epsilon_i &\sim \text{Normal}(0, \sigma^2)
\end{aligned}$$

**Where:**

- $y_i \mid \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma)$  represents the outcome variable  $y_i$  (comprehension: number of words understood by child  $i$ ), which is modeled as a normal distribution with mean  $\mu_i$  and standard deviation  $\sigma$ .
- $\beta_0$  is the intercept term.
- $\beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7, \beta_8$  are the coefficients for each predictor.
- $\sigma^2$  is the variance of the error term.

The model is executed in **R** (R Core Team 2023) using the `rstanarm` package (`rstanarm?`). Default priors from `rstanarm` (`rstanarm?`) are used, with the priors set to have a mean of zero and a moderate standard deviation to ensure a reasonable level of regularization.

### 3.3 Model justification

Existing research in developmental psychology and linguistics suggests that factors such as age, production (vocabulary), birth order, caregiver’s level of education, and other demographic variables can have a great impact on children’s comprehension. For example, older children generally comprehend more words due to longer exposure to language and cognitive development. Since these two aspects of language acquisition are closely related, higher language expression tends to correlate with better comprehension. Caregiver education is associated with a child’s vocabulary since more educated caregivers provide a richer language environment. Birth order may also play a role, with research suggesting that first-born children typically receive more one-on-one attention, which may lead to differences in comprehension levels.

Bayesian linear regression modeling was chosen to predict comprehension because the outcome variables are continuous and can reasonably be assumed to follow a normal distribution. Linear regression is suitable for quantifying the relationship between multiple predictors and continuous outcomes, providing interpretable coefficients for each independent variable. Bayesian

methods allow us to incorporate a priori knowledge into the model and are particularly useful in situations where existing research provides information on expected effects. In addition, the Bayesian inference method provides a way to quantify the uncertainty of parameter estimates, which provides us with a deeper understanding of the modeled relationships.

## 4 Results

Section 4 examines the relationship between age, production, is\_norming, birth order, caregiver education, race, sex, and monolingual status with respect to early childhood word comprehension. Using a dataset containing observations of children’s linguistic abilities and various influencing factors, we apply a Bayesian linear regression model to assess which key predictors have the most significant impact on comprehension. Below, we present the results of our model and discuss the implications of our findings.

We predicted comprehension scores using our test dataset, but due to missing data for certain groups such as some levels of caregiver education or race categories, our model could not make predictions for all children. To address these gaps, we examined demographic trends from similar children in the dataset; if a group showed consistent patterns in linguistic development, we assumed that trend would continue and used it to estimate missing comprehension values. This limitation arises because groups with missing values are not represented in the model, preventing accurate forecasts for those individuals. Further discussion on this limitation can be found in Section 5.

### 4.1 Model results and interpretation

The bayesian linear regression model built using our training dataset, which consists of 2,457 data points, estimated the factors that influence comprehension levels in children. For brevity, Table 1 shows only the first ten rows of the model’s coefficients, while the full model summary is provided in Appendix C. The intercept is estimated at 78.515, representing the baseline comprehension level when all predictors are at their reference levels. The model achieved an  $R^2$  value of 0.936, indicating that 93.6% of the variance in comprehension outcomes is explained by the predictors included. The adjusted  $R^2$  value of 0.936 further confirms that the model captures the relationships between variables effectively without overfitting.

Table 1: Summary of key coefficients from the Bayesian linear regression model predict comprehension

	coefficient
(Intercept)	78.515
age	-3.479
production	0.983

Table 1: Summary of key coefficients from the Bayesian linear regression model predict comprehension

	coefficient
<code>is_norming</code>	26.115
<code>birth_order</code>	-1.922
<code>caregiver_educationGraduate</code>	1.577
<code>caregiver_educationPrimary</code>	10.381
<code>caregiver_educationSecondary</code>	6.403
<code>caregiver_educationSome College</code>	-2.204
<code>caregiver_educationSome Graduate</code>	-0.086

Certain predictors have a more pronounced impact than others. For instance, `production` (0.983) shows that a one-unit increase in the number of words produced by a child corresponds to a 0.983 increase in predicted comprehension outcomes. The variable `is_norming` (26.115) also has a positive effect, indicating that children in the norming status tend to have higher comprehension levels. Other predictors, such as `age` (-3.481) and `birth_order` (-1.922), have negative coefficients, suggesting that as age increases (within the studied age range), there may be a slight decrease in comprehension, and children born later in the birth order tend to have lower comprehension scores.

The caregiver’s education level also contributes to variations in comprehension. For example, `caregiver_educationPrimary` (10.381) indicates that children whose caregivers have a primary education level tend to have higher comprehension compared to those whose caregivers have no formal education. However, `caregiver_educationSome College` (-2.204) and `caregiver_educationSome Graduate` (-0.086) show negative coefficients, suggesting a potential decrease in comprehension associated with these categories. A complete overview of all predictors is available in Appendix C.

Based on Table 3 in Appendix C, Figure 6 shows the 90% credible intervals for selected coefficients from the Bayesian linear regression model used to predict comprehension by using package `broom.mixed` (Bolker and Robinson 2024). The variables in Figure 6 are chosen based on their coefficient magnitude in the regression results table, specifically for those coefficients that have an absolute value greater than 3, which indicate a significant influence on the model.

The plot visualizes the estimates of key parameter and their confidence intervals, which provide a deeper understanding of the direction and uncertainty associated with these estimates. The intervals provide an understanding of the possible range of where the impact of each variable.

`raceWhite`, `raceOther` and `raceBlack` show a relatively large credible interval, indicating high uncertainty.

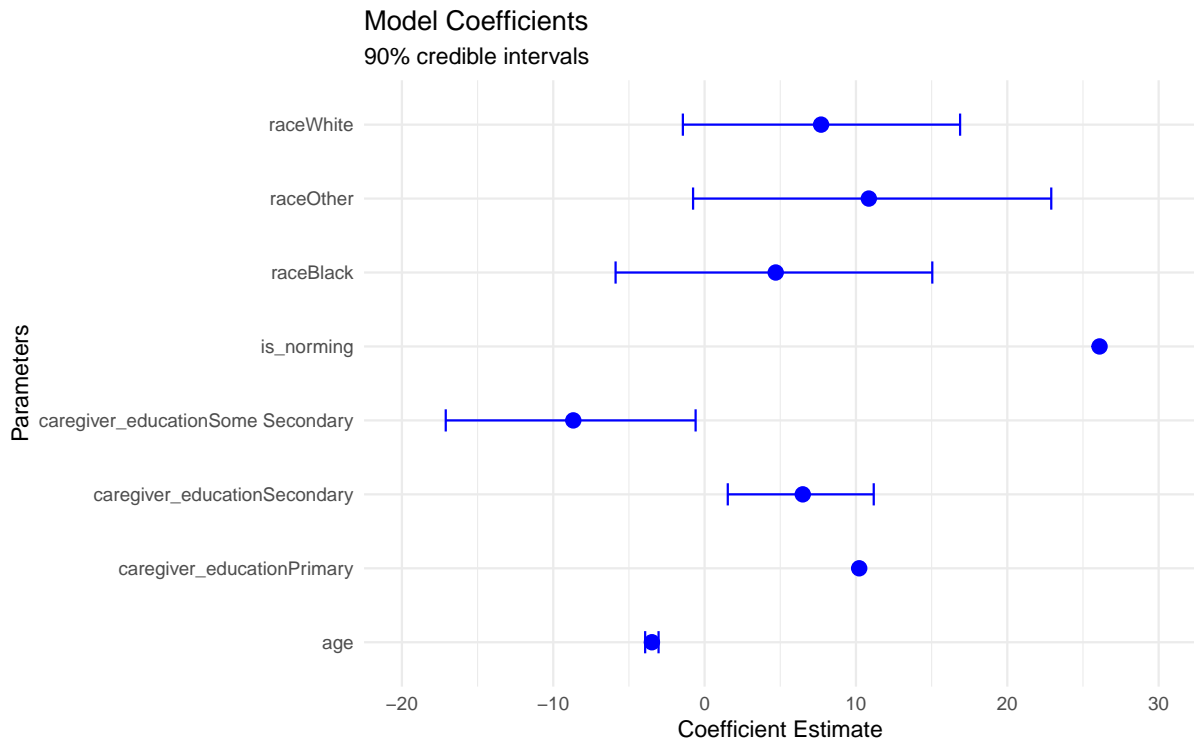


Figure 6: Coefficient Estimates with Confidence Intervals from Bayesian linear regression model

`caregiver_educationPrimary` has a negative estimate with its credible interval crossing zero, suggesting its effect might not be significant. The credible interval for `caregiver_educationSecondary` does not extend very far into the negative side, which might indicate that the effect could be significant, but there is still some uncertainty. And for `caregiver_educationSome Secondary`, the credible interval is quite wide and crosses zero, indicating a high degree of uncertainty regarding the effect of this parameter. This implies that it might not be significantly different from zero, meaning that it could have little or no effect on the `comprehension`

`age` has a very small credible interval, which means it has a more accurate estimate of its `comprehension`.

The parameter `is_norming` does not have an error bar, which suggests that this estimate has high uncertainty or insufficient information.

## 5 Discussion

This study provides insights into early childhood language comprehension by examining various factors, such as age, caregiver education, word production, and birth order, through a Bayesian linear regression model. The findings underscore the significant influence of these predictors, contributing to a broader understanding of the mechanisms that underpin language acquisition during early childhood development.

### 5.1 The Role of Caregiver Education in Language Development

An important finding of this study is that caregiver education plays an important role in early childhood language comprehension. Children show better comprehension with caregivers who are more educated, which highlights the great importance of a rich home environment in early language development. This suggests that caregiver education can provide an important representation of language exposure, socioeconomic status, and access to resources that work together to create a favourable learning environment for children.

### 5.2 Effect of birth order on language comprehension

Another finding was the subtle effect of birth order on comprehension. The results showed that first-born children generally had better comprehension skills compared to their siblings. This may be because they receive more one-on-one attention from their caregivers before subsequent siblings are born. The effect of birth order suggests that parental dynamics and family structure play a key role in shaping early language skills, which suggests that attention needs to be more equitably distributed among children in a multi-child family.



### 5.3 Model Limitations and Uncertainties

While this study provides valuable observations, there are many limitations. A major limitation is the high degree of uncertainty in some of the model's coefficient estimates, particularly those related to race and the educational level of some caregivers. The wide confidence intervals indicate variability and inconsistency in the data, which may prevent the generalization of the findings. In addition, using samples may introduce bias and reduce the representativeness of the results across populations. In addition, missing data for some population groups limits the ability of the model to make predictions for all groups, which may affect the overall robustness of the study results.

### 5.4 Recommendations for Future Research

Future research should address these limitations by collecting more representative data across different demographic groups. This would reduce uncertainty in model estimates and provide a better understanding of the relationships between different predictors and language comprehension. Longitudinal studies would also be valuable in assessing these relationships' stability over time, providing a deeper understanding of the long-term effects of early intervention. In addition, further exploration of the role of caregiver education, particularly in families with multiple children, could provide more effective strategies to support equitable language development. The expansion of the dataset to include more vocabulary categories and contextual factors such as socioeconomic status, bilingual ability, and parental language practices also enhances our understanding of the mechanisms of early language learning and supports the development of more nuanced educational policies.

## Appendix

### A CDI methodology overview and assessment

#### A.1 Overview

The Communicative Development Inventory (CDI) is a standardized assessment tool used to study early language development in children aged 8 to 30 months. Originally developed by Fenson et al. (1994), the CDI provides a reliable and comprehensive method for evaluating both vocabulary comprehension and production based on parent reports. This chapter revisits the strengths and limitations of CDI, its psychometric properties, and the reliability of the instrument in longitudinal studies.

#### A.2 Target population, frame, and sample

- **Target Population:** Infants and toddlers aged 8 to 30 months whose parents are willing to report on their child's language development.
- **Sample Frame:** The target sample includes children from different demographic backgrounds to capture a diverse representation of language development. Recruitment was conducted across pediatric clinics, childcare centers, and online forums targeting parents, ensuring participants come from various socioeconomic, educational, and geographic contexts.
- **Sample Size:** The study collected data from approximately 2,000 parents. A sample size of this scale ensures generalizability while maintaining a margin of error within  $\pm 3\%$  at a 95% confidence level, thereby enhancing the robustness of findings.

#### A.3 Sample recruitment

##### Recruitment Methods:

- **Parents were recruited via two primary methods:** direct outreach at pediatric clinics and digital advertisement on online parenting platforms such as forums and social media. This double approach ensured coverage of typically underrepresented groups such as rural families. The recruitment was conducted with bilingual support (English and Spanish) to improve inclusivity and response rates.
- **Online Panel Use:** An additional recruitment method involved using an online panel to increase the diversity of the sample. This panel was instrumental in reaching parents from areas with limited in-person outreach capabilities.

## **A.4 Sampling approach and trade-offs**

### **Sampling Approach:**

- The approach involves collecting CDI (Communicative Development Inventory) data across various languages and dialects, using both cross-sectional and longitudinal sampling. This data is sourced from multiple studies conducted over the years, sometimes gathered through paper forms, mail-in surveys, in-person sessions, or electronically. Samples are typically convenience samples, with limited control over representativeness.
- Types of Data: The data includes both “Words & Gestures” (WG, for infants) and “Words & Sentences” (WS, for toddlers), gathered from children ranging from infancy to around 36 months. Instruments include checklists and questionnaires aimed at assessing language comprehension and production.

### **Trade-offs:**

- Advantages: The flexibility of using both WG and WS instruments allows for capturing a broad spectrum of early language development stages. The data provides valuable insight into different aspects of language development, including gestures, morphology, and grammar.
- Disadvantages: The sampling method is not necessarily representative due to reliance on convenience samples, which leads to biases depending on how researchers choose or access participants. Moreover, inconsistent administration methods such as electronic, paper and in-person lead to potential variability in data quality, affecting comparability between studies.

## **A.5 Non-response handling**

Non-response Issues: The variability of the data is mentioned in the CDI methodology and is sometimes described as a “difficult dataset” where inconsistent or surprising results occur. This may indicate that non-response issues or participant misunderstanding of the questions affected data quality.

### **Strategies for Handling Non-response:**

- Multiple Attempts: The research team attempted to handle potential non-response issues by improving the clarity of instructions such as simplifying written instructions for electronic forms, thereby reducing misinterpretations that could lead to “floor and ceiling” effects.

- **Data Inclusion:** Although recognizing the variability caused by nonresponse, the research team decided not to exclude any data sets. They instead accept the confusion and variability inherent in it, and recognize that exclusion of these data may lead to biased conclusions. This method of inclusion was intended to avoid circular arguments and preserve the natural variability in the dataset.

## **A.6 Longitudinal Stability of CDI Measurements**

**Stability Across Time:** Longitudinal data from CDI demonstrate high stability in early language scores. Studies like De Houwer et al. (2005) reveal significant correlations in CDI measurements across different ages, indicating that early language acquisition remains relatively stable over time.

## **A.7 Psychometric Analysis and Item Response Theory (IRT)**

- **Measurement Properties of CDI Items:** Using IRT, each CDI item was evaluated for difficulty and discrimination. Items like common nouns showed high discrimination, indicating their effectiveness in distinguishing between different levels of vocabulary ability.
- **Psychometric Weaknesses:** Items such as “mommy” and “daddy” have low discrimination because they are almost universally known by children, providing limited information about individual differences in vocabulary. However, these items are retained due to their cultural and developmental relevance.

### **IRT Evaluation Results:**

- **Difficulty:** Items varied in difficulty, with abstract words proving harder for parents to report on than concrete items.
- **Discrimination:** High discrimination scores were associated with verbs and adjectives, suggesting these word types are better indicators of a child’s linguistic development level compared to function words.

## **A.8 Simulation of CDI Sampling Approach**

To further improve the reliability and representativeness of CDI findings, Bayesian inference can be used to assess different sampling strategies. Bayesian methods can integrate a priori information about sampling bias and provide a probability framework to assess the uncertainty of sampling outcomes. By incorporating a priori information related to demographic characteristics such as socioeconomic status and bilingual families, researchers can better model the effects of sampling bias and improve the robustness of their conclusions.

Simulations can also be extended to a variety of sampling methods using Bayesian modelling to estimate the impact of sampling bias, such as over-representation of certain demographic characteristics or under-sampling of key groups. Generating synthetic datasets containing these underrepresented groups would illustrate how these populations might affect language development patterns, and updates to the Bayesian model would adjust the model to reflect observed trends in the data over time.

## A.9 Conclusion

The CDI method, while not without its challenges, provides a generalized approach to understanding early language development through parental reports. The stratified sampling methodology used in CDI studies allows for a more representative understanding of language development across demographic contexts, thus mitigating potential biases in traditional parent-report measures. The CDI has been used in some studies to assess language development in children and adolescents. Although there are some limitations, such as potential bias in parent reports and subjectivity in assessing comprehension, the CDI has demonstrated strong reliability and utility in psychometric assessments, making it an important tool for researchers of child language development.

## B Additional data details

### B.1 Data manipulation and cleaning

In the data cleaning process, we prepared the raw data for analysis by applying transformations, filtering, and restructuring using several R packages, including `dplyr` (Wickham et al. 2023), `arrow` (Richardson et al. 2024), and `rsample` (Frick et al. 2024).

1. **Removing columns with missing values:** Using the `select(where(~ !all(is.na(.))))` function from `dplyr` to remove any columns that were entirely NA, thus retaining only those columns that contain meaningful data.
2. **Removing unnecessary columns:** We further removed columns deemed unnecessary for analysis, they were `downloaded`, `language`, `form`, `dataset_name`, `child_id`, `ethnicity`, `language_exposures`, `health_conditions`, `typically_developing` by using `dplyrselect()` function. This helped to streamline the dataset by eliminating redundant information.
3. **Handling missing values:** Removing rows with any missing values by using the `na.omit()` function, ensuring that our dataset contained only complete cases. This step was important to maintain data integrity and ensure accurate model training

4. **Transforming birth order:** The `birth_order` column was converted to a numeric sequence using `mutate()` and `factor()`. We mapped levels like “First”, “Second”, “Third” to a corresponding numeric value to facilitate easier use in statistical models.
5. **Converting categorical variables to binary:** Several categorical variables were converted into binary values to prepare the data for modeling. We used `mutate()` and `ifelse()` to recode variables:
  - `is_norming` was converted to 1 for “TRUE” and 0 for “False”.
  - `sex` was coded as 1 for “Male” and 0 for “False”.
  - `monolingual` was similarly converted to 1 for “TRUE” and 0 for “False”.
6. **Splitting the data:** We used the `rsample` package to perform a stratified split of the cleaned dataset based on the `race` variable. This ensured that all levels of `race` were well-represented in both the training (70%) and testing (30%) sets. We set a seed (`set.seed(123)`) for reproducibility.
7. **Saving the cleaned data:** Finally, the cleaned dataset was saved in different formats (CSV and Parquet) using `arrow` to facilitate efficient storage and accessibility. We saved both the full cleaned dataset and the resulting training/testing splits, enabling their use in further analyses and modeling.

## B.2 Descriptions of Each Predictor Variable

Table 2

Variable	Description
<code>age</code>	The age of the child in months.
<code>production</code>	The number of words a child can produce.
<code>is_norming</code>	A binary variable indicating if the child is included in the norming sample (1 = Yes, 0 = No).
<code>birth_order</code>	The birth order of the child (e.g., 1 = first-born, 2 = second-born).
<code>caregiver_education</code>	The education level of the primary caregiver (e.g., Graduate, Some College).
<code>race</code>	The race of the child (e.g., White, Asian).
<code>sex</code>	The sex of the child (0 = Female, 1 = Male).
<code>monolingual</code>	A binary variable indicating if the child is monolingual (1 = Yes, 0 = No).

Note: This table provides descriptions of each predictor variable used in the study, explaining their relevance and values.

Descriptions of Each Predictor Variable

## C Model details

### C.1 Model summary

Table 3 presents the coefficients from our model analyzing factors that influence comprehension in early childhood. The `intercept` value is 78.522, represents the baseline level of comprehension when all other predictors are at their reference levels. Key variables include `is_norming` with a positive coefficient indicating that children included in the norming sample tend to have higher comprehension scores. Different caregiver education levels show varying effects. For example, `"caregiver_educationPrimary"` has a positive coefficient of 10.214, suggesting that having a caregiver with primary education is associated with increased comprehension. On the other hand, `"caregiver_educationSome Secondary"` has a negative coefficient, indicating a potential decrease in comprehension scores for children whose caregivers have some secondary education.

### C.2 Posterior predictive check

In Table 4 we conduct a posterior predictive check, showing the overlap between the observed data (denoted by  $y$ ) and the replicated data generated from the model (denoted by  $y_{rep}$ ). In a well-fitting model, the distribution of the replicated data should match the observed data closely. In this case, the replicated lines generally follow the shape of the observed data's density, especially around the peak, showing that the model is effectively capturing the main structure of the data. Any discrepancies between the replicated and observed lines would highlight areas where the model fails to adequately capture the observed data characteristics.

### C.3 Variance Inflation Factors for Each Predictor in the Model

Table 5 presents the Variance Inflation Factors (VIF) for each predictor in the model. VIF values are used to assess multicollinearity among the predictors, with values greater than 10 typically indicating a high level of multicollinearity. In this model, all predictors have VIF values below 3, suggesting that multicollinearity is not a significant concern. For example, `age` and `production` have VIF values of 2.559 and 2.500, respectively, indicating a moderate relationship with other predictors. Other variables, such as `is_norming` and `monolingual`, show lower VIF values, which indicates minimal multicollinearity.

Table 5: VIF for Each Predictor in the Model

	GVIF	Df	GVIF..1..2.Df..
age	2.559	1	1.600
production	2.500	1	1.581
is_norming	1.282	1	1.132

Table 5: VIF for Each Predictor in the Model

	GVIF	Df	GVIF <sup>1/2</sup>
birth_order	1.092	1	1.045
caregiver_education	1.326	6	1.024
race	1.344	3	1.050
sex	1.026	1	1.013
monolingual	1.196	1	1.094

## C.4 Diagnostics

Figure 7 is a visualization of the Gelman-Rubin diagnosis (often denoted as R-hat) for evaluating the convergence of Markov chain Monte Carlo (MCMC) samples in Bayesian models. The plot illustrates the R-hat values for various parameters. An R-hat value near 1 indicates that the MCMC chains have mixed well and are effectively sampling from the same posterior distribution. In Figure 7, all parameters exhibit an R-hat value at or below 1.05, suggesting that the model has reached adequate convergence, and the resulting parameter estimates are reliable for interpretation.

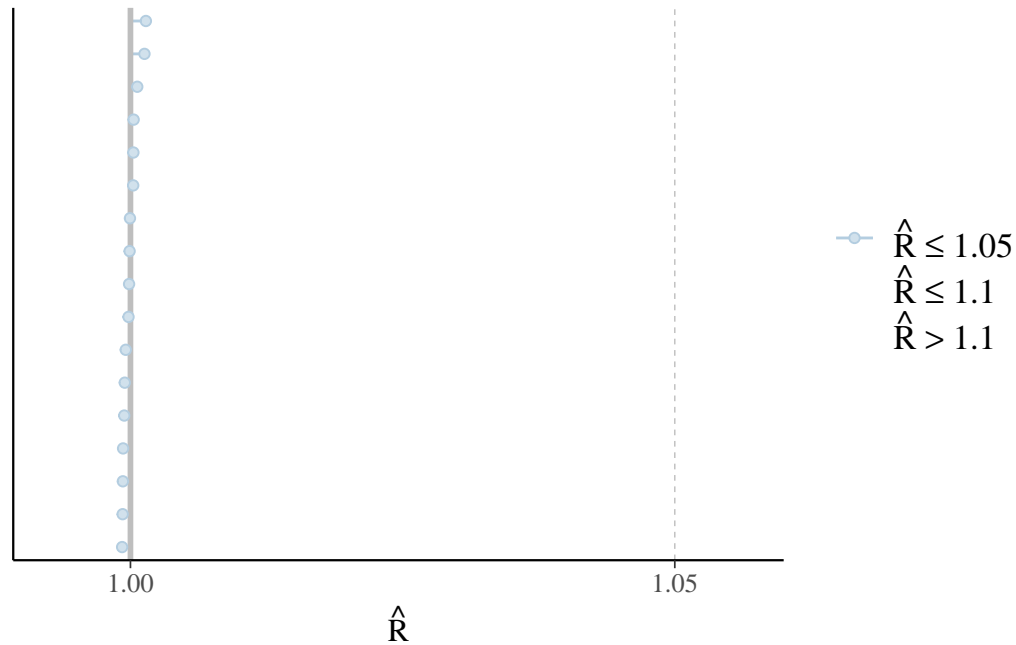


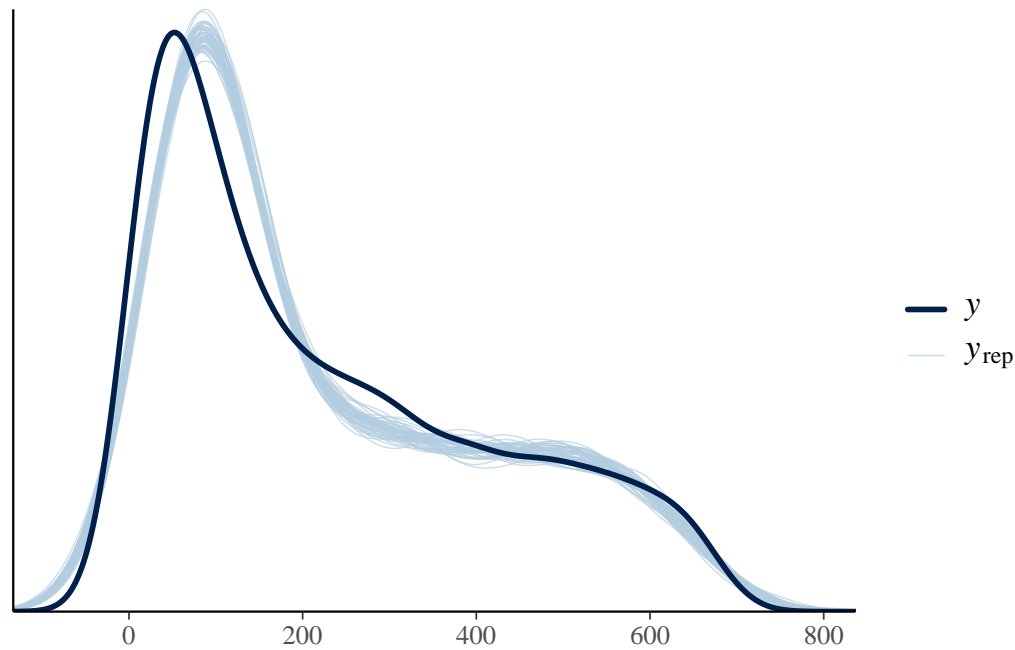
Figure 7: All parameters have R-hat values at or below 1.05, which indicates strong convergence of the MCMC chains and reliable parameter estimates.



Table 3: Coefficients from a regression model examining factors influencing comprehension

	coefficient
(Intercept)	78.522
age	−3.481
production	0.983
is_norming	26.095
birth_order	−1.921
caregiver_educationGraduate	1.516
caregiver_educationPrimary	10.214
caregiver_educationSecondary	6.486
caregiver_educationSome College	−2.216
caregiver_educationSome Graduate	−0.135
caregiver_educationSome Secondary	−8.677
raceBlack	4.699
raceOther	10.851
raceWhite	7.697
sex	0.034
monolingual	−1.091
Num.Obs.	2457
R2	0.936
R2 Adj.	0.935
Log.Lik.	−13 004.701
ELPD	−13 020.2
ELPD s.e.	75.0
LOOIC	26 040.4
LOOIC s.e.	150.0
WAIC	26 040.3
RMSE	48.09

Table 4



## References

- Alexander, Rohan. 2023. *Telling Stories with Data*. Chapman; Hall/CRC. <https://tellingstorieswithdata.com/>.
- Arel-Bundock, Vincent. 2022. “modelssummary: Data and Model Summaries in R.” *Journal of Statistical Software* 103 (1): 1–23. <https://doi.org/10.18637/jss.v103.i01>.
- Bolker, Ben, and David Robinson. 2024. *Broom.mixed: Tidying Methods for Mixed Models*. <https://CRAN.R-project.org/package=broom.mixed>.
- Fox, John, and Sanford Weisberg. 2019. *An R Companion to Applied Regression*. Third. Thousand Oaks CA: Sage. <https://www.john-fox.ca/Companion/>.
- Frank, Mika et al., Braginsky. n.d. “Wordbank Database.” *Wordbank*. [https://wordbank.stanford.edu/data/?name=admin\\_data](https://wordbank.stanford.edu/data/?name=admin_data).
- Frick, Hannah, Fanny Chow, Max Kuhn, Michael Mahoney, Julia Silge, and Hadley Wickham. 2024. *rsample: General Resampling Infrastructure*. <https://rsample.tidymodels.org>.
- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2022. “rstanarm: Bayesian applied regression modeling via Stan.” <https://mc-stan.org/rstanarm/>.
- Lawton, Will, Ozzy Araujo, and Yousif Kufaishi. 2023. “Language Environment and Infants’ Brain Structure.” *Journal of Neuroscience* 43 (28): 5129–31. <https://doi.org/10.1523/JNEUROSCI.0787-23.2023>.
- Müller, Kirill. 2020. *Here: A Simpler Way to Find Your Files*. <https://CRAN.R-project.org/package=here>.

- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Richardson, Neal, Ian Cook, Nic Crane, Dewey Dunnington, Romain François, Jonathan Keane, Dragoş Moldovan-Grünfeld, Jeroen Ooms, Jacob Wujciak-Jens, and Apache Arrow. 2024. *arrow: Integration to 'Apache' 'Arrow'*. <https://github.com/apache/arrow/>.
- Wickham, Hadley. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *dplyr: A Grammar of Data Manipulation*. <https://dplyr.tidyverse.org>.
- Wilke, Claus O. 2024. *Ggridges: Ridgeline Plots in 'Ggplot2'*. <https://CRAN.R-project.org/package=ggridges>.
- Xie, Yihui. 2014. “Knitr: A Comprehensive Tool for Reproducible Research in R.” In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman; Hall/CRC.
- Zhu, Hao. 2024. *kableExtra: Construct Complex Table with 'Kable' and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.
- Zonarich, Elizabeth. 2024. “Why Do Some Kids Learn to Talk Earlier Than Others?” *Harvard Gazette*, October. [https://news.harvard.edu/gazette/story/2024/01/why-do-some-kids-learn-to-talk-earlier-than-others-childhood-development-linguistics/?utm\\_source=chatgpt.com](https://news.harvard.edu/gazette/story/2024/01/why-do-some-kids-learn-to-talk-earlier-than-others-childhood-development-linguistics/?utm_source=chatgpt.com).
- Zubrick, Stephen R., Catherine L. Taylor, and Daniel Christensen. 2015. “Patterns and Predictors of Language and Literacy Abilities 4-10 Years in the Longitudinal Study of Australian Children.” *PLOS ONE*, September. [https://journals.plos.org/plosone/article?id=10.1371%2Fjournal.pone.0135612&utm\\_source=chatgpt.com#abstract0](https://journals.plos.org/plosone/article?id=10.1371%2Fjournal.pone.0135612&utm_source=chatgpt.com#abstract0).